# Regular Object Types

Vladimir Gapeyev      Benjamin C. Pierce

University of Pennsylvania

### Abstract

Regular expression types have been proposed as a foundation for statically typed processing of XML and similar forms of tree-structured data. To date, however, regular expression types have been explored mainly in the setting of special-purpose languages (e.g., XDuce, CDuce, and XQuery) whose type systems were designed around regular expression types "from the ground up." The goal of the Xtatic language is to bring regular expression types to a broader audience by presenting them as a lightweight extension of a popular object-oriented language, $C^{\#}$.

We develop here the formal core of the Xtatic design—a combination of the tree-structured data model of XDuce with the classes-and-objects data model of a conventional object-oriented language. Our tool for this investigation is a tiny language called FX, whose features are drawn from Featherweight Java (FJ) and from the core of XDuce. Points of interest include a smooth interleaving of the two value spaces, in which XDuce's tree structures are grafted into of FJ's class hierarchy while objects and object types play the role of XDuce's label values and label types; a "semantic" definition of the subtype relation, inherited from XDuce and extended to objects; and a natural encoding of XML documents and their schemas using a simple form of singleton classes.

## 1   Introduction

The popularity of XML is due in part to the existence of a number of formalisms for specifying the structures of XML documents. By supporting dynamic consistency checking, ensuring that information being exchanged (e.g., between modules in an application or nodes in a distributed system) has the expected structure, these schema languages significantly increase the robustness of complex XML-based information systems.

However, the exploitation of schema languages by current XML technologies falls far short of what is possible. In particular, schemas now play little part in the static analysis of programs that operate on XML structures: they are not used for checking code for inconsistencies at compile time, or for optimization—in short, they are not used as types in the usual programming-language sense of the term. Taking advantage of this missed opportunity, and thereby improving both the robustness and the efficiency of XML-based information systems, is the long-range goal of the Xtatic project at the University of Pennsylvania.

The key technology for this project is *regular expression types*. Regular expression types are based on well-known constructions from automata theory (they are a mild generalization of nondeterministic tree automata). Their basic constructors (union, concatenation, repetition, etc.) are similar to those found in existing XML schema formalisms such as DTDs [35] and XML-Schema [36]. The difference from these schema languages is that, in Xtatic, XML trees are built-in values of the language, and static analysis of the shapes of trees that may appear at run time (as values of variables, parameters to methods, results of complex expressions, etc.) is part of the ordinary behavior of the typechecker.

Past work on regular expression types at Penn led to a language prototype called XDuce [16, 18, 15, 17]. XDuce is a statically typed language for writing recursive tree transformers—roughly, a statically typed fragment of the popular XSLT language [37]. Beyond regular expression types, its main innovation is a powerful form of *regular expression pattern matching*—a statically typed "tree grep" primitive that arises naturally from types [15]. The XDuce implementation incorporates efficient algorithms for subtyping and typechecking [18].

XDuce has had a significant impact in parts of the XML world; in particular, its influence can be seen in the type system of the XML Query Algebra [11], the core of the W3C standard query language for

XML, as well as newer schema languages such as TREX [8] and Relax NG [9]. However, significant work remains before the benefits of regular expression types can be made available to the vast majority of XML programmers. In particular, the simple tree-data model of XDUCE must be enriched to include objects.

We have begun a new phase of the XDUCE project—a redesign and re-implementation along more ambitious lines, dubbed XTATIC, whose main focus is on inter-operability both at source level and at run time with an established, object-oriented host language. We have chosen compatibility with $C^{\#}$ as our immediate target. (A similar exercise could of course be carried out for similar languages such as Java.) The goal is to make XTATIC as lightweight an extension of $C^{\#}$ as possible, smoothly merging the tree values and types of XDUCE with the familiar object model of $C^{\#}$ and re-using existing $C^{\#}$ features wherever possible in the design, rather than introducing new, XML-specific mechanisms.

This paper develops the formal core of the XTATIC design—a combination of the tree-structured data model of XDUCE with the classes-and-objects model of $C^{\#}$. Our tool for this investigation is a tiny language called FX, which combines Featherweight Java [19] with the core features of XDUCE. The main points of interest may be summarized as follows.

- The two original data models are tightly interwoven in FX. On one hand, the subtype hierarchy of tree types is grafted into the class hierarchy, allowing tree values to be passed to generic library facilities (collection classes, etc.), stored in fields of objects, etc. Conversely, the role of labels and label types in XDUCE is played by objects and classes in FX.

- Subtyping in FX is a natural extension of both the object-oriented subclass relation and the richer subtype relation of regular expression types. XDUCE's simple "semantic" definition of subtyping (sans inference rules) is extended to objects and classes.

- FX enriches XDUCE's regular expression pattern matching construct with a natural form of type-based pattern-matching on objects.

The rest of the paper is organized as follows. Section 2 gives a brief illustrative example of XTATIC code. In Section 3, we review some details of XDUCE and FJ's data models. Section 4—the heart of the paper—combines these to produce the data model of FX. The remainder of the the FX language is informally described in Section 5; standard soundness properties are sketched in Section 6. In Section 7, we show how the FX data model encodes XML types and values. Section 8 overviews relevant work, and Section 9 sketches our plans for the future development of XTATIC.

## 2   Example

XTATIC aims to provide a general mechanism for constructing tree- (and sequence-)structured data. This mechanism can be used to provide an encoding (indeed, multiple possible encodings) for XML data. We give here a sketch of the idea of a possible encoding which is sufficient to get a feel of the language, postponing the details till Section 7.

Assume there is a class `Tag` whose objects are supposed to encode XML tags. Let also classes `Person`, `Name`, `Email` and `Phone` be descendants of `Tag` aimed to encode XML tags `<person>`, `<name>`, `<email>`, `<phone>`, and let `person`, `name`, `email`, `phone` be variables that contain objects of the corresponding classes. Then an expression

```
[ <person>[
      <name>[<"Queen Elisabeth">[]]
      <email>[<"queen@buckingham.uk">[]] ]
   <person>[
      <name>[<"Tony Blair">[]]
      <phone>[<"+49 34 3456">[]] ]
]
```

can be thought of as representing an XML document of a similar structure. The type of this expression is the sequence type

```
[ <Person>[
      <Name>[<String>[]]
      (<Email>[<String>[]] | <Phone>[<String>[]])
   ]* ]
```

The outer parentheses [, ] in both cases delimit regular expressions and types, keeping them distinct from expressions of the host language ($C^\#$ or FJ). A tree is constructed using the form <...>[...] where <...> contains the tree's label, and [...] contains the sequence of children trees. A sequence is built by placing trees adjacently to each other. The type constructor "|" is type alternation (union), and "*" is repetition. Note that native $C^\#$ values, like the tag objects person, name, or strings, occur only inside tree labels.

Sequence values can be examined using type-based pattern matching. For example, assuming the variables list and phonebook each contain a sequence of the type given above and that the variable spamlist holds a string, the code fragment

```
match (list) {
   case [ <Person>[<Name>[String]  <Email>[String e] ]
          //...
        ]:
      spamlist = spamlist + "," + e;
      //...
   case [ (<Person>[String] p)  //...
        ]:
      phonebook = [ phonebook  p ]
      //...
   case []:
      //...
}
```

inspects the first tree in the string list and, if the corresponding person has an email, extracts the address into a string variable e and uses it to extend the string spamlist; otherwise, the person must have a phone, and the second case branch handles this case by binding the whole person's entry to the variable p and using it to extend phonebook sequence. (Note that Java comments //... elide parts of code above.)

# 3   Technical Background

The *data model* of a language is the collection of *values* that programs in the language manipulate, their *types*, and fundamental relations such as *value typing* and *subtyping*. The data model is the bedrock on which the full language definition (the syntax, typing rules, and evaluation rules for expressions) rests. Because the primary topic of this paper is the combination of trees and objects (and their types), the data model of FX is where we need to concentrate our most careful attention.

As background for what follows, this section briefly—and informally—sketches the data models found in XDuce and FJ.

## 3.1   The XDuce Data Model

The data model of XDuce is parameterized on a language of *labels*. The details of these labels can vary (and do vary, across the several published XDuce papers and implementations), but all the variations contain the following common structure:

- a set $L$ of *label values*, ranged over by l,

- a set of *label types*, ranged over by L,

- a denotation function $[\![\cdot]\!]$ giving the set $[\![L]\!] \subseteq L$ of label values that are members of each label type L.

The subtyping relation on label types, written $L_1 \sqsubseteq: L_2$, is generated by $[\![\cdot]\!]$—that is $L_1 \sqsubseteq: L_2$ iff $[\![L_1]\!] \subseteq [\![L_2]\!]$.

For example, one simple choice of label language is to select an arbitrary set of identifiers and designate them to be the set $L$ of label values; for each value $l \in L$, we consider $l$ to be a label type as well (i.e., $l$ is the *singleton type* whose denotation contains just $l$); we also introduce the wildcard label type ~, denoting the whole set $L$. Of course, a yet simpler choice would be to omit ~, but having a maximal label type turns out to be extremely useful in pattern matching, where it functions as a "don't care" pattern.

