

XDuce: A Typed XML Processing Language (Preliminary Report)

Haruo Hosoya¹ and Benjamin Pierce²

¹ Department of CIS, University of Pennsylvania
hahosoya@cis.upenn.edu

² Department of CIS, University of Pennsylvania
bcpierce@cis.upenn.edu

1 Introduction

Among the reasons for the popularity of XML is the hope that the static typing provided by DTDs [XML] (or more sophisticated mechanisms such as XML-Schema [XS00]) will improve the safety of data exchange and processing. But, in order to make the best use of such typing mechanisms, we need to go beyond types for *documents* and exploit type information in static checking of *programs* for XML processing.

In this paper, we present a preliminary design for a statically typed programming language, *XDuce* (pronounced “transduce”). XDuce is a tree transformation language, similar in spirit to mainstream functional languages but specialized to the domain of XML processing. Its novel features are *regular expression types* and a corresponding mechanism for *regular expression pattern matching*. Regular expression types are a natural generalization of DTDs, describing, as DTDs do, structures in XML documents using regular expression operators (i.e., *, ?, |, etc.). Moreover, regular expression types support a simple but powerful notion of *subtyping*, yielding a substantial degree of flexibility in programming. Regular expression pattern matching is similar to ML pattern matching except that regular expression types can be embedded in patterns, which allows even more flexible matching.

In this preliminary report, we show by example the role of these features in writing robust and flexible programs for XML processing. After discussing the relationship of our work to other work, we briefly sketch some larger applications that we have written in XDuce, and close with remarks on future work. Appendices give a formal definition of the core language.

2 Programming in XDuce

We develop a series of examples of programming in XDuce, using regular expression types and regular expression pattern matching.

2.1 Regular Expression Types

Values and Types XDuce's values are XML documents. A XDuce program may read in an XML document as a value and write out a value as an XML document. Even values for intermediate results during the execution of the program have a one-to-one correspondance to XML documents (besides some trivial differences).

As concrete syntax, the user has two choices: XML syntax or XDuce's native syntax. We can either write the following XDuce value (we assign it to the variable `mybook` for later explanation)

```
val mybook = addrbook[
    name["Haruo Hosoya"],
    addr["Tokyo"],
    name["ABC"],
    addr["Def"],
    tel["123-456-789"],
    name["Benjamin Pierce"],
    addr["Philadelphia"]]
```

in the native syntax, or the following corresponding document in standard XML syntax:

```
<addrbook>
  <name>Haruo Hosoya</name>
  <addr>Tokyo</addr>
  <name>ABC</name>
  <addr>Def</addr>
  <tel>123-456-789</tel>
  <name>Benjamin Pierce</name>
  <addr>Philadelphia</addr>
</addrbook>
```

XDuce provides term constructors of the form `label[...]`, where `...` is a sequence of other values. This corresponds to `<label>...</label>` in XML notation. We enclose strings in double-quotes, unlike XML.

Observe the sequence contained in `addrbook`. It is natural to impose a structure on the seven children so that they can be regarded as three “entries,” each of which consists of fields tagged `name`, `addr` and optional `tel`. We can capture this structure by defining the following regular expression types.

```
type Addrbook = addrbook[(Name,Addr,Tel?)*]
type Name = name[String]
type Addr = addr[String]
type Tel = tel[String]
```

These XDuce definitions roughly correspond to the following DTD:

```
<!ELEMENT addrbook (name,addr,tel?)*>
<!ELEMENT name #PCDATA>
<!ELEMENT addr #PCDATA>
<!ELEMENT tel #PCDATA>
```

(Just as XDuce can read standard XML documents, we also provide a construct to import DTDs as regular expression types.) Type constructors `label[...]` have the same form as the term constructors that they classify. In addition, types may be formed using the regular expression operators `*` for repetition, `|` for alternation, and `?` for optional elements. (We will show examples of alternations later.) For instance, the type `(Name,Addr,Tel?)*` stands for zero or more repetitions of the sequence of a `Name`, an `Addr`, and an optional `Tel`.

The notion of *subtyping* will play a crucial role in the calculation that justifies assigning the type `Addrbook` to the value `mybook`.

Subtyping Before showing the subtyping relation, we need to clearly state this: the elements of every type in XDuce are *sequences*. For example, the type `Tel*` contains the following sequences.

<code>()</code>	the empty sequence
<code>tel["123"]</code>	sequence with one <code>Tel</code>
<code>tel["123"],tel["234"]</code>	sequence with two <code>Tel</code> 's
...	

In the type language, comma is the type constructor for concatenation of sequences. For example, the type `(Name,Tel*,Addr)` contains

```
name["abc"],addr["ABC"]
name["abc"],tel["123"],addr["ABC"]
name["abc"],tel["123"],tel["234"],addr["ABC"]
...
```

i.e., sequences with one `Name` value, followed by zero or more `Tel` values, then followed by one `Addr` value. The comma operator on types is associative: the types `((Name,Tel*),Addr)` and `(Name,(Tel*,Addr))` have exactly the same set of elements.

The *subtype* relation between two types is simply inclusion of the sets denoted by types (see Appendix A for the formal definition).

