

# A General Framework for Method Chaining Style Embedding of Domain Specific Languages

Hao Xu

University of North Carolina at Chapel Hill  
xuh@cs.unc.edu

## Abstract

Embedded Domain Specific Languages(DSL) are used in various programming tasks, such as SQL in database applications and regular expressions in pattern matching. Embedding DSL programs as strings in the host language that are interpreted at runtime make it extremely difficult to statically check for errors in the embedded programs. One solution is using specialized checkers, but the programming and integration of these checkers into the software development process is not easy. Another emerging paradigm is to extend the host language with support for the DSL, such as LINQ[8]. However, this requires modification of the compiler which is sometimes even more difficult, especially when the source code of the compiler is not available. Embedding a DSL into a host language using nested functions or constructors or combinators has been studied extensively[6]. However, method chaining style embedding seems to be largely ignored, probably because most of the research is done in a functional setting. As the object-oriented programming languages become more and more mature in their type systems, method chaining style may become a viable alternative as it seems to be very natural in object-oriented programming languages and widely accepted in practice[7]. In this paper, we explore the idea of type-level programming to embed DSLs into a strongly typed host language in the method chaining style. We use Scala[1, 4] to demonstrate our ideas, but these ideas should be easily applicable in other object-oriented programming languages such as Java. The ideas are partially implemented in an ongoing project – the EriLex parser generator.

**Keywords** Doman Specific Language, Method Chaining Style

## 1. Introduction

Embedded Domain Specific Languages (DSL) are used in various programming tasks, such as SQL in database applications and regular expressions in pattern matching. However, code written in these DSLs is usually embedded as a string in the host language and are interpreted at runtime. Consequently, the errors in the embedded code are not detected until runtime. One solution is checking statically for errors in the embedded code using specialized checkers. However, programming and integration of these checkers into the

software development process is not easy for programmers. Another emerging paradigm is to extend the host language with support for the DSL, such as LINQ. However, this requires modification of the compiler which is sometimes even more difficult, especially when the source code of the compiler is not available. In this paper, we explore the idea of method chaining style embedding of DSLs into a strongly typed object-oriented programming language.

The idea of embedding a DSL into the host language using constructs of the host language is not new. Libraries that use this technique include various monadic parser libraries that embed BNF-like DSLs, software engineering libraries such as jMock and Hibernate Criteria Query, etc. The contribution of this paper is that we present a general framework for method chaining style embedding of DSLs into a strongly typed object-oriented programming language. We consider both the grammar and typing of the embedded language. There are several advantages to this idea in general: first, it uses the type checker of host language to check the grammar and typing of the embedded language without extra tools, which makes it less error-prone; second, given that some advanced integrated development environments such as Netbeans and Eclipse provide semantics based assistance, the embedded program can take advantage of these tools provided for the host language.

There are at least two styles of embedding a DSL. One of them embeds a DSL into a host language using nested functions or constructors or combinators, which we call the functional nesting style(FNS). FNS has been studied extensively[6]. In FNS, most terms in the object language are represented by terms in the host language, such as combinators or inductively defined data types. In contrast, the method chaining style(MCS) encodes most terms of the object language as methods in the host language, which are sometimes partial or second-class constructs in the host language, such as in Java. However there are a number of advantages in the MCS. First, sometimes the embedding of MCS is much more compact than that of FNS in an object-oriented programming language. For example to embed

```
if x then 1 else 0
```

in Java using FNS, we would write code that looks like:

```
new If(new Var("X"), new Con(1), new Con(0))
```

which is very difficult to read because of all the parentheses and keywords. Even using case classes in Scala the encoding is still a little difficult to read, because the role of `Var("X")`, `Con(1)`, and `Con(0)` are not clear from the representation.

```
If(Var("X"), Con(1), Con(0))
```

Adding constructors for `then` and `else` would further "clog" the code.

On the other hand, in MCS, we could write

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

If ( . var ("X") . Then ( ) . con ( 1 ) . Else ( ) . con ( 0 )

in Java and

If ( ) varx ("X") Then ( ) con ( 1 ) Else ( ) con ( 0 )

in Scala<sup>1</sup>, which is much easier to read.

Another advantage is that MCS makes only limited use of host language delimiters in the object language compared to FNS. for example, if we want to encode

**app f app f X**

in FNS, the encoding may look like

App ( Var (" f " ) , App ( Var (" f " ) , Var (" X " ) )

in Scala, which uses parenthesis in the host language to delimit the pair, even when in the object language it may be clear that **app** associates to the right. In the functional programming style, usually something like

app ( var f ) ( app ( var f ) ( var X ) )

is necessary because unless given special rules for the **app** function in the host language,

app ( var f ) app ( var f ) ( var X )

would be a totally different expression. In MCS, we can encode it in Scala as

app ( ) varx (" f " ) app ( ) var (" f " ) var (" X " )

which is much less relying on the host language.

The idea presented in this paper is implemented in an ongoing project – the EriLex<sup>2</sup> parser generator.

The paper is organized as follows: Section 2 reviews various concepts related to pushdown automata and introduces an embedding of a special class of automata into the types of Scala; Section 3 discusses practical aspects of implementing these types; Section 4 extends the embedding to deal with typing; Section 5 discusses related work; Section 6 discusses extensions to this framework.

## 2. Embedding Realtime Deterministic Pushdown Automata in Scala

Most of the definitions of grammars and automata in this section are based on [11].

**Definition 2.1.** A (context-free) grammar (CFG)  $G$  is a quadruple  $(A_G, Z_G, z_G, P_G)$ , where  $A_G$  is a finite set of terminals,  $Z_G$  is a finite set of nonterminals,  $P_G$  a finite set of productions, and  $z_G$  is the start symbol. The language of a context-free grammar, written  $L(G)$ , is the set of all finite length strings generated by  $z_G$ , which is usually called a context-free language (CFL).

Let  $a, a_1, a_2 \dots$  range over  $A_G$ ,  $z, z_1, z_2 \dots$  range over  $Z_G$ ,  $a^n = a_1^n \dots a_n^n, i^n = i_1^n \dots i_n^n, \dots$  range over  $A_G^n$ , variables  $z^n = z_1^n \dots z_n^n, \dots$  range over  $Z_G^n$ . Most CFG can be converted to certain equivalent normal forms. One of them is the Greibach normal form (GNF). There are usually two different definitions of GNF. One of them includes the production  $z_G \rightarrow \epsilon$ . We use the following definition.

**Definition 2.2.** A CFG  $G$  is in the Greibach normal form if all productions in  $P_G$  are of the form  $z \rightarrow az^n$  where  $n \in \mathbb{N}$ .

**Theorem 2.3.** (Greibach) Every CFL  $L$  where the empty string  $\epsilon \notin L$  can be generated by a CFG in the Greibach normal form.

<sup>1</sup> Here we use `varx` instead of `var` because the latter is a keyword in Scala.

<sup>2</sup> <http://erilex.kenai.com>

Given a CFL  $G$ , a parse tree  $\text{tree}(as_1 \dots s_n)$  of a string  $as_1 \dots s_n$  is a tree defined inductively

$$\text{tree}(as_1 \dots s_n) = z(a, \text{tree}(s_1), \dots, \text{tree}(s_n))$$

where  $z \rightarrow az_1 \dots z_n \in P_G$  and  $s_k$  is generated by  $z_k$  for all  $k \in \{1, \dots, n\}$ . The strings  $s_k$  are called substrings. Let  $t^n = t_1^n \dots t_n^n$  range over  $n$  subtrees and a pattern matching on trees is written  $z(at^n)$ . CFLs can be accepted by pushdown automata. Languages of CFG in GNF can be accepted by a special class of pushdown automata.

**Definition 2.4.** A realtime pda over an input alphabet  $A$  is  $(Q, Z, i, K, A, T)$ , where  $Q$  is a finite set of states,  $Z$  a stack alphabet,  $T$  a finite subset of  $A \times Q \times Z \times Z^* \times Q$  which is called the transition rules,  $i$  an initial configuration, and  $K$  a set of accepting configurations. a realtime pda is deterministic, or a dpda, if and only if  $T$  is a function  $A \times Q \times Z \rightarrow Z^* \times Q$ .

The different between a realtime pda and a pda is that realtime pdas do not have  $\epsilon$ -rules. Because any proper context-free language that does not contains  $\epsilon$  can be generated by a CFG in GNF without  $\epsilon$ -productions, the languages accepted by realtime pdas differ from those accepted by pdas only in whether  $\epsilon$  may be included, which is not a significant restriction in DSL design.

What is interesting is the difference between realtime deterministic pda and deterministic pda. Although dpda is strictly less powerful than pda, it has been shown that languages generated by dpda are exactly  $LR(k)$ , which contains most programming languages. Therefore, dpda is powerful enough for our purpose. But realtime dpda is strictly less powerful than dpda. There are languages that are accepted by some dpda but no realtime dpda. For example, the language  $L = \{ \langle^p \{^n \rangle^p \mid p, n > 0 \} \cup \{ \langle^p \{^n \}^n \mid p, n > 0 \}$ , which is a modification of an example in [11], can be accepted by a dpda but no realtime dpda. The problem is that in order to make sure that at least one of “>” or “}” matches with “<” or “{”, we need to push state symbols on the stack to match the number of “<”s and “{”s. If “>” is encountered, we need to remove those symbols that are generated by “{”. However, this is a very rare situation in language design, because here we want to design a language that does not keep all brackets matched on purpose. If the object language has a universal hierarchy of programming constructs in which programming constructs in the lower level are strictly and properly contained (i.e. no overlapping boundary) by programming constructs in the higher level, as in Java where the delimiter such as “{” and “}” are used verbosely, then realtime dpda is expressive enough because the number of pops is always bounded. We work with dpda in the following sections.

Next, we look at how to embed a realtime dpda in Scala. Before proceeding to define the embedding, we first present an equivalent form of definition of a realtime dpda.

**Definition 2.5.** A stateless realtime dpda (slrpd)  $\mathcal{A}$  over an input alphabet  $A$  is  $(Z, z_A, A, T)$ , where  $Z$  is a stack alphabet,  $T$  a function  $A \times Z \rightarrow Z^*$  which is called the transition rules, and  $z_A$  an initial configuration, with  $\epsilon$  as an accepting configuration. The language accepted by  $\mathcal{A}$  is the set of all strings over  $A$  accepted by it.

