

# Extracting Typed Values from XML Data

Richard Connor<sup>‡</sup>   David Lievens<sup>‡</sup>   Paolo Manghi<sup>‡</sup>  
Steve Neely<sup>‡</sup>   Fabio Simeoni<sup>‡</sup>

<sup>‡</sup>Department of Computer Science  
University of Strathclyde, Glasgow (Scotland)

<sup>†</sup>Dipartimento di Informatica  
Università di Pisa, Pisa (Italy)

28th September 2001

## Abstract

Values of existing typed programming languages are increasingly generated and manipulated outside the language jurisdiction. Instead, they often occur as fragments of XML documents, where they are uniformly interpreted as labelled trees in spite of their domain-specific semantics. In particular, the values are divorced from the *high-level* type with which they are conveniently, safely, and efficiently manipulated within the language.

We propose language-specific mechanisms which extract language values from arbitrary XML documents and inject them in the language. In particular, we provide a general framework for the formal interpretation of extraction mechanisms and then instantiate it to the definition of a mechanism for a sample language core **L**. We prove that such mechanism can be built by giving a sound and complete algorithm that implements it.

The values, types, and type semantics of **L** are sufficiently general to show that extraction mechanisms can be defined for many existing typed languages, including object-oriented languages. In fact, extraction mechanisms for a large class of existing languages can be directly derived from **L**'s. As a proof of this, we introduce the *SNAQue* prototype system, which transforms XML fragments into CORBA objects and exposes them across the ORB framework to any CORBA-compliant language.

## 1 Introduction

Values of existing typed programming languages are increasingly generated and manipulated outside the language jurisdiction. Instead, they often occur as fragments of XML documents (cf. [8])<sup>1</sup>.

This may be because the containing documents are *semistructured*, i.e. their structure is too irregular or unstable to be effectively handled by traditional programming languages or DBMSs (cf. [1, 9]). It may also occur when the document is more disciplined, but needs to be exchanged across proprietary boundaries in a standard and self-describing format.

<sup>1</sup>Here and throughout the paper, XML is implicitly intended as a format for arbitrary data, not only document-oriented data.

As an example, consider the following XML document **d**, where some irregularities have been intentionally added to the data for sake of illustration.

```
<staff>
  <member code = "123517">
    <name>Richard Connor</name>
    <home>www.cis.strath.ac.uk/~richard</home>
  </member>
  <member code = "123345">
    <name>Steve Neely</name>
    <ext>4565</ext>
    <project>
      <name>SNAQue</name>
    </project>
  </member>
  <member code = "175417">
    <ext>4566</ext>
    <name>Fabio Simeoni</name>
  </member>
</staff>
```

On a much larger scale, this irregularity would prevent the document from being conveniently managed within the typed framework of conventional technology. While union types and object-oriented features may accommodate some of the irregularity, their abuse would soon degrade the performance of the system and complicate program specification and maintenance.

Consider instead the fragment **d'** of **d** shown next:

```
<staff>
  <member code = "123517">
    <name>Richard Connor</name>
  </member>
  <member code = "123345">
    <name>Steve Neely</name>
  </member>
  <member code = "175417">
    <name>Fabio Simeoni</name>
  </member>
</staff>
```

For most object-oriented languages, **d'** may be an XML encoding of an object **staff** of class **Staff**, where

```
class Staff {
    private Member[] members;
    Member[] getMembers() {...}
    void setMembers(Member[] members) {...}
    ...}

```

where `Member` is the class:

```
class Member {
    private String name;
    private int code;
    String getName() {...}
    void setName (String n){...}
    int getCode() {...}
    void setCode (int c){...}
    ...}

```

This simple observation raises the expectation that programming over `d'` be as simple, safe, and efficient as programming over `staff` with existing programming languages. In particular, we require these good properties to scale, i.e. hold for generalised computations over XML fragments considerably larger than `d'`. Unfortunately, we believe that none of the current approaches fully satisfies such requirement.

## 1.1 Background

To date, computations over XML data can be specified in a variety of paradigms, models and languages. Two kinds of approaches, however, appear to prevail: dedicated query languages and bindings to programming languages, typically object-oriented ones.

In query languages such as [11, 3, 14], queries have a familiar SQL-like structure, but contain powerful path expressions specified against the tree topology of the data. This gives the languages the flexibility required to compute over data with irregular or partially known structure. It makes them also more succinct than full-fledged programming languages for most operations of data filtering and transformation.

However, query languages are not computationally complete, nor well suited to complex programming tasks over large datasets, possibly involving recursion. Furthermore, they are essentially untyped, except with respect to the tree structure of the data. While this is justified in the general case by the typeless nature of the format, potential regularity in (subsets of) the data could and should be exploited for program verification and optimisation.

Language bindings are instead defined by implementing programming interfaces to one of two possible in-memory representations of the data. In the Document Object Model interface (cf. [15]), the data is organised and manipulated as a labelled tree. In the Simple API for XML (cf. [23]), the data is a string of characters organised and processed along parsing events.

Beside performance-related differences, both solutions impose an interpretation of the data which generalises their structural relationships (e.g. nodes of a

tree), but conveys only indirectly and too concretely their intended meaning (e.g. staff of a university department). When computations explicitly address the structural properties of the data (e.g. adding or removing a node, searching for a string in the data), this interpretation is adequate. In most cases, however, it complicates program specification, making it tedious, error-prone and hard to maintain.

Consider, for example, any computation over the names and codes of the staff members in `d`. In a pure implementation of the DOM interface for a Java-like language, the code may include something like:

```
int code;
String name=null;
Element staff=d.getDocumentElement();
NodeList members =
    staff.getElementsByTagName("member");
int memberCount = members.getLength();
for (int i=0;i<memberCount,i++) {
    Element member = (Element) members.item(i);
    code = Integer.parseInt(member.getAttribute("code"));
    NodeList children = member.getChildNodes();
    int length = children.getLength();
    for (int j=0;j<length,j++) {
        Node child = children.item(j);
        if (child.getNodeType()==Node.ELEMENT_NODE) {
            String tagName= ((Element) child).getTagName();
            if (tagName.equals("name")) name=
                ((CharacterData) child.getFirstChild()).getData(); }
        ...do something with name and code ...}

```

Even for a simple task, the code is highly convoluted and inefficient. Partly, this is due to the document-oriented nature of any XML programming interface. For instance, semantically related data must be accessed with the different algebras of elements and attributes. Similarly, manipulating atomic data requires an implicit or explicit cast from the type of strings, the only available. More generally, the logic of the computation is unnecessarily expressed in an algebra of trees, while domain-specific concepts (e.g. names and codes) are relegated to the role of run-time parameters.