Having selected the language of labels, the XDuce data model can be defined in a uniform way. First, a *tree value* t consists of a label value as its node and a sequence of children tree values:

$$t \quad ::= \quad l[t_1, \ldots, t_n] \qquad \text{where } n \geq 0$$

Now, a *sequence value* is a sequence $t_1, \ldots, t_n$ of zero or more tree values combined using the comma ','. (We use the shorthand notation $\bar{t}$ throughout the paper for sequences, and () to denote the empty sequence. We write $\bar{s}, \bar{t}$ for the concatenation of the sequences $\bar{s}$ and $\bar{t}$.)

XDuce types—*regular expression types*—are built from tree types and references to a collection of globally defined type names X:

$$
\begin{array}{lll}
T & ::= & \\
& \texttt{<L>[T]} & \text{tree} \\
& \texttt{()} & \text{empty sequence} \\
& \texttt{T T} & \text{concatenation} \\
& \texttt{T|T} & \text{union} \\
& \texttt{T*} & \text{repetition} \\
& \texttt{X} & \text{type name}
\end{array}
$$

The meanings of types are given with respect to a global collection of type definitions, each of the form X = T. We write *Typenames* for the set of of type names X appearing on the left of one of these definitions, *def* for the function that maps each X to the corresponding body type T. The global definitions may be recursive or mutually recursive, but (to limit the power of the type language to just regular, rather than context-free, sets of trees), we impose the condition that all "loops" from a variable X back to itself must pass through the body of at least one `<L>[T]` construct—i.e., "top-level" recursion is not allowed.

Next, the denotation function $[\![\cdot]\!]$ mapping types T to sets of sequence values $\bar{t}$ is defined as the least solution of the following equations:

$$
\begin{array}{lll}
[\![\texttt{<L>[T]}]\!] & = & \{\, \texttt{<l>}[\bar{t}] \mid l \in [\![L]\!] \text{ and } \bar{t} \in [\![T]\!] \,\} \\
[\![\texttt{()}]\!] & = & \{\, () \,\} \\
[\![T_1\ T_2]\!] & = & \{\, \bar{t}_1\ \bar{t}_2 \mid \bar{t}_1 \in [\![T_1]\!], \bar{t}_2 \in [\![T_2]\!] \,\} \\
[\![T_1\,|\,T_2]\!] & = & [\![T_1]\!] \cup [\![T_2]\!] \\
[\![T*]\!] & = & \{\, t_1, \ldots, t_n \mid n \geq 0 \wedge \forall k \in 1 \ldots n.\ t_k \in [\![T]\!] \,\} \\
[\![X]\!] & = & [\![def(X)]\!]
\end{array}
$$

The subtyping relation for regular expression types is defined in the simplest imaginable way:

$$T_1 <: T_2 \quad \text{iff} \quad [\![T_1]\!] \subseteq [\![T_2]\!].$$

The fact that subtyping can be defined semantically is actually quite important in XDuce. The alternative—writing down a collection of inference rules characterizing the relation inductively—would be much heavier and harder to understand than the subtype relations of most languages, since the regular expression type constructors satisfy many algebraic laws arising from the associativity of comma and the associativity, commutativity, and distributivity (over comma and `<L>[...]`) of the (*non*-disjoint!) union. An inference-rule presentation of the subtyping relation can certainly be given—indeed, it must be: it is the basis for the algorithm for subtype checking [18]—but it is not pretty.

## 3.2   The FJ Data Model

Featherweight Java, or FJ, is a tiny calculus designed to capture the essential typing mechanisms of class-based object-oriented type systems in programming languages such as Java and $C^{\#}$. It was first used by

Igarashi, Pierce, and Wadler [19] to formalize the GJ [5] type system, and has since formed the basis of numerous formal studies of Java and related languages [20, 30, 2, 3, 38, 1, 28, 22, etc.]. FJ embodies the core mechanisms of object creation, field access, method invocation, and inheritance (and—in the most common presentation, though not here—casting) in exactly the same form as they are found in Java, while omitting everything else... from reflection and concurrency to interfaces, overloading, static members, and even assignment.

An FJ program consists of a collection of class declarations plus a single expression to be evaluated. The types in an FJ program are just class names C. FJ values are objects, which (since FJ is a declarative language, so the only thing that distinguishes one object from another are its class and the arguments passed to its constructor) are simply identified with new expressions.

$$
\begin{aligned}
\mathtt{o} \quad ::= \quad & \\
& \mathtt{new\ C(o_1,\ \ldots o_n)}
\end{aligned}
$$

The constructor arguments $o_1, \ldots, o_n$ (usually written just $\overline{o}$) are required to correspond exactly to the fields of the class C. For example, if C has fields a and b and its immediate subclass D has fields e and f, then an instance of D will always have the form $\mathtt{new\ D(o_1,o_2,o_3,o_4)}$, where $o_1$ is the value for the a field of the new object, $o_2$ is the value of the b field, $o_3$ of the c field, and $o_4$ of the d field.

The global set of class definitions in an FJ program is formalized as a *static context*, which provides several functions for examining different aspects of the class definitions: the set of all defined classes (which always includes the special class Object); the immediate-subclass relation, which must be tree-structured with Object at the root; the list of field names and types in each class; the method names and signatures in each class; and the method bodies for each class. This static context is used to define the typing and evaluation relations. For purposes of discussing the FJ data model, we can restrict attention to the part of the static context comprising just the set of class names and the immediate-subclass relation, which we call the *static data context*.

The subtype relation in FJ, written $\mathtt{C_1}\ \sqsubseteq: \mathtt{C_2}$, is the reflexive and transitive closure of the immediate-subclass relation. Like XDUCE's, this definition of subtyping is pleasingly simple; however, it has a completely different—more syntactic—character. In order to combine the two data models, we need to look for a more "semantic" presentation of this one (as we remarked above, a syntactic presentation of XDUCE subtyping is an unattractive alternative). This can be achieved as follows.

We say that a value $\mathtt{new\ C(\overline{o})}$ is an *instance* of the class C. That is, an object is an instance of the class from which it was created. The *denotation* of a class C is then the set of all instances of this class and all its subclasses:

$$
[\![\mathtt{C}]\!] \quad = \quad \{\, \mathtt{o} \mid \mathtt{o} \text{ is an instance of } \mathtt{D}, \text{ for some } \mathtt{D} \sqsubseteq: \mathtt{C} \,\}
$$

Note that this does *not* require that the constructor arguments $\overline{o}$ belong to the types of the fields in class C. This may appear overly permissive, but it has some useful implications:

1. It is obvious from the definition that the "semantic" subtyping relation derived from it coincides exactly with the syntactic subclass relation:

$$
\mathtt{C_1} \sqsubseteq: \mathtt{C_2} \quad \text{iff} \quad [\![\mathtt{C_1}]\!] \subseteq [\![\mathtt{C_2}]\!].
$$

2. This definition requires no changes if we enrich the language with imperative features. A more precise definition ("the values in type C are objects of the form $\mathtt{new\ D(\overline{o})}$, where $\mathtt{D} \sqsubseteq: \mathtt{C}$ and, $o_i \in [\![\mathtt{F_i}]\!]$, where $\mathtt{F_i}$ is the type of the $i^{th}$ field of class D") would require a co-inductive reading (cf. [33]) to make sense in this setting.

Intuitively, the reason we can get away with this "loose" interpretation of types is that later, e.g., in the proof of soundness, we will never deal with arbitrary elements of $[\![\mathtt{C}]\!]$, but only with elements that we also know are well typed (according to the expression typing relation, which *does* ensure that constructor arguments have the right types).

| Values | | | Types | | |
|---|---|---|---|---|---|
| | | *Full FX language* | | | |
| a | ::= | FX value | A | ::= | FX type |
| | new C(ā) | object | | C | class type |
| | [t̄] | delimited sequence | | [T] | delimited RE type |
| | | *Regular expression sublanguage* | | | |
| t | ::= | | T | ::= | RE type |
| | <a>[t̄] | tree value | | X | RE type variable |
| | | | | <C>[T] | tree type |
| | | | | () | empty sequence |
| | | | | T T | concatenation |
| | | | | T\|T | alternation |
| | | | | T* | repetition |

Figure 1: FX values and types

# 4   The FX Data Model

The interweaving of XDUCE's and FJ's data models in FX is founded on two observations.

1. We can treat sequences of trees as objects simply by "grafting" the whole collection of regular expression types into the class hierarchy, inventing a special *class* Seq whose sub*types* are all the regular expression types. This grafting is justified by a compilation model—reminiscent of GJ's homogeneous translation [5, 19]—in which all regular expression types in an FX program are "erased" to the single class type Seq and all tree values are translated into objects of class Seq.

2. The data model of objects and classes qualifies as a "label language" in the sense discussed in Section 3.1, so we can use arbitrary objects as the labels in XDUCE trees and classes as label types.

Formally, the data model is defined in three steps. First, we give the syntax of values and types. Next, we give the notion of a static context, which summarizes the type-related information defined in a program. Finally, fixing a particular static context, we define the membership relation for values in types.