It is easy to see that we can construct an equivalent slrpd from every realtime dpda by encoding the current state in the top stack symbol.

Let  $\mathcal{A}$  be a slrpd. We assume that a name embedding function  $e : A \cup Z \rightarrow \text{ScalaName}$  that maps injectively slrpd symbols to valid Scala names is predefined. We write  $e(a)$  as  $e_a$ , and  $e(z)$  as  $e_z$ .

To simplify the presentation, we may write a Scala class

```
class c{def1 def2 ... defn}
```

as

```
class c{def1}, class c{def2}, ..., class c{defn}
```

where  $c$  is a class signature and  $def_1, def_2, \dots, def_n$  are method definitions. When combining `class c1{def1}`, `class c2{def2}`, ..., `class cn{defn}` where  $c_1, \dots, c_n$  all have the same type constructor, back to `class c{def1 def2 ... defn}`, we may need to rename the type parameters in  $c_1, \dots, c_n$ , which can be done using simultaneous unification. For example, `class A[τ1]{def1}`, ..., `class A[τn]{defn}` are combined to `class A[τ]{def1, ..., defn}`. For succinctness, we write type  $Null$  as  $\perp$ .

**Definition 2.6.** The embedding of an slrpd  $\mathcal{A} = (Z, z_{\mathcal{A}}, A, \mathcal{T})$ , written  $en_{\perp}^{\mathcal{A}}(\mathcal{A}, e)$ , is a set of classes with the name of  $e_z$  for all  $z \in Z$ .

$$en_{\perp}^{\mathcal{A}}(\mathcal{A}, e) = \{ \text{abstract class } e_z[\kappa] \{ \\ \text{def } e_a() : e_{z_1}[\dots e_{z_n}[\kappa]] \\ |(a, z, z_1 \dots z_n) \in \mathcal{T} \} \}$$

The embedding of an input string  $a^n$ , written  $en_{\mathcal{S}}^{\mathcal{A}}(a^n)$ , is the method chain

$$en_{\mathcal{S}}^{\mathcal{A}}(a^n, e) = e_{a_1}^n() \dots e_{a_n}^n()$$

The well-formness is summarized in the following theorem.

**Theorem 2.7.** *If  $\mathcal{A}$  is a slrpd and  $imp_{z_{\mathcal{T}}}$  is an implementation of  $e_{z_{\mathcal{T}}}$ , then  $a^n \in L(\mathcal{A})$  if and only if `new impzT[⊥]()`ens( $a^n, e$ ) has type  $\perp$  in Scala.*

*Proof.* (Sketch) [4] did not deal with generics, but we can still use the (METHOD) rule, by viewing the type metavariables as also ranging over constructed types. Given a sequence of input  $a^n$ , we construct a sequence of configurations of the slrpd. Also, given the method chain, we construct a sequence of types by instantiating the (METHOD) rule. The initial empty stack is represented by  $\perp$ . By induction, we can prove that the each configuration is represented by the corresponding type. Because the type of the Scala expression is  $\perp$ , the input is accepted by the slrpd.  $\square$

A grammar in the Greibach normal form (GNF) and LL(1) (i.e. no common prefix) can be converted directly to an slrpd. The slrpd is quite straightforward to construct, by just using nonterminals as stack symbols. The embedding  $en_{\perp}^G(G)$  of a LL(1) GNF grammar  $G$  is defined as

**Definition 2.8.** The embedding of  $G$  in LL(1) GNF w.r.t name embedding function  $e$ , written  $en_{\perp}^G(G, e)$ , is a set of classes with the name of  $e_z$  for all  $z \in Z$ .

$$en_{\perp}^G(G, e) = \{ \text{abstract class } e_z[\kappa] \{ \\ \text{def } e_a() : e_{z_1}^n[\dots e_{z_n}^n[\kappa]] \\ |z \rightarrow az^n \in P \} \}$$

The embedding of an input string  $a^n$ , written  $en_{\mathcal{S}}^G(a^n, e)$ , is the method chain

$$en_{\mathcal{S}}^G(a^n, e) = e_{a_1}^n() \dots e_{a_n}^n()$$

**Theorem 2.9.** *If  $G$  is a grammar in LL(1) GNF and  $imp_{z_G}$  is an implementation of  $e_{z_G}$ , then  $a^n \in L(G)$  if and only if `new impzG[⊥]()`ens( $a^n$ ) has type  $\perp$  in Scala.*

*Proof.* (Sketch) Prove by converting the grammar to an equivalent slrpd.  $\square$

Next, let's look at an example.

```
abstract class es[κ]{
  def ec() : er[κ]
}
abstract class er[κ]{
  def ea() : er'[κ]
  def ec() : er[er'[κ]]
}
abstract class er'[κ]{
  def ey() : κ
  def e*() : er'[κ]
  def ea() : er'[κ]
  def ec() : er[er'[κ]]
}
```

Figure 1. Embedding Scala Classes

**Example 2.10.** Let  $\Sigma$  be a finite alphabet. We can embed in Scala the regular expression DSL given by  $s \rightarrow (r), r \rightarrow a|r * |rr|r'|r(r)$ , where  $a \in \Sigma$ . We added an additional symbol  $s$ , which makes the productions much easier to work with. First, we convert it to an grammar in LL(1) GNF.

$$\begin{array}{ll} r & \rightarrow ar' & r' & \rightarrow *r' \\ r & \rightarrow (rr' & r' & \rightarrow '|r \\ r' & \rightarrow ) & r' & \rightarrow (rr' \\ r' & \rightarrow ar' & s & \rightarrow (r \end{array}$$

We can generate the embedding directly from the GNF. For example, for the first production, we have `class er[κ]{def a() : er'[κ]}`. When combined together, the embedding classes are in the Figure 1.

If  $imp_s$  is an implementation of  $e_s$ , then for a regular expression  $(a^*)$ ,

$$\text{new imp}_s[\perp]().e_c().e_a().e_*(e_y())$$

has type  $\perp$ ; and for a string  $|a|$ ,

$$\text{new imp}_s[\perp]().e_c().e_a().e_l()$$

is rejected by Scala.

The intuition behind this embedding is that a type  $e_z[e_{z_1}[\dots e_{z_n}[\perp]]]$  represents the configuration of the parser as a stack of the type (constructors). The current parser state is parsing nonterminal  $z$ . When  $z$  is successfully parsed, the parser continues to parse  $z_1$ , and so forth. When parsing for some  $z$  in the grammar, the leading terminals decide which productions to enter and which nonterminals to push onto the stack that represents the next configuration of the parser. The choice is always deterministic because the grammar is LL(1).

**Definition 2.11.** We call the stack of type (constructors) representing the current configuration of the slrpd the continuation stack. The top of continuations stack refers to the outermost type (constructor).

Suppose that the continuation stack is  $e_z[\perp]$ . If the selected production is of the form  $z \rightarrow a$ , then it just removes the current symbol from the top of the continuation stack by setting the return type of the method the type  $\kappa$ , and the new stack is represented by type  $\kappa$  as shown in Figure 2(a). If the selected production of the form  $z \rightarrow az_1z_2$ , then it pushes the type constructors representing  $z_1, z_2$  onto the continuation stack as shown in Figure 2(b).

### 3. Untyped MCS

In the previous section, we defined an embedding function  $en^G$  that embeds a grammar in LL(1) GNF into abstract classes of Scala. But in this embedding the embedding classes can not be instantiated because they are abstract classes. In this section, we discuss

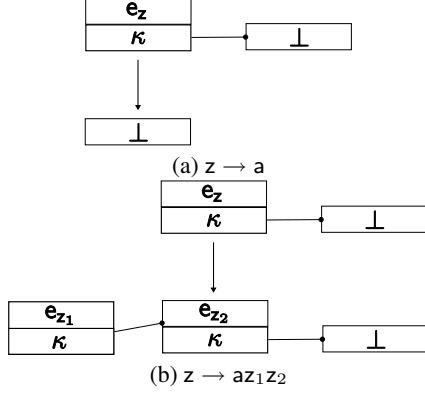


Figure 2. Examples

systematically the methodology of embedding that generates concrete classes.

### 3.1 Methodology

The methodology of embedding is that each embedding  $en$  of a grammar  $G$  is a pair of embedding functions  $(en_L, en_S)$  parameterized over  $G$  and a name embedding function  $e$ , where  $en_L$  is a *language embedding* that generates support classes from  $G$ , and  $en_S$  is a *string embedding* that generates an expression for a string  $s \in L(G)$ . And the soundness theorem will be in the form that if  $G$  satisfies certain properties, then with embedding  $en_L(G)$ , a string  $s \in L(G)$  if and only if the embedding  $en_S(s)$  of  $s$  is well-typed in Scala.

### 3.2 Workaround for Restrictions in Scala

Before we look at the embedding, let us look at a restriction of generics in Scala and how to work around it. One of the restrictions in languages that use erasure such as Scala and Java is that type parameters can not be used to create an object because they are erased during compilation. For example, `new  $\tau$ ()` where  $\tau$  is a type parameter is not allowed. This is in contrast to C++ templates and C# generics, which expand the generic types to instances of these types during compilation or at runtime.

**Example 3.1.** We implement the abstract classes in the previous example. In a first attempt, we may be tempted to implement the class  $e_{r'}[\kappa]$  as

```
class e_{r'}[\kappa]{
  def e_s() = new \kappa()
  ...
}
```

but this is not a valid Scala program! This problem occurs for methods that represents productions of the form  $z \rightarrow a$ . We can work around this restriction by letting the creator of this class pass in a continuation that knows how to create an object of type  $\kappa$ . In this example, we can simply pass in an object of type  $\kappa$ .

```
class e_{r'}[\kappa](k : \kappa){
  def e_s() = k
  ...
}
```

The implementation of other methods are straightforward. The implementation of  $e_s, e_r, e_{r'}$  is shown in Figure 3. In methods that add more than one type constructors the type of the continuation, such as  $e_c()$ , we need to rewire the continuation properly. Now, in order to initialize the parser, we set the type of the initial continuation to  $\perp$  and pass in an initial continuation *null* which is the only object of that type .

```
class e_s[\kappa](k : \kappa){
  def e_c() = new e_r[\tau](k)
}
class e_r[\kappa](k : \kappa){
  def e_a() = new e_{r'}[\kappa](k)
  def e_c() = new e_r[e_{r'}[\kappa]](new e_{r'}[\kappa](k))
}
class e_{r'}[\kappa](k : \kappa){
  def e_a() = new e_{r'}[\kappa](k)
  def e_s() = new e_{r'}[\kappa](k)
  def e_c() = new e_r[e_{r'}[\kappa]](new e_{r'}[\kappa](k))
  def e_s() = k
}
```

Figure 3. Implementation

$$\begin{aligned} \kappa_0^n &\triangleq \kappa \\ \kappa_{i+1}^n &\triangleq e_{z_{n-i}}^n[\kappa_i^n] \\ k_0^n &\triangleq k \\ k_{i+1}^n &\triangleq \text{new } e_{z_{n-i}}^n(k_i^n) \end{aligned} \quad (a)$$

$$\begin{aligned} en_L^1(G, e) &= \{en_L^1(p) \mid p \in P_G\} \\ en_L^1(z \rightarrow az^n, e) &= \text{class } e_z[\kappa](k : \kappa) \{en_L^1(az^n)\} \\ en_L^1(az^n, e) &= \text{def } e_a() : \kappa_n^n = k_n^n \end{aligned} \quad (b)$$

$$en_S^1(G, a^n, e) = \text{new } e_{z_G}[\perp](null)e_{a_1}^n() \dots e_{a_n}^n() : \perp \quad (c)$$

Figure 4. Embedding Rules

### 3.3 The General Embedding Functions

The embedding  $en^1$  embeds a grammar  $G$  as concrete classes. The definition of the language embedding  $en_L^1$  is given in Figure 4(a) and Figure 4(b). The rules represents a production of the form  $z \rightarrow a$  by a method that removes the type constructor at the top of the continuation stack by returning the continuation  $k$  of type  $\kappa$  and represents a production of the form  $z \rightarrow az_1 \dots z_n$  by a method that pushes the type constructors representing  $z_1, z_2, \dots, z_n$  onto the continuation stack by constructing objects `new  $e_{z_1}(\dots e_{z_n}(k))$` . If written separately, they are

$$\begin{aligned} en_L^{1'}(a) &= \text{def } e_a() : \tau = \kappa \\ en_L^{1'}(az_1 \dots z_n) &= \text{def } e_a() : e_{z_1}[\dots e_{z_n}[\kappa]] = \text{new } e_{z_1}(\dots e_{z_n}(k)) \end{aligned}$$

The string embedding  $en_S^1$  is given in Figure 4(c).

**Corollary 3.2.** *If  $G$  is a grammar in LL(1) GNF,  $e$  is a name embedding function, then  $a^n \in L(G)$  if and only if  $en_L^1(G, e) \cup \{en_S^1(G, a^n, e)\}$  is well-typed in Scala.*

### 3.4 Parametrized Grammar

In this section, we introduce an extension to grammars, named parametrized grammars, that will be used in this paper.

**Definition 3.3.** A parametrized grammar is a grammar  $G = (Z_G, z_G, A_G, P_G)$  in the GNF with a mapping  $arity : P_G \rightarrow \mathbb{N}$ .

For each production  $z \rightarrow az^n$  of the grammar, arity assigns an integer  $k$  such that  $0 \leq k \leq n$ . The first  $k$  nonterminals on the right hand side of the production are called parameters of  $a$ . Furthermore, we require that  $Z_G$  can be divided into two disjoint subsets  $Z_{G_0}, Z_{G_1}$  such that in all productions for  $z \in Z_{G_1}$ , any  $z_0 \in Z_{G_0}$  only appears as parameters and for all productions for  $z_0 \in Z_{G_0}$ , no  $z \in Z_{G_1}$  or parameter appears.

Effectively, we divided  $Z_G$  into two levels. The terminals in productions of top level nonterminals may take parameters from the second level nonterminals. And the second level can be written as a set of inductively defined data types, with the terminals as constructors and nonterminals as types. We use  $\rightarrow$  to mark top level symbols and  $=$  to mark second level symbols. Given a parametrized grammar  $G$ , we write the productions for  $Z_{G_0}$  as  $P_{G_0}$  and productions for  $Z_{G_1}$  as  $P_{G_1}$ .

**Example 3.4.** The untyped lambda calculus with de Bruijn index:

$$\begin{aligned} e &\rightarrow \text{var}(i)|\text{app } e \text{ } e|\text{abs } e \\ i &= z|s \text{ } i \end{aligned}$$

$e$  is a top level nonterminal, while  $i$  is a second level nonterminal. The terminal  $\text{var}$  takes  $i$  as a parameter.

We will encode second level nonterminals as method parameters in the embedding, which is summarized in the following definition.

To disambiguate with name embedding function  $e$  which are used to encode top level symbols, we use another name embedding functions  $d$  to encode second level symbols. The embedding  $\text{en}_L^0$  of second level symbols of a parametrized grammar  $G$  with regard to the name embedding function  $d$  is a set of case classes (or inductive data types, as sometimes they are called) defined in Figure 5(a). The embedding  $\text{en}_L^1$  of top level symbols of a parametrized grammar  $G$  with regard to the name embedding functions  $e$  and  $d$  is a set of classes with the name  $e_z$  for all  $z \in Z_{G_1}$  as defined in Figure 5(b). The string embedding of a string  $a_1 \dots a_n$  generated by some  $z \in Z_{G_0}$  is inductively defined on its parse tree where  $t_k^n$  is a subtree for  $k \in \{1, \dots, n\}$  as shown in Figure 5(c). The string embedding of string  $a_1 i^{k_1} \dots a_n i^{k_n}$ , where  $i^{k_i}$  for  $i \in \{1, \dots, n\}$  are strings generated by some  $z \in Z_0$ , is a method chain as shown in Figure 5(d).

### 3.5 Building Abstract Syntax Trees

In this section we discuss how to modify the parser to build an abstract syntax tree when the parser is run. We first look at another example,

**Example 3.5.** Since the grammar in Example 3.4 is already in GNF, we can embed it as classes in Scala using  $\text{en}^1$ . We use the following name embedding:  $e$  map  $\text{var}$  to  $\text{Var}$ ,  $\text{app}$  to  $\text{App}$ ,  $\text{abs}$  to  $\text{Abs}$ , and  $e$  to  $E$ ;  $d$  maps  $e$  to  $\text{Term}$ ,  $\text{var}$  to  $\text{Var}$ ,  $\text{app}$  to  $\text{App}$ ,  $\text{abs}$  to  $\text{Abs}$ ,  $i$  to  $\text{Ref}$ ,  $z$  to  $\text{Zero}$ , and  $s$  to  $\text{Succ}$ .

```
class E[κ](k : κ){
  def var(x : Ref) : κ = k
  def app() : E[E[κ]] = new E(E(k))
  def abs() : E[κ] = new E(k)
}
```

We extend embedding  $\text{en}^0$  in Section 3.4 to top level symbols to generate data structures for storing the abstract syntax tree. A string generated by  $e$  is represented by  $\text{Term}$ , and a string generated by  $i$  is represented by  $\text{Ref}$ , as shown in Figure 6(a). Next, we add a builder function to the implementation as shown in Figure 6(b).

We made three modifications here. First, a new type parameter  $\sigma$  is added to the class indicating the type of the value built by the builder. Second, a parameter is added to the class  $b : \text{Term} \Rightarrow \sigma$ , which is the builder. Third, the type of  $k$  is modified to  $\sigma \Rightarrow \tau$ . The

$$\begin{aligned} \text{en}_L^0(G, d) &= \{\text{en}_L^0(p, d) | p \in P_{G_0}\} \cup \{\text{en}_L^0(z, d) | z \in Z_{G_0}\} \\ \text{en}_L^0(z \rightarrow az^n, d) &= \text{case class } d_a(i_1 : d_{z_1^n}, \dots, i_n : d_{z_n^n}) \\ &\quad \text{extends } d_z() \\ \text{en}_L^0(z, d) &= \text{abstract case class } d_z() \end{aligned} \tag{a}$$

$$\begin{aligned} \text{en}_L^1(G, d, e) &= \{\text{en}_L^1(p, d, e) | p \in P_{G_1}\} \\ \text{en}_L^1(z \rightarrow az^p z^n, d, e) &= \text{class } e_z[\kappa]\{ \\ &\quad \text{def } e_a(i_1 : d_{z_1^p}, \dots, i_p : d_{z_p^p}) : \\ &\quad \quad \kappa^n = \kappa^n \} \end{aligned} \tag{b}$$

$$\begin{aligned} \text{en}_S^0(a^n, G, d) &= \text{en}_S^{0'}(\text{tree}(a^n), G, d) \\ \text{en}_S^{0'}(z(\text{at}^n), G, d) &= d_a(\text{en}_S^{0'}(t_1^n, G, d), \dots, \text{en}_S^{0'}(t_n^n, G, d)) \end{aligned} \tag{c}$$

$$\begin{aligned} \text{en}_S^1(a_1 i^{k_1} \dots a_n i^{k_n}, G, d, e) &= \text{new } e_{z_G}[\perp]().s_1 \dots s_n : \perp \\ &\quad \text{where} \\ s_i &= e_{a_i}(\text{en}_S^0(i_1^{k_i}, G, d), \\ &\quad \dots, \text{en}_S^0(i_{k_i}^{k_i}, G, d)) \end{aligned} \tag{d}$$

Figure 5. Embedding Rules

```
abstract case class Term()
case class Var(x : Ref) extends Term()
case class App(t1 : Term, t2 : Term) extends Term()
case class Abs(t : Term) extends Term()

abstract case class Ref()
case class Zero() extends Ref()
case class Succ(r : Ref) extends Ref()

class E[κ, σ](b : Term ⇒ σ, k : σ ⇒ κ){
  def var(x : Ref) : κ = k(b(Var(x)))
  def app() : E[E[κ, σ], Term ⇒ σ] = new E(
    (t1 : Term) ⇒ (t2 : Term) ⇒ b(App(t1, t2)),
    (b1 : Term ⇒ σ) ⇒ E(b1, k))
  def abs() : E[κ, σ] = new E((t : Term) ⇒ b(Abs(t)), k)
}
```

Figure 6. Implementation with a Builder Function

reason is that we need to pass the builder to the continuation, so that the currently building tree is part of the continuation constructed by  $k$  during runtime.