The same task could have been specified directly against the object `staff` of class `Staff` defined above, and as simply as:

```
Member[] members = staff.getMembers();
for (int i=0;i<members.length,i++) {
    int code = members[i].getCode();
    String name = members[i].getName();
    ...do something with name and code ...}

```

The code is now aligned to the semantics of the application. It is also more succinct and less redundant, for generic operations on staff and staff members do not have to be repeated within specific computations, but can be factored out in class definitions, thoroughly tested, and then reused.

Inadequate data abstractions also compromise static checking of computations. Correctness can be guaranteed for operations on trees and strings, but

not staff members. For example, the following invocation:

```
NodeList members =
    staff.getElementsByTagName("mebmer");
```

where "mebmer" is a typo for "member", would silently compile and return a `null` value only at run-time. Safety is thus responsibility of the programmer, not the system. Programmatic checks worsen readability and maintainability of the code, and are not always sufficient to guarantee correct behaviour. In the lack of some description of the data (e.g. a DTD), the typo may be interpreted as the absence of required data and thus trigger unintended behaviour. Even assuming some data description, the typo may accidentally identify some other data or, in the best case, be simply signalled at run-time.

For similar reasons, the system can optimise resources only within the limits of its static knowledge of the data. For instance, it ignores the fact that all staff members have names and codes.

## 1.2 Extraction Mechanisms

Motivated by the previous observations, we aim at defining *high-level bindings* between XML and existing programming languages, which preserve the intended semantics of the data.

Specifically, we propose language-specific mechanisms that *extract* self-describing representations of language values from arbitrary XML documents, and transform them into their counterparts within the language. Thereafter, the extracted data are computed over in a familiar, expressive, and robust environment.

To achieve this for a given language, we interpret the extraction of a value as the projection of its language type over the containing XML document. Following the previous example, the projection of class `Staff` over `d` would result in the extraction of the object `staff`.

More formally, let  $\mathbf{D}$  be the set of XML documents,  $\mathbf{L}$  a typed programming language, and  $\mathbf{V}$  and  $\mathbf{T}$  the value and type spaces of  $\mathbf{L}$ , respectively (see Figure 1). Let also  $sd : \mathbf{V} \rightarrow \mathbf{D}$  be a self-describing interpretation of  $\mathbf{L}$ 's values in  $\mathbf{D}$ , and  $\preceq$  in  $\mathbf{D} \times \mathbf{D}$  a relation of containment between XML documents.

**Definition 1.1** *Let  $v \in \mathbf{V}$ .  $v$  is extractable from  $d \in \mathbf{D}$  according to  $T \in \mathbf{T}$  if: (i)  $v$  has type  $T$ , and (ii) there exists  $d' \preceq d$  such that  $sd(v) = d'$ .*

Finally, an *extraction mechanism* for  $\mathbf{L}$  takes both a document  $d$  and a type  $T$ , and returns a value  $v$  extractable from  $d$  according to  $T$ , if one exists.

An extraction mechanism is invoked by a programmer, via some interface to the language. If no extractable value can be returned, the programmer is

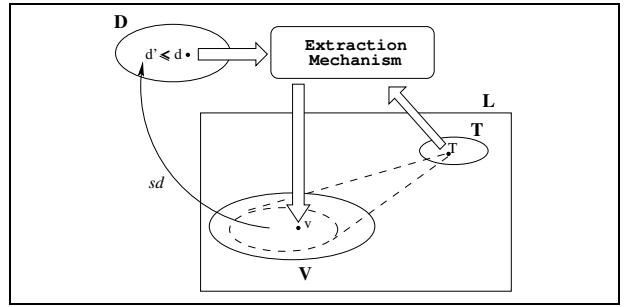


Figure 1: An extraction mechanism for  $\mathbf{L}$

notified of the failure. In case of successful extraction, the returned value may still not fully satisfy the programmer's requirements.

The reason is that, when passing a type to the mechanism, the programmer may not be aware of the exact structure of the target document. Such a type will probably be defined after an eye-inspection of the document or, if available, a description of its structure. As a result, the type may be an under-specification or an over-specification of the data that are potentially relevant to the programmer.

These sort of misjudgements may be very frequent, especially in correspondence with large documents in which data have been inserted at different times and possibly by different users with a different cognition of data representation.

In general, we require that some quantification of relevance be always returned along with the value extracted by the mechanism, and Section 6 will present one quantification scheme in more detail. By interpreting the quantification, the programmer may conclude the inadequacy of the proposed type, refine it, and then re-invoke the mechanism in a prototyping fashion.

## 1.3 Outline

In the rest of the paper, we concentrate on the formal definition of an extraction mechanism. This is done in Section 4, after the definition of a simplified syntax for XML data and a typed language core in Section 2 and Section 3, respectively.

Section 5 completes the definition of the extraction mechanism by presenting an algorithm that implements it. The algorithm is then used to illustrate a sample quantification scheme in Section 6. Section 7 presents a prototype implementation of the mechanism, while Section 8 and Section 9 examine related work, draw conclusions, and outline further work.

Due to lack of space, auxiliary definitions and proofs have been moved to Appendix A, or else only informally included in the paper.

## 2 A Document Syntax

In this section, we define a syntax for XML documents that isolates the data-oriented features of the format (e.g. naming and nesting) from its document-oriented features (e.g. ordering, attributes, processing instructions, etc.)<sup>2</sup>.

**Definition 2.1** *Let  $\mathbf{D}$  be the language of documents defined by the following grammar:*

$$d, d_1, d_2 ::= \langle \rangle \mid s \mid \langle l_1 \rightarrow d \rangle \mid d_1 \wedge d_2$$

where  $l_i \in \text{Lbl}$ ,  $s \in \text{Str}$ , and the sets  $\text{Lbl}$  of labels and  $\text{Str}$  of strings are pre-defined languages over the same alphabet of characters.

A document  $d$  is atomic or complex. An atomic document is either the *empty document*  $\langle \rangle$  or else a string. A complex document is either the singleton document  $\langle l \rightarrow d \rangle$  or the concatenation  $d_1 \wedge d_2$  of two complex or empty documents. Finally, we shall consider equal two documents that differ only in the ordering of components (e.g.  $d_1 \wedge d_2 = d_2 \wedge d_1$ ).<sup>3</sup>

In the rest of the paper, we shall abbreviate  $d = \langle l_1 \rightarrow d_1 \rangle \wedge \dots \wedge \langle l_p \rightarrow d_p \rangle$  with  $\langle l_1 \rightarrow d_1, \dots, l_p \rightarrow d_p \rangle$  and refer to  $l \rightarrow d$  as to an  $l$ -field with name  $l$  and value  $d$ . We will also use the function  $\text{FV}_l$  which takes a document and returns the set of its  $l$ -field values.