Figure 1 defines the syntax. An FX value a can have one of two forms: it is either an object new C(ā) or a sequence [t̄] delimited by special brackets. Observe that, inside an object new C(ā), the values of fields may be arbitrary FX values ā; in particular, they can be sequences. The organization of FX types A is similar, combining class types C and regular types [T]. Regular values t̄ and regular types T are essentially those of XDUCE, where any FX value can be used as a label in a tree value and any class type C can be used as a label in a tree type.[1]

A *static data context* is a tuple

$$DatCtx = \langle Classes, \sqsubseteq:, Typenames, def \rangle$$

where

- *Classes* is a set of class names, ranged over by C and containing special names Object and Seq;

- $\sqsubseteq$: is a binary relation on *Classes*, generated as a reflexive and transitive closure from a relation corresponding to an "immediate predecessor" function *Parent* : *Classes* \ {Object} → *Classes*;

- *Typenames* is a set of type names, ranged over by X;

---

[1] The careful reader may note a certain discrepancy here: a sequence can be used as a label in another tree, as in <[s̄]>[t̄], but a regular expression *type* cannot be similarly used as a label. This raises the question of what type can be given to a value of the above form. As we shall see later, the type would have the form <Seq>[T] where Seq is the special class type containing all sequences.

$$
\begin{array}{ll}
\textit{Instances} : \texttt{C} \mapsto \mathcal{P}(\texttt{a}) \\[4pt]
\quad \textit{Instances}(\texttt{C}) \;=\; \left\{ \begin{array}{ll} \{\, [\overline{\texttt{t}}] \mid \overline{\texttt{t}} \text{ arbitrary } \} & \text{if } \texttt{C} = \texttt{Seq} \\ \{\, \texttt{new C}(\overline{\texttt{a}}) \mid \overline{\texttt{a}} \text{ arbitrary} \} & \text{otherwise} \end{array} \right. \\[16pt]
\llbracket \_ \rrbracket : \texttt{A} \mapsto \mathcal{P}(\texttt{a}) \\[2pt]
\quad \llbracket \texttt{C} \rrbracket & = \quad \bigcup \{\, \textit{Instances}(\texttt{D}) \mid \texttt{D} \sqsubseteq : \texttt{C} \,\} \\[2pt]
\quad \llbracket \texttt{X} \rrbracket & = \quad \llbracket\, \textit{def}(\texttt{X}) \rrbracket \\[2pt]
\quad \llbracket \texttt{<C>[T]} \rrbracket & = \quad \{\, [\texttt{<a>}[\overline{\texttt{t}}]] \mid \texttt{a} \in \llbracket \texttt{C} \rrbracket \text{ and } [\overline{\texttt{t}}] \in \llbracket \texttt{T} \rrbracket \,\} \\[2pt]
\quad \llbracket \texttt{()} \rrbracket & = \quad \{\, [] \,\} \\[2pt]
\quad \llbracket \texttt{T}_1 \ \texttt{T}_2 \rrbracket & = \quad \{\, [\overline{\texttt{t}}_1 \ \overline{\texttt{t}}_2] \mid [\overline{\texttt{t}}_1] \in \llbracket \texttt{T}_1 \rrbracket, [\overline{\texttt{t}}_2] \in \llbracket \texttt{T}_2 \rrbracket \,\} \\[2pt]
\quad \llbracket \texttt{T}_1 \mid \texttt{T}_2 \rrbracket & = \quad \llbracket \texttt{T}_1 \rrbracket \cup \llbracket \texttt{T}_2 \rrbracket \\[2pt]
\quad \llbracket \texttt{T*} \rrbracket & = \quad \{\, [\overline{\texttt{t}}_1 \ldots \overline{\texttt{t}}_n] \mid \forall k \in 1 \ldots n. \ \overline{\texttt{t}}_k \in \llbracket \texttt{T} \rrbracket, \text{ for some } n \geq 0 \,\}
\end{array}
$$

Figure 2: Type denotations.

- *def* is a function from *Typenames* to types, that maps each type name X to a regular expression type expression T (its *definition*);

and such that

1. if a type name X$'$ appears in *def*(X), then X$' \in$ *Typenames*; and

2. a grammar obtained from *def* by considering variables from *Typenames* as non-terminals generates a regular language (see [18] for a formal syntactic restriction on the grammar that guarantees this condition).

The semantics of types is given by the denotation function $\llbracket \ \rrbracket$, which maps each type A to its set of inhabitants a. This function is the least solution of the equations in Figure 2. Note the special role of the class Seq, whose denotation does not contain objects (new Seq($\overline{\texttt{a}}$) is not in the denotation of any type), but instead contains all sequence values.

Subtyping on FX types is defined semantically:

$$
\texttt{A}_1 <: \texttt{A}_2 \qquad \Longleftrightarrow \qquad \llbracket \texttt{A}_1 \rrbracket \subseteq \llbracket \texttt{A}_2 \rrbracket.
$$

The XDUCE subtyping algorithm [18] can be used to decide this relation, since it is parameterized by the subyping relation for tree label types (called there "subtagging"), which corresponds in FX to the subclass relation $\texttt{C}_1 \sqsubseteq : \texttt{C}_2$.

# 5   The FX Language

The FX data model described in the previous section establishes a skeleton, on which a full-blown programming language can be constructed—providing ways of interrogating and destructing values, as well as abstraction mechanisms and all the other usual apparatus. Naturally, FX's value-destruction mechanisms are contributed by the corresponding sublanguages: FJ provides field projection on objects and XDUCE brings in regular-expression pattern matching on sequences and trees. The abstraction mechanisms of FX—classes, methods, and inheritance—are taken entirely from FJ.

Figure 3 gives the syntax of FX expressions and their constituent patterns. The behavior of most of these constructs is standard; therefore we discuss the language semantics mostly informally, commenting in more detail on the issues that are novel in FX. Full definitions can be found in Appendix A.

We do not describe concrete syntax for class and method declarations: for the present discussion it is more convenient to think about an FX program as an abstract *static context Ctx* defined along the lines,

```
            e, d   ::=                          expression
                   x                               value variable
                   new C(e̅)                        new object creation
                   e.f                             field access
                   e.m(e̅)                          method call
                   <e>[e̅]                          tree
                   [e̅]                             sequence
                   match(e){case [P̅]:  e̅}          pattern match
```

| R | ::= | FX pattern | | P | ::= | RE pattern |
|---|---|---|---|---|---|---|
| | Q | class pattern | | | X | RE type name |
| | [P] | delimited RE pattern | | | <Q>[P] | tree |
| | | | | | () | empty sequence |
| Q | ::= | class pattern | | | P P | concatenation |
| | C | class | | | P\|P | alternative |
| | C x | object binding | | | T* | type repetition |
| | | | | | P x | RE value binding |

Figure 3: FX language syntax.

and as an extension of, the static data context *DatCtx* of Section 4. Namely, in addition to the items from *DatCtx*, the full context *Ctx* associates with each class a collection of methods available for calling on the objects of the class. For each method, *Ctx* provides its signature (types of the arguments and the return type), the list of argument variables, and the expression of the body. Additionally, *Ctx* must obey constraints on method types in subclasses, stemming from the $C^{\#}$ inheritance rules.

The only significant difference of an FX context *Ctx* from the information provided by an FJ program is that the types appearing in method signatures are arbitrary FX types, i.e. they can be regular types as well as classes. Consequently, the subtyping relation used for checking the inheritance constraints (as part of the process of checking that a class is well formed) is the semantic subtyping relation $<:$. Similarly, FX variables x (which can only originate in FX as method argument names or as binders in patterns) can hold any FX values, either objects or sequences. As in FJ, there is a special variable this that can be used in expressions to refer to the current object. The typing and evaluation rules treat this variable specially.

The FX data model permits only tree values to be members of sequences. That is, something like [ [t̅] (new C(a̅)) [s̅] ] is not a well-formed value. The syntax of *expressions*, however, does allow nested sequences. The reason is that we want an expression like

$$[ \ db.getPapers("POPL") \ db.getPapers("ICFP") \ ]$$

to be legal—provided the method getPapers() returns values of a sequence type—and to mean the concatenation of the sequences returned by the two calls. But, as long as this expression is legal, the expression obtained by replacing the method calls by their results—which is a nested sequence—should also be legal. Therefore, a nested sequence [ [e̅] [d̅] ] is a valid FX expression, which evaluates to the same value as [ e̅ d̅ ]. Of course, for the latter to be type safe, the types of expressions [e̅] and [d̅] must both be valid regular types. The FX typing rules ensure that this is indeed the case for nested sequence expressions.

On the other hand, an object is never legal as a member of a sequence and, symmetrically, a tree expression <e>[d̅] is never allowed outside the sequence parentheses [ ... ]. Since both are permitted syntactically, this condition is checked by the typing rules.

Deconstruction of sequence values is done by matching them against patterns using the match construct, which syntactically resembles $C^{\#}$ switch statement but behaves more like XDUCE's match. That is, the behavior of an expression

```
match (d) {
    case [P₁]:   e₁;
    case [P₂]:   e₂;
...
    case [Pₙ]:   eₙ;
}
```

is to evaluate d and match the result against each of the patterns in turn until the first one, say [$P_i$], matching the value is encountered. The successful match produces an environment that maps variables declared in [$P_i$] to the appropriate portions of the value computed from d. The result of the whole expression is the result of evaluating $e_i$, assuming variable mappings from the environment. So, the case bodies do not have the "fall through" behavior of switch. The value of match' input d must be a sequence, and all case patterns must be sequence patterns.