Unlike in FNS, sometimes we need to construct a partial AST whose subtrees are filled in later. Partial trees are built by constructing new builders. For example, the  $\text{app}$  function constructs the builder  $(t_1 : \text{Term}) \Rightarrow (t_2 : \text{Term}) \Rightarrow b(\text{App}(t_1, t_2))$  and passes it to the continuation  $k$ , which encodes how to build the  $\text{App}$  object. But the actual construction of the  $\text{App}$  node happens when all parameters to the constructor are built.

To use this class, we need to pass in an initial builder to the constructor of the class where we create an object of the class.

Using the new embedding, we can encode terms of  $L(G)$  in Scala. We define an auxiliary Scala method.

```
def term() = new E((t : Term) => null, (s : Term) => s)
```

Using this method, the following term

```
app abs varx(z) con(1)
```

is encoded in Scala as

```
term() app() abs() varx(Zero()) con(1) : ⊥
```

and the following term

```
app app abs abs con(2) app abs varx(z) con(1) con(1)
```

is encoded in Scala as

```
term()
  app()
    app()
      abs()
        abs()
          con(2)
        app()
          abs()
            varx(Zero())
          con(1) : ⊥
```

Next, we formalize the ideas of Example 3.5 in the embedding  $\text{en}^B$ .  $\text{en}_B$  is defined as a pair of embeddings  $(\text{en}^d, \text{en}^b)$  which are parametrized over grammar  $G$ , name embedding function  $e$  which maps all symbols in  $G$  to a unique valid Scala names, and name embedding function  $d$  which maps all symbols in  $G$  to unique valid Scala names such that  $d_z \neq e_z$  for all  $z \in A_G \cup Z_G$ .  $e$  is used in  $\text{en}^b$  while  $d$  is used both in  $\text{en}^b$  and  $\text{en}^d$ .  $\text{en}^d$  is an extension of  $\text{en}^0$  to all top level productions  $p$  and for all top level nonterminals  $z$ , as shown in Figure 7(a) and Figure 7(b).  $\text{en}^b$  is an extension of  $\text{en}^1$  to include the builder function. We use the auxiliary definitions as shown in Figure 7(c).  $\sigma_n^n$  is the type for the  $n$ -ary builders that build values of type  $\sigma$ .  $p$  is the number of parameters of  $a$ .  $B_n^{n,p}$  is an  $n$ -ary builder and its type is always  $\sigma_n^n$ . The degenerated builder  $B_0^{0,p}$  is actually the value built by the builder.  $\kappa_n^n$  is the continuation stack.  $k_n^n$  is the continuation whose type is always  $\sigma_n^n \Rightarrow \kappa_n^n$ .  $z^p$  are parameters. The rules are shown in Figure 7(d). The string embedding  $\text{en}_S^b$  of the object language is defined in Figure 7(e), where  $a_k$  is a top level symbols for  $k \in \{1, \dots, n\}$  and  $i^{k_i}$  are second level strings for  $i \in \{1, \dots, n\}$ . The embedding  $\text{en}^B$  is defined in Figure 7(f).

We need the following theorem to complete our discussion.

**Theorem 3.6.** *If  $G$  is a grammar of the terms of the object language in  $LL(1)$  GNF,  $e$  maps  $A_G \cup Z_G$  in  $G$  to unique Scala names,  $d$  maps  $A_G \cup Z_G$  to unique Scala name such that  $d_z \neq e_z$  for all  $z \in Z_G$ . If  $a^n$  is a term of the object language, then  $a^n \in L(G)$  if and only if  $\text{en}_L^B(G, d, e) \cup \{\text{en}_S^B(a^n, G, d, e)\}$  is well-typed in Scala.*

*Proof.* (Sketch) Note that any encoded term is automatically well-formed because the constructor  $d_a$  only encodes well-formed terms. First, we prove that all encoded terms are well-formed. Then, by the type safety, the encoding always constructs encoding of some well-formed term. Finally, we prove that the term constructed is the same term that is encoded in the method chain.  $\square$

## 4. Typing MCS

In this section we discuss how to ensure well-typedness of an object language program embedded in MCS. One of the key observations of the previous section is that we can use the tree built to enforce well-formedness of the program. In this section, we explore how to

$$\begin{aligned} \text{en}_L^d(G, d) &= \{\text{en}_L^d(p, d) \mid p \in P_G\} \cup \{\text{en}_L^d(z, d) \mid z \in Z_G\} \\ \text{en}_L^d(z \rightarrow az^n, d) &= \text{case class } d_a(i_1 : d_{z_1^n}, \dots, i_n : d_{z_n^n}) \\ &\quad \text{extends } d_z() \end{aligned}$$

$$\text{en}_L^d(z, d) = \text{abstract case class } d_z() \quad (\text{a})$$

$$\begin{aligned} \text{en}_S^d(a^n, d) &= \text{en}_S^{d'}(\text{tree}(a^n), d) \\ \text{en}_S^d(z(\text{at}^n), d) &= d_a(\text{en}_S^{d'}(t_1^n, d), \dots, \text{en}_S^{d'}(t_n^n, d)) \end{aligned} \quad (\text{b})$$

$$\begin{aligned} \sigma_0^n &\triangleq \sigma \\ \sigma_{i+1}^n &\triangleq d_{z_{n-i}^n} \Rightarrow \sigma_i^n \\ \kappa_0^n &\triangleq \kappa \\ \kappa_{i+1}^n &\triangleq e_{z_{n-i}^n}[\kappa_i^n, \sigma_i^n] \\ k_0^n &\triangleq k \\ k_{i+1}^n &\triangleq (b_{i+1} : \sigma_{i+1}^n) \Rightarrow \text{new } e_{z_{n-i}^n}(b_{i+1}, k_i^n) \\ B_0^{n,p} &\triangleq b(e_a(i_1, \dots, i_p, t_1, \dots, t_n)) \\ B_{i+1}^{n,p} &\triangleq (t_{n-i} : d_{z_{n-i}^n}) \Rightarrow B_i^{n,p} \end{aligned} \quad (\text{c})$$

$$\begin{aligned} \text{en}_L^b(G, d, e) &= \{\text{en}_L^b(p, d, e) \mid p \in P_{G_1}\} \\ \text{en}_L^b(z \rightarrow az^p z^n, d, e) &= \text{class } e_z[\kappa] \{ \\ &\quad \text{def } e_a(i_1 : d_{z_1^p}, \dots, i_p : d_{z_p^p}) : \\ &\quad \quad \kappa_n^n = k_n(B_n^{n,p}) \} \end{aligned} \quad (\text{d})$$

$$\begin{aligned} \text{en}_S^b(a_1^{i^{k_1}} \dots a_n^{i^{k_n}}, G, d, e) &= \text{new term()}s_1 \dots s_n : \perp \\ &\quad \text{where} \\ s_i &= e_{a_i}(\text{en}_S^d(i_1^{k_i}, d), \\ &\quad \dots, \text{en}_S^d(i_{k_i}^{k_i}, d)) \\ \text{def term()} &= \text{new } e_{z_G}[\perp, d_{z_G}]( \\ &\quad (t : d_{z_G}) \Rightarrow \text{null}, \\ &\quad (s : d_{z_G}) \Rightarrow s \\ &\quad ) \end{aligned} \quad (\text{e})$$

$$\begin{aligned} \text{en}_L^B(G, d, e) &= \text{en}_L^d(G, d, e) \cup \text{en}_L^b(G, d, e) \\ \text{en}_S^B(a^n, G, d, e) &= \text{en}_S^b(a^n, G, d, e) \end{aligned} \quad (\text{f})$$

**Figure 7.** Implementation Rules

use that to guarantee well-typedness of the program. The key is to use the type checker of the host language to check the type of the object language.

### 4.1 An Example

First, let's look at an example.

**Example 4.1.** The simply typed lambda calculus with de Bruijn indices, booleans, and integers:

$$\begin{aligned}
e &\rightarrow \text{var}(i)|\text{con}(n|b)|\text{app } e_1 e_2|\text{abs } e \\
i &= z|s \ i \\
n &= \text{integers} \\
b &= \text{booleans} \\
t &= N|B|\text{fun } t_1 t_2 \\
E &= \text{nil}|\{t, E\}
\end{aligned}$$

where  $e$  are terms,  $t$  are types, and  $E$  are typing environments. And the following typing rules:

$$\begin{aligned}
\frac{}{E \vdash \text{con}(n) : N} \quad (\text{TNAT}) \quad & \frac{}{E \vdash \text{con}(b) : B} \quad (\text{TBOOL}) \\
\frac{E \vdash \text{var}(i) : t_1}{\{t_2, E\} \vdash \text{var}(s \ i) : t_1} \quad (\text{TSUCC}) & \\
\frac{}{\{t, E\} \vdash \text{var } x(z) : t} \quad (\text{TZERO}) & \\
\frac{E \vdash e_1 : \text{Fun}(t_1, t_2) \quad E \vdash e_2 : t_1}{E \vdash \text{app } e_1 e_2 : t_2} \quad (\text{TAPP}) & \\
\frac{\{t_1, E\} \vdash e : t_2}{E \vdash \text{abs}(t_1) e : \text{fun}(t_1, t_2)} \quad (\text{TABS}) &
\end{aligned}$$

where  $t, t_1, t_2$  are type metavariables.