**Definition 2.2** *Let  $l \in \text{Lbl}$ . Let  $\text{FV}_l : \mathbf{D} \rightarrow \wp(\mathbf{D})$  be the function:*

$$\text{FV}_l(d) = \begin{cases} \{d'\} & d = [l \rightarrow d'] \\ \text{FV}_l(d_1) \cup \text{FV}_l(d_2) & d = d_1 \wedge d_2 \\ \emptyset & \text{otherwise} \end{cases}$$

As an example, the following document  $d$

$\langle \text{member} \rightarrow \langle \text{name} \rightarrow \text{Richard}, \text{age} \rightarrow \langle \rangle \rangle \rangle$

is a rewriting of the XML syntax

```
<member>
  <name>Richard</name>
  <age/>
</member>
```

while  $\text{FV}_{\text{member}}(d) = \{\langle \text{name} \rightarrow \text{Richard}, \text{age} \rightarrow \langle \rangle \rangle\}$ .

## 3 A Language Core

In this Section, we present the language  $\mathbf{L}$  for our sample extraction mechanism of Section 4. For our purposes, it suffices a language core defined around a value notation, a language of structural types, and a relationship of typing between the two. Extensions to

<sup>2</sup>In practice, element attributes may be replaced by subelements.

<sup>3</sup>Essentially, a document is an edge-labelled tree, in slight contrast with the standard interpretation of XML documents as node-labelled trees. This choice allows us to simplify the formal treatment but has no impact on the applicability of our results.

full-fledged languages with value operators or object-oriented types do not present particular problems.

Along the way, we follow a principle of generality that allows extraction mechanisms for other typed languages to be derived from  $\mathbf{L}$ 's.

We have chosen for  $\mathbf{L}$  a selection of the type constructs commonly found in existing programming languages. Available types are constructed from a range of atomic types  $B_1, B_2, \dots, B_N$ . They include record, set, and union types, possibly recursively defined. Record constructors are the singleton record type  $[l : T]$  and the concatenation  $T_1 \wedge T_2$  of two disjoint record types<sup>4</sup>. Set types are denoted by  $\text{set}(T)$ , where  $T$  is the member type of the set<sup>5</sup>, recursive types by  $\mu X.T$ , where  $T$  is the body of the type and  $X$  a type variable, and union types by  $T_1 \vee T_2$ , where  $T_1$  and  $T_2$  are the branch types of the union. Contrary to most languages, union types are untagged, as we do not need tags to describe alternatives in the structure of the data<sup>6</sup>.

**Definition 3.1** *Let  $\mathbf{T}$  be the language of types generated by the following grammar:*

$$T, T_1, T_2 ::= B_k \mid X \mid [l : T] \mid T_1 \wedge T_2 \mid T_1 \vee T_2 \mid \text{set}(T) \mid \mu X.T$$

where  $l \in \text{Lbl}$ ,  $X \in \text{Var}$ ,  $k \in [1, n]$ .

As for documents, two types are equal if they differ only in the ordering of the components (e.g.  $T_1 \wedge T_2 = T_2 \wedge T_1$  and  $T_1 \vee T_2 = T_2 \vee T_1$ ). Similarly, we shall implicitly extend to types the abbreviations, conventions, and auxiliary functions introduced for documents.

The values of  $\mathbf{L}$  include the elements of the atomic types, singleton records  $[l = v]$ , disjoint concatenations  $v_1 \wedge v_2$  of two record values,<sup>7</sup> sets  $\{v_1, \dots, v_n\}$ , and the empty set  $\{\}$ .<sup>8</sup> Finally, value equality follows the same rules as document and type equality.

**Definition 3.2** *Let  $\mathbf{V}$  be the language of values generated by the following grammar:*

$$v, v_1, \dots, v_n ::= b_k \mid [l = v] \mid v_1 \wedge v_2 \mid \{\} \mid \{v_1, \dots, v_n\}$$

where  $l \in \text{Lbl}$  and  $b_k \in B_k$ .

The relation of typing between values and types is standard and can be defined as follows:

<sup>4</sup>Two record types are disjoint when they have no field name in common.

<sup>5</sup>We could have chosen bag types that, differently from set types, can describe repeated sub-documents of a given document. The extension does not present particular problems but it slightly complicates the formal treatment and has been avoided here.

<sup>6</sup>Of course, this does not exclude the use of tagged unions to increase static checking of programs. In this case, tags would be simply ignored for extraction purposes.

<sup>7</sup>To improve readability, we do not syntactically distinguish the operations of concatenation of documents, types, and values. The context shall clarify the domains of definition.

<sup>8</sup>In the following examples, we shall assume that  $\mathbf{L}$  include integer numbers  $n$  of type  $\text{int}$  and strings "s" of type  $\text{string}$  among its atomic values and types.

**Definition 3.3** Let  $d \in \mathbf{D}, T \in \mathbf{T}$ .  $d$  has type  $T$  if  $d : T$ , where  $: \subseteq \mathbf{D} \times \mathbf{T}$  is the typing relation inductively defined by the following rules:

$$\begin{array}{lcl}
b_k : B_k & (\text{ATM}) & \{ \} : \text{set}(T) \quad (\text{ESET}) \\
\frac{v_1 : T_1, \dots, v_n : T}{\{v_1, \dots, v_n\} : \text{set}(T)} & (\text{SET}) & \frac{v : T}{[l = v] : [l : T]} \quad (\text{SREC}) \\
\frac{v_1 : T_1 \quad v_2 : T_2}{v_1 \wedge v_2 : T_1 \wedge T_2} & (\text{REC}) & \frac{v : T [\mu X. T/X]}{v : \mu X.T} \quad (\text{RCS}) \\
\frac{v : T_1}{v : T_1 \vee T_2} & (\text{ULFT}) & \frac{v : T_2}{v : T_1 \vee T_2} \quad (\text{URGT})
\end{array}$$

and where  $T [\mu X. T/X]$  denotes the standard operation of (capture-avoiding) variable substitution.

## 4 An Extraction Mechanism

In this Section, we provide a self-describing interpretation of language values as well as an inclusion relation between documents.

**Definition 4.1** Let  $sd : \mathbf{V} \rightarrow \mathbf{D}$  be the function defined as:

$$sd(b_k) = \text{textify}(b_k)$$

$$sd([l = v]) = \begin{cases} \langle \rangle & v = \{ \} \\ \bigwedge_{i=1}^n \langle l \rightarrow sd(v_i) \rangle & v = \{v_1, \dots, v_n\} \\ \langle l \rightarrow sd(v) \rangle & \text{otherwise} \end{cases}$$

$$sd(v_1 \wedge v_2) = sd(v_1) \wedge sd(v_2)$$

where *textify* is any function that returns string representations of atomic values.

The interpretation is straightforward, except perhaps for the case of set values, which are interpreted only within record values. The reason is that set values are not directly supported in  $\mathbf{D}$  and must be interpreted in correspondence with repeated field names within complex documents.

For example, the document

```
<member → < name → Steve >,
  member → < name → Fabio >>
```

interprets the value

$$[\text{member} = \{[\text{name} = \text{"Steve"}], [\text{name} = \text{"Fabio"}]\}].$$

In particular, values such as  $\{1, 2\}$ ,  $v \vee \{1, 2\}$ , or  $[a : \{\{1, 2\}, 3\}]$  cannot be interpreted in  $\mathbf{D}$  for no field label is available for their interpretation.

Inclusion of documents is susceptible of different interpretations. Here, we have followed the simple intuition according to which  $d'$  is included in  $d$  if  $d'$  is syntactically contained in  $d$ , with the exception that the empty document is included in any document.

**Definition 4.2** Let  $d, d' \in \mathbf{D}$ .  $d'$  is contained in  $d$  if  $d' \preceq d$ , where  $\preceq : \subseteq \mathbf{D} \times \mathbf{D}$  is the relation inductively defined by the following rules:

$$\begin{array}{lcl}
\langle \rangle \preceq d & (\text{EMP}) \\
s \preceq s & (\text{STR}) \\
\frac{d_1 \preceq d_2}{\langle l \rightarrow d_1 \rangle \preceq \langle l \rightarrow d_2 \rangle} & (\text{SDOC}) \\
\frac{d_1 \preceq d_3 \quad d_2 \preceq d_4}{d_1 \wedge d_2 \preceq d_3 \wedge d_4} & (\text{DOC})
\end{array}$$

## 5 Extraction Algorithm

In this Section, we show the algorithm **Ext** that implements the extraction mechanism defined in Section 4.

**Definition 5.1** Let  $\mathbf{V}_\perp = \mathbf{V} \cup \{\perp\}$ . Let  $\text{Ext} : \mathbf{D} \times \mathbf{T} \rightarrow \mathbf{V}_\perp$  be the algorithm defined as

```
Ext(d, T) =
1  if unf(T)={T}
2  case of T
3    Bk:
4      { if ∃bk s.t. textify(bk)=d
          return bk
5      return ⊥ }
6  T1 ∧ T2: return (Ext(d, T1) ∧⊥ Ext(d, T2))
7  [l : T'], T' = set(T''):
8    {if d ∈ Str return ⊥
9    A := ∅
10   for each d' ∈ FVl(d) {
11     v = Ext(d', T'')
12     if v ≠ ⊥ A = A ∪ {v}
13   return [l=collect(A)] }
14  [l : T'], T' ≠ set(T''):
15   { for each d' ∈ FVl(d) {
16     v = Ext(d', T')
17     if v ≠ ⊥ return [l=v] }
18   return ⊥ }
19  otherwise: return ⊥
20 else
21   { for each T' ∈ unf(T) {
22     v = Ext(d, T')
23     if v ≠ ⊥ return v }
24   return ⊥ }
```

where  $\wedge_\perp : \mathbf{V}_\perp \times \mathbf{V}_\perp \rightarrow \mathbf{V}_\perp$  is the function defined as:

$$v_1 \wedge_\perp v_2 = \begin{cases} \perp & v_1 = \perp \text{ or } v_2 = \perp \\ v_1 \wedge v_2 & \text{otherwise} \end{cases}$$

Given  $d \in \mathbf{D}$  and  $T \in \mathbf{T}$ , **Ext** performs a recursive analysis of both type and document and it either fails (i.e. returns  $\perp$ ) or else derives a value  $v \in \mathbf{V}$  extractable from  $d$  according to  $T$ . To achieve this, **Ext** solves two main problems.

The first is that set values must be extracted along with record values, for the same reason underlying

Definition 4.1. This explains why `Ext` processes set types only when processing record types and fails with types such as  $\text{set}(T) \vee T$  or  $\text{set}(\text{set}(T))$ .

The situation is complicated further by the possibility that set types do occur within record types but are ‘protected’ by union or recursive types. Given the type  $[a : \text{set}(T) \vee T]$ , for example, we cannot recursively delegate to the union case the extraction of a value of type  $\text{set}(T)$ , for we would lose the label  $a$  necessary to extract the value set from the document.

`Ext` solves the problem by extracting values only according to record types that are ‘flattened’, i.e. contain no union or recursive type fields. This is achieved with the preliminary check on line 1, which invokes an implementation of the function `unf`. `unf` is a generalisation of the standard operation of one-step unfolding of recursive types. In particular, it one-steps unfolds a union type into the set of its branches, and a record type into the set of records obtained by one-step unfolding all its non-record field values.

**Definition 5.2** Let  $\text{unf} : T \rightarrow T$  be the function defined as:

$$\text{unf}(T) = \begin{cases} \{T' [\mu X. T/X]\} & T = \mu X. T' \\ \{T_1, T_2\} & T = T_1 \vee T_2 \\ \{[l : T'] \mid T' \in \text{unf}_R(T)\} & T = [l : T] \\ \{T'_1 \wedge T'_2 \mid T'_1 \in \text{unf}(T_1), T'_2 \in \text{unf}(T_2)\} & T = T_1 \wedge T_2 \\ \{T\} & \text{otherwise} \end{cases}$$

where  $\text{unf}_R : T \rightarrow T$  is the function

$$\text{unf}_R(T) = \begin{cases} \{T\} & T = [l : T] \\ \text{unf}(T) & \text{otherwise} \end{cases}$$

For example,

$$\text{unf}([a : \text{int} \vee \text{set}(\text{string}), b : [c : \text{int} \vee \text{set}(\text{int})]]) = \{ [a : \text{int}, b : [c : \text{int} \vee \text{set}(\text{int})]], [a : \text{set}(\text{string}), b : [c : \text{int} \vee \text{set}(\text{int})]] \}.$$

A type  $T$  is then *unfolded* if  $\text{unf}(T) = \{T\}$ , otherwise is *folded*.

The second problem occurs when multiple values of type  $T$  can be extracted from  $d$ . Due to the presence of set and union types, this possibility is in fact the norm. For example, consider the document  $d = \langle a \rightarrow 1, a \rightarrow 2, b \rightarrow 3, c \rightarrow \text{four} \rangle$  and the type  $T = T_1 \vee T_2$ , where  $T_1 = [a : \text{set}(\text{int}), b : \text{int}]$  and  $T_2 = [a : \text{set}(\text{int}), c : \text{string}]$ .

From Definition 1.1 and Definition 4.1, it is easy to see that the values  $[a = \{1, 2\}, b = 3]$  and  $[a = \{1, 2\}, c = \text{"four"}]$  are extractable from  $d$  according to  $T_1$  and  $T_2$ , respectively, and thus according to  $T$ . The same is true of the values  $[a = \{\}, b = 3]$ ,  $[a = \{1\}, b = 3]$ ,  $[a = \{2\}, b = 3]$ , etc.

`Ext` returns one extractable value on the basis of both a *best-attempt* and a *first-attempt* policy. Specifically, it returns: (i) the largest value extractable from  $d$  according to a set type, and (ii) the first value extractable from  $d$  according to one of the branches

of union types, when these are ordered from left to right.<sup>9</sup> In the previous example, `Ext` derives the value  $[a = \{1, 2\}, b = 3]$ .

The best-attempt policy is justified by an immediate *principle of maximisation* of the data contained in a correctly extracted value. The first-attempt policy is instead more arbitrary but it simplifies the definition of the algorithm and is thus adequate for the purpose of a proof of concept.

In lines 21-23, `Ext` implements its first-attempt policy and either fails or returns the first extractable value returned by a recursive execution of `Ext` with  $d$  and an unfolding of  $T$ .

In lines 2-19, `Ext` processes basic, singleton record, and concatenated record types. For a basic type  $B_k$ , `Ext` is successful only if  $d$  ‘textifies’ a value  $v$  of type  $B_k$  (lines 3-5).

For singleton record types  $[l : T']$ , `Ext` tries to derive a singleton record value  $[l : v]$ , where the shape of  $v$  depends on whether  $T'$  is a set type.

If  $T'$  is a set type  $\text{set}(T'')$  (lines 7-13), `Ext` implements its best-attempt policy and tries to extract a value of type  $T''$  from each  $l$ -field value of  $d$ . The extracted values are then memorised and eventually grouped into a set value  $v$  by the sub-routine `collect`. Notice that, in this case, `Ext` fails only when the document is a string (line 8). Otherwise, it returns at worst the singleton record  $[l = \{\}]$ .

If  $T'$  is not a set type (lines 14-18), `Ext` returns the first value of type  $T'$  extractable from an  $l$ -field value of  $d$ . If such value does not exist, `Ext` fails.

If  $T$  is the concatenation of two record types  $T_1$  and  $T_2$ , `Ext` concatenates the values extracted from  $d$  according to  $T_1$  and  $T_2$  respectively. The operation of concatenation  $\wedge_{\perp}$  refines standard concatenation by returning  $\perp$  any time one of the operands is  $\perp$ . This ensures that failing the extraction according to either singleton record type fails the entire process.

`Ext` is clearly terminating, for it recursively operates on subdocuments of the input document and because `unf`, `textify`, and `collect` are trivially terminating.<sup>10</sup> In Appendix A, we shall also give a formal proof that `Ext` is sound, i.e. returns a value extractable from  $d$  according to  $T$ , and complete, i.e. it fails only when no value of type  $T$  can be extracted from  $d$ .

## 6 Relevance

In this Section, we present a simple relevance quantification scheme for the extraction mechanism de-

<sup>9</sup>Accordingly, left-to-right is the branch ordering followed by the sub-routine `unf`.

<sup>10</sup>Note that termination holds under the standard assumption that contractive recursive types, such as  $\mu X.X$ , are not part of the type language (cf. [5]).

defined in Section 4. The scheme can be easily embedded in **Ext**, but we give it here a separate specification to improve readability.

Let us start with a motivating example. Consider the following document  $d$ :

```
< staff
  < member →
    < name → David >
    < project →
      < name → SNAQue, ... > ... >
    < project →
      < name → GLOSS, ... > ... > ... >
  < member →
    < name → Paolo >
    < project →
      < name → Tequyla, ... > ... >
    < project →
      < name → TQL, ... > ... > ... > ... >
```

and assume that the programmer requires to compute over named projects of staff members.

From an initial analysis of  $d$ , the programmer proposes the type  $T$ :

```
staff : [member : set([project : [name : string]])]
```

and the mechanism returns the value  $v$ :

```
[staff = [member = {[project = [name = "SNAQue"],
                    [project = [name = "Tequyla"]]}]]]
```

The programmer did not notice that some staff members have more than one project, i.e.  $d$  contains more relevant data than  $T$  makes possible to extract.

To inform the user, we quantify the *precision* with which  $T$  describes the data in  $d$  that are relevant to the programmer. To return a readable measure, we distribute it along the singleton record types occurring in  $T$ , i.e. where loss of relevant data may actually occur. The result is a set of annotations for  $T$  that may help the user to refine the type and improve extraction.

In particular, the precision of a singleton record type  $[l : T']$  is measured with respect to all the documents that are processed with  $[l : T']$  on a successful execution path of **Ext**( $d, T$ ).<sup>11</sup>

Let thus  $D_{[l:T']}$  be the set of all such documents, and let

$$P_{[l:T']} = \{(d', [l = v']) \mid d' \in D_{[l:T]}, [l = v'] = \mathbf{Ext}(d', [l : T'])\}.$$

The precision  $\text{prec}_{[l:T']}$  of  $[l : T']$  is then calculated as:

$$\text{prec}_{[l:T']} = \frac{\sum_{p \in P_{[l:T']}} \text{vprec}(p)}{|P_{[l:T']}|}$$

where the *value precision*  $\text{vprec}$  of a pair in  $P_{[l:T']}$  is defined as:

$$\text{vprec}(d', [l = v']) = \begin{cases} \frac{|v'|}{\sigma} & \sigma > 0 \\ 1 & \sigma = 0 \end{cases}$$

<sup>11</sup>By successful execution path of **Ext**( $d, T$ ), we intend a path of the execution tree of **Ext**( $d, T$ ) along which **Ext** never fails.

and, in turn,  $\sigma = |FV_l(d')|$  and

$$|v'| = \begin{cases} 0 & v' = \{\} \\ n & v' = \{v_1, \dots, v_n\} \\ 1 & \text{otherwise} \end{cases}$$

Informally,  $\text{prec}_{[l:T']}$  is the average of the precisions calculated at each pair  $(d', [l = v']) \in P_{[l:T']}$ . Each of these is in turn the ratio between the number of  $l$ -field values from which **Ext** extracted a value and the number of those it did not.

In particular, low precision for a singleton  $T_1 = [l : \text{set}(T')]$  in  $T$  suggests that **Ext** did not extract values of type  $T'$  from many  $l$ -field values in  $d$ . A renewed analysis of the data may then reveal that the singleton  $T_2 = [l : \text{set}(T' \vee T'')]$  allows to extract more relevant data from  $d$  and should thus replace  $T_1$ . Similarly,  $[l : \text{set}(T')]$  may do better than a low-precision singleton  $[l : T']$  if  $T' \neq \text{set}(T'')$ .

For example, the precision of type  $T$  in the previous example may be returned as the following annotation:

```
1[staff 1[member : set(1/2 [project : 1 [name : string]])]]
```

The user may then improve extraction by refining  $T$  into the type:

```
[staff : [member : set([project : set([name : string]]))]]
```

which has precision 1 on all its singleton record types.

The problem of relevance quantification is certainly complex and identifies an interesting research topic per se. The scheme presented is fairly simple, and we have introduced it as a proof of concept. Although we have not yet gathered experimental results, we believe that the scheme can be useful with large datasets, where the exact structure of relevant data is not known when the mechanisms is first invoked.

## 7 Snaque

Based on the algorithm **Ext**, we have built a distributed system prototype. The system, maintained at the Computer Science Department of the University of Strathclyde, is called **SNAQue** - *the Strathclyde Novel Architecture for Querying extensible mark-up language* (cf. [29]). Although some parts are still under development, the system is currently being tested in a number of biodiversity projects by the Palaeobiology Research Group at the University of Glasgow.

SNAQue is a CORBA application that implements an extraction mechanism for a subset of the *CORBA Interface Definition Language* (cf. [28]), and therefore for any CORBA-compliant language (e.g. C, C++, Smalltalk, Java, Ada95, etc.).

The correctness of the system relies directly on the correctness of the extraction mechanism defined in Section 4. In particular, SNAQue receives an XML document and an IDL type description from a remote

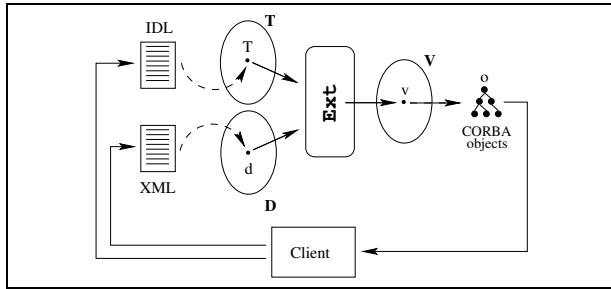


Figure 2: The SNAQue architecture

client, and maps them onto a document  $d \in \mathbf{D}$  and a type  $T \in \mathbf{T}$ , respectively. It then invokes `Ext` and transforms the output  $v \in \mathbf{V}$  of type  $T$  in a number of interrelated CORBA objects. By virtue of the mapping from IDL to  $\mathbf{T}$ , the objects expose interfaces corresponding to the initial IDL description.

According to the XML specification, there is always one entry point to the CORBA objects: the object  $o$  that corresponds to the root of the XML document. Finally, a reference to  $o$  is returned to the client for local binding in programs written in any CORBA-compliant language of choice. Optionally, SNAQue may register  $o$  with a public alias provided by the client. Then any client informed of the alias can come along and gain remote access to  $o$  (see Figure 2).

The details of the mappings from XML to  $\mathbf{D}$ , IDL to  $\mathbf{T}$ , and  $\mathbf{V}$  to the corresponding CORBA objects are out of the scope of this paper. Roughly, the mapping from XML concentrates on the data-oriented features of the format, while the mapping from IDL converts interfaces, sequences, and tagged union types, into records, sets, and untagged union types in  $\mathbf{T}$ .

As an example, consider the document  $d$  and its fragment  $d'$  introduced in Section 1. Using SNAQue, it is straightforward to compute over  $d'$  with a code analogous to that shown in Section 1. One has only to provide SNAQue with the following IDL type description:

```
interface Staff {
  typedef sequence <Member> MemberSeq;
  attribute MemberSeq members; }

interface Member {
  attribute String name;
  attribute long code; }
```

With respect to these inputs, SNAQue will create four CORBA objects: one conforming to the `Staff` interface and three conforming to the `Member` interface. It will then return a reference to the first object to which clients can bind in their programs.

SNAQue chooses Java to implement the extracted CORBA objects. In particular, the Java classes generated by the system for the objects derived above are

exactly the `Staff` and `Member` classes shown in Section 1. Similarly, the Java-like code suggested there could be immediately used to compute over the data.

The choice of CORBA IDL as the type language is an obvious one. It increases the applicability of the extraction mechanism for  $\mathbf{L}$  and makes it distributed. However, the distributed nature raises performance issues, such as those related with every read or write operation performed on the data across the network.

SNAQue offers two ways of tackling such problems. On the one hand it allows the use of *value types*<sup>12</sup>. This allows clients to *pull* the generated CORBA objects over the network and inject them in the local environment. On the other hand, SNAQue gives the possibility of pushing computations to the server, by allowing the clients to specify additional methods in the IDL interfaces. For example, the `Staff` interface could be extended with a new method:

```
interface Staff {
  typedef sequence<Member> MemberSeq;
  attribute MemberSeq members;
  Member getMember(in String name) }
```

SNAQue can not automatically generate the implementation for this method, which has to be provided by the client. Due to the lack of space, we cannot discuss this facility in more detail.

## 8 Related Work

The differences between extraction mechanisms and existing approaches has been largely discussed in Section 1. It is worth noticing here that extraction mechanisms operate on arbitrary XML documents and can thus be easily coupled with untyped query languages. In an integrated environment, the convenience of the first would complement the flexibility of the second for data with a varying degree of structural regularity.

Other high-level bindings between XML and existing programming languages have been recently presented (cf. [31, 16, 22]). They all map some form of data description (usually a DTD or an XML Schema) onto language types that capture directly the semantics intended for the data. For this reason, they operate on fairly regular XML documents and do not provide facilities for extracting regular subsets from arbitrary documents. In addition, they have been developed for specific languages (e.g. Haskell, Java) and do not generalise.

The idea of exploiting regularity in XML, and more generally, semistructured data, has also motivated a number of approaches.

Early proposals were for extending standard database technology to accommodate some degree of

<sup>12</sup>Value types have been recently introduced into CORBA to allow objects to be passed by value, rather than by reference.



irregularity in the data, typically via the provision of union types. Although similar in motivation, such approaches differ from ours in their attempt to provide a total description of the data. As mentioned in Section 1, significantly irregular data lead to an uncontrolled use of union types, thereby progressively decreasing system performance and complicating program specification.

Later proposals assume a dedicated query language as a starting point, but differ in their data-first or type-first strategy.

Approaches of the first kind infer type information from existing datasets. In this case, type inference can be performed by the system for the entire database, automatically or semi-automatically (cf. [27, 25, 26, 20, 19]). The resulting types are mainly for users to understand the data and, to some extent, for query optimisers to improve execution (cf. [18]). Partial inference can be also performed by users, and the results then fed to the system as hints to reduce the scope of a search (cf. [10]). Overall, inference-based approaches exploit typing for resource optimisation, while computations remain essentially untyped.

Approaches of the second kind exploit static knowledge to guarantee computational safety (cf. [21, 4, 12, 13]). To achieve this against a tree-based model, they resort to low-level types for XML documents. Due to the support of regular expressions, such *tree types* are more flexible than high-level types in capturing irregularities in the data (cf. [17, 30, 7, 6]).

To the best of our knowledge, Ozone (cf. [2]) is the only attempt to seamlessly integrate structured and semistructured data in the same typed environment. The system extends the ODMG model to include semistructured data, and allows structured objects to be queried with semistructured primitives. Interestingly, it also supports a function for coercing semistructured data to structured objects according to a type and, as such, implements a simple extraction mechanism for ODMG. However, our mechanism is proved correct and returns values of a larger set of types.

## 9 Conclusions and Future issues

We have presented a novel approach to programming over XML data based on language bindings. The bindings are defined as mechanisms that identify and derive language values from subsets of arbitrary XML documents. When programming over such subsets, the approach delivers the computational advantages associated with the host language. Furthermore, the derived values preserve the semantics intended for the data, and thus facilitate program specification.

These mechanisms can be formally defined and correctly implemented, and we have done it for a sample but representative core language. In particular, we have proven the generality of the sample mechanism by deriving extraction mechanisms for all CORBA-compliant languages directly from it.

Future research directions concern both theoretical and practical aspects of the investigation. Beyond XML, we have already extended our results to more general forms of semistructured data. In particular, we are able to extract language values from graph-structured data, i.e. in the presence of cycles and sharing. The interested reader is referred to [24] for the full treatment.

Another interesting direction relates to the definition of inclusion between documents. The one we proposed follows first intuitions, but alternative definitions could be considered. As a first example example, inclusion checks may start from arbitrary elements of the target document, not necessarily the root element. This would save the user the often tedious task of describing the structure that leads from the root of the document to the data of interest. It would also give a hint of the flexibility achieved by navigational query paradigms without reducing the advantages of the approach.

The extraction algorithm has been proved sound but the belief that it is also tractable in pragmatic terms has not been supported by a formal analysis of its complexity. Although tests on large data samples have shown acceptable performance even on desktop machines, the impact of a considerable use of union types remains to be measured. Furthermore, we are currently working with back-tracking techniques towards algorithms with stronger properties of completeness.

Relevance quantification could be certainly improved over the sample scheme proposed in Section 6. In particular, an extraction mechanism could customise a general scheme to the programmer's specifications.

Finally, SNAQue is under continuous development, and a web interface to the system is being published at the time of writing. The research agenda is currently focusing on whether values can be virtually injected in the value space of the target language rather than materialised. Single or multiple indexes to regular subsets suggest the possibility to dynamically synchronise the interface between the language and the database under updates. At the same time, they raise the opportunity for incremental extractions.

In addition, the client/server scenario raises a number of questions related to the efficiency of the system and to the possibility of integrating data from distributed XML servers.

Investigation is needed to identify the cases in which it is more convenient to pull extracted values

at the client side or else push client computations to the server. Completeness quantification could here be used by both client and server to make intelligent decisions about data or code migration.

In addition, the possibility of storing and publishing typed interfaces over the data at the server side suggests interesting data protection and data evolution policies. For example, the usage and volume of XML data referenced through the interfaces could be gathered into statistical information that may be used to assess the impact of changes to the data.

## References

- [1] S. Abiteboul. Querying semi-structured data. *Database Theory - ICDT '97, 6th International Conference, Delphi, Greece*, pages 1–18, January 8-10 1997.
- [2] S. Abiteboul, R. Goldman, T. Lahiri, J. McHugh, and J. Widom. Ozone: Integrating structured and semistructured data. Technical report, Stanford University Database Group, 1998.
- [3] S. Abiteboul, D. Quass, J. McHugh, J. Widom, and J.L. Wiener. The Lorel Query Language for Semistructured Data. *Journal of Digital Libraries*, 1(1), pages 68–88, April 1997.
- [4] Antonio Albano, Dario Colazzo, Giorgio Ghelli, Paolo Manghi, and Carlo Sartiani. A type system for querying XML documents. In *Proceedings of ACM SIGIR 2000 Workshop On XML and Information Retrieval*, Athens, Greece, July 2000.
- [5] G. Amadio and L. Cardelli. Subtyping recursive types. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 15(4), pages 575–631, 1993.
- [6] Arnaud Sahuguet. Everything You Ever Wanted to Know About DTDs, But Were Afraid to Ask. In *WebDB-2000*, 2000.
- [7] Paul V. Biron and Ashok Malhotra. XML Schema Part 2: Datatypes. Technical report, World Wide Web Consortium, December 1999. W3C Working Draft.
- [8] T. Bray, J. Paoli, and C. Sperberg-McQueen. Extendible Markup Language (XML) 1.0. Technical report, World Wide Web Consortium, 1998. W3C Recommendation.
- [9] P. Buneman. Semi-structured data. In *Proceedings of the Sixteenth ACM SIGACT – SIGMOD – SIGART Symposium on Principles of Database Systems*, pages 117–121, 1997.
- [10] P. Buneman, S. Davidson, M. Fernandez, and D. Suciu. Adding structure to unstructured data. *Lecture Notes in Computer Science*, 1186:336–350, 1997.
- [11] Peter Buneman, Mary Fernandez, and Dan Suciu. UnQL: A Query Language and Algebra for Semistructured Data Based on Structural Recursion. *VLDB Journal*, 9(1):76–110, 2000.
- [12] Don Chamberlin, Daniela Florescu, and Jonathan Robie. Quilt: an XML query language for heterogeneous data sources. In *Proceedings of WebDB*, Dallas, TX, May 2000.
- [13] Sophie Cluet, Claude Delobel, Jérôme Siméon, and Katarzyna Smaga. Your mediators need data conversion! In *Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD-98)*, volume 27,2 of *ACM SIGMOD Record*, pages 177–188, New York, June 1998. ACM Press.
- [14] Alin Deutsch, Mary Fernandez, Daniela Florescu, Alon Levy, and Dan Suciu. XML-QL: A Query Language for XML. Technical report, World Wide Web Consortium, August 1998. Submission to the World Wide Web Consortium.
- [15] Document object model (DOM). <http://www.w3.org/DOM>.
- [16] The Breeze Factor. Xml data binding and breeze xml studio (white paper). 2001.
- [17] David C. Fallside. XML Schema Part 0: Primer. Technical report, World Wide Web Consortium, October 2000. W3C Candidate Recommendation.
- [18] Mary Fernandez and Dan Suciu. Optimizing regular path expressions using graph schemas (full version), February 1997. Manuscript available from [http://www.research.att.com/~{mff,suciu}](http://www.research.att.com/~mff,suciu).
- [19] R. Goldman and J. Widom. Approximate DataGuides. In *Proceedings of the second International Workshop WebDB '99, Pennsylvania*, June 1999.
- [20] Roy Goldman and Jennifer Widom. DataGuides: Enabling query formulation and optimization in semistructured databases. In *VLDB'97, Proceedings of 23rd International Conference on Very Large Data Bases, August 25-29, 1997, Athens, Greece*, pages 436–445. Morgan Kaufmann, 1997.
- [21] Haruo Hosoya and Benjamin C. Pierce. XDuce: A typed XML processing language (preliminary report), May 2000. WebDB workshop.
- [22] Sun Microsystems Inc. Web services made easier: The java apis for xml<sup>®</sup>. 2001.
- [23] Megginson Technologies Ltd. SAX 2.0: The Simple API for XML, 2000. <http://megginson.com/SAX/>.
- [24] Paolo Manghi. *Extracting Typed Values from Semistructured Databases*. PhD thesis, Dipartimento di Informatica, Università di Pisa - Italy, 2001.
- [25] S. Nestorov, S. Abiteboul, and R. Motwani. Inferring structure in semistructured data. *Workshop on Management of Semistructured Data, in conjunction with PODS/SIGMOD, Tucson - Arizona*, 1997.
- [26] S. Nestorov, S. Abiteboul, and R. Motwani. Extracting schema from semistructured data. *ACM SIGMOD*, pages 295–306, 1998.
- [27] Svetlozar Nestorov, Jeffrey D. Ullman, Janet L. Wiener, and Sudarshan S. Chawathe. Representative objects: Concise representations of semistructured, hierarchical data. In Alex Gray and Per-Åke Larson, editors, *Proceedings of the Thirteenth International Conference on Data Engineering, April 7-11, 1997 Birmingham U.K.*, pages 79–90. IEEE Computer Society, 1997.
- [28] OMG. CORBA OMG IDL text file - The Object Management Group, 1999. <ftp://ftp.omg.org/pub/docs/formal/99-04-01.txt>.
- [29] Anna Stavrianou. Querying XML through a CORBA gateway. Master's thesis, University of Glasgow (UK), September 1999.
- [30] Henry S. Thompson, David Beech, Murray Maloney, and Noah Mendelson. XML Schema Part 1: Structures. Technical report, World Wide Web Consortium, December 1999. W3C Working Draft.
- [31] Malcolm Wallace and Colin Ranciman. Haskell and XML: Generic combinators or type-based translation? In *Proceedings of the Fourth ACM SIGPLAN International Conference on Functional Programming (ICFP'99)*, volume 34-9 of *ACM Sigplan Notices*, pages 148–159, N.Y., September 27–29 1999. ACM Press.

# A Appendix

**Lemma A.1** *Let  $v \in \mathbf{V}, T \in \mathbf{T}$ .  $v : T$  if and only if there exists  $T' \in \text{unf}(T)$  such that  $v : T'$ .*

**Proof.** By structural induction on  $T$ . The proof is immediate and we shall here discuss only the cases  $T = T_1 \wedge T_2$ .

Assume  $v : T$ . For typing scheme (REC),  $v = v_1 \wedge v_2$ , with  $v_1 : T_1$ , and  $v_2 : T_2$ . For the inductive hypothesis,  $v_1 : T'_1$  and  $v_2 : T'_2$ , for some  $T'_1 \in \text{unf}(T_1)$  and  $T'_2 \in \text{unf}(T_2)$ . For typing scheme (REC) and Definition 5.2,  $v = v_1 \wedge v_2 : T'_1 \wedge T'_2 \in \text{unf}(T_1 \wedge T_2) = \text{unf}(T)$ .

Vice versa, assume  $v : T'$ , with  $T' \in \text{unf}(T)$ . For Definition 5.2,  $T' = T'_1 \wedge T'_2$ , with  $T'_1 \in \text{unf}(T_1)$  and  $T'_2 \in \text{unf}(T_2)$ . For type scheme (REC),  $v = v_1 \wedge v_2, v_1 : T'_1, v_2 : T'_2$ . For the inductive hypothesis and typing scheme (REC),  $v = v_1 \wedge v_2 : T_1 \wedge T_2 = T$ . ■

**Proposition A.2** *Let  $d \in \mathbf{D}, T \in \mathbf{T}$ . If  $\text{Ext}(d, T) \neq \perp$  then  $\text{Ext}(d, T)$  is extractable from  $d$  according to  $T$ .*

**Proof.** By induction on the height  $h$  of the execution tree of  $\text{Ext}(d, T)$ .

For the hypothesis, the case  $h = 1$  corresponds to one of the cases  $T = B_k$  and  $T' = [l = T']$ , with  $T = \text{set}(T')$ . In the first case,  $\text{Ext}(d, T) = b_k$  and  $\text{sd}(b_k) = d$ . The thesis follows then from typing scheme (ATM), and the reflexivity of document inclusion.<sup>13</sup> In the second case,  $v = [l = \{\}]$ . For Definition 4.1 and inclusion scheme (EMP),  $\text{sd}(v) = \text{sd}([l = \{\}]) = \langle \rangle \preceq d$ . For typing schemes (SREC) and (ESET),  $v : T$  and the thesis is proven.

Assume now an execution tree of height  $h > 1$  and that the thesis is proven for all execution trees of height  $h - 1$ . Let us distinