

**Universitetet i Oslo
Institutt for informatikk**

**Towards integration
of XML in the Creol
object-oriented
language**

Arild Torjusen, Olaf
Owe, and Gerardo
Schneider

**Research report 365
ISBN 82-7368-323-0
(Revised version)**

**October 2007
Revised February 2008**



Towards integration of XML in the Creol object-oriented language

Arild Torjusen, Olaf Owe, and Gerardo Schneider

Department of Informatics, University of Oslo
PO Box 1080 Blindern, NO-0316 Oslo, Norway
email: {aribraat,olaf,gerardo}@ifi.uio.no

Abstract

The integration of XML documents in object-oriented programming languages is becoming paramount with the advent of the use of Internet in new applications like web services. Such integration is not easy in general and demands a careful language design. In this paper we propose an extension to Creol, a high level object-oriented modeling language for distributed systems, for handling XML documents.

1 Introduction

XML (eXtensible Markup Language) [7] is a flexible and generic format for structured data aimed at being shared on the World Wide Web and intranets. The need for XML documents as first-class citizens is acknowledged by academic as well as by business-oriented communities [21].

XML *documents* are ordered labeled tree structures containing *markup* symbols describing their content. The document structure is described by a document type -or *schema*- written in a schema language. Many such languages have been proposed, among them DTD (Document Type Definition) [7] and XML-Schema [10]. Unlike other markup languages (like HTML), XML has no restrictions on the tags or attributes used to mark up a document. One remarkable feature of XML is its plain-text-based nature. The advantage is that there is no problem with proprietary nor deciphering data. The disadvantages are the large bandwidth needed for transmission of documents and the need of encryption because of security issues. Part of the manipulation of XML documents includes the retrieval of information through queries. XQuery [5] provides a sound foundation for XML query, based on infosets. The situation is not ideal for developers since they need to know one language for analyzing the tuples e.g., SQL, another language for the Infoset e.g., XQuery, and a third one for operating on objects e.g., Java. Some attempts have been done to combine object-oriented languages and XML, but this turned out to be a complex task; this problem is known as the *impedance mismatch* [25], which arises when trying to combine object-oriented programming languages and (relational) databases.

The integration of XML on current object-oriented languages is far from trivial. The initial approach has been to treat XML through APIs which uses strings for representing literals. One problem of this approach is that it limits the use of static checking tools. Furthermore, the representation of programs as text involves potential security risks. See [21] for a more detailed description of the main problems arising with the integration of XML in object-oriented languages.

In addition to the integration of XML documents within OOP languages, another question is what to do with these data, i.e., how easy it is to make queries, getting useful information from such XML-documents.

1.1 Creol

Our research project concerns integration of XML into the object-oriented language Creol [16, 15, 18]. The main features of Creol are:

- It supports both object-oriented classes, with late binding and multiple inheritance, as well as user defined data types and functions. This gives flexibility in our choices when representing XML.
- It is oriented towards open distributed systems. Exchange of XML documents fits naturally in this context.
- It supports concurrency and method calls based on asynchronous communication. We wish to explore the processing and sharing of XML documents in this setting.
- It is strongly typed, supporting subtypes and subinterfaces, with a type hierarchy including both by means of the universal type, *Data*.
- It has a formal operational semantics, defined in rewriting logics. This enables us to formalize the extension to XML by reuse of the operational semantics.
- It has a small kernel with an operational semantics consisting of only 11 rewrite rules. This makes it easy to extend and modify the language and the semantics.
- Creol has an executable interpreter defined in the Maude language. This provides a useful framework for implementation and testing of our XML representations.

1.2 Related Work

The list of languages for processing XML documents is extensive, so it is not possible to be exhaustive here. We briefly discuss below some of the most influential works, namely XDuce, CDuce and C_ω . We mention other related work as reference for further reading, without entering into detail.

XDuce XDuce [13] is a functional programming language for XML processing. Its basic data values are XML documents and its types—called *regular expressions types*—correspond to document schemas. The language is statically typed but it also provides dynamic type-checking. Other interesting feature of XDuce is regular expression pattern matching which includes tag checking, subtree extraction and conditional branching.

An XML document in XDuce is represented as a sequence of nodes, and types use similar constructs as string regular expressions like “*” for representing that zero or more occurrences may happen, “?” for indicating an item may be omitted, “+” for one or more time repetition, “|” for alternation and “,” for concatenation. The main difference with string regular expressions, is that regular expression types describe sequences of tree nodes instead of sequences of characters.

The type-checking algorithm is based on the following subtype relationship: one type is a subtype of another if and only if the former denotes a subset of the latter. The subtype checker may be used both for checking that the actual type of a function's body is a subtype of the programmer-declared result type and for verifying function call arguments against parameter types given by the programmer. Although the theoretical complexity of the corresponding problem to subtype checking on tree automata is exponential, it is claimed in [13] that it works well in practice.

CDuce CDuce [3] is a typed functional language born from an attempt to solve some of the limitations of XDuce [13]. It extends XDuce on three areas:

Type system In addition to regular expression types and type-based patterns, CDuce adds recursive types and other less XML specific constructs: products, records (open and closed), general Boolean connectives (intersection, union and difference) and arrow types. This extension takes care of not breaking down the nice subtype relation of XDuce.

Language design The following language constructions are included in CDuce: overloaded functions (useful for code sharing and reuse), iterators on sequences and trees and other extensions of the pattern algebra. Besides, XML tags are first-class citizens and strings are simple sequence of characters. The language support higher-order programming, so all functions are first-class citizens.