As in the previous example, we look at how to embed both the grammar and the typing as case classes. In the previous example, we used an embedding  $\text{en}^d$  to map terms of the object language to terms of the host language. Here, we also need a mapping that maps types of the object language to the host language. A mapping function  $\varepsilon$  maps the types of the object language to Scala.

$$\begin{aligned}
\varepsilon_t(N) &= \text{Int} \\
\varepsilon_t(B) &= \text{Boolean} \\
\varepsilon_t(\text{fun } t_1 t_2) &= \text{Fun}[\varepsilon_t(t_1), \varepsilon_t(t_2)]
\end{aligned}$$

Here we use the metavariable  $t_1, t_2$  to disambiguate with type metavariables  $t_1, t_2$ . For booleans and integers, we use Scala's types, and for function type, we defined a new type constructor to disambiguate with the function type in Scala.

Let us first focus on the rules (TNAT), (TBOOL), (TAPP), and partial (TABS) without the typing environment. We really can not claim well-typedness without the rest of the rules, but we will consider those rules later. We represent types and terms of the object language using case classes as shown in Figure 8.

Next we implement the embedding classes. Since the grammar is already in GNF, we can encode it directly in Scala, this time parametrize  $e_e$  also by  $\chi$ , the type of the term that is currently being constructed. In a first attempt, we tried to implement it with a builder function as

```

class E[κ, σ, χ](b : Term[χ] ⇒ σ, k : σ ⇒ κ){
  ...
  def abs[χ1, χ2]() : E[κ, σ, χ2] =
    new E((t : Term[χ2]) ⇒ b(Abs[χ1, χ2](t)), k)
}

```

But this program does not compile, because in the definition of function  $\text{abs}$  the builder  $b$  expects some value of type  $\text{Term}[\chi]$ , but what we supplied is of type  $\text{Term}[\text{Fun}[\chi_1, \chi_2]]$ . To solve the type mismatch, we need to test the definitional equality of types, which

```

case class Fun[τ1, τ2]()

abstract case class Term[τ]()
case class Var(x : Ref[τ]) extends Term[τ]()
case class Con(x : τ) extends Term[τ]()
case class App[τ1, τ2](t1 : Term[Fun[τ1, τ2]], t2 : Term[τ1])
  extends Term[τ2]()
case class Abs[τ1, τ2](t : Term[τ1])
  extends Term[Fun[τ1, τ2]()]

abstract case class Ref[τ]()
case class Zero[τ]() extends Ref[τ]()
case class Succ[τ](r : Ref[τ]) extends Ref[τ]()

```

**Figure 8.** Data structure for Implementation with a Builder Function and Typing, Partial Version

```

class E[κ, σ, χ](b : Term[χ] ⇒ σ, k : σ ⇒ κ){
  def var(x : Ref[χ]) : κ = k(b(Var[χ](x)))
  def con(c : χ) : κ = k(b(Con[χ](c)))
  def app[χ1]() :
    E[E[κ, σ, χ1], Term[χ1] ⇒ σ, Fun[χ1, χ]] =
    new E(
      (t1 : Term[Fun[χ1, χ2]]) ⇒
      (t2 : Term[χ1]) ⇒ b(App(t1, t2)),
      (b1 : Term ⇒ σ) ⇒ E(b1, k))
  def abs[χ1, χ2](cast : Term[χ1, χ2] ⇒ Term[χ]) :
    E[κ, σ, χ2] =
    new E((t : Term[χ2]) ⇒
      b(cast(Abs[χ1, χ2](t))), k)
}

```

**Figure 9.** Implementation with a Builder Function and Typing, Partial Version

can be done by passing in a cast function:  $\text{cast} : \text{Term}[\chi_1, \chi_2] \Rightarrow \text{Term}[\chi]$ . The modified code is shown in Figure 9.

We take a simple approach here – we assume that when the method  $\text{app}$  is called we pass the following generic function

$$\text{def } \text{id}[\tau](x : \tau) = x$$

and initiate the type parameter to the desired type of the abstraction. For example,  $\text{abs}(\text{id}[\text{Fun}[\text{Int}, \text{Int}]])$ . Note that if the type is incorrect, then Scala will report type errors, which means that the only "requirement" for the user is ensuring that other functions are not passed in as the cast function. We postpone the discussion of a more sophisticated way of automatically enforcing definitional equality of types that removes this "requirement" to Section 6.

We have now implemented (TNAT), (TBOOL), (TAPP), and partial (TABS). Next, we add typing environments to complete the implementation. First, we define a mapping for typing environments.

$$\begin{aligned}
\varepsilon_e(\text{nil}) &= \text{Nothing} \\
\varepsilon_e(\{t, E\}) &= (\varepsilon_t(t), \varepsilon_e(E))
\end{aligned}$$

We add a type parameter that represents the typing environment to all classes that represent terms, The  $\text{Ref}$  class and its subclasses represent the typing rules (TZERO), (TSUCC) as in [2]. Now the classes together guarantee that only well-formed and well-typed terms can be constructed. The modification to the classes are shown in Figure 10.

We made two changes here. First, we added a new type parameter  $\eta$ . Second, in the type of  $\text{abs}$ , we modified it to implement the modification of the environment in (TABS). Note that here we only

```

abstract case class Term[τ, η]()
case class Var[τ, η](x : Ref[τ, η]) extends Term[τ, η]()
case class Con[τ, η](x : τ) extends Term[τ, η]()
case class App[τ1, τ2, η](
  t1 : Term[Fun[τ1, τ2, η], η], t2 : Term[τ1, η])
  extends Term[τ2, η]()
case class Abs[τ1, τ2, η](t : Term[τ2, (τ1, η)])
  extends Term[Fun[τ1, τ2, η], η]()

```

```

abstract case class Ref[τ, η]()
case class Zero[τ, η]() extends Ref[τ, η]()
case class Succ[τ, η](r : Ref[τ, η])
  extends Ref[τ, (r, η)]()

```

```

class E[κ, σ, χ, η](b : Term[χ] ⇒ σ, k : σ ⇒ κ){
  def varx(i : Ref[τ, η]) : κ = k(b(Var[χ, η](i)))
  def con(c : χ) : κ = k(b(Con[χ, η](c)))
  def app[χ1, η]() :
    E[E[τ, σ, χ1, η], Term[χ1, η] ⇒ σ, Fun[χ1, χ], η] =
    new E(
      (t1 : Term[Fun[χ1, χ2, η]) ⇒
        (t2 : Term[χ1, η]) ⇒ b(App[χ1, χ, η](t1, t2)),
      (b1 : Term[χ, η] ⇒ σ) ⇒ E(b1, k))
  def abs[χ1, χ2, η](cast : Term[χ1, χ2, η] ⇒ Term[χ, η]) :
    E[κ, σ, χ2, (χ1, η)] =
    new E((t : Term[χ2, (χ1, η)] ⇒
      b(cast(Abs[χ1, χ2, η](t))), k)
}

```

**Figure 10.** Implementation with a Builder Function and Typing, Complete Version

need to change the environment in the continuation because, by the typing rule (TABS), the new environment is used only in the term immediately following *abs*.

We can encode terms in  $L(G)$  in Scala. To use this embedding, not only do we need to pass in an initial builder to the *E* class where we create an object of the class, but also specify the type parameters representing the initial typing environment and the type of the program. We define auxiliary methods.

```

def term[χ]() = new E[
  ⊥,
  Term[χ, Nothing],
  χ,
  Nothing](
  (t : Term[χ, Nothing]) ⇒ null,
  (s : Term[χ, Nothing]) ⇒ s
)
def id[τ, η](x : Term[τ, η]) = x

```

Using these methods, the following term

```
app abs varx(z) con(1)
```

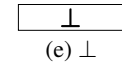
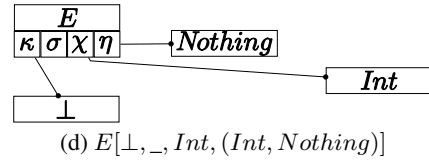
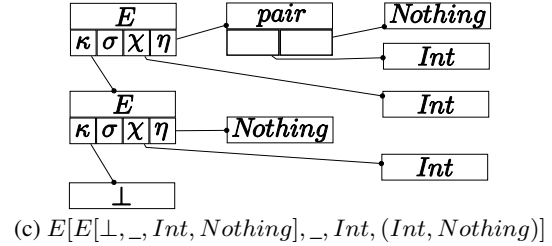
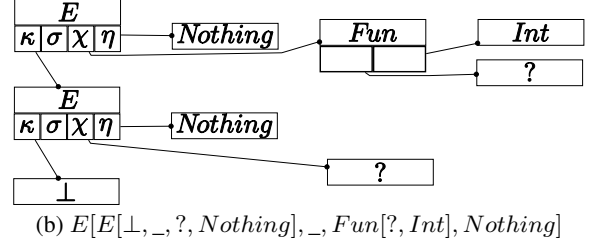
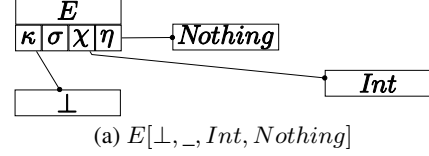
is encoded in Scala as

```
term[Int]() app()
abs(id[Fun[Int, Int], Nothing]) varx(Zero()) con(1) : ⊥
```

The types of the partial method chains are shown in Figure 11. The type of initial *term()* is shown in Figure 11(a); The type of *term()* *app()* is shown in Figure 11(b); and so forth. We do not care about the type  $\sigma$ , therefore, we put “ $\perp$ ”. The question marks stand for type variables. The type of the whole method chain is shown in Figure 11(e).