The syntax of FX sequence patterns [P] is essentially that of XDuce:[2] a pattern is just a type annotated with variable binders. This intuition is also extended to class types in FX. A *class pattern* has the form C x, where C is a class name and x is a variable to be bound. At run time, an object value new D($\bar{a}$) matches the pattern whenever D $\sqsubseteq$ C, what agrees with the denotational relationship $[\![D]\!] \subseteq [\![C]\!]$. As it is the case for the class Seq, the pattern Seq x is treated specially: it matches any sequence value [$\bar{t}$]. (And, naturally, [$\bar{t}$] is matched by any ancestor class of Seq, including Object.)

Since classes are types of labels in tree types, it is natural to use a class pattern in the label position in a tree pattern. This allows one to extract a label from a tree as an object for later use in the program. It is worthwhile noticing that this is a benefit of our goal to use, whenever possible, C$^{\#}$ features for the needs of regular types. To compare, in XDuce a label is an integral part of a tree and cannot be extracted from it as a first-class value.

The typing of match depends on the type inference for variables bound in its patterns. In XDuce, the type inference, formalized as the judgment T ▷ P ⇒ Γ, is precise, meaning that for each type Γ(x), each value from its denotation can be possibly bound to x at run time as a result of matching some value from T's denotation against P, and $[\![\Gamma(x)]\!]$ does not contain values that cannot be thus obtained. The precision is achieved thanks to the availability of unrestricted union operation on XDuce types. In FX, however, we cannot have union for class types, and have to use an upper bound instead, sacrificing the precision of type inference. For an example, suppose class D has A, B, and C as its direct subclasses, and consider matching values of type [<D>[]] against pattern [<A x>[] | <B x>[]]. Since there is no a class whose denotation is an exact union of the denotations of A and B, the only reasonable type assignment for x is D, which is not precise—x can never be bound to an object of C, another D's descendant. Therefore, we decided to formalize FX type assignment for pattern variables by a simpler relation, ▷R ⇒ Γ, which does not take the input type into account and assigns Γ(x) to be the type on which x appears as the annotation in pattern R (and, in the case of the alternation pattern like in the above example—the *join* of the types of the alternatives, defined as their smallest common computable upper bound).

In XDuce, precise pattern type inference extends to the whole collection of patterns of a match: the input type for the $i^{\text{th}}$ pattern $P_i$ is not the type of input to the whole match, but the input restricted to those values that could not be matched by any of the previous patterns. Implementation of this feature depends on availability of type difference. Since difference is not available for classes, we had to give up on this feature as well.

We are able, however, to check for exhaustiveness of patterns—by checking if the input type is a subtype of the join of (types of) all the patterns, and provide a restricted form of pattern redundancy checking—by comparing, for each prefix of the pattern list, the join of the prefix's types and the input type.

# 6   Properties

We will now formulate for FX the standard results of soundness and progress. The proofs of them follow standard induction techniques and are omitted.

---

[2]We also demand that each pattern P satisfy the same regularity constraint as for types, and that it be *linear*. Intuitively, linearity means that no variable is bound in P twice, except in the alternation subpatterns, where alternative branches must bind exactly the same variables (see [15], appendix A.2, for the formal definition).

All the results are stated assuming there is a well-formed static context corresponding to an FX program.

Value environments are mappings $\Sigma : x \mapsto a$, and typing environments are mappings $\Gamma : x \mapsto A$. An environment with the empty domain is written $\bullet$ (both for value and type environments).

The following three relations formalizing FX operational semantics can be obtained by adopting the corresponding relations from FJ and XDuce, tacking into account the comments in Section 5 (we have to omit them from the main body of the paper for the lack of space):

- $\Gamma \vdash e \in A$, "in the typing environment $\Gamma$, expression $e$ gets type $A$",

- $\Sigma \vdash e \Downarrow a$, "in the value environment $\Sigma$, expression $e$ evaluates to the value $a$",

- $\Sigma \vdash e \not\Downarrow$ "evaluation of $e$ gets stuck in the finite number of steps" (this relation is specific to big-step semantics, the analogous property for small-step semantics says that $e$ gets reduced to a non-value expression to which none of the evaluation rules is applicable).

Write $a \in: A$ to mean that $\bullet \vdash a \in A'$ for some $A' <: A$.

A value environment $\Sigma$ *conforms* to a typing environment $\Gamma$, written $\Gamma \vdash \Sigma$, if $dom(\Sigma) = dom(\Gamma)$ and $\Sigma(x) \in: \Gamma(x)$, for all $x$.

**6.1 Theorem [Soundness]:** For $\Gamma \vdash \Sigma$, if $\Gamma \vdash e \in A$ and $\Sigma \vdash e \Downarrow a$, then $a \in: A$.

**6.2 Theorem [Progress]:** If $\bullet \vdash e \in A$, then not $\bullet \vdash e \not\Downarrow$.

Both of the standard theorems depend (in the parts of their proofs corresponding to the `match` construct) on the following property of pattern matching, which is interesting in itself. Recall that the object-against-pattern case in our pattern matching relation $a \triangleright R \Rightarrow \Sigma$ does not check for the well-typedness of object fields. The property says that, despite of this, if pattern matching is done against a well typed value $a$, any binding in the resulting environment is also well-typed.

**6.3 Proposition [Pattern matching preserves well-typedness]:** Let $a$ and $A$ be such that $\bullet \vdash a \in A$. If $a \in R \Rightarrow \Sigma$ and $\triangleright R \Rightarrow \Gamma$, then $A <: tyof(R)$ and, for all $x$, $\bullet \vdash \Sigma(x) \in B$ for some $B <: \Gamma(x)$.

# 7   XML in FX

So far, none of the mechanisms we have described have been especially tied to XML—we have simply established a generic foundation for representing and manipulating ordered, labeled tree structures in an object-oriented setting. Our final job is to show how this foundation supports a natural encoding of (most of) XML itself, based on a simple form of singleton types and a modicum of syntactic sugar.

We begin by explaining how the textual "leaf data" of XML documents, known as PCDATA (parsed character data), can be treated. Our first step is to extend, conceptually, the $C^{\#}$ data model by introducing singleton classes for individual characters. We assume that the data context *DatCtx* provides a class `Char`, corresponding to the standard $C^{\#}$ character class, plus, for each character c, a class $\mathtt{Char_c}$ extending `Char`. All these classes have no instance variables and nullary constructors—thus, each class $\mathtt{Char_c}$ contains only a single object, `new` $\mathtt{Char_c}$`()`, which we can identify with the character c itself. Now, a $C^{\#}$ character literal, say $'a'$, is considered as syntactic sugar for either the object `new` $\mathtt{Char_a}$`()`, when used in an expression, or for the class $\mathtt{Char_a}$, when used in a type.

We can now define a regular expression type `PCDATA` for representing XML character data:

$$def(\mathtt{PCDATA}) \quad = \quad (\ \mathtt{<Char>[]}\ )*$$

That is, an XML text value is represented by a sequence of trees, where each tree has no children and has a character object as its label. The type `PCDATA` contains arbitrary text strings, so we can write patterns like

$$\mathtt{<Object>[PCDATA],}$$

which matches a tree whose body contains only character data.

Why we did not adopt the more obvious choice of using C$^{\#}$'s `String` class to hold XML character data? One reason is that the `PCDATA` representation opens the way to interesting uses of pattern matching for string regular expression processing. Since each `Char`$_a$ is a subtype of `Char`, we can write types that restrict text to a particular form. For example, all character sequences starting with `'a'` and ending with `'b'` belong to the type

$$\texttt{<'a'>[] PCDATA <'b'>[]}.$$

This type, like any XTATIC type, can be annotated with variable binders to obtain a pattern. The general pattern-matching facility, then, offers functionality somewhat similar to that of Perl's regular expression string patterns, but with static typing support. (See [32] for a deeper exploration of this idea.)

Another reason for using `PCDATA` instead of `String` is that, in XML, two character sequences following each other are indistinguishable from a single larger character sequence. The `PCDATA` type satisfies this requirement,

$$[\![\texttt{PCDATA PCDATA}]\!] = [\![\texttt{<Char>[]* <Char>[]*}]\!] = [\![\texttt{<Char>[]*}]\!] = [\![\texttt{PCDATA}]\!]$$

but a `String`-based representation does not, since $[\![\texttt{<String>[] <String>[]}]\!] \neq [\![\texttt{<String>[]}]\!]$.

The encoding of XML documents in Xtatic now follows naturally—all we need is an encoding for XML tags, and this can be obtained by following exactly the same intuitions that we used for characters. We assume that the data context *DatCtx* contains a special class `Tag` and, for each XML tag `<g>`, a singleton class `Tag`$_{\texttt{<g>}}$ (with the object `new Tag`$_{\texttt{<g>}}$`()` as its only inhabitant) as an immediate subclass of `Tag`. Then, for an XML fragment

        `<basket> <apple/> <banana/> </basket>`

the corresponding Xtatic value is

        `<new Tag`$_{\texttt{<basket>}}$`()>[ <new Tag`$_{\texttt{<apple>}}$`()>[] <new Tag`$_{\texttt{<banana>}}$`()>[] ]`

and the corresponding type is

        `<Tag`$_{\texttt{<basket>}}$`>[ <Tag`$_{\texttt{<apple>}}$`>[] <Tag`$_{\texttt{<banana>}}$`>[] ]`

Of course, an implementation needs special syntax that makes these values and types readable (and even writable!). The concrete syntax in our current prototype implementation looks very close to standard XML.

Together, the encodings of character data and tags allow a good-size fragment of XML to be represented very directly in FX. (There are still important parts missing, though. Most urgently, we still lack a good treatment of attributes which, until very recently [14], was also lacking in XDUCE.)