We now show the sequence of steps involved in verifying that `mybook` has type `Addrbook`. First, from the intuition that `?` means “optional,” we would expect the following relations:

$$\begin{aligned} \text{Name,Addr} &<: \text{Name,Addr,Tel?} \\ \text{Name,Addr,Tel} &<: \text{Name,Addr,Tel?} \end{aligned}$$

Notice that each right hand side has more possibilities than the left hand side. Similarly, `*` means “zero or more” intuitively, so in particular it could be three:

$$T,T,T <: T^*$$

Combining these relations, we obtain

$$(\text{Name,Addr}), (\text{Name,Addr,Tel}), (\text{Name,Addr}) <: (\text{Name,Addr,Tel?})^*$$

Since comma is associative, we can get rid of parentheses:

$$\begin{aligned} & \text{Name, Addr, Name, Addr, Tel, Name, Addr} \\ &= (\text{Name, Addr}), (\text{Name, Addr, Tel}), (\text{Name, Addr}) \end{aligned}$$

(Here, we mean by $T = U$ that both $T <: U$ and $U <: T$.) Finally, combining these two relations and enclosing both sides by `addrbook` constructor, we obtain

$$\begin{aligned} & \text{addrbook}[\text{Name, Addr, Name, Addr, Tel, Name, Addr}] \\ & <: \text{addrbook}[(\text{Name, Addr, Tel?})^*] \\ & \stackrel{\text{def}}{=} \text{Addrbook}. \end{aligned}$$

Since the `mybook` value trivially has the type on the left hand side, it has also the type on the right hand side.

Union Types XDuce also provides a union (or alternation) type constructor `|`. For example, we write `(Name|Tel)` to mean “either `Name` or `Tel`”; the basic subtyping relations for union types are the following.

$$\begin{aligned} \text{Name} & <: \text{Name} \mid \text{Tel} \\ \text{Tel} & <: \text{Name} \mid \text{Tel} \end{aligned}$$

Notice that each right hand side offers more possibilities, and so describes a larger set of sequences.

Union types substantially increase our flexibility in programming. In particular, union types yield two interesting relations: “forget ordering” subtyping and “distributivity.” These are the distinguishing points in union types, as opposed to conventional tagged sum types (as found, say, in ML). To illustrate these, let us consider the following scenario of a “database evolution.”

Suppose we begin with a trivial database consisting of just a list of names, with type `Name*`. At some point, this database is copied to two different sites and maintained and evolved separately. At one site, address information is added to each name and the type of the database becomes `(Name,Addr)*`, while at the other telephone numbers are added and it becomes `(Name,Tel)*`.

Now, suppose we want to re-integrate these databases—that is, combine the copies `src1`, whose type is `(Name,Addr)*`, and `src2` of type `(Name,Tel)*` by concatenating them: `src1, src2`. The type of this merged database is, of course, `(Name, Addr)*, (Name,Tel)*`.

Next, suppose we want to do something with our new database that involves extracting the common part (i.e., the name) from each record. Since we have two repetitions in the type, we might expect to need two loops in the program. (We do not show such a program explicitly, but it is easy to write.) However, we can do better by making the two loops into one, using the following “forget ordering” subtype relation:

$$(\text{Name, Addr})^*, (\text{Name, Tel})^* <: ((\text{Name, Addr}) \mid (\text{Name, Tel}))^*$$

The intuition behind this relation is that the ordering information of the left hand side is forgotten on the right hand side. That is, on the left hand side, any $(\text{Name}, \text{Tel})$ pairs must occur after any $(\text{Name}, \text{Addr})$ pairs, while on the right hand side, these pairs can appear in any order.

Finally, since we have two alternatives joined by $|$ in the new type, we might expect to need two branches in our inner loop, to extract the Name field from each of them. But we don't: we can use the following distributive subtyping law

$$(\text{Name}, \text{Addr}) \mid (\text{Name}, \text{Tel}) = \text{Name}, (\text{Addr} \mid \text{Tel})$$

to reorganize the type so that the Name field can be accessed directly.

2.2 Regular Expression Pattern Matching

Our term language is based on a powerful form of pattern matching. Our pattern matching is similar in spirit to ML's (or Haskell's, etc.), but somewhat more powerful, since it includes the use of regular expression types to dynamically match values of those types. Our patterns also require a different treatment of the usual checks for exhaustiveness and ambiguity of patterns.

The body of a XDuce program is a series of function definitions. As an example, the following function converts an address book into a telephone list.

```
fun mkTelList : (Name,Addr,Tel?)* → (Name,Tel)* =
  name[n:String], addr[a:String], tel[t:String],
  rest:(Name,Addr,Tel?)*
  → name[n], tel[t], mkTelList(rest)
| name[n:String], addr[a:String], rest:(Name,Addr,Tel?)*
  → mkTelList(rest)
| ()
  → ()
```

This function takes a value of type $(\text{Name}, \text{Addr}, \text{Tel}?)^*$ and returns a value of type $(\text{Name}, \text{Tel})^*$. The body is a pattern match that breaks up the possibilities on the input values into three cases. The first case matches when the input value is a sequence beginning with name , addr , and tel labels, followed by some further sequence of type $(\text{Name}, \text{Addr}, \text{Tel}?)^*$. In this case, we pick out the name and tel elements and prepend them to the result of calling mkTelList recursively on the remainder. The second case matches when we cannot find tel after addr , and simply calls mkTelList recursively. The third case matches the empty sequence, and returns the empty sequence.