Run-time system A new approach for avoiding unnecessary computation at runtime is added in CDuce, allowing the programmer to use a more declarative style when writing patterns, without degrading performance. The underlying theory is based on a new kind of tree automata.

CDuce provides also a tool for translating DTDs into CDuce's types.

C_ω C_ω [24] is a programming language developed at Microsoft Research, combining features from two other research languages: (a) Polyphonic C# [2]: a control flow extension with asynchronous wide-area concurrency, and (b) Xen [21]: a data type extension for processing XML and table manipulation. Besides other interesting features, C_ω allows the construction of objects using XML syntax.

The C_ω type systems combines the following three data models: relational, object and XML data-access, and it is more oriented to XML constrained using W3C XML Schema. The language covers the following XML and XML Schema features: document order, distinction between elements and attributes, multiplicity of fields with equal name but different values and content models for specifying choice (union) types for fields.

One of the nice features of the C_ω type system are *streams*. It is possible to invoke methods on streams, which are applied to all the elements of the stream; XPath-style queries over objects graphs are easily written in this way. It also includes the concept of *apply-to-all* expressions construct. Choice (union) types allow the programmer to specify one of different possible values for a certain field. Moreover, *null* is a valid value for a type, which have been proved useful in XML and relational databases. Document order and multiplicity of equal names for child elements, are solved through the use of *anonymous structs*. In C_ω DTDs (and XML Schemas) are represented by *content classes*.

Other languages The following languages try to extend Java with XML processing: XJ [12], XACT [20], XOBJE [19], BPELJ [4].

XL [11] is a language whose only type system is the XML type system, and not a language whose syntax is described using XML vocabulary. It is specially designed for the implementation of Web services. XL is portable and fully compliant with all W3C standards such as XQuery, XML Protocol, and XML Schema.

PiDuce¹ is CDuce-like language based on the π -calculus. ECMAScript for XML (E4X) is a set of programming language extensions adding native XML support to ECMAScript. E4X is standardized by Ecma International in ECMA-357 standard.²

See [22] for a good survey on static type-checking for XML transformation languages.

1.3 Our Agenda

In order to integrate XML documents in Creol, we intend to follow the following agenda:

1. *Parsing and well-formedness checking.* We will enhance the language as to be able to take a given XML document as input and generate some internal data structure from it.
2. *Internal representation of XML in Creol.* We aim at extending Creol for supporting XML documents with the least possible changes to the existing framework. One of the key features we would like to preserve is Creol static type-safety. In order to make a lightweight integration of XML into Creol and keep static type safety we will restrict type checking of XML in this implementation to only *well-formedness* of XML values, i.e. that some value of type `XMLDoc` (the Creol type for XML documents) checks out as an `XMLDoc`.
3. *Simple validity-checking of XML data-structures.* We will validate XML data-structures against some schema. Schema is here taken in a broad sense, meaning a formal description of the type of an XML document, without regards to any specific schema language as e.g. DTD, XML-Schema or RELAX NG (cf. Sec 3). Validity checking will be done by functions “on top” of the type system and not within the type system itself.
4. *More complex validity-checking of XML data-structures.* We will perform more complex validity checking after enhancing the Creol language with *regular expression types*, following the work of Hosoya et.al. [14].
5. *Queries.* We will also demonstrate how to perform queries and data extraction from XML document instances.³
6. *Transformations.* We will perform more complex operations such as construction and transformations on XML documents.

In this paper, however, we will concentrate on items 2 and 3 above. In the next section we show how XML documents are integrated in Creol. In Section 3 we show how schemas are represented in Creol after a short discussion on existing schema languages. Section 4 is concerned with the validation of XML documents. In Section 6 we conclude and present further work.

¹<http://www.cs.unibo.it/~laneve/PiDuce/>

²See <http://www.ecma-international.org/publications/standards/Ecma-357.htm>.

³Cf. e.g. <http://www.w3.org/TR/2005/WD-xquery-use-cases-20050915/> for test use cases.

2 A model for XML in Creol

Different XML documents may vary in physical representation due to syntactic changes permitted by the XML standard. W3C has issued a recommendation which describes how any XML document can be normalized into a canonical form [6]. The data model defined in the XPath 1.0 Recommendation [26] is the basis for canonical XML and we will use this as the point of departure for the internal representation of XML in Creol.

2.1 The XPath Data model

XPath models an XML document as an ordered tree containing nodes of seven different types:

- **root:** The root node is the root of the tree and will correspond to an XML document instance. It contains a list of processing instructions, a list of comment nodes, and exactly one element which is the root element of the document.
- **element:** The element node has a name (corresponding to the XML tag for the element) and may have as its children element nodes, comment nodes, processing instruction (PI) nodes and text nodes. It is also associated to a set of attribute nodes and a set of namespace nodes.
- **text:** A text node contains a string, representing character data in the XML document.
- **attribute:** An attribute node contains a name and a value.
- **namespace:** A namespace node contains a string value for the namespace prefix and a value for the namespace URI.
- **processing instructions:** A PI node has a name identifying the target application and a string which is to be passed to the application.
- **comment:** A comment node contains a string.

To simplify the initial XML implementation for Creol we will leave out the last three kinds of nodes from our model. According to [7], comments “are not part of the document’s character data; an XML processor MAY, but need not, make it possible for an application to retrieve the text of comments.”, we choose not to retain comments in the Creol representation of XML. Processing instructions are not relevant for our purpose of demonstrating lightweight integration of XML in Creol and can also be left out. As will be explained later we will adopt the DTD language for specification of schemas; since the DTD does not support namespaces it is natural not to represent namespace nodes in the model. These design choices also simplifies the definition of element and root nodes.

2.2 The Creol representation of XML

Given the two-tiered type-system of Creol where objects are typed by interfaces and local computations on terms occur in a functional language, we introduce XML into Creol by adding type constructors for a new `XMLDoc` type, as a subtype of the universal type `Data`, as well as functions on this type.

Creol has an operational semantics defined in rewriting logic, which is executable with Maude [9] and provides an interpreter and analysis platform for system models. So to accommodate XML we extend the operational semantics with some Maude sorts (type names) and constructors (Creol definitions would be very similar):

```

sorts    XMLName ElemNd TextNd AttNd ContentNd XMLDoc .
subsort ElemNd TextNd < ContentNd .
subsort String < XMLName .

```

introducing sorts for XML names, element, attribute, text and content nodes, letting `ElemNd` and `TextNd` be subsorts of `ContentNd`. The sort `XMLName` includes `String`, the predefined sort for strings.

To simplify the writing of XML values in a program we use mix-fix notation (indicating argument positions by underline symbols) to provide a compact syntax by adding the following constructors for attributes, text nodes and elements (with and without attributes).

```

op (=__)      : XMLName String                -> AttNd [ctor] .
op _(_)[_]    : XMLName AttNdList ContentNdList -> ElemNd [ctor] .
op _[_]       : XMLName ContentNdList         -> ElemNd [ctor] .
op tx         : String                        -> TextNd [ctor] .

```

where the clause `[ctor]` after an operator (`op`) indicates that it is a constructor, and where `ContentNdList` and `AttNdList` represent lists of `ContentNd` and `AttNd`, respectively, defined as conventional in Maude.

Note that there is no specific constructor for root nodes. Since we leave out processing instructions and comments, the root node is just the element node occurring at the root of an XML document tree. Thus, the XML document constructor is

```

op xmlDoc     : ElemNd XMLSchema -> XMLDoc [ctor] .

```

We define the operator

```

op noSchema   :                               -> XMLSchema [ctor] .

```

for XML documents with no `XMLSchema`. Other `XMLSchema` constructors are defined further below.

Example The following simple XML fragment

```

<email>
  <head>
    <sender>Arild</sender>
    <rcp addr="vera@foo.com">Vera</rcp>
    <subject>Test</subject></head>
  <body>
    <message>Hello there, you wrote in an earlier message:
    <quote>We'll meet again</quote> See you later</message>
  </body>
</email>

```


can be represented as an `ElemNd` term with the Creol/Maude syntax:

```
"email" [
  ("head" [
    ("sender" [tx("Arild")])
    ("rcp" ("addr"="vera@foo.com") [tx("Vera")])
    ("subject" [tx("Test")]))
  ("body" [
    ("message" [
      tx("Hello there, you wrote in an earlier message:")
      ("quote" [tx("We'll meet again")]) tx("See you later")])))] .
```

As conventional in Maude, the list constructor (concatenation) is here denoted by white space (blank).

3 Schemas and type checking

3.1 Regular expression types vs. schema types

Static type checking of XML documents in a programming language can be achieved by introducing types for XML fragments in the language. Xduce and CDuce mentioned earlier are examples of projects going in this direction, by introducing regular expressions types for XML schemas and letting the type system handle the validation.

For the current integration of XML in Creol we will take a less involved approach by introducing one data type for XML schema, together with functions to validate documents against schema. We may then specify a type for an XML document as a value of type `XMLSchema` and thus the validation takes place within the existing type system and does not constitute an addition to the type system itself. The advantage of this approach is that we do not need significant modifications to the type system, the disadvantage is that we get a less fine grained tool for working with XML schema.

3.2 Expressive power of schema languages

There exists several generally adopted XML schema languages with different expressive power. Murata, Lee, and Mani [23] suggest a taxonomy of schema languages based on the formal theory of regular tree grammars. Some of the most common schema languages can be ranked in order of increasing expressivity thus: The DTD language, The W3Cs XML Schema, The RELAX NG specification. Validation of the first two can be done by simple adaptations of word automata, while the last requires a more complicated tree automaton. However the DTD language is sufficiently expressive for our purpose which is to demonstrate how XML can be integrated in the object oriented modeling framework of Creol. Therefore in our model for XML schema values in Creol we adapt the restrictions inherent in the DTD language to achieve simple validation, (i.e. only deterministic regular expressions is allowed in the definition of an element as explained below).⁴

⁴Roughly corresponding to “Local Tree Grammars” in [23].

3.3 The schema type for Creol

A DTD is a list of markup declarations where markup declarations are either element type declarations, attribute-list declarations, entity declarations, or notation declarations.

For our purpose we only consider element type declarations and attribute-list declarations. Entity declarations may be considered as a kind of macro notation for strings that may appear in a DTD or an XML document, since our focus is on internal processing we will assume that these already are expanded by the parser and will abstract away from them in our model. Notation declarations are similarly a kind of shorthand for notations and are also left out. Accordingly the XML Schema constructor is:

```
op xmlSchema : XMLName ElemDeclList AttDeclList -> XMLSchema .
```

Element type declarations consist of a name referring to an element and a specification of the legal content. There are four kinds of specifications: either one of the designated keywords “EMPTY” or “ANY”, or the specification of a *content model*. A content model is a context free grammar governing the allowed types of the child elements and the order in which they are allowed to appear. The fourth kind of content specification is the *Mixed-content Declaration* which is of the form:

$$(\#PCDATA | e_1 | e_2 | \dots | e_n)^*$$

Where each e_i is an element name and n may be 0 in which case the ‘*’ is optional.

Example A DTD for the XML fragment given above could be:

```
<!ELEMENT email (head, body, foot*) >
<!ELEMENT head (sender, rcp, subject?)>
<!ELEMENT body (message)*>
<!ELEMENT foot (#PCDATA)>
<!ELEMENT sender (#PCDATA)>
<!ELEMENT rcp (#PCDATA)>
<!ELEMENT subject (#PCDATA)>
<!ELEMENT message (#PCDATA|quote)*>
<!ELEMENT quote (#PCDATA)>
```

The first three element declarations specify content models and the rest are instances of mixed-content declarations. We model the content models as *regular expressions*. Let Σ be an alphabet over element names, including the reserved name PCDATA. By including PCDATA in Σ we can model a mixed-content declaration as a special kind of a content model specification. The set of regular expressions over Σ^* are obtained in the standard way: The empty string ϵ and each member of Σ are regular expressions. If α is a regular expression, then so are (α) , $\alpha?$, α^* and α^+ . If α and β are regular expressions, then so is $\alpha\beta$, and $\alpha|\beta$. The operators $?$, $*$, and $+$ has higher precedence than concatenation. Concatenation has higher precedence than union ($|$). The regular expression combinators have the expected semantics. We model element declarations as follows:

```
subsort XMLName < ReToken < RegExp .
op elemDecl    : XMLName ContentModel -> ElemDecl [ctor] .
op empty      :                               -> ContentModel [ctor] .
```

```

op any      :                               -> ContentModel [ctor] .
op elemCt   : RegExp                        -> ContentModel [ctor] .
op PCDATA  :                               -> ReToken .
op _?      : RegExp                        -> RegExp [ctor prec 40] .
op _*      : RegExp                        -> RegExp [ctor prec 40] .
op _+      : RegExp                        -> RegExp [ctor prec 40] .
op _@_     : RegExp RegExp                 -> RegExp [ctor assoc prec 42]5
op _|_     : RegExp RegExp                 -> RegExp [ctor prec 44]

```

The XML specification adds the requirement that the content models must be deterministic [7, Appendix E], i.e. a content model must not allow an element to match more than one occurrence of an element name in the content model. This ensures that when matching an element name σ with a schema we do not have to look ahead beyond the σ in the input string to decide which regular expression in the content model matches σ . This requirement is included in the XML specification to ensure compatibility with SGML. For a detailed discussion see e.g. [8].

Example The Maude syntax for a document type declaration containing the DTD given above is:

```

xmlSchema("email", (
  elemDecl("email", elemCt("head"@("body"@("foot"*)))
  elemDecl("head", elemCt("sender"@("rcp"@("subject"?)))
  elemDecl("body", elemCt("message"*))
  elemDecl("foot", elemCt(PCDATA))
  elemDecl("sender", elemCt(PCDATA))
  elemDecl("rcp", elemCt(PCDATA))
  elemDecl("subject", elemCt(PCDATA))
  elemDecl("message", elemCt((PCDATA|"quote"*))
  elemDecl("quote", elemCt(PCDATA)), noAttDecl6) .

```

4 Validating XML in Creol

Well-formedness of any value of type XMLDoc is ensured by Maude type checking. The XML specification defines an XML document to be *valid* “if it has an associated document type declaration and if the document complies with the constraints expressed in it” [7].

The XML document constructor associates the root element of a document with a schema, (which may also be the special value `noSchema`). Hence, an XML document is validated by first checking for existence of a schema and by checking that the root node element name matches that schema name. Secondly we check that each element node in the tree is valid with respect to the element declarations in the schema.

Validation of a document is performed by the function

```

op validate : XMLDoc -> ValResult .

```

⁵We here use ‘@’ as the concatenation operator to avoid problems with overloading of ‘,’ or whitespace.

⁶Attribute declarations are not yet supported.

```

op res : Bool String -> ValResult .
eq collate( res(b,s) , res(b',s')) = res((b and b') , (s + s')) .

eq validate( xmlDoc(nm(atts)[cts], noSchema )) = res(false,"No Schema") .

eq validate( xmlDoc(nm(atts)[cts] , xmlSchema(nm',elDs,attDs) ) ) =
  if ( nm /= nm' ) then res(false,("Document root-element: " + nm +
    ", must match schema type: " + nm' + " \n"))
  else val(nm(atts)[cts] , elDs) fi .

eq val(tx(str),elDs) = res(true,"") .
eq val(emp,elDs) = res(true,"") .
eq val((ct cts) , elDs) = collate( val(ct,elDs) , val(cts,elDs)) [owise] .

ceq val(nm(atts)[cts],elDs) =
  if cm == undefined then
    res(false,("Element-type : " + nm + " must be declared.\n"))
  else check(nm(atts)[cts],cm,elDs) fi      if cm := getCM(nm,elDs) .

eq check(nm(atts)[cts],empty,elDs) =
  if (cts == emp) then res(true,"Empty elem: " + nm + "\n")
  else
    res(false,"Elem: " + nm + " declared as EMPTY, but has content.\n")
  fi .

eq check(nm(atts)[cts],any,elDs) =
  collate(res(true,"Elem: " + nm + " defined as ANY.\n"),val(cts,elDs)) .

eq check(nm(atts)[cts], elemCt(regex) ,elDs) =
  if match(getTokens(cts), regex) then
    collate (res(true, nm + ": (" + ctToS(cts) + ")
      matches [" + reToS(regex) + "]\n") , val(cts,elDs))
  else
    collate (res(false, nm + ": (" + ctToS(cts) + ")
      does NOT match [" + reToS(regex) + "]\n"), val(cts,elDs) ) fi .

```

Figure 1: Maude code for validation of XML documents.

where a `ValResult` is a boolean/string pair with the boolean value indicating validity and the string containing an error message or a record of the processing. `validate` checks whether there is a schema with a name matching the document root node associated with the document, in which case the recursive function `val` is called, otherwise validation ends with a negative result. The helper function `collate` builds the final validation result for a document from validation of its parts. The relevant parts of the Maude code are given in fig. 1. The function:

```

op val : ContentNdList ElemDeclList -> ValResult .

```

validates a content node list against the element declaration list defined by the schema. For a list of nodes, `val` is called recursively on each node in the list. For a single node, the element

type declaration corresponding to the node is retrieved (by name) from the list of element declarations and the node is checked against the retrieved declaration by a call to the function

```
op check : ContentNd ContentModel ElemDeclList -> ValResult .
```

If there is no `ContentModel` for some node the document is invalid. Note also that according to [7] an element type must not be declared more than once so uniqueness of declarations may be assumed.

In the call to `check`, the complete list of element declarations is passed on as a parameter since any child nodes to the node currently being processed must also be validated.

For a `ContentNd` to be valid relative to a `ContentModel` we need to consider three cases: The `ContentModel` is `empty` and the element should have no content or the `ContentModel` is `any` and the element can consist of any sequence of (declared) elements intermixed with character data. These two cases are easy to check. The third case is where the `ContentModel` specifies a regular expression, in this case the function `match` will be called to determine whether the list of actual children elements matches the regular expression specified in the corresponding element declaration, in addition `val` is called on the list of children elements.

The function

```
op getTokens : ContentNdList -> TokenList .
```

builds a list of tokens from the element content, i.e. it builds a list consisting of; element names for content nodes of sort `ElemNd`, and the special token `'PCDATA'` for content nodes of sort `TextNd`. As tokens we use the Maude built-in sort `Qid`. The token list and the regular expression from the element type declaration are then processed by the `match` function:

```
op match : TokenList RegExp -> Bool .
```

Matching of a list of element names from Σ against a regular expression is implemented by constructing a deterministic finite automaton from the regular expression and test whether the automaton accepts the string corresponding to the list of names. The implementation details are left out here, but see e.g. [1,17] for a description of how this is done in Maude. Our implementation is based on the work done in [1]. `ctToS` and `reToS` are just string conversion functions for content nodes and regular expressions for logging purposes.

Example Validation of the sample document with the DTD specified above gives the following result:

```
reduce in XML-VALIDATE-TEST : validate(xmlDoc(email, emailSchema)) .
rewrites: 9454 in 8ms cpu (8ms real) (1050561 rewrites/second)
result ValResult:
res(true, "email: (head ,body) matches [head @ body @ (foot*)]
  head: (sender ,rcp ,subject) matches [sender @ rcp @ (subject?)]
  sender: (PCDATA) matches [PCDATA]
  rcp: (PCDATA) matches [PCDATA]
  subject: (PCDATA) matches [PCDATA]
  body: (message) matches [(message*)]
  message: (PCDATA , quote ,PCDATA) matches [(PCDATA | quote*)]
  quote: (PCDATA) matches [PCDATA]")
```

5 A Creol example

The current implementation of XML in Creol and the validation algorithm enable Creol applications to use XML documents as a data storage and exchange format.

As an example we look at a simple system consisting of a `LibraryServer` and a `LibraryClient` which exchange messages on XML format. The example shows how XML elements and XML documents can be used as parameters in method calls and returns and it also illustrates validation of the XML Data against DTDs, the Creol code is given in fig. 2. The library server keeps a catalogue of books in an XML document and a client may call the method

```
getEntries(in query:ElemNd ; out result:ElemNd)
```

```
class LibraryClient implements LibraryCl
begin
  var res1 : ElemNd var res2 : ElemNd
  var res3 : ElemNd var res4 : ElemNd

  op run == var server : LibraryServ ;
  server := new LibraryServer();
  server.getEntries("query"[("title"[tx("TCP/IP Illustrated")]))] ; res1) ;
  server.getEntries("query"[("price"[tx("65.90")]))] ; res2) ;
  server.getEntries("query"[("publisher"[tx("Addison-Wesley")]))];res3) ;
  // Initialize a new server with a invalid catalogue.
  server := new LibraryServer(xmlDoc("bib"[tx("A non valid library catalogue")],noSchema)) ;
  server.getEntries("query"[("title"[tx("TCP/IP Illustrated")]))] ; res4)
end

class LibraryServer(catalogue:XMLDoc) implements LibraryServ
begin
  var queryType : XMLSchema := <queryTypeValue>
  var defCat : XMLDoc := <catalogueValue>
  var status : String

  //Use a default catalogue if no parameter is given to class.
  op init == if catalogue = null then catalogue := defCat end

  with LibraryCl
  op getEntries(in query:ElemNd ; out result:ElemNd ) ==
    //Make sure that library catalogue is valid before accepting calls.
    if xmlValid(catalogue) then
      status := "Valid, accepting calls";
      if xmlValid(query , queryType) then
        result := subElemQuery(query,catalogue)
      else
        result := "result"[("err_result"[tx("Invalid query type")])]
      end
    else
      status := "Invalid catalogue";
      result := "result"[("err_result"[tx("Library catalogue invalid")])]
    end
  end
end
```

Figure 2: Creol program

on the server. The server will first ensure that its own catalogue is valid w.r.t. a DTD for the library catalogue with a call to the function

```
xmlValid(catalogue:XMLDoc) .
```

The server then checks that the query conforms to the specified DTD for queries with a call to the function

```
xmlValid(query:ElemNd,queryType:XMLSchema),
```

and executes the query by calling the function

```
subElemQuery(query:ElemNd,catalogue:XMLDoc) .
```

If the query is valid, the server will perform the query given as a parameter and return an `ElemNd` as response, if it is not valid it will return an `ElemNd` containing an error message as a response, the same will happen if the catalogue is invalid. All responses conform to a query result DTD.

To execute the program we extend the Creol language with the three functions mentioned above. For the two validation functions this amounts to extending the Maude interpreter by specifying equations which map the Creol syntax above to our previously defined `validate` function in Maude to enable the interpreter to execute the program as a Creol program. With `D` and `D'` being of sort `Data` we add the following two equations to the interpreter:

```
eq "xmlValid"(D) = bool(getB(validate(D))) .
eq "xmlValid"(D # D') = bool(getB(validate(xmlDoc(D,D')))) .
```

The `subElemQuery` function is another addition to the Creol API and is likewise implemented in Maude. We leave out the details of the implementation since the function is tailored for this specific case for the purpose of the example. To extend the Creol API with more general XML operations we would need to consider carefully which operations to add to the API to ensure that we chose a minimal set of useful basic functions. This is left for further work.

In the program text above, the `<queryTypeValue>` is a Creol value representing the DTD for queries. the `<catalogueValue>` is a Creol value representing an XML document instance. See the appendix for an example of an execution of the program with some specific values.

6 Conclusion

Integrating XML documents in object-oriented languages is not easy in general as witnessed by the extensive research conducted in this area, and nicely presented in the survey [21]. We have shown here how to integrate XML documents into Creol, an object-oriented language with formal semantics in rewriting logic. We have also presented an algorithm for validating XML documents against XML schemas, to show that the former are instances of the latter.

This paper is a first step towards a full integration of XML into Creol, and we intend to pursue our work as to complete our agenda described in Section 1.3. In particular, we find it extremely interesting to be able to manipulate and reason about XML documents, to include regular expression types, and to adapt the semantic sub-typing algorithm from CDuce and XDuce discussed in the introduction.

References

- [1] E. W. Axelsen. A meta-level framework for recording and utilizing communication histories in Maude. Master's thesis, Dept. of Informatics, Univ. of Oslo, Norway, Aug. 2004.
- [2] N. Benton, L. Cardelli, and C. Fournet. Modern concurrency abstractions for c#. *ACM Trans. Program. Lang. Syst.*, 26(5):769–804, 2004.
- [3] V. Benzaken, G. Castagna, and A. Frisch. CDuce: an XML-centric general-purpose language. *SIGPLAN Not.*, 38(9):51–63, 2003.
- [4] M. Blow, Y. Goland, M. Kloppmann, F. Leymann, G. Pfau, D. Roller, and M. Rowley. BPELJ: BPEL for java. <http://ftpn2.bea.com/pub/downloads/ws-bpelj.pdf>.
- [5] S. Boag, D. Chamberlin, M. Fernández, D. Florescu, J. Robie, and J. Siméon. *XQuery 1.0: An XML Query language*, November 2005. <http://www.w3.org/TR/2004/WD-xquery-20040723/>.
- [6] J. Boyer. *Canonical XML Version 1.0*. W3C, 2001. W3C Recommendation 15 March 2001. <http://www.w3.org/TR/xml-c14n>.
- [7] T. Bray, J. Paoli, C. Sperberg-McQueen, E. Maler, and F. Yergeau. *Extensible Markup Language (XML) 1.0*, third edition, February 2004. <http://www.w3.org/TR/REC-xml/>.
- [8] A. Brüggemann-Klein and D. Wood. Deterministic regular languages. In *STACS*, pages 173–184, 1992.
- [9] M. Clavel, F. Durán, S. Eker, P. Lincoln, N. Martí-Oliet, J. Meseguer, and J. F. Quesada. Maude: Specification and programming in rewriting logic. *Theoretical Comput. Sci.*, 285:187–243, Aug. 2002.
- [10] D. Fallside and P. Walmsley. *XML Schema Part 0: Primer*, second edition, October 2004. <http://www.w3.org/TR/xmlschema-0/>.
- [11] D. Florescu, A. Grünhagen, and D. Kossmann. Xl: an xml programming language for web service specification and composition. *Comput. Networks*, 42(5):641–660, 2003.
- [12] M. Harren, M. Raghavachari, O. Shmueli, M. Burke, R. Bordawekar, I. Pechtchanski, and V. Sarkar. XJ: Facilitating XML processing in Java. In *WWW'05*, pages 278–287. ACM Press, May 2005.
- [13] H. Hosoya and B. C. Pierce. XDuce: A statically typed XML processing language. *ACM Trans. Inter. Tech.*, 3(2):117–148, 2003.
- [14] H. Hosoya, J. Vouillon, and B. C. Pierce. Regular expression types for XML. *ACM SIGPLAN Notices*, 35(9):11–22, 2000.
- [15] E. B. Johnsen and O. Owe. An asynchronous communication model for distributed concurrent objects. *Software and Systems Modeling*, 6(1):35–58, Mar. 2007.

- [16] E. B. Johnsen, O. Owe, and M. Arnestad. Combining active and reactive behavior in concurrent objects. In *Proc. of the Norwegian Informatics Conference (NIK'03)*, pages 193–204. Tapir, Nov. 2003.
- [17] E. B. Johnsen, O. Owe, and E. W. Axelsen. A run-time environment for concurrent objects with asynchronous method calls. In N. Martí-Oliet, editor, *Proc. 5th International Workshop on Rewriting Logic and its Applications (WRLA'04), Mar. 2004*, volume 117 of *Electronic Notes in Theoretical Computer Science*, pages 375–392. Elsevier Science Publishers, Jan. 2005.
- [18] E. B. Johnsen, O. Owe, and I. C. Yu. Creol: A type-safe object-oriented model for distributed concurrent systems. *Theoretical Computer Science*, 365(1–2):23–66, Nov. 2006.
- [19] M. Kempa and V. Linnemann. On XML Objects. In *Informal Proceedings of the Workshop on Programming Language Technologies for XML (PLAN-X 2002), PLI 2002, Pittsburgh, USA*, pages 44–54, 3.-8. October 2002.
- [20] C. Kirkegaard, A. Møller, and M. I. Schwartzbach. Static analysis of XML transformations in Java. *IEEE Transactions on Software Engineering*, 30(3):181–192, March 2004.
- [21] E. Meijer, W. Schulte, and G. Bierman. Programming with circles, triangles and rectangles. In *Proceedings of the XML Conference*, 2003.
- [22] A. Møller and M. Schwartzbach. The design space of type checkers for XML transformation languages. In *ICDT'05*, volume 3363 of *LNCS*, pages 17–36. Springer-Verlag, January 2005.
- [23] M. Murata, D. Lee, and M. Mani. Taxonomy of XML schema languages using formal language theory. In *Extreme Markup Languages*, Montreal, Canada, 2001.
- [24] D. Obasanjo. Overview of C_ω . <http://msdn.microsoft.com/library/default.asp?url=/library/en-us/dnexxml/html/xml01142005.asp>, January 2005.
- [25] D. A. Thomas. The impedance imperative - tuples + objects + infosets = too much stuff! *Journal of Object Technology*, 2(5):7–12, 2003.
- [26] W3C (World Wide Web Consortium). *XML Path Language (XPath) Version 1.0*, 1999. W3C Recommendation 16 November 1999. <http://www.w3.org/TR/1999/REC-xpath-19991116>.
- [27] W3C (World Wide Web Consortium). *XML Query Use Cases*, 2007. W3C Working Group Note 23 March 2007. <http://www.w3.org/TR/2007/NOTE-xquery-use-cases-20070323/>.

APPENDIX

A Execution of the Creol example

We here give an example of executing the program in Sec 5 with the following values:

The query DTD

```
<!ELEMENT query (title | author | editor | publisher)>
<!ELEMENT title (#PCDATA )>
<!ELEMENT author (#PCDATA )>
<!ELEMENT editor (#PCDATA )>
<!ELEMENT publisher (#PCDATA )>
```

in Creol syntax:

```
queryType =
xmlSchema("query",
  elemDecl("query" , elemCt("title" | ("author" | ("editor" | "publisher"))))
  elemDecl("title", elemCt(PCDATA))
  elemDecl("author", elemCt(PCDATA))
  elemDecl("editor", elemCt(PCDATA))
  elemDecl("publisher", elemCt(PCDATA)), noAttDecl) .
```

The query result DTD

This value is never used in the program but all query results conform to this DTD.

```
<!ELEMENT result (book | err_result)>
<!ELEMENT err_result (#PCDATA )>
<!ELEMENT book (title, (author+ | editor+ ), publisher, price )>
<!ELEMENT title (#PCDATA )>
<!ELEMENT author (#PCDATA )>
<!ELEMENT editor (#PCDATA )>
<!ELEMENT publisher (#PCDATA )>
<!ELEMENT price (#PCDATA )>
```

in Creol syntax:

```
resultType =
xmlSchema("result",
  elemDecl("result" , elemCt("book" | "err_result"))
  elemDecl("err_result", elemCt(PCDATA))
  elemDecl("book", elemCt("title" @ ("author"+ |"editor"+) @ "publisher" @ "price"))
  elemDecl("author", elemCt(PCDATA))
  elemDecl("editor", elemCt(PCDATA))
  elemDecl("title", elemCt(PCDATA))
  elemDecl("publisher", elemCt(PCDATA))
  elemDecl("price", elemCt(PCDATA)) .
```

The library catalogue

The catalogue of books is a modified version of the bibliography document used as an example in [27]:

```
<!DOCTYPE bib [  
  <!ELEMENT bib (book* )>  
  <!ELEMENT book (title, (author+ | editor+ ), publisher, price )>  
  <!ATTLIST book >  
  <!ELEMENT author (#PCDATA )>  
  <!ELEMENT editor (#PCDATA )>  
  <!ELEMENT title (#PCDATA )>  
  <!ELEMENT publisher (#PCDATA )>  
  <!ELEMENT price (#PCDATA )>  
  ]>  
<bib>  
  <book>  
    <title>TCP/IP Illustrated</title>  
    <author>Stevens, W.</author>  
    <publisher>Addison-Wesley</publisher>  
    <price>65.95</price>  
  </book>  
  <book>  
    <title>Advanced Programming in the Unix environment</title>  
    <author>Stevens, W.</author>  
    <publisher>Addison-Wesley</publisher>  
    <price>65.95</price>  
  </book>  
  <book>  
    <title>Data on the Web</title>  
    <author>Abiteboul, Serge</author>  
    <author>Buneman, Peter</author>  
    <author>Suciu, Dan</author>  
    <publisher>Morgan Kaufmann Publishers</publisher>  
    <price>39.95</price>  
  </book>  
  <book>  
    <title>The Economics of Technology and Content for Digital TV</title>  
    <editor>Gerbarg, Darcy</editor>  
    <publisher>Kluwer Academic Publishers</publisher>  
    <price>129.95</price>  
  </book>  
</bib>
```

The library catalogue as a Creol value is:

```
eq catalogue =  
xmlDoc("bib" [  
  ("book" [  
    ("title"[tx("TCP/IP Illustrated"))]  
    ("author"[tx("Stevens, W.")])  
    ("publisher"[tx("Addison-Wesley")])  
    ("price"[tx("65.95")]))])
```

```

("book" [
  ("title" [tx("Advanced Programming in the Unix environment")])
  ("editor" [tx("Stevens, W.")])
  ("publisher" [tx("Addison-Wesley")])
  ("price" [tx("65.95")])])
("book" [
  ("title" [tx("Data on the Web")])
  ("author" [tx("Abiteboul, Serge")])
  ("author" [tx("Buneman, Peter")])
  ("author" [tx("Suciu, Dan")])
  ("publisher" [tx("Morgan Kaufmann Publishers")])
  ("price" [tx("39.59")])])
("book" [
  ("title" [tx("The Economics of Technology and Content for Digital TV")])
  ("author" [tx("Gerbarg, Darcy")])
  ("publisher" [tx("Kluwer Academic Publishers")])
  ("price" [tx("129.95")])]),
xmlSchema("bib",
  elemDecl("bib", elemCt("book" * ))
  elemDecl("book", elemCt("title" @ ("author"+ |"editor"+) @ "publisher" @ "price"))
  elemDecl("author", elemCt(PCDATA))
  elemDecl("editor", elemCt(PCDATA))
  elemDecl("title", elemCt(PCDATA))
  elemDecl("publisher", elemCt(PCDATA))
  elemDecl("price", elemCt(PCDATA),noAttDecl)) .

```

The client sends four queries to the server. In XML notation the queries and responses would look like this:

Query 1

```
<query><title>TCP/IP Illustrated</title><query>
```

```

<result>
  <book>
    <title>TCP/IP Illustrated</title>
    <author>Stevens, W.</author>
    <publisher>Addison-Wesley</publisher>
    <price>65.95</price>
  </book>
</result>

```

Query 2

```
<query><price>65.95</price><query>
```

```

<result>
  <err_result>Invalid query type</err_result>
</result>

```

Query 3

```
<query><publisher>Addison-Wesley</publisher><query>
```

```

<result>
  <book>

```



```

< 'getEntries : Mtdname | Param: 'query,Latt: 'query |-> null, 'result |-> null,Code: if
  'xmlValid[['catalogue]] th 'status := str("Valid, accepting calls") el ('status := str(
    "Invalid catalogue")) ; return("result"(emp)["err_result"(emp)[tx(
      "Library catalogue invalid"])] fi ; if 'xmlValid[['query,, 'queryType]] th 'result := (
      'subElemQuery[['query,, 'catalogue]]) el 'result := "result"(emp)["err_result"(emp)[tx(
        "Invalid query type"])] fi ; return('result) > *
< 'init : Mtdname | Param: emp,Latt: noSub,Code: if 'equal[['catalogue,,null]] th 'catalogue
  := 'defCat el nil fi ; return(emp) >,0cnt: 3 >

< ob('LibraryClient1) : 'LibraryClient | Att: 'res1 |-> "result"(emp)[
  "book"(emp)[("title"(emp)[tx("TCP/IP Illustrated")]) ("author"(emp)[tx("Stevens, W.")]) (
    "publisher"(emp)[tx("Addison-Wesley")]) ("price"(emp)[tx("65.95")])]], 'res2 |->
  "result"(emp)["err_result"(emp)[tx("Invalid query type")]], 'res3 |-> "result"(emp)[(
    "book"(emp)[("title"(emp)[tx("TCP/IP Illustrated")]) ("author"(emp)[tx("Stevens, W.")]) (
      "publisher"(emp)[tx("Addison-Wesley")]) ("price"(emp)[tx("65.95")])]) ("book"(emp)[(
        "title"(emp)[tx("Advanced Programming in the Unix environment")]) ("editor"(emp)[tx(
          "Stevens, W.")]) ("publisher"(emp)[tx("Addison-Wesley")]) ("price"(emp)[tx("65.95")])])]],
  'res4 |-> "result"(emp)["err_result"(emp)[tx("Library catalogue invalid")]], 'this |->
  ob('LibraryClient1),Pr: idle,PrQ: noProc,Lcnt: 7 >

< ob('LibraryServer1) : 'LibraryServer | Att: 'catalogue |-> <catalogue> , 'defCat |->
  <catalogue> , 'status |-> str("Valid, accepting calls"), 'this |-> ob('LibraryServer1),
  Pr: idle,PrQ: noProc,Lcnt: 3 >

< ob('LibraryServer2) : 'LibraryServer | Att: 'catalogue |-> xmlDoc("bib"(emp)[tx(
  "A non valid library catalogue")], noSchema), 'defCat |-> <catalogue>, 'status |-> str(
  "Invalid catalogue"), 'this |-> ob('LibraryServer2),Pr: idle,PrQ: noProc,Lcnt: 3 >

< ob('main) : 'class | Att: 'var |-> ob('LibraryClient1),Pr: idle,PrQ: noProc,Lcnt: 0 >

```