The only basic data type provided by the XML standard is character sequences. Some schema formalisms, however, introduce *datatypes*—a set of conventions by which a schema can specify that a particular textual fragment in an XML document is supposed to represent a non-textual value, e.g. a float or a date. Some of these datatype descriptions can be captured using subtypes of `PCDATA` built from regular expression operators to mimic the string regular expressions that describe particular datatype formats. Alternatively, the FX framework could accomodate a Schema-datatype-aware encoding of XML, when a text representing a Schema datatype value gets translated directly into a value of an appropriate C$^{\#}$ type (placed as a label of a childless tree), bypassing the `PCDATA` representation.

# 8   Related Work

There is a substantial literature and many formalisms and tools for dynamic validation of XML documents against expected schemas, either by stand-alone processors or during document construction, as has been proposed for DOM Level 3 [10]. While XTATIC shares some formal background with these techniques, its central goal—to support *static* checking of XML-manipulating code—falls completely outside their purview.

Among static approaches, there are two, not entirely distinct, kinds of work that are directly relevant to ours: work on providing XML processing capabilities in a pre-existing programming language with static guarantees of correctness, and work on combining object-orientation with other XML-like data models.

A popular direction for work of the first kind is to provide a translation that generates type definitions (and value constructors) in the original language corresponding to XML types of interest. Examples include JAXB [31], Relaxer [27], HaXML [34], and XMλ [23, 29]. One disadvantage of these translations is that they tend to introduce "spurious structure," destroying some useful flexibility in the subtype relation. This point is discussed in detail in [17] and [13].

There can be varying degrees of integration of a "foreign" data model into the OO data model. One is creating a combined data model that incorporates the features of both on the equal level. A successful example is the ODMG data model [6], an accepted standard for object-oriented databases, which offers a class-based object-oriented type system analogous to that of programming languages like $C^{\#}$, together with a few other built-in type constructors: records, sets, bags, lists and arrays (all of them typed).

A greater degree of integration can be achieved by taking the object-oriented data model as primary and the other data model as somewhat subsidiary, in the sense that its values can also be viewed as objects. This approach has the advantage of better integration with "legacy" software written entirely under the original object-oriented model. Examples of this approach can be found in both the programming language and database communities.

Pizza [25] project extended Java with parametric polymorphism, higher-order functions and tagged union types with pattern matching. (The polymorhism component was superseded by the GJ [5] proposal of generic types for Java.) All these features where implemented by translation into pure Java in such a way that the extended data model is used to type check Pizza source, while run time representations (if any) of the additional features are objects of either pre-determined Java classes (for their *homogeneous* translation), or of classes generated w.r.t. the Pizza source (for the *heterogeneous* translation). We plan to use a scheme analogous to the homogeneous translation in the final implementation of XTATIC. Currently, the programming language Scala [24] is developed to incorporate many on the same ideas in a larger context aimed for programming Web services, including XML processing.

Even before XML became popular, the database community was actively investigating the management of semistructural data, with Object Exchange Model (OEM) [26] being one of the formalizations. An OEM data value is a directed graph (often just a tree) with edges labeled by tags, internal nodes containing unique identifiers, and leaf nodes containing atomic values (integers, strings, etc.).

As has also been argued in the programming languages world, the combination of ordinary algebraic types and objects within the ODMG data model proved too inflexible for working with semistructural data, as it involved encodings within the structural ODMG model, which are usually complex and difficult to manage and evolve. The Ozone project [21] approached this problem by integrating the OEM data model into the ODMG data model. Their solution is similar to ours at the level of values: first, the OEM model is generalized to allow arbitrary ODMG values, including objects and structural values as leaves; second, a special ODMG class, `OEM`, is designated to hold all OEM values. The OEM values are ultimately implemented as objects of `OEM` subclasses. The OEM data model, however, is not statically typed. The motivation for Ozone was to allow convenient manipulation of semistructural data in an object-oriented database while avoiding the overly strict ODMG typing restrictions. The XDUCE type system, on the other hand, is indicating that it might be the right choice for typing OEM-like data (XML), without sacrificing its flexibility. Our contribution can be seen as an observation that an Ozone-like integration of objects and semistructural data can be carried out in a fully typed way, as long as the appropriate alternative to algebraic types, namely regular expression types, is chosen.

Two ongoing language design efforts that are very close to Xtatic in intentions and approach are CDuce and Jwig. CDuce [4] starts from the XDUCE type system and extends it in several significant directions: it provides a full set of boolean operations (including intersection and difference) on types, as well as higher-order functional types and overloading in the style of λ-&. Like XTATIC, CDuce extends XDUCE's semantic interpretation of subtyping to the types of the whole extended language [12]. Jwig [7] is an extension of Java intended for programming interactive sessions between Web servers and clients. Although quite different in style from XTATIC (it uses data flow analysis to check well formedness of XML expressions constructed by filling in "templates", rather than a conventional type system and tree expression language), the basic expressive power of Jwig's analysis is close to that of XDUCE's type system (see [7] for a detailed discussion of this point).

# 9    Future Work

We currently have a prototype interpreter for a fragment of XTATIC, including the extended data model described here. Though it still lacks most of the features of C$^{\#}$, the language implemented by the interpreter goes quite a bit beyond the simple FX core—in particular, it includes imperative features; we have used it to experiment with a number of small demos. Our short-term goals include handling a larger fragment of C$^{\#}$, building more ambitious demos, and replacing the simple interpreter by a back end targeting the .NET Common Language Runtime.

Another important near-term goal is to extend the type system to encompass a larger part of XML—most urgently, attributes. Hosoya and Murata [14] have recently proposed a typing mechanism and corresponding algorithms based on the attribute-element constraints of Relax NG; we hope to be able to adapt this proposal to XTATIC. We also plan to implement translators from standard XML schema languages (in particular, the XML-Schema standard) into XTATIC.

Our longer-term goals concentrate in two major areas: improving the efficiency of the underlying algorithms and run-time representations, and refining and extending the design of the core language. On the efficiency side, the main development currently in the works is high-performance compilation of pattern matching. We also need to come up with better run-time representations for certain special cases, while keeping compliance with the basic data model. One case in point is the PCDATA type. The typing and pattern-matching properties of the PCDATA definition given in Section 7 are attractive, but the naive representation that we sketched is clearly too heavy to perform well; something more clever will be needed. At the level of the core language design, there are also numerous questions to be considered. Can objects and trees be further unified? E.g., could pattern matching be used to extract object fields? Could attributes and fields be unified? Can we offer other kinds of pattern matching primitives, e.g. support for XPath? And, last but not least, can the XTATIC design be extended to cope with parametric polymorphism ("generics" in C$^{\#}$ parlance)?

# Acknowledgements

# References

[1] J. Aldrich, V. Kostadinov, and C. Chambers. Alias annotations for program understanding. In *ACM Symposium on Object Oriented Programming: Systems, Languages, and Applications (OOPSLA)*, Nov. 2002.

[2] D. Ancona, G. Lagorio, and E. Zucca. A core calculus for java exceptions. In *ACM Symposium on Object Oriented Programming: Systems, Languages, and Applications (OOPSLA)*, pages 16–30, 2001.

[3] D. Ancona and E. Zucca. True modules for java-like languages: Design and foundations, Aug. 2000. Technical Report DISI-TR-00-12, Dipartimento di Informatica e Scienze dellInformazione, Università di Genova.

[4] V. Benzaken, G. Castagna, and A. Frisch. CDuce: a white paper. `ftp://ftp.ens.fr/pub/di/users/castagna/cduce-wp.ps.gz`, 2002. Workshop on Programming Language Technologies for XML (PLAN-X).

[5] G. Bracha, M. Odersky, D. Stoutamire, and P. Wadler. Making the future safe for the past: Adding genericity to the Java programming language. In C. Chambers, editor, *ACM Symposium on Object Oriented Programming: Systems, Languages, and Applications (OOPSLA)*, ACM SIGPLAN Notices volume 33 number 10, pages 183–200, Vancouver, BC, Oct. 1998.

[6] R. Catell, editor. *The Object Database Standard: ODMG-93*. Morgan Kaufmann, 1994.

[7] A. S. Christensen, A. Moller, and M. I. Schwartzbach. Extending Java for high-level web service construction. `http://www.brics.dk/~mis/jwig.ps`, 2002.

[8] J. Clark. TREX: Tree Regular Expressions for XML. `http://www.thaiopensource.com/trex/`, 2001.

[9] J. Clark and M. Murata. RELAX NG. `http://www.relaxng.org`, 2001.

[10] Document object model (dom) level 3 validation specification, w3c working draft. `http://www.w3.org/TR/DOM-Level-3-Val`, 2002.

[11] M. F. Fernández, J. Siméon, and P. Wadler. A semi-monad for semi-structured data. In J. V. den Bussche and V. Vianu, editors, *Proceedings of 8th International Conference on Database Theory (ICDT 2001)*, volume 1973 of *Lecture Notes in Computer Science*, pages 263–300. Springer, 2001.

[12] A. Frisch, G. Castagna, and V. Benzaken. Semantic subtyping. In *LICS*, 2002.

[13] H. Hosoya. *Regular Expression Types for XML*. PhD thesis, The University of Tokyo, Japan, 2000.