As another example, consider the following function firstTriple , which takes out the first entry with a tel element.

```
fun firstTriple : (Name,Addr,Tel?)* → (Name,Addr,Tel)? =
  ps:(Name,Addr)*, t:(Name,Addr,Tel), rest:(Name,Addr,Tel?)* → t
| whole:(Name,Addr,Tel?)* → ()
```

The function firstTriple has a pattern matching with two cases. In the first case, we skip all "pair" entries (i.e., $(\text{Name}, \text{Addr})$) from the beginning and then

pick out the first “triple” entry (i.e., `(Name, Addr, Tel)`) if such an entry exists. The second case matches otherwise and returns the empty.

The second example is more interesting in that the use of regular expression types is more critical there than in the first example. In the first case in the first example, the pattern matcher will walk over the first three elements of label `name`, `addr`, and `tel`, and then try to match the rest value against the pattern `rest: (Name, Addr, Tel?)*`. However, any value should already have this type. Therefore such a matching would not be meaningful. This is not true in the first case in the second example. When the pattern matcher looks at the first pattern `ps: (Name, Addr)*` in the first case, there is no hint about how many entries are “pairs.” Therefore the matcher must walk through the input value to find where the chain of pairs ends. This matching for a *variable* length sequence is beyond ML pattern matching.

In these examples, pattern matchings are exhaustive. That is, all the values of type `(Name, Addr, Tel?)*` are covered by these patterns. In order to check exhaustiveness, we again use subtyping. For instance, in the first example, we check the following subtype relation

```
(Name, Addr, Tel?)*
<: name[String], addr[String], tel[String], (Name, Addr, Tel?)*
   | name[String], addr[String], (Name, Addr, Tel?)*
   | ()
```

where the left hand side is the parameter type on the annotation and the right hand side is the type constructed from the patterns (i.e., the union of the patterns with all the term variables `n`, `a`, etc. removed). (Although in these examples subtyping of the other way around also holds, we do not check this since allowing this sometimes makes programming easier. Such a situation typically occurs when a pattern contains variables whose type information is useless in the body.)

Our pattern matchings can have two kinds of ambiguity. The first ambiguity occurs when multiple patterns match the same input value. For example, the patterns in `firstTriple` function above are ambiguous since any value that matches the first pattern also matches the second pattern. In such a case, we simply take the first matching pattern (“first matching policy”). The second ambiguity occurs when a single pattern can have multiple ways for variable bindings. This is intrinsic in regular expression pattern matching. For example, suppose we replace the first case in `firstTriple` with the following:

```
es: (Name, Addr, Tel?)*, t: (Name, Addr, Tel), rest: (Name, Addr, Tel?)* → t
```

since we now skip both pair and triple entries at the beginning using the pattern `es: (Name, Addr, Tel?)*`, it is not clear which triple entry the variable `t` is bound to. We resolve this ambiguity by adopting a “longest match” policy where patterns appearing earlier have higher priority. In this example, the first `(Name, Addr, Tel?)*` matches as a long sequence as possible and therefore `t` is bound to the last triple entry. See Section Appendix C for a more formal treatment.

Another possible approach to resolving this ambiguity issue would be to simply disallow ambiguity. However, when we want to write a “default case” in a pattern matching, this restriction would force to write a somewhat cumbersome pattern that captures the “negation” of the other cases.

2.3 More Complex Example: Folder Manipulation

Up to now, the types that we have seen looked like regular expressions on strings. More interesting programs involve regular expressions on *trees*. Consider the following program.

```
type Folder = Record*
type Record = name[String], folder[Folder]
             | name[String], url[String], exists[Bool]

fun tidyFolder : Folder→Folder =
  record:Record, folder:Folder
  → tidyRecord(record), tidyFolder(folder)
| () → ()

fun tidyRecord : Record→Record? =
  name[nm:String], folder[f1:Folder]
  → name[nm], folder[tidyFolder(f1)]
| name[nm:String], url[s:String], exists[false[]]
  → ()
| name[nm:String], url[s:String], exists[true[]]
  → name[nm], url[s], exists[true[]]
```

The mutually recursive types `Folder` and `Record` define a simple template for storing structured lists of bookmarks, as might be found in a web browser: a folder is a list of records, while a record is either a named folder or a named URL plus a boolean indicating whether the link is good or broken. The functions `tidyFolder` and `tidyRecord` traverse a bookmark list recursively, preserving leaves with good links and dropping ones with bad links.

3 Related Work

Mainstream XML-specific languages can be divided into query languages such as XML-QL [DFF⁺] and Lorel [AQM⁺97] and programming languages such as XSLT [XSL]. In general, when one is interested in rather simple information extraction from XML databases, programs in programming languages are less succinct than the same programs in a suitable query language. On the other hand, programming languages tend to be more suitable for writing complicated transformations like conversion to a display format. XDuce is categorized as a programming language.

Static typing of programs for XML processing has been approached from several different angles. One popular approach is to embed a type system for

XML in an existing typed language. The advantage is that we can enjoy not only the static safety and typechecking, but also all the other features provided by the host language. The cost is that XML values and their corresponding DTDs must be some how “injected” into the value and type space of the host language; this usually involves adding more layers of tagging than were present in the original XML documents, which inhibits subtyping. The lack of subtyping (or availability of only restricted forms of subtyping) is not a serious problem for simple traversal of tree structures; it becomes a stumbling block, though, in tasks like the “database evolution” that we discussed in Section 2, where forget-ordering subtyping and distributivity were critically needed.