As a more complex example, the following term

```
app app abs abs con(2) app abs varx(z) con(1) con(1)
```



**Figure 11.** Examples

is encoded in Scala as

```

term[Int]()
app()
app()
abs(id[Fun[Int, Fun[Int, Int]], Nothing])
abs(id[Fun[Int, Int], (Int, Nothing)])
con(2)
app()
abs(id[Fun[Int, Int], Nothing])
varx(Zero())
con(1) : ⊥

```

## 4.2 Embedding Typing

Before we proceed to define the general metafunction  $\text{en}^T$ , we would like to summarize what we have learned from the example. First of all, we must take typing rules into account when defining this metafunction. It is very difficult to define a general embedding function for all type systems, because it is very difficult to define a small language that is expressive enough for all type systems. However, if we focus on languages built around simply typed lambda calculus, then there are a few assumptions that we can make. The most important assumption is that

**Assumption 4.2.** *There is exactly one typing rule for each production in the grammar.*



Each typing rule

$$\frac{E_{s_1} \vdash e_{s_1} : t_{s_1} \quad \dots \quad E_{s_m} \vdash e_{s_m} : t_{s_m}}{E \vdash ae_1 \dots e_n : t} \quad (\mathcal{TR}^*)$$

can be written as

$$\frac{E_{a,s_1}(E, t) \vdash e_{s_1} : T_{a,s_1}(t) \quad \dots \quad E_{a,s_m}(e, t) \vdash e_{s_m} : T_{a,s_m}(t)}{E \vdash ae_1 \dots e_n : t} \quad (\mathcal{TR})$$

where  $\{s_1, \dots, s_m\} \subset \{1, \dots, n\}$ . For each typing rule  $\mathcal{TR}$ , we define two groups of mappings. The first group is  $T_{a,s_i}$ , which maps the type  $t$  of  $e_1 \dots e_n$  to the types of  $e_{s_1}, \dots, e_{s_m}$ . Another is  $E_{a,s_i}$  which maps the environment  $E$  and type  $t$  of  $e_1 \dots e_n$  to the environments of  $e_{s_1}, \dots, e_{s_m}$ . The second assumption is that

**Assumption 4.3.** *Both the group  $T_{a,s_i}$  and the group  $E_{a,s_i}$  are definable with either construction or pattern matching on the parameters.*

It turns out that many syntactic sugars fall into this framework.

**Example 4.4.** For example, consider extending our language with the rules for `if-then-else` expressions and pairs.

$$\frac{E \vdash e_1 : B \quad E \vdash e_2 : t \quad E \vdash e_3 : t}{E \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 : t} \quad (\text{TIFTHENELSE})$$

$$\frac{E \vdash e_1 : t_1 \quad E \vdash e_2 : t_2}{E \vdash \text{pair } e_1 e_2 : (t_1, t_2)} \quad (\text{TPAIR})$$

$$\frac{E \vdash e : (t_1, t_2)}{E \vdash \text{fst } e : t_1} \quad (\text{TFST}) \quad \frac{E \vdash e : (t_1, t_2)}{E \vdash \text{snd } e : t_2} \quad (\text{TSND})$$

The (TIFTHENELSE) rules can be written as

$$\frac{E_{if,1}(E, t) \vdash e_1 : T_{if,1}(t) \quad E_{if,3}(E, t) \vdash e_2 : T_{if,3}(t) \quad E_{if,5}(E, t) \vdash e_3 : T_{if,5}(t)}{E \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 : t} \quad (\text{TIFTHENELSE})$$

where

$$\begin{aligned} E_{if,1} | E_{if,3} | E_{if,5}(E, t) &= E \\ T_{if,1}(t) &= B \\ T_{if,3} | T_{if,5}(t) &= t \end{aligned}$$

The (TPAIR) rule can be written as

$$\frac{E_{pair,1}(E, (t_1, t_2)) \vdash e_1 : T_{pair,1}(t) \quad E_{pair,2}(E, (t_1, t_2)) \vdash e_2 : T_{pair,2}(t)}{E \vdash \text{pair } e_1 e_2 : (t_1, t_2)} \quad (\text{TPAIR})$$

where

$$\begin{aligned} E_{pair,1} | E_{pair,2}(E, t) &= E \\ T_{pair,1}((t, \_)) &= t \\ T_{pair,2}(\_, t) &= t \end{aligned}$$

The (TABS) rule can be written as

$$\frac{E_{abs,1}(E, \text{fun}(t_1, t_2)) \vdash e : T_{abs,1}(\text{fun}(t_1, t_2))}{E \vdash \text{abse} : \text{fun}(t_1, t_2)} \quad (\text{TABS})$$

where

$$\begin{aligned} E_{abs,1}(E, \text{fun}(t, \_)) &= \{t, E\} \\ T_{abs,1}(\text{fun}(\_, t)) &= t \end{aligned}$$

The (TAPP) rule can be written as

$$\frac{E_{app,1}(E, t) \vdash e_1 : T_{app,1}(t) \quad E_{app,2}(E, t) \vdash e_2 : T_{app,2}(t)}{E \vdash e_1 e_2 : t} \quad (\text{TAPP})$$

where

$$\begin{aligned} E_{app,1} | E_{app,2}(E, t) &= E \\ T_{app,1}(t) &= \text{fun}(t_1, t) \\ T_{app,2}(t) &= t_1 \end{aligned}$$

where  $t_1$  is a fresh type metavariable.

First, if type string matching or environment string pattern matching is required in the consequence of a typing rule for  $z \rightarrow az_1 \dots z_n$  (which means that the type of the consequence is not a type metavariable), such as in (TABS) and (TPAIR), then we need to pass in a proof of definitional equality by a *cast* function.

Second, if type string pattern matching or environment string pattern matching is required in the antecedent of a typing rule for  $z \rightarrow az_1 \dots z_n$  (which means that construction is required in  $T_{a,k}$  or  $E_{a,k}$ , for some  $k \in \{1, \dots, n\}$ ), then we need to construct a new type for parameter  $\chi$  or  $\eta$ , which represents the type and the environment, respectively, either for type constructors that are pushed onto the continuation stack in the return type of  $e_a$  or the types of the parameters of  $e_a$ .

Third, all type variables in the rules other than that represented by  $\chi$  become the type variables of the method definitions such as the *app* and *abs* method.

By viewing the grammar as parametrized, we actually embedded part of the grammar not in the method chaining style but in the functional nesting style. Therefore, it is possible to use some of the existing techniques to avoid requiring the *cast* function, such as in (TSUC). Conversely, because an initial object with necessary information is needed to initialize a method chain, the method chaining style would induce unnecessary overhead and clog to the code if used within the functional nesting style. This implies that we can use method chaining style to embed a language while using functional nesting style to deal with variables.

One potential problem that scrutinous readers may have noticed is that in Example 4.1, the typing rules for constants are only partially represented because the *con* method does not receive a cast function as an argument. As a result, the universe of types may be expanded. In general, the cast function is necessary. However, this is not a problem in our example since all types of constants are either specified in the type argument of *term*, or by the type argument of cast functions for *abs*. Even if we add a cast function as a parameter of the *con* method, we still rely on the user to provide a type argument within the defined types for the cast function. We call these types that are reused in the object language "native" types. We need to deal with these types carefully, when generalizing our method exemplified.

### 4.3 The General Embedding Function for Typed Languages

Now, let's define the embedding  $\text{en}^\top$ .  $\text{en}^\top$  is defined by a quadruple of embedding functions ( $\text{en}^t, \text{en}^e, \text{en}^d, \text{en}^b$ ). These embedding functions are parametrized over  $G, T, E, R, e$ , and  $d$ .  $G$  is a parametrized grammar of the object language in LL(1) GNF.  $T$  is a grammar for types of the object language in LL(1) GNF such that  $Z_T = \{z_T\}$  (which may be a subgrammar of  $G$ ).  $E$  is a grammar for typing environments of the object language in LL(1) GNF such that  $Z_E = \{z_E\}$ .  $R$  is a function that maps each production in  $G$  to a typing rule of the form  $\mathcal{TR}$ , so that there is exactly one typing rule for each production in  $G$ . The name embedding function  $e$  maps terminals and nonterminals in  $G$  to unique Scala names, and the name embedding function  $d$  maps terminals and nonterminals of

$$\begin{aligned} \text{en}_L^t(\mathbb{T}, d) &= \{\text{en}_L^t(p, d) \mid p \in P_T\} \\ \text{en}_L^t(z \rightarrow az^n, d) &= \text{case class } d_a[v_1, \dots, v_n] \end{aligned} \quad (a)$$

$$\begin{aligned} \text{en}_S^t(a^n, d) &= \text{en}_S^{t'}(\text{tree}(a^n), d) \\ \text{en}_S^{t'}(z(\text{at}^n), d) &= d_a[\text{en}_S^{t'}(t_1^n, d), \dots, \text{en}_S^{t'}(t_n^n, d)] \end{aligned} \quad (b)$$

**Figure 12.** Embedding of Types in the Object Language

$$\begin{aligned} \text{en}_L^e(\mathbb{T}, d) &= \{\text{en}_L^e(p, d) \mid p \in P_E\} \\ \text{en}_L^e(z \rightarrow az^n, d) &= \text{case class } d_a[v_1, \dots, v_n] \end{aligned} \quad (a)$$

$$\begin{aligned} \text{en}_S^e(a^n, d) &= \text{en}_S^{e'}(\text{tree}(a^n), d) \\ \text{en}_S^{e'}(z(\text{at}^n), d) &= d_a[\text{en}_S^{e'}(t_1^n, d), \dots, \text{en}_S^{e'}(t_n^n, d)] \end{aligned} \quad (b)$$

**Figure 13.** Embedding of Typing Environments in the Object Language

$G$ ,  $T$ , and  $E$ , and metavariables to unique Scala names such that  $d_z \neq e_z$  for all  $z \in Z_G$ .

We write the embedding of metavariable  $t$  as  $d_t$ .

**Definition 4.5.** A extension  $\tilde{G}$  of grammar  $G$  with metavariables  $t_1, t_2, \dots$  is  $(Z_G, z_G, A_G \cup \{t_1, t_2, \dots\}, P_G \cup \{z \rightarrow t_k \mid z \in Z_G, k \in \{1, 2, \dots\}\})$ .

We start by defining the embedding  $\text{en}^t$ , which is shown in Figure 12, where  $v_1, \dots, v_n$  are type parameters,  $a^n$  is a string in  $\tilde{T}$ , and  $t^n$  are subtrees. Each terminal in  $\mathbb{T}$  is mapped to a type constructor. This embedding generates the Scala classes that represent types in the object language. This embedding actually allows any type of the host language to be a type of the object language. We do not need to restrict those types here because possible types are restricted by the typing rules  $R$ . One advantage of this is that we avoid casting for many types that we can reuse from the host language.

The string embedding  $\text{en}_S^t$  is defined on  $\tilde{T}$  which maps type metavariables  $t$  to  $d_t$ . For example,

$$\begin{aligned} \text{en}_S^t(N, d) &= \text{Int} \\ \text{en}_S^t(B, d) &= \text{Boolean} \\ \text{en}_S^t(\text{fun } t_1 \ t_2, d) &= \text{Fun}[\text{en}_S^t(t_1, d), \text{en}_S^t(t_2, d)] \\ \text{en}_S^t(t, d) &= d_t \end{aligned}$$

The embedding  $\text{en}^e$  of  $E$  is shown in Figure 13, where  $v_1, \dots, v_n$  are type parameters,  $a^n$  is a string in  $\tilde{E}$ , and  $t^n$  are subtrees. Because  $E$  usually contains  $T$ , the embedding of  $E$  may share the same classes as the embedding of  $T$ . For example, if the environment is  $\{t_1, \dots, \{t_n, \text{nil}\}\}$ , then  $\text{en}_S^e(\{t_1, \dots, \{t_n, \text{nil}\}\}, d) = d_{\{[\text{en}_S^t(t_1, d), d_{\{[\text{en}_S^t(t_{n-1}, d), d_{\text{nil}}]\}}]$ .

The typing rules in  $R$  are encoded by both  $\text{en}^d$  and  $\text{en}^b$ .

As the name suggests,  $\text{en}^d$  is an extension of the embedding  $\text{en}^d$  defined in Section 3 for productions in  $G$ . Each embedding class is parametrized over type parameters  $\chi, \eta$  which represent types and typing environments of the object language, respectively. The  $\text{fv}$  metavariables are replaced by lists of free type variables that occur

$$\begin{aligned} \text{nat}(a, b, c) &= \begin{cases} b & b \text{ native} \\ a[b, c] & \text{otherwise} \end{cases} \\ \varphi &\triangleq \text{nat}(d_z, \text{en}_S^t(t, d), \text{en}_S^e(E, d)) \\ \varphi_i^n &\triangleq \text{nat}(d_{z_i^n}, \text{en}_S^t(t_i, d), \text{en}_S^e(E_i, d)) \end{aligned} \quad (a)$$

$$\begin{aligned} \text{en}_L^d(G, R, d) &= \{\text{en}_L^d(p) \mid p \in P_G\} \cup \{\text{en}_L^d(z) \mid z \in Z_G\} \\ \text{en}_L^d(z \rightarrow az^n, R, d) &= \text{case class } d_a[\text{fv}( \\ &\quad i_1 : \varphi_1^n, \dots, i_n : \varphi_n^n) \\ &\quad \text{extends } \varphi() \\ &\quad \text{where } R(z \rightarrow az^n) = \\ &\quad \frac{E_1 \vdash e_1 : t_1 \quad \dots \quad E_n \vdash e_n : t_n}{E \vdash ae_1 \dots e_n : t} \\ \text{en}_L^d(z, d) &= \text{abstract case class } d_z[\chi, \eta]() \end{aligned} \quad (b)$$

$$\begin{aligned} \text{en}_S^d(a^n, d) &= \text{en}_S^{d'}(\text{tree}(a^n), d) \\ \text{en}_S^{d'}(z(\text{at}^n), d) &= \begin{cases} d_a & d_a \text{ native} \\ d_a(\text{en}_S^{d'}(t_1^n, d), \dots, \text{en}_S^{d'}(t_n^n, d)) & \text{otherwise} \end{cases} \end{aligned} \quad (c)$$

**Figure 14.** Embedding of Terms in the Object Language

in the constructor definition. Each constructor represents a typing rule. The techniques for defining these data structures are typing rule directed, as in [2]. Using auxiliary definitions in Figure 14(a), where "native" means that the type is a predefined Scala type such as *Int*, or the name is a name for constants such as integers, the embedding is defined in Figure 14(b) and Figure 14(c), where  $a^n$  is a string in  $G$ , and  $t^n$  are subtrees. Note that the string embedding is exactly the same as in Section 3, because Scala automatically infers the type parameters.

Next, we define  $\text{en}^b$ . The rules are shown in Figure 15. We redefined the auxiliary symbols from Section 3, adding type variables that represent types and typing environments to  $\kappa_k^n, B_k^{n,p}$  for  $k, n, p \in \mathbb{N}$ . Also, in the definition of  $B_0^{n,p}$ , we apply the *cast* function after applying the builder.

Each method defined in class  $e_z$  corresponds to a production, and its typing rule. The class  $e_z$  is parametrized over  $\kappa, \sigma, \chi, \eta$ , which represents the continuation stack, the type of the builder, the type of the expected string that is generated by  $z$ , and the typing environment, respectively. In the embedding, a *cast* :  $\varphi \Rightarrow d_z[\text{tvt}, \text{tve}]$  function is passed in as a parameter for each method. In the typing rule, the metavariables  $t$  and  $t_k$  for  $k \in \{1, \dots, n\}$  are replaced with the actual types and similarly for  $E$  and  $E_k$  for  $k \in \{1, \dots, n\}$ . If the actual type is a type metavariable  $t$  and the actual environment is a metavariable  $E$ , then  $t = t, E = E$  and neither  $T_{a_i, k}$  nor  $E_{a_i, k}$  require pattern matching, in which case we can ignore the *cast* function because it has a trivial type  $d_z[d_t, d_e] \Rightarrow d_z[d_t, d_e]$ , and can be removed in a code generator<sup>3</sup>; but we still list it in the parameters for uniformity. Similarly for  $\text{en}_S^b$  defined in Figure 16, the code generator can remove the *id* argument for trivial cast function parameters to generate more succinct code.

<sup>3</sup>EriLex does this by testing the definitional equality when generating code.

$$\begin{aligned}
p &= \text{arity}(z \rightarrow az^p z^n) \\
\text{nat}(a, b, c) &= \begin{cases} b & \text{b native} \\ a[b, c] & \text{otherwise} \end{cases} \\
\varphi &\triangleq \text{nat}(d_z, \text{en}_S^t(t, d), \text{en}_S^e(E, d)) \\
\varphi_i^n &\triangleq \text{nat}(d_{z_i^n}, \text{en}_S^t(t_i^n, d), \text{en}_S^e(E_i^n, d)) \\
\varphi_i^p &\triangleq \text{nat}(d_{z_i^p}, \text{en}_S^t(t_i^p, d), \text{en}_S^e(E_i^p, d)) \\
\sigma_0^n &\triangleq \sigma \\
\sigma_{i+1}^n &\triangleq \varphi_{n-i}^n \Rightarrow \sigma_i^n \\
\kappa_0^n &\triangleq \kappa \\
\kappa_{i+1}^n &\triangleq e_{z_{n-i}^n}[\kappa_i^n, \sigma_i^n, \text{en}_S^t(t_{n-i}^n, d), \text{en}_S^e(E_{n-i}^n, d)] \\
\kappa_0^n &\triangleq \kappa \\
\kappa_{i+1}^n &\triangleq (b_{i+1} : \sigma_{i+1}^n) \Rightarrow \text{new } e_{z_{n-i}^n}(b_{i+1}, \kappa_i^n) \\
B_0^{n,p} &\triangleq b(\text{cast}(d_a(i_1, \dots, i_p, t_1, \dots, t_n))) \\
B_{i+1}^{n,p} &\triangleq (t_{n-i} : \varphi_{n-i}^n) \Rightarrow B_i^{n,p}
\end{aligned} \tag{a}$$

$$\begin{aligned}
\text{en}_L^b(G, R, d, e) &= \{\text{en}_L^b(p, R, d, e) \mid p \in P_{G_1}\} \\
\text{en}_L^b(z \rightarrow az^p z^n, R, d, e) &= \text{class } e_z[\kappa, \sigma, \text{tvt}, \text{tve}] ( \\
&\quad b : d_z[\text{tvt}, \text{tve}] \Rightarrow \sigma, \\
&\quad k : \sigma \Rightarrow \kappa) \{\text{en}_L^b(az^n, d, e)\} \\
&\quad \text{where } R(z \rightarrow az^n) = \\
&\quad \frac{E_1^p \vdash e_1^p : t_1^p \quad \dots \quad E_n^n \vdash e_n^n : t_n^n}{E \vdash ae^p e^n : t} \\
\text{en}_L^b(az^p z^n, d, e) &= \text{def } e_a[\text{fv}] ( \\
&\quad i_1 : \varphi_1^p, \dots, i_p : \varphi_p^p, \\
&\quad \text{cast} : \varphi \Rightarrow d_z[\text{tvt}, \text{tve}] : \\
&\quad \kappa_n^n = \kappa_n^n(B_n^{n,p}) \\
&\quad \text{where} \\
\text{tvt} &= \begin{cases} \chi & \text{t is not a metavariable} \\ d_t & \text{otherwise} \end{cases} \\
\text{tve} &= \begin{cases} \eta & E \text{ is not a metavariable} \\ d_E & \text{otherwise} \end{cases} \\
\text{fv} &= \text{free type variables in the method definition}
\end{aligned} \tag{b}$$

**Figure 15.** Embedding of Terms in MCS with a Builder Function and Typing

The source embedding of programs in the object language defined using the following auxiliary method definitions:  $\text{term}()$  and  $\text{id}()$ , is shown in Figure 16.

The embedding in  $\text{en}_S^b$  for a program

$$a_1 i^{k_1} \dots a_n i^{k_n}$$

where  $a_k$  are top level symbols for  $k \in \{1, \dots, n\}$  and  $i^{k_i}$  are second level strings for  $i \in \{1, \dots, n\}$ , is defined in Figure 16, where  $t_k^n, E_k^n$  are the types and typing environments of the substrings of program at  $a_k$ , for  $k \in \{1, \dots, n\}$ .

The embedding of the object language is defined in Figure 17.

$$\begin{aligned}
\text{en}_S^b(a_1 i^{k_1} \dots a_n i^{k_n}, \\
G, R, d, e, t^n, E^n) &= \text{new term}() s_1 \dots s_n : \perp \\
&\quad \text{where} \\
s_i &= e_{a_i}(\text{en}_S^d(i_1^{k_i}, d), \dots, \text{en}_S^d(i_{k_i}^{k_i}, d) \\
&\quad \quad \text{id}[\text{en}_S^t(t_i^n, d), \text{en}_S^e(E_i^n, d)]) \\
\varphi_G &\triangleq d_{z_G}[\text{en}_S^t(t_1^n, d), \text{en}_S^e(E_1^n, d)] \\
\text{def term}() &= \text{new } e_{z_G}[\perp, \varphi_G, \text{en}_S^t(t_1^n, d), \text{en}_S^e(E_1^n, d)] ( \\
&\quad (t : \varphi_G) \Rightarrow \text{null}, \\
&\quad (s : \varphi_G) \Rightarrow s \\
&\quad ) \\
\text{def id}[\tau, \eta] &= (x : \text{Term}[\tau, \eta]) \Rightarrow x
\end{aligned}$$

**Figure 16.** String Embedding

$$\begin{aligned}
\text{en}_L^T(T, E, G, R, d, e) &= \text{en}_L^t(T, d) \cup \text{en}_L^e(E, d) \\
&\quad \cup \text{en}_L^d(G, R, d) \cup \text{en}_L^b(G, R, d, e) \\
\text{en}_S^T(a^n, G, d, e, t^n, E^n) &= \text{en}_S^b(a^n, G, d, e, t^n, E^n)
\end{aligned}$$

**Figure 17.** The Embedding  $\text{en}_L^T$

We need the following theorem to complete our discussion.

**Theorem 4.6.** *If  $G, T$ , and  $E$  are grammars of the terms, types, and typing environments of the object language in  $LL(1)$  GNF,  $R$  is a function that maps each production in  $G$  to a typing rule of the form  $T\mathcal{R}$ ,  $e$  maps  $A_G \cup Z_G$  to unique Scala names,  $d$  maps  $A_T \cup Z_T \cup A_E \cup Z_E \cup A_G \cup Z_G$  to unique Scala names such that  $d_z \neq e_z$  for all  $z \in Z_G, s = a_1 i^{k_1} \dots a_n i^{k_n} \in A_G^*, t^n$  are strings of  $T, E^n$  are strings of  $E$ , then  $s \in L(G)$  and  $E^n \vdash s_k : t_k^n$ , where  $s_k$  is the substring of  $s$  at  $a_k$ , for all  $k \in \{1, \dots, n\}$  if and only if  $\text{en}_L^T(T, E, G, R, d, e) \cup \{\text{en}_S^T(s, G, d, e, t^n, E^n)\}$  is well-typed in Scala.*

*Proof.* (Sketch) First prove that any encoded term is automatically well-formed and well-typed because the constructor  $d_a$  only encodes well-formed and well-typed terms. Next, by type safety, the encoding always constructs encoding of some well-formed and well-typed term. Finally, we prove that the term constructed is the same term as that is encoded in the method chain.  $\square$

It is clear that with the addition of any nontrivial typing (i.e. everything does not have the same type), our embedding generally can no longer be modeled by an  $\text{slrpd}$  with finite stack symbols and input symbols. However, if we allow the  $\text{slrpd}$  to have countably infinite number of stack symbols and input symbols, our embedding can be easily modeled, by viewing types constructed from embedding classes as just a complex stack symbol, thus providing an upper bound of the expressiveness of the embedding framework.

## 5. Related Work

The problems of embedding a typed DSL into a host language in the FNS has been studied extensively, see [6] for a list of references. However, MCS embedding seems to be largely ignored, probably because most of the research are done in a functional setting, and embedding MCS in a functional programming language without object-oriented features is awkward. As the object-oriented programming languages become more and more mature in their type

systems, MCS may be a viable alternative to FNS since it seems to be very natural in object-oriented programming languages.

Scala [1, 4, 9, 12] is well known for its combination of flexible syntax, functional programming features, and object-oriented programming features, and consequently the ease to define DSLs. The Scala Parser Combinators is a parser library that embeds a BNF-like parser DSLs in the Scala programming language in a very succinct way. Combinator based parser libraries are also widely implemented and used [5]. We would like to point out that the monadic parsers provide a very concise way of embedding parser DSLs. However, as pointed out in [3], there is no easy way to transform the resulting parser to make optimizations. In our approach, the builder function are used to build a typed abstract syntax tree so that the techniques introduced in [3] can be applied. Also, syntactically, which is probably less essential, in combinator based parser libraries, the sequencing combinator (which usually does not appear in BNF) usually can not be omitted unless there are special language support in the host language (such as the `do` notation in Haskell), whereas in our approach as shown in Section 2, the sequencing operator is not needed, which results in a more succinct embedding. Also, it is not clear in an object-oriented programming language how the advantages of parser combinator would translate to embedding DSLs other than parser DSLs, such as simply typed lambda calculus based languages.

Two case studies of type safe DSLs and many practical aspects of implementing embedded DSL are discussed in [7], in which type safety is dealt with in a case by case manner. Typing environments are partially encoded in class names which results in that the number of classes required is linear to the maximum size of the typing environment allowed. In our approach, this problem is solved based on the embedding from [2] and using the two groups of functions  $E$ ,  $T$  defined in Section 4.2. In this respect, our embedding can be viewed as an extension of embedding from [2].

## 6. Discussion

It may seem that the difference between FNS and MCS are purely syntactical because one can be encoded in another, but we believe that MCS is a very important tool for object-oriented programming languages. On one hand, the syntax is an important aspect of embedding a DSL, for the expressiveness of an embedded language can not exceed that of the host language. On the other hand, if encoding is concerned only, then most common features in object-oriented programming languages can be encoded in a functional programming language as demonstrated in [10]. Conversely, functional programming can also be encoded in an object-oriented programming language with parametric types and closure. But alternatives are still useful because some encoding may make the code "bulky" at best, and inefficient at worst. After all every programming language is more polished at the task it is intended for. For mainstream object-oriented programming languages, it is creating objects and invoking methods, but not dealing with higher order functions.

Our approach may not be the most efficient because the builder function in the embedding actually builds an AST. However, because the constructors  $d_s$  are not required to be constructors, they can be any functions with the same type. Therefore, they can be easily interchanged with functions that actually act as the interpreter, thus improving the performance. Practically, if we want to obtain the AST for transformation, we can simply change the initial continuation  $k$  to  $(t : Term) \Rightarrow t$  to make it return the tree built instead of `null`. Sometimes if the grammar is complex, the resulting AST of GNF may not be the desired AST. For example in Example 2.10, we may discard the nodes generated by  $r'$  and move its subnodes one level up, as it is apparent that only nodes generated by  $r$  are needed in a more compact form of the grammar.

To ensure that the cast function is always the identity function, we may use the technique in [2], but this requires higher-rank polymorphism of kind  $* \rightarrow *$  (i.e. universally quantified type constructor in the domain type of a function), which is not directly supported but can be encoded in Scala. For practical reasons we choose not to use this encoding in this framework. However, if direct support for higher-rank polymorphism of higher kind is implemented in Scala, then we will definitely use that because that seems to be a more elegant solution.

Our framework actually allows the object language to have not just "terms", but multiple classes of programming constructs, which is useful for embedding practical DSLs such as SQL in which different classes of programming constructs are not interchangeable.

To conclude, we demonstrated how to embed DSLs in Scala in this paper in a general framework. We tried to keep to the object-oriented programming paradigm as close as possible, so that our framework can be used in a wide range of object-oriented programming languages and are easily assessable to the majority of programmers. In fact, most of the framework can be directly translated to Java, albeit in a more verbose manner, using the classical inner class encoding. There are many other features in Scala that we have not used, such as abstract types, existential types, and explicit bounded polymorphism. Exploring these features for embedding DSLs is part of future work for seeking more elegant solutions.

## References

- [1] K.U.Leuven Adriaan Moors, Frank Piessens and Martin Odersky. Safe type-level abstraction in scala. In *FOOL '08: International Workshop on Foundations of Object-Oriented Languages*, 2008.
- [2] Arthur I. Baars and S. Doaitse Swierstra. Typing dynamic typing. In *ICFP '02: Proceedings of the seventh ACM SIGPLAN international conference on Functional programming*, pages 157–166, New York, NY, USA, 2002. ACM.
- [3] Arthur I. Baars, S. Doaitse Swierstra, and Marcos Viera. Typed transformations of typed abstract syntax. In *TLDI '09: Proceedings of the 4th international workshop on Types in language design and implementation*, pages 15–26, New York, NY, USA, 2009. ACM.
- [4] Vincent Cremet, Francois Garillot, Serguei Lenglet, Martin Odersky, Ecole Normale Supérieure, and Ecole Normale Supérieure De Lyon. A core calculus for scala type checking. In *Proceedings of the 31st International Symposium on Mathematical Foundations of Computer Science, Springer LNCS*, pages 1–23. Springer, 2006.
- [5] Graham Hutton and Erik Meijer. Monadic Parsing in Haskell. *Journal of Functional Programming*, 8(4):437–444, July 1998.
- [6] Oleg Kiselyov Jacques Carette and Chung-Chieh Shan. Finally tagless, partially evaluated: Tagless staged interpreters for simpler typed languages. *Journal of Functional Programming*, 19:509–543, September 2009.
- [7] Jevgeni Kabanov and Rein Raudjärv. Embedded typesafe domain specific languages for java. In *PPPJ '08: Proceedings of the 6th international symposium on Principles and practice of programming in Java*, pages 189–197, New York, NY, USA, 2008. ACM.
- [8] Fabrice Marguerie, Steve Eichert, and Jim Wooley. *LINQ in Action*. February 2008.
- [9] Adriaan Moors, Frank Piessens, and Martin Odersky. Generics of a higher kind. In *OOPSLA '08: Proceedings of the 23rd ACM SIGPLAN conference on Object-oriented programming systems languages and applications*, pages 423–438, New York, NY, USA, 2008. ACM.
- [10] Benjamin C. Pierce. *Types and Programming Languages*. 2002.
- [11] Grzegorz Rozenberg and Arto Salomaa, editors. *Handbook of formal languages, vol. 1: word, language, grammar*. Springer-Verlag New York, Inc., New York, NY, USA, 1997.
- [12] Dean Wampler and Alex Payne. *Programming Scala*. 2009.