[14] H. Hosoya and M. Murata. Validation and boolean operations for attribute-element constraints. In *Workshop on Programming Language Technologies for XML (PLAN-X)*, 2002.

[15] H. Hosoya and B. Pierce. Regular expression pattern matching. In *ACM Symposium on Principles of Programming Languages (POPL), London, England*, 2001. Full version to appear in *Journal of Functional Programming*.

[16] H. Hosoya and B. C. Pierce. XDuce: A typed XML processing language (preliminary report). In D. Suciu and G. Vossen, editors, *International Workshop on the Web and Databases (WebDB)*, May 2000. Reprinted in *The Web and Databases, Selected Papers*, Springer LNCS volume 1997, 2001.

[17] H. Hosoya and B. C. Pierce. Xduce: A statically typed xml processing language. *ACM Transactions on Internet Technology*, 2002. Submitted for publication.

[18] H. Hosoya, J. Vouillon, and B. C. Pierce. Regular expression types for XML. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 2001. To appear; short version in ICFP 2000.

[19] A. Igarashi, B. Pierce, and P. Wadler. Featherweight Java: A minimal core calculus for Java and GJ. In *ACM Symposium on Object Oriented Programming: Systems, Languages, and Applications (OOPSLA)*, Oct. 1999. Full version in ACM Transactions on Programming Languages and Systems (TOPLAS), 23(3), May 2001.

[20] A. Igarashi and B. C. Pierce. On inner classes. In *European Conference on Object-Oriented Programming (ECOOP)*, 2000. Also in informal proceedings of the Seventh International Workshop on Foundations of Object-Oriented Languages (FOOL). To appear in *Information and Computation*.

[21] T. Lahiri, S. Abiteboul, and J. Widom. Ozone: integrating structured and semistructured data. In *International Workshop on Database Programming Languages*. 1999.

[22] C. League, Z. Shao, and V. Trifonov. Type-preserving compilation of Featherweight Java. *ACM Transactions on Programming Languages and Systems*, 24(2):112–152, 2002.

[23] E. Meijer and M. Shields. XM$\lambda$: A functional programming language for constructing and manipulating XML documents. Submitted for publication, 1999.

[24] M. Odersky. Report on the programming language scala. `http://lamp.epfl.ch/~odersky/scala/reference.ps`, 2002.

[25] M. Odersky and P. Wadler. Pizza into Java: Translating theory into practice. In *ACM Symposium on Principles of Programming Languages (POPL), Paris, France*, 1997.

[26] Y. Papaconstantinou, H. Garcia-Molina, and J. Widom. Object exchange across heterogeneous information sources. In *International Conference on Data Engineering*, Mar. 1995.

[27] Relaxer. `http://www.asahi-net.or.jp/~dp8t-asm/java/tools/Relaxer/index.html`.

[28] U. P. Schultz. Partial evaluation for class-based object-oriented languages. In *Programs as Data Objects (PADO), Aarhus, Denmark*, volume 2053 of *Lecture Notes in Computer Science*, pages 173–197, 2001.

[29] M. Shields and E. Meijer. Type-indexed rows. In *ACM Symposium on Principles of Programming Languages (POPL), London, England*, Jan 2001.

[30] T. Studer. Constructive foundations for featherweight java. In R. Kahle, P. Schroeder-Heister, and R. Stärk, editors, *Proof Theory in Computer Science*. Springer-Verlag, 2001. Lecture Notes in Computer Science, volume 2183.

[31] I. Sun Microsystems. The Java architecture for XML binding (JAXB). `http://java.sun.com/xml/jaxb`, 2001.

[32] N. Tabuchi, E. Sumii, and A. Yonezawa. Regular expression types for strings in a text processing language. In J. V. den Bussche and V. Vianu, editors, *Proceedings of Workshop on Types in Programming (TIP)*, pages 1–18, July 2002.

[33] M. Tofte. Type inference for polymorphic references. *Information and Computation*, 89(1), Nov. 1990.

[34] M. Wallace and C. Runciman. 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., Sept. 27–29 1999. ACM Press.

[35] Extensible markup language (XML™), Feb. 1998. XML 1.0, W3C Recommendation, `http://www.w3.org/XML/`.

[36] XML Schema Part 0: Primer, W3C Working Draft. `http://www.w3.org/TR/xmlschema-0/`, 2000.

[37] XSL Transformations (XSLT), 1999. `http://www.w3.org/TR/xslt`.

[38] M. Zenger. Type-safe prototype-based component evolution. In *European Conference on Object-Oriented Programming (ECOOP), Malaga, Spain*, June 2002.

$$\frac{\text{D} \sqsubseteq: \text{C}}{\text{new } \text{D}(\overline{\text{a}}) \in \text{C} \Rightarrow \bullet} \quad \text{(P-ObjClass)}$$

$$\frac{\text{Seq} \sqsubseteq: \text{C}}{[\overline{\text{t}}] \in \text{C} \Rightarrow \bullet} \quad \text{(P-RecClass)}$$

$$\frac{\text{a} \in \text{C} \Rightarrow \Sigma}{\text{a} \in \text{C} \ \text{x} \Rightarrow \text{x}:\text{a}, \Sigma} \quad \text{(P-BindClass)}$$

$$\frac{\begin{array}{c} REtype(\text{X}) = \text{T} \\ [\overline{\text{t}}] \in [\text{T}] \Rightarrow \Sigma \end{array}}{[\overline{\text{t}}] \in [\text{X}] \Rightarrow \Sigma} \quad \text{(P-REtype)}$$

$$\frac{\begin{array}{c} \text{a} \in \text{Q} \Rightarrow \Sigma_1 \\ [\overline{\text{t}}] \in [\text{T}] \Rightarrow \Sigma_2 \end{array}}{[\langle\text{a}\rangle[\overline{\text{t}}]] \in [\langle\text{Q}\rangle[\text{T}]] \Rightarrow \Sigma_1, \Sigma_2} \quad \text{(P-Tree)}$$

$$[\,] \in [()] \Rightarrow \bullet \quad \text{(P-Eps)}$$

$$\frac{\begin{array}{c} [\overline{\text{t}}_1] \in [\text{P}_1] \Rightarrow \Sigma_1 \\ [\overline{\text{t}}_2] \in [\text{P}_2] \Rightarrow \Sigma_2 \end{array}}{[\overline{\text{t}}_1 \ \overline{\text{t}}_2] \in [\text{P}_1 \ \text{P}_2] \Rightarrow \Sigma_1, \Sigma_2} \quad \text{(P-Cat)}$$

$$\frac{[\overline{\text{t}}] \in [\text{P}_1] \Rightarrow \Sigma}{[\overline{\text{t}}] \in [\text{P}_1|\text{P}_2] \Rightarrow \Sigma} \quad \text{(P-Alt1)}$$

$$\frac{\begin{array}{c} [\overline{\text{t}}] \notin [\text{P}_1] \\ [\overline{\text{t}}] \in [\text{P}_2] \Rightarrow \Sigma \end{array}}{[\overline{\text{t}}] \in [\text{P}_1|\text{P}_2] \Rightarrow \Sigma} \quad \text{(P-Alt2)}$$

$$\frac{\forall k \in 1 \dots n. \ [\overline{\text{t}}_k] \in [\text{T}]}{[\overline{\text{t}}_1 \dots \overline{\text{t}}_n] \in [\text{T*}] \Rightarrow \bullet} \quad \text{(P-Rep)}$$

$$\frac{[\overline{\text{t}}] \in [\text{P}] \Rightarrow \Sigma}{[\overline{\text{t}}] \in [\text{P} \ \text{x}] \Rightarrow \text{x} : [\overline{\text{t}}], \Sigma} \quad \text{(P-BindRE)}$$

Figure 4: Pattern-matching relation $\text{a} \in \text{R} \Rightarrow \Sigma$.

# A    Appendix: Additional Definitions

This appendix presents in full some parts of the FX language definition that were just sketched in the body. These definitions are included only for the sake of completeness (and in case the referees may be curious) — the main points in the body of the paper do not depend on them.

## A.1    Patterns

Figure 4 defines the pattern-matching relation $\text{a} \in \text{R} \Rightarrow \Sigma$, read "an FX value expression $\text{a}$ matches a pattern R, yielding an environment $\Sigma$".

The simple type inference algorithm that we use in this definition of XTATIC is shown in Figure 5 in the form of the relation $\triangleright \text{R} \Rightarrow \Gamma$. According to this definition, a type associated by $\Gamma$ to a variable $\text{x}$ is exactly the type to which $\text{x}$ is bound in R. The *type join* operation $\sqcup$ used in the rule [PI-Alt] is given in Figure 6).

## A.2    Static semantics

A *static typing context* conforming to a static data context *DatCtx* is a tuple

$$TypCtx = \langle fields, mtype \rangle$$

where, assuming $DatCtx = \langle Classes, \sqsubseteq:, Typenames, def \rangle$,

- *fields* : $\text{C} \mapsto \overline{\text{F}} \ \overline{\text{f}}$ is a function defined on *Classes* with *fields*(C) being the list $\text{F}_1 \ \text{f}_1, \dots, \text{F}_n \ \text{f}_n$ of field types and field names such that

  - *fields*(Object) and *fields*(Seq) are empty;
  - if $\text{C} \sqsubseteq: \text{D}$, then the list *fields*(D) is a prefix of the list *fields*(C);

- *mtype* : $\text{C} \mapsto \text{m} \mapsto (\overline{\text{A}} \rightarrow \text{A})$ is a function defined on *Classes* with *mtype*(C) being a partial mapping from method names $\text{m}$ to method signatures $(\overline{\text{A}} \rightarrow \text{A})$ such that

  - if $\text{C} \sqsubseteq: \text{D}$, then $dom(\text{D}) \subseteq dom(\text{C})$ and, for any $\text{m} \in dom(\text{D})$, $mtype(\text{C})(\text{m}) = mtype(\text{D})(\text{m})$.

Figure 7 defines the relation $\Gamma \vdash \text{e} \in \text{A}$, read "in the typing environment $\Gamma$, expression $\text{e}$ has type $\text{A}$". Auxiliary operations used in the rules are defined in Figure 6.

$$\rhd \mathtt{C} \Rightarrow \bullet \qquad \text{(PI-Class)}$$

$$\frac{\rhd \mathtt{C} \Rightarrow \Gamma}{\rhd \mathtt{C}\ \mathtt{x} \Rightarrow \mathtt{x} \colon \mathtt{C}, \Gamma} \qquad \text{(PI-BindClass)}$$

$$\frac{REtype(\mathtt{X}) = \mathtt{T} \qquad \rhd [\mathtt{T}] \Rightarrow \Gamma}{\rhd [\mathtt{X}] \Rightarrow \Gamma} \qquad \text{(PI-REtype)}$$

$$\frac{\rhd \mathtt{Q} \Rightarrow \Gamma_1 \qquad \rhd [\mathtt{T}] \Rightarrow \Gamma_2}{\rhd [\texttt{<Q>}[\mathtt{T}]] \Rightarrow \Gamma_1, \Gamma_2} \qquad \text{(PI-Tree)}$$

$$\rhd [] \Rightarrow \bullet \qquad \text{(PI-Eps)}$$

$$\frac{\rhd [\mathtt{P_1}] \Rightarrow \Gamma_1 \qquad \rhd [\mathtt{P_2}] \Rightarrow \Gamma_2}{wrt[\mathtt{P_1\ P_2}] \Rightarrow \Gamma_1, \Gamma_2} \qquad \text{(PI-Cat)}$$

$$\frac{\rhd [\mathtt{P_1}] \Rightarrow \Gamma_1 \qquad \rhd [\mathtt{P_2}] \Rightarrow \Gamma_2}{\rhd [\mathtt{P_1 \mid P_2}] \Rightarrow \Gamma_1 \mid \Gamma_2} \qquad \text{(PI-Alt)}$$

$$\rhd [\mathtt{T*}] \Rightarrow \bullet \qquad \text{(PI-Rep)}$$

$$\frac{\rhd [\mathtt{P}] \Rightarrow \Gamma}{\rhd [\mathtt{P}\ \mathtt{x}] \Rightarrow \mathtt{x} \colon tyof(\mathtt{P}), \Gamma} \qquad \text{(PI-BindRE)}$$

Figure 5: Declared type inference $\rhd \mathtt{R} \Rightarrow \Gamma$.

$$
\begin{aligned}
\mathtt{C} \sqcup \mathtt{D} &= \sup{}_{\sqsubseteq:}\{\mathtt{C}, \mathtt{D}\} \\
\mathtt{C} \sqcup [\mathtt{T}] &= \sup{}_{\sqsubseteq:}\{\mathtt{C}, \mathtt{Seq}\} \\
[\mathtt{T}] \sqcup \mathtt{C} &= \sup{}_{\sqsubseteq:}\{\mathtt{Seq}, \mathtt{C}\} \\
[\mathtt{T_1}] \sqcup [\mathtt{T_2}] &= [\mathtt{T_1 \mid T_2}] \\
\Gamma_1 \sqcup \Gamma_2 &= \lambda \mathtt{x}.\Gamma_1(\mathtt{x}) \sqcup \Gamma_2(\mathtt{x}) \\
\\
ClassOf(\mathtt{C}) &= \mathtt{C} \\
ClassOf([\mathtt{T}]) &= \mathtt{Seq}
\end{aligned}
$$

Figure 6: Auxiliary typing definitions.

$$\frac{\Gamma(\mathtt{x}) = \mathtt{A}}{\Gamma \vdash \mathtt{x} \in \mathtt{A}} \qquad\qquad (\text{T-VAR})$$

$$\frac{\begin{array}{cc} \mathtt{C} \neq \mathtt{Seq} & \mathit{fields}(\mathtt{C}) = \overline{\mathtt{F}}\ \overline{\mathtt{f}} \\ \Gamma \vdash \overline{\mathtt{e}} \in \overline{\mathtt{A}} & \overline{\mathtt{A}} <: \overline{\mathtt{F}} \end{array}}{\Gamma \vdash \mathtt{new}\ \mathtt{C}(\overline{\mathtt{e}}) \in \mathtt{C}} \qquad\qquad (\text{T-NEW})$$

$$\frac{\Gamma \vdash \mathtt{e} \in \mathtt{C} \qquad \mathit{fields}(\mathtt{C}) = \overline{\mathtt{F}}\ \overline{\mathtt{f}}}{\Gamma \vdash \mathtt{e}.\mathtt{f}_k \in \mathtt{F}_k} \qquad\qquad (\text{T-FIELD})$$

$$\frac{\begin{array}{cc} \Gamma \vdash \mathtt{e} \in \mathtt{C} & \mathit{mtype}(\mathtt{C})(\mathtt{m}) = \overline{\mathtt{F}} \rightarrow \mathtt{B} \\ \Gamma \vdash \overline{\mathtt{d}} \in \overline{\mathtt{A}} & \overline{\mathtt{A}} <: \overline{\mathtt{F}} \end{array}}{\Gamma \vdash \mathtt{e}.\mathtt{m}(\overline{\mathtt{d}}) \in \mathtt{B}} \qquad\qquad (\text{T-INVK})$$

$$\Gamma \vdash \texttt{[]} \in \texttt{[()]} \qquad\qquad (\text{T-EPS})$$

$$\frac{n \geq 2 \qquad \Gamma \vdash \texttt{[}\mathtt{e}_k\texttt{]} \in \texttt{[}\mathtt{T}_k\texttt{]}}{\Gamma \vdash \texttt{[}\overline{\mathtt{e}}\texttt{]} \in \texttt{[}\mathtt{T}_1\ \ldots\ \mathtt{T}_n\texttt{]}} \qquad\qquad (\text{T-SEQ})$$

$$\frac{\begin{array}{cc} \Gamma \vdash \mathtt{d} \in \mathtt{A} & \mathit{ClassOf}(\mathtt{A}) = \mathtt{C} \\ \multicolumn{2}{c}{\Gamma \vdash \texttt{[}\overline{\mathtt{e}}\texttt{]} \in \texttt{[}\mathtt{T}\texttt{]}} \end{array}}{\Gamma \vdash \texttt{[<}\mathtt{d}\texttt{>[}\overline{\mathtt{e}}\texttt{]]} \in \texttt{[<}\mathtt{C}\texttt{>[}\mathtt{T}\texttt{]]}} \qquad\qquad (\text{T-TREE})$$

$$\frac{\Gamma \vdash \texttt{[}\mathtt{e}\texttt{]} \in \texttt{[}\mathtt{T}\texttt{]}}{\Gamma \vdash \texttt{[[}\mathtt{e}\texttt{]]} \in \texttt{[}\mathtt{T}\texttt{]}} \qquad\qquad (\text{T-COLLAPSE})$$

$$\frac{\begin{array}{cc} \Gamma \vdash \mathtt{d} \in \texttt{[}\mathtt{T}\texttt{]} & \mathtt{T} \triangleright \overline{\mathtt{P}}\ \mathit{ok} \\ \mathtt{T} \triangleright \texttt{[}\mathtt{P}_k\texttt{]} \Rightarrow \Gamma_k & \Gamma, \Gamma_k \vdash \mathtt{e}_k \in \mathtt{A}_k \end{array}}{\Gamma \vdash \mathtt{match}(\mathtt{d})\{\overline{\mathtt{case}\ \texttt{[}\overline{\mathtt{P}}\texttt{]}:\overline{\mathtt{e}}}\} \in \mathtt{A}_1 \sqcup \mathtt{A}_2 \sqcup \ldots \sqcup \mathtt{A}_n} \qquad\qquad (\text{T-MATCH})$$

Figure 7: Expression typing relation $\Gamma \vdash \mathtt{e} \in \mathtt{A}$.

$$\frac{\Sigma(\mathtt{x}) = \mathtt{a}}{\Sigma \vdash \mathtt{x} \Downarrow \mathtt{a}} \qquad\qquad (\text{E-Var})$$

$$\frac{\Sigma \vdash \overline{\mathtt{e}} \Downarrow \overline{\mathtt{a}}}{\Sigma \vdash \mathtt{new\ C(\overline{e})} \Downarrow \mathtt{new\ C(\overline{a})}} \qquad\qquad (\text{E-New})$$

$$\frac{\Sigma \vdash \mathtt{e} \Downarrow \mathtt{new\ C(\overline{a})} \qquad \mathit{fields}(\mathtt{C}) = \overline{\mathtt{F}}\ \overline{\mathtt{f}}}{\Sigma \vdash \mathtt{e.f}_k \Downarrow \mathtt{a}_k} \qquad\qquad (\text{E-Field})$$

$$\frac{\begin{array}{c}\Sigma \vdash \mathtt{e}_0 \Downarrow \mathtt{new\ C(\overline{a})} \qquad \mathit{mbody}(\mathtt{C,m}) = (\overline{\mathtt{x}}, \mathtt{e}) \\ \Sigma \vdash \overline{\mathtt{d}} \Downarrow \overline{\mathtt{b}} \qquad \overline{\mathtt{x}} : \overline{\mathtt{b}}, \mathtt{this} : \mathtt{new\ C(\overline{a})} \vdash \mathtt{e} \Downarrow \mathtt{a}\end{array}}{\Sigma \vdash \mathtt{e}_0.\mathtt{m}(\overline{\mathtt{d}}) \Downarrow \mathtt{a}} \qquad\qquad (\text{E-J-Invk})$$

$$\frac{n \geq 2 \qquad \Sigma \vdash [\mathtt{e}_k\} \Downarrow [\mathtt{t}_k]}{\Sigma \vdash [\overline{\mathtt{e}}] \Downarrow [\overline{\mathtt{t}}_1, \ldots, \overline{\mathtt{t}}_n]} \qquad\qquad (\text{E-Seq})$$

$$\frac{\Sigma \vdash \mathtt{e} \Downarrow \mathtt{a} \qquad \Sigma \vdash [\overline{\mathtt{d}}] \Downarrow [\overline{\mathtt{t}}]}{\Sigma \vdash [\mathtt{<e>}[\overline{\mathtt{d}}]] \Downarrow [\mathtt{<a>}[\overline{\mathtt{t}}]]} \qquad\qquad (\text{E-Tree})$$

$$\frac{\Sigma \vdash [\mathtt{e}] \Downarrow [\overline{\mathtt{t}}]}{\Sigma \vdash [[\mathtt{e}]] \Downarrow [\overline{\mathtt{t}}]} \qquad\qquad (\text{E-Collapse})$$

$$\frac{\begin{array}{c}\Sigma \vdash \mathtt{d} \Downarrow [\overline{\mathtt{t}}] \\ [\overline{\mathtt{t}}] \notin [\mathtt{P}_1] \quad \ldots \quad [\overline{\mathtt{t}}] \notin [\mathtt{P}_{i-1}] \qquad [\overline{\mathtt{t}}] \in [\mathtt{P}_k] \Rightarrow \Sigma' \\ \Sigma, \Sigma' \vdash \mathtt{e}_k \Downarrow \mathtt{a}\end{array}}{\Sigma \vdash \{\mathtt{match(d)\{case\ }[\overline{\mathtt{P}}] : \quad \overline{\mathtt{e}}\} \Downarrow \mathtt{a}} \qquad\qquad (\text{E-Match})$$

Figure 8: Expression evaluation relation $\Sigma \vdash \mathtt{e} \Downarrow \mathtt{a}$.

## A.3   Dynamic semantics

A *static execution context* conforming to static typing and data contexts *TypCtx* and *DatCtx* (with *TypCtx* conforming to *DatCtx*) is a tuple

$$ExeCtx = \langle mbody \rangle$$

where, assuming $DatCtx = \langle Classes, \sqsubseteq:, Typenames, def \rangle$ and $TypCtx = \langle fields, mtype \rangle$,

- $mbody : \mathtt{C} \mapsto \mathtt{m} \mapsto (\mathtt{D}, \overline{\mathtt{x}}, \mathtt{e})$ is a function defined on *Classes* with $mbody(\mathtt{C})$ being a partial mapping that associates a type name $\mathtt{m}$ with the triple $(\mathtt{D}, \overline{\mathtt{x}}, \mathtt{e})$ of class $\mathtt{D}$ where the method definition is located, method parameters $\mathtt{x}$, and method body $\mathtt{e}$, such that

  - $dom(mbody(\mathtt{C})) = dom(mtype(\mathtt{C}))$ for all $\mathtt{C} \in Classes$,
  - for any $\mathtt{C}$ and $\mathtt{m}$, if $mtype(\mathtt{C})(\mathtt{m}) = \overline{\mathtt{A}} \rightarrow \mathtt{A}$ and $mbody(\mathtt{C})(\mathtt{m}) = (\mathtt{D}, \overline{\mathtt{x}}, \mathtt{e})$, then $\mathtt{this} : \mathtt{D}, \overline{\mathtt{x}} : \overline{\mathtt{A}} \vdash \mathtt{e} \in \mathtt{B}$ for some $\mathtt{B} <: \mathtt{A}$.

Figure 8 defines the evaluation relation $\Sigma \vdash \mathtt{e} \Downarrow \mathtt{a}$, read "in environment $\Sigma$, expression $\mathtt{e}$ evaluates to the value $\mathtt{a}$".

## A.4   The "stuck" relation

The notion "evaluation of $\mathtt{e}$ gets stuck in the finite number of steps" is formalized by the following *stuck evaluation* relation $\Sigma \vdash \mathtt{e} \not\Downarrow$:

- $\Sigma \vdash \mathtt{new\ C(\overline{e})} \not\Downarrow$ if $\Sigma \vdash \mathtt{e}_i \not\Downarrow$ for either one of $\mathtt{e}_i$ from $\overline{\mathtt{e}}$;

- $\Sigma \vdash$ `d.f` $\Downarrow\!\!\!/$ if either

  1. $\Sigma \vdash$ `d` $\Downarrow\!\!\!/$ or
  2. $\Sigma \vdash$ `d` $\Downarrow$ `[`$\overline{\texttt{t}}$`]`, or
  3. $\Sigma \vdash$ `d` $\Downarrow$ `new C(`$\overline{\texttt{b}}$`)`, but `f` is not among *fields*(`C`);

- $\Sigma \vdash$ `d.m(`$\overline{\texttt{e}}$`)` $\Downarrow\!\!\!/$ if either

  1. $\Sigma \vdash$ `d` $\Downarrow\!\!\!/$ or
  2. $\Sigma \vdash$ `d` $\Downarrow$ `[`$\overline{\texttt{t}}$`]`, or
  3. $\Sigma \vdash$ `d` $\Downarrow$ `new C(`$\overline{\texttt{b}}$`)` but *mbody*(`C`)(`m`) is not defined, or
  4. $\Sigma \vdash$ `d` $\Downarrow$ `new C(`$\overline{\texttt{b}}$`)`, *mbody*(`C`)(`m`) is defined, but $\Sigma \vdash$ `e`$_i$ $\Downarrow\!\!\!/$ for either one of `e`$_i$ from $\overline{\texttt{e}}$;
  5. $\Sigma \vdash$ `d` $\Downarrow$ `new C(`$\overline{\texttt{b}}$`)`, *mbody*(`C`)(`m`) $=$ (`$\overline{\texttt{x}}$`,e)`, $\Sigma \vdash \overline{\texttt{e}} \Downarrow \overline{\texttt{a}}$, but `this : new C(`$\overline{\texttt{b}}$`),`$\overline{\texttt{x}} : \overline{\texttt{a}} \vdash$ `e` $\Downarrow\!\!\!/$

- $\Sigma \vdash$ `[<d>[`$\overline{\texttt{e}}$`]]` $\Downarrow\!\!\!/$ if either

  1. $\Sigma \vdash$ `d` $\Downarrow\!\!\!/$, or
  2. $\Sigma \vdash$ `d` $\Downarrow$ `a`, but $\Sigma \vdash \mid e \mid \Downarrow\!\!\!/$;

- $\Sigma \vdash$ `[`$\overline{\texttt{e}}$`]` $\Downarrow\!\!\!/$ if $|\overline{\texttt{e}}| \geq 2$ and $\Sigma \vdash$ `e`$_i$ $\Downarrow\!\!\!/$ for either one of `e`$_i$ from $\overline{\texttt{e}}$;

- $\Sigma \vdash$ `[[`$\overline{\texttt{e}}$`]]` $\Downarrow\!\!\!/$ if $\Sigma \vdash$ `[`$\overline{\texttt{e}}$`]` $\Downarrow\!\!\!/$

- $\Sigma \vdash$ `match(d){case [`$\overline{\texttt{P}}$`]:` $\overline{\texttt{e}}$`}` $\Downarrow\!\!\!/$ if either

  1. $\Sigma \vdash$ `d` $\Downarrow\!\!\!/$, or
  2. $\Sigma \vdash$ `d` $\Downarrow$ `new C(`$\overline{\texttt{b}}$`)`, or
  3. $\Sigma \vdash$ `d` $\Downarrow$ `[`$\overline{\texttt{t}}$`]`, but `[`$\overline{\texttt{t}}$`]` $\in$ `[P`$_i$`]` $\not\Rightarrow$ for all $i$, or
  4. $\Sigma \vdash$ `d` $\Downarrow$ `[`$\overline{\texttt{t}}$`]`, and, for some $i$, `[`$\overline{\texttt{t}}$`]` $\in$ `[P`$_1$`]` $\not\Rightarrow$, `[`$\overline{\texttt{t}}$`]` $\in$ `[P`$_{i-1}$`]` $\not\Rightarrow$, `[`$\overline{\texttt{t}}$`]` $\in$ `[P`$_i$`]` $\Rightarrow \Sigma'$, but $\Sigma, \Sigma' \vdash$ `e`$_i$ $\Downarrow\!\!\!/$.