A recent example of the embedding approach is Wallace and Runciman proposal to use Haskell as a host language [WR99] for XML processing. The only thing they add to Haskell is a mapping from DTDs into Haskell datatypes. This allows their programs to make use of other mechanisms standard in functional programming languages, such as higher-order functions, parametric polymorphism, and pattern matching. However, they do not have any notion of subtyping. Moreover, pattern matching in XDuce is more powerful than Haskell’s in some cases. For instance, as shown in Section 2.2, we can concisely write patterns that skip a variable length sequence by using regular expression types. A difference in the other direction is that XDuce does not currently support higher-order functions or parametric polymorphism. (We are working on both of these extensions.)

Another piece of work along similar lines is the functional language $X\lambda$ for XML processing, proposed by Meijer and Shields [MS99]. Their type system is not described in detail in this paper, but seems to be close to Haskell’s, except that they incorporate *Glushkov automata* in type checking, resulting in a more flexible type system.

A closer relative to XDuce is the query language YAT [CS98], which allows optional use of types similar to DTDs. The notion of subtyping between these types is somewhat weaker than ours (lacking, in particular, the distributivity laws used in the “database evolution” example in Section 2.1).

Milo, Suciu, and Vianu have studied a typechecking problem for their general framework called *k-pebble tree transducers*, which can capture a wide range of query languages for XML [MSV00]. The types there are based on tree automata and conceptually identical to those of XDuce. Papakonstantinou and Vianu present a typechecking algorithm for their query language *loto-ql* by using extensions to DTDs [PV00]. One of their extensions is equivalent to tree automata. The type checking algorithms presented in both papers are *semantically complete*, while XDuce’s is not (since such typechecking would be undecidable in general for Turing-complete languages).

The type system of XDuce was originally motivated by the observation by Buneman and Pierce [BP99] that untagged union types corresponded naturally to forms of variation found in semistructured databases.

4 Conclusions and Future Work

We have presented several examples of XDuce programming and shown how we can write flexible and robust programs for processing XML by combining regular expression types and regular expression pattern matching.

We consider XDuce suitable for applications involving rather complicated tree transformation. Moreover, for such applications, our static typing mechanism would help in reducing development periods.

In this view, we have built a prototype implementation of XDuce and used it to develop some small but non-trivial applications:

Bookmarks can be viewed as a simple database query. It takes as input an Netscape bookmarks file of type `Bookmarks`, which is a subset of the (much larger) type `HTML`. It extracts a certain folder named “Public”, formats it as a free-standing document, adds a table of contents at the front, and inserts links between the contents and the body. The result has type the full `HTML` type. (Total: 224 lines)

Html2Latex takes an `HTML` file (of type `HTML`) and converts it into `LaTeX` format (of type `String`). (264 lines)

Diff implements the “tree diff” algorithm proposed by Chawathe [Cha99]. It takes a pair of XML files of generic `Xml` type and returns a tree with annotations indicating whether each subtree has been retained, inserted, deleted, or changed between the two inputs. (300 lines)

The first two applications are written in the way that traverses the input tree by several simple recursive functions. The third one is more complex. The first phase is a dynamic programming algorithm, where regular expression types are used for representing the internal data structures; the second phase combines two input trees and inserts annotations at each node, where types ensure that the annotations and the actual trees are never confused. In the course of writing these applications, our type checker gave us tremendous help in finding silly mistakes.

The implementation of XDuce raises many algorithmic issues. The primary source of complication is that types and patterns in XDuce are essentially tree automata and therefore we need to use operations on tree automata [CDG⁺], which are in general expensive. For instance, decision for subtyping is inclusion of tree automata, which is known to be EXPTIME-complete [Sei90]. We have addressed this problem and obtained an algorithm that runs efficiently in practice [HVP00]. In particular, the `HTML` type¹ used in the above applications is generally considered to be one of the largest XML DTDs, yet type checking of our programs involving it takes a fraction of a second on stock hardware. As other implementation issues, we are working on a type inference algorithm to eliminate spurious type annotations in patterns, and a pattern compilation scheme to improve run-time efficiency [HP00].

¹ More precisely, we use `XHTML`, which is an XML implementation of `HTML`.

XDuce's language design is far from finished. We plan to add standard features from functional programming, such as higher-order functions and parametric polymorphism. We also consider a support for object-oriented features found in XML-Schema specifications [XS00]. The combination of these features with regular expression types raises some subtle issues, which we are now seeking to solve.

Our prototype implementation is written in O'Cam1 (6500 lines excluding libraries such as the XML parser). Interested readers are invited to visit our home page:

<http://www.cis.upenn.edu/~hahosoya/xduce.html>

5 Acknowledgments

Our main collaborators in the XDuce project are Peter Buneman, Jérôme Vouillon, and Phil Wadler. We have also learned a great deal from discussions with Nils Klarlund and Volker Renneberg, with the DB Group and the PL Club at Penn, and with members of Professor Yonezawa's group at Tokyo.

This work was supported by the University of Pennsylvania and by an NSF Career grant, CCR-9701826. Haruo Hosoya is supported by Japan Society for the Promotion of Science.

References

- [AQM⁺97] Serge Abiteboul, Dallon Quass, Jason McHugh, Jennifer Widom, and Janet L. Wiener. The Lorel query language for semistructured data. *International Journal on Digital Libraries*, 1(1):68–88, 1997.
- [BP99] Peter Buneman and Benjamin Pierce. Union types for semistructured data. In *Proceedings of the International Database Programming Languages Workshop*, September 1999. Also available as University of Pennsylvania Dept. of CIS technical report MS-CIS-99-09.
- [CDG⁺] Hubert Common, Max Dauchet, Rémy Gilleron, Florent Jacquemard, Denis Lugiez, Sophie Tison, and Marc Tommasi. Tree automata techniques and applications. Draft book; available electronically on <http://www.grappa.univ-lille3.fr/tata>.
- [Cha99] Sudarshan S. Chawathe. Comparing hierarchical data in external memory. In *Proceedings of the Twenty-fifth International Conference on Very Large Data Bases*, pages 90–101, Edinburgh, Scotland, U.K., September 1999.
- [CS98] Sophie Cluet and Jerome Simeon. Using YAT to build a web server. In *Intl. Workshop on the Web and Databases (WebDB)*, 1998.
- [DFF⁺] Alin Deutsch, Mary Fernandez, Daniela Florescu, Alon Levy, and Dan Suciu. XML-QL: A Query Language for XML. <http://www.w3.org/TR/NOTE-xml-ql>.
- [HP00] Haruo Hosoya and Benjamin Pierce. Tree automata and pattern matching, July 2000. Available through <http://www.cis.upenn.edu/~hahosoya/papers/tapat-full.ps>.

- [HU79] John E. Hopcroft and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, 1979.
- [HVP00] Haruo Hosoya, Jérôme Vouillon, and Benjamin C. Pierce. Regular expression types for XML. In *Proceedings of the International Conference on Functional Programming (ICFP)*, 2000.
- [MS99] Erik Meijer and Mark Shields. XML: A functional programming language for constructing and manipulating XML documents. Submitted to USENIX 2000 Technical Conference, 1999.
- [MSV00] Tova Milo, Dan Suciu, and Victor Vianu. Typechecking for XML transformers. In *Proceedings of the Nineteenth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, pages 11–22. ACM, May 2000.
- [PV00] Yannis Papakonstantinou and Victor Vianu. DTD Inference for Views of XML Data. In *Proceedings of the Nineteenth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, pages 35–46, Dallas, Texas, May 2000.
- [Sei90] Hermut Seidl. Deciding equivalence of finite tree automata. *SIAM Journal of Computing*, 19(3):424–437, June 1990.
- [WR99] Malcolm Wallace and Colin Ranciman. Haskell and XML: Generic combinators or type-based translation? In *Proceedings of the Fourth ACM SIGPLAN International Conference on Functional Programming (ICFP'99)*, volume 34-9 of *ACM Sigplan Notices*, pages 148–159, N.Y., September 27–29 1999. ACM Press.
- [XML] Extensible markup language (XMLTM). <http://www.w3.org/XML/>.
- [XS00] XML Schema Part 0: Primer, W3C Working Draft. <http://www.w3.org/TR/xmlschema-0/>, 2000.
- [XSL] XSL Transformations (XSLT). <http://www.w3.org/TR/xslt>.

A Types

This section gives the formal definitions for types and subtyping.

A.1 Syntax

We assume a fixed set of base types, ranged over by B , and a countably infinite set of labels, ranged over by l . Each base type B denotes a set V_B of base values, ranged over by b_B . We assume that all V_B are disjoint.

Values are defined as follows.

$$v ::= M_1, \dots, M_n \text{ sequence } (n \geq 0)$$

$$\begin{array}{ll} M ::= b_B & \text{base element} \\ & l[v] \quad \text{label element} \end{array}$$

We write $()$ for the empty value (i.e., the empty sequence). We write v, w for the concatenation of values (i.e., sequences) v and w .

A XDuce program consists of a set of type definitions, a set of function definitions, and an expression to be evaluated. Type definitions have the form

$$\text{type } X = T$$

where X ranges over variables and T ranges over type expressions. Type expressions are defined as follows.

$$\begin{array}{ll} T ::= X & \text{variable} \\ & B \quad \text{base type} \\ & () \quad \text{empty sequence} \\ & l[T] \quad \text{label} \\ & T, T \quad \text{concatenation} \\ & T | T \quad \text{union} \end{array}$$

We sometimes write $\text{type } \bar{X} = \bar{T}$ to abbreviate a set of type definitions $\text{type } X_1 = T_1, \dots, \text{type } X_n = T_n$. Since the set of type definitions is fixed for a given XDuce program, we lighten the notation below by assuming a fixed set of definitions $E = \text{type } \bar{X} = \bar{T}$.

A.2 Semantics of Types

The denotation of each type is a set of values, defined as follows.

First, we define a function $\llbracket \cdot \rrbracket$, which takes a type and an environment mapping type variables to sets of values, and returns a set of values.

$$\begin{array}{l} \llbracket B \rrbracket \rho = V_B \\ \llbracket () \rrbracket \rho = \{ () \} \\ \llbracket l[T] \rrbracket \rho = \{ l[f] \mid f \in \llbracket T \rrbracket \rho \} \\ \llbracket T_1, T_2 \rrbracket \rho = \{ f, g \mid f \in \llbracket T_1 \rrbracket \rho \wedge g \in \llbracket T_2 \rrbracket \rho \} \\ \llbracket T_1 | T_2 \rrbracket \rho = \llbracket T_1 \rrbracket \rho \cup \llbracket T_2 \rrbracket \rho \\ \llbracket X \rrbracket \rho = \rho(X) \end{array}$$

Let μ be the smallest mapping satisfying $\mu(X_i) = \llbracket T_i \rrbracket \mu$ for all the definitions `type Xi=Ti` in E . Finally, write $\llbracket T \rrbracket$ for $\llbracket T \rrbracket \mu$.

A.3 Derived Forms

The regular expression type constructors are derivable as combinations of the above constructs. We represent the Kleene closure T^* by a variable X that is recursively defined as follows (cf. lists as a datatype in ML).

```
type X = T, X | ()
```

The other regular expression constructors are defined as follows.

```
T+ ≡ T , T*
T? ≡ T | ()
```

A.4 Regularity

As we have defined them so far, types correspond to arbitrary context-free grammars—for example, we can write definitions like:

```
type X = Int, X, String | ()
```

Since the decision problem for inclusion between context free languages is undecidable [HU79], we need to impose an additional restriction to reduce the power of the system so that types correspond to regular tree languages. Deciding whether an arbitrary context-free grammar is regular is also undecidable [HU79], so we adopt a simple syntactic condition, called *well-formedness*, that ensures regularity. Intuitively, well-formedness allows recursive uses of variables to occur only in tail positions. For example, we allow the following type definitions:

```
type X = Int, Y
type Y = String, X | ()
```

More precisely, we define well-formedness in terms of a “right-linearity” judgment of the form $\sigma \vdash T : rl(X)$, where σ is a set of variables. It should be read “ T is right-linear in X , assuming that all bodies of variables in σ are right-linear in X .” This judgment uses an auxiliary “disconnectedness” judgment of the form $\sigma \vdash T : dc(X)$, read “ T is disconnected from X (i.e., X does not occur in the top level of T), assuming that all bodies of variables in σ are disconnected from X .”

These two judgments are defined by the following rules (where $X \neq Y$):

$$\begin{array}{ll}
\sigma \vdash T : rl(X) & \text{for } T = () \text{ or } 1[T] \text{ or } X \\
\sigma \vdash Y : rl(X) & \text{if } Y \in \sigma \\
\sigma \vdash Y : rl(X) & \text{if } Y \notin \sigma \text{ and } \sigma \cup \{Y\} \vdash E(Y) : rl(X) \\
\sigma \vdash T|U : rl(X) & \text{if } \sigma \vdash T : rl(X) \text{ and } \sigma \vdash U : rl(X) \\
\sigma \vdash T, U : rl(X) & \text{if } \emptyset \vdash T : dc(X) \text{ and } \sigma \vdash U : rl(X) \\
\sigma \vdash T : dc(X) & \text{for } T = () \text{ or } 1[T] \\
\sigma \vdash Y : dc(X) & \text{if } Y \in \sigma \\
\sigma \vdash Y : dc(X) & \text{if } Y \notin \sigma \text{ and } \sigma \cup \{Y\} \vdash E(Y) : dc(X) \\
\sigma \vdash T|U : dc(X) & \text{if } \sigma \vdash T : dc(X) \text{ and } \sigma \vdash U : dc(X) \\
\sigma \vdash T, U : dc(X) & \text{if } \sigma \vdash T : dc(X) \text{ and } \sigma \vdash U : dc(X)
\end{array}$$

The empty sequence, a label, and the variable X are right-linear in X . For variables Y other than X , we recursively check the right-linearity of the body of Y . To ensure termination, we keep track in σ of variables that have already been checked. For $(T|U)$, both T and U should be right-linear in X . For (T,U) , we check if U is right-linear in X , while T is disconnected from X . The disconnectedness judgment is defined similarly, except for the first rule, in which X is *not* disconnected from X . Now, the set of type definitions E is said to be *well-formed* if

$$\emptyset \vdash E(X) : rl(X) \text{ for all } X \in \text{dom}(E)$$

A.5 Subtyping

We define the subtyping relation in the simplest possible way:

$$S <: T \text{ iff } \llbracket S \rrbracket \subseteq \llbracket T \rrbracket.$$

A.1 Theorem: There is an algorithm to decide the subtyping relation.

The concrete algorithm and its soundness, completeness, and termination theorems are given in our another paper [HVP00].

B Terms

B.1 Syntax

At the top level, we have a set of declarations of mutually recursive functions, each of the following form.

```

fun f : S→T =
  p1 → e1
  | p2 → e2
  | ...
  | pn → en

```

Each function has type annotations for both the parameter and the result. The body consists of one or more pattern matching clauses. Each clause has a pattern and a term for the body. (We write $\text{fun } f : S \rightarrow T = \bar{p} \rightarrow \bar{e}$ to abbreviate the above form.) In addition to the fixed set E of type definitions, we fix a set F of function definitions for the remainder of the document.

The syntax of terms and patterns is:

$e ::= x$ variable
 b_B base value
 $()$ empty sequence
 $l[e]$ label
 e, e concatenation
 $f(e)$ application

$p ::= x:T$ variable
 $()$ empty sequence
 $l[p]$ label
 p, p concatenation

B.2 Typing Rules

A context Γ is a mapping from variables to types, written $x_1:T_1, \dots, x_n:T_n$. We have three judgments:

$\Gamma \vdash e \in T$ e has type T
 $\vdash p \in T \Rightarrow \Gamma'$ p accepts type T and yields context Γ'
 $\vdash \text{fun } f : S \rightarrow T = \bar{p} \rightarrow \bar{e}$ f is well-typed

The typing relation for terms is:

$$\frac{\Gamma(x) = T}{\Gamma \vdash x \in T} \quad (\text{TE-VAR})$$

$$\Gamma \vdash b_B \in B \quad (\text{TE-BASE})$$

$$\Gamma \vdash () \in () \quad (\text{TE-EMP})$$

$$\frac{\Gamma \vdash e \in T}{\Gamma \vdash l[e] \in l[T]} \quad (\text{TE-LAB})$$

$$\frac{\Gamma \vdash e_1 \in T_1 \quad \Gamma \vdash e_2 \in T_2}{\Gamma \vdash e_1, e_2 \in T_1, T_2} \quad (\text{TE-CAT})$$

$$\frac{\text{fun } f : S \rightarrow T = \dots \in F \quad \Gamma \vdash e \in S'}{\Gamma \vdash f(e) \in T} \quad (\text{TE-APP})$$

$$\frac{\Gamma \vdash e \in S \quad S <: T}{\Gamma \vdash e \in T} \quad (\text{TE-SUB})$$

For patterns:

$$\vdash x : T \in T \Rightarrow x : T \quad (\text{TP-VAR})$$

$$\vdash () \in () \Rightarrow \cdot \quad (\text{TP-EMP})$$

$$\frac{\vdash p \in T \Rightarrow \Gamma}{\vdash \mathbb{1}[p] \in \mathbb{1}[T] \Rightarrow \Gamma} \quad (\text{TP-LAB})$$

$$\frac{\vdash p_1 \in T_1 \Rightarrow \Gamma_1 \quad \vdash p_2 \in T_2 \Rightarrow \Gamma_2}{\vdash p_1, p_2 \in T_1, T_2 \Rightarrow \Gamma_1, \Gamma_2} \quad (\text{TP-CAT})$$

For functions:

$$\frac{\vdash p_i \in S_i \Rightarrow \Gamma_i \quad \Gamma_i \vdash e_i \in T_i \quad S <: S_1 | \dots | S_n}{\vdash \text{fun } f : S \rightarrow T = \bar{p} \rightarrow \bar{e}} \quad (\text{TF})$$

The rule TF needs a little explanation: each pattern p_i accepts type S_i and yields the context Γ_i , under which the body e_i can be given the result type T (as annotated). Also, the parameter type S is required to be a subtype of the union of all S_i 's. This subtype checks exhaustiveness of patterns.

C Operational Semantics

The operational semantics consists of evaluation relations for terms and pattern matching. The former is fairly standard, the latter a little unusual due to “ambiguity” in patterns.

There are two kinds of ambiguity. One appears when there are multiple clauses whose patterns match the input value. For example, in the function

```
fun is_single : Int* → Bool =
  x: Int → true
  | x: Int* → false
```

we have two possibilities when the input is a single integer.

The other source of ambiguity is when there are multiple ways in which a single clause can match a given value. For example, in the function

```
fun f : Int* → Int* =
  x: Int*, y: Int* → x
```


it is unclear how many integers x takes from the beginning of the input. We address the first source of ambiguity by adopting a “first match” policy, as in ML. For the second one, we adopt a “longest match” policy as in Emacs.

In fact, the longest-match policy turns out to be just a special case of the first-match policy. Recall that Int^* is a derived form meaning a variable X defined as follows:

```
type X = Int,X | ()
```

The ordering of the branches in the definition is significant: Int,X comes before $()$. Now, in the above example, we first process the pattern $x:\text{Int}^*$, traversing the input value and the definition of X in parallel and taking the first clause (Int,X) as often as possible. When the input value is exhausted and the second clause $()$ is taken, we move on to the pattern $y:\text{Int}^*$, where we can now only take the second clause because the remaining of the input value is the empty sequence.

To describe the first-match policy, we use a notion of *choice sequence*. During pattern matching, we remember the index of the branch we take at each choice point. A choice sequence is a sequence of such indices, listed according to the order of traversal—from left to right and from outer to inner. Finally, we take the *smallest* choice sequence in the dictionary order (written \prec).

The evaluation relations are defined with respect to an environments V —a mapping from variables to values. There are three judgments.

$$\begin{array}{ll} V \vdash e \Downarrow v & e \text{ evaluates to } v \\ \vdash p \triangleright v \Rightarrow V / \alpha & p \text{ matches } v \text{ and yields } V \text{ and choice sequence } \alpha \\ \vdash T \triangleright v \Rightarrow \alpha & T \text{ matches } v \text{ and yields choice sequence } \alpha \end{array}$$

Evaluation of terms is defined as follows:

$$V \vdash x \Downarrow V(x) \quad (\text{EE-VAR})$$

$$V \vdash b_B \Downarrow b_B \quad (\text{EE-BASE})$$

$$V \vdash () \Downarrow () \quad (\text{EE-EMP})$$

$$\frac{V \vdash e \Downarrow v}{V \vdash \mathbb{1}[e] \Downarrow \mathbb{1}[v]} \quad (\text{EE-LAB})$$

$$\frac{V \vdash e_1 \Downarrow v \quad V \vdash e_2 \Downarrow w}{V \vdash e_1, e_2 \Downarrow v, w} \quad (\text{EE-CAT})$$

$$\frac{\begin{array}{l} V \vdash e \Downarrow v \\ \text{fun } f : S \rightarrow T = \bar{p} \rightarrow \bar{e} \in F \\ \vdash p_i \triangleright v \Rightarrow W / \alpha \\ \forall j, \beta. (\vdash p_j \triangleright v \Rightarrow U / \beta \implies (i, \alpha) \prec (j, \beta)) \\ W \vdash e_i \Downarrow w \end{array}}{V \vdash f(e) \Downarrow w} \quad (\text{EE-APP})$$

Pattern matching:

$$\frac{\vdash T \triangleright v \Rightarrow \alpha}{\vdash x:T \triangleright v \Rightarrow x \mapsto v / \alpha} \quad (\text{EP-VAR})$$

$$\vdash () \triangleright () \Rightarrow \cdot / \cdot \quad (\text{EP-EMP})$$

$$\frac{\vdash p \triangleright v \Rightarrow V / \alpha}{\vdash 1[p] \triangleright 1[v] \Rightarrow V / \alpha} \quad (\text{EP-LAB})$$

$$\frac{\begin{array}{l} \vdash p \triangleright v \Rightarrow V / \alpha \\ \vdash q \triangleright w \Rightarrow W / \beta \end{array}}{\vdash p, q \triangleright v, w \Rightarrow V, W / \alpha, \beta} \quad (\text{EP-CAT})$$

Type matching:

$$\frac{\begin{array}{l} \text{type } X=T \in E \\ \vdash T \triangleright v \Rightarrow \alpha \end{array}}{\vdash X \triangleright v \Rightarrow \alpha} \quad (\text{ET-VAR})$$

$$\vdash B \triangleright b_B \Rightarrow \cdot \quad (\text{ET-BASE})$$

$$\vdash () \triangleright () \Rightarrow \cdot \quad (\text{ET-EMP})$$

$$\frac{\vdash T \triangleright v \Rightarrow \alpha}{\vdash 1[T] \triangleright 1[v] \Rightarrow \alpha} \quad (\text{ET-LAB})$$

$$\frac{\begin{array}{l} \vdash T \triangleright v \Rightarrow \alpha \\ \vdash U \triangleright w \Rightarrow \beta \end{array}}{\vdash T, U \triangleright v, w \Rightarrow \alpha, \beta} \quad (\text{ET-CAT})$$

$$\frac{\vdash T \triangleright v \Rightarrow \alpha}{\vdash T|U \triangleright v \Rightarrow 1, \alpha} \quad (\text{ET-OR1})$$

$$\frac{\vdash U \triangleright v \Rightarrow \alpha}{\vdash T|U \triangleright v \Rightarrow 2, \alpha} \quad (\text{ET-OR2})$$

The third and fourth premises in EE-APP say that we choose the match that yields the smallest choice sequence in the dictionary order. Also, notice that in EP-CAT and ET-CAT we concatenate the choice sequences left to right, and that in ET-OR1 and ET-OR2 we adjoin the present choice number to the front. These reflect our policy that the priority of choice is from left to right and from outer to inner.

We can prove the following type soundness: a well-typed term evaluates to a value inhabiting the expected type.

C.1 Theorem [Type Soundness]: Suppose $\vdash \text{fun } f : S \rightarrow T = \bar{p} \rightarrow \bar{e}$ for all functions in F . Then, $\emptyset \vdash e \Downarrow v$ and $\emptyset \vdash e \in T$ imply $v \in \llbracket T \rrbracket$.

Proof: We prove the result by showing the following stronger statements:

- If $\vdash T \triangleright v$, then $v \in T$.
- If $\vdash p \triangleright v \Rightarrow V / \alpha$ and $\vdash p \in T \Rightarrow \Gamma$, then $\Gamma \vdash V$.
- If $V \vdash e \Downarrow v$ and $\Gamma \vdash e \in T$ where $\Gamma \vdash V$, then $v \in \llbracket T \rrbracket$.

(Here, $\Gamma \vdash V$ means that $\text{dom}(\Gamma) = \text{dom}(V)$ and $V(x) \in \llbracket \Gamma(x) \rrbracket$ for each $x \in \text{dom}(\Gamma)$.) Each statement can be proved by induction on the derivation of the type matching, pattern matching, or term evaluation relation. For the last statement, we need the following inversion properties.

- $\Gamma \vdash x \in T$ implies $\Gamma(x) <: T$.
- $\Gamma \vdash \mathbf{b}_B \in T$ implies $B <: T$.
- $\Gamma \vdash () \in T$ implies $() <: T$.
- $\Gamma \vdash \mathbf{1}[e] \in T$ implies $\Gamma \vdash e \in S$ with $\mathbf{1}[S] <: T$ for some S .
- $\Gamma \vdash e_1, e_2 \in T$ implies $\Gamma \vdash e_1 \in S_1$ and $\Gamma \vdash e_2 \in S_2$ with $S_1, S_2 <: T$ for some S_1 and S_2 .
- $\Gamma \vdash f(e) \in T$ implies $\Gamma \vdash e \in S$ where $\text{fun } f : S \rightarrow T' = \dots \in F$ and $T' <: T$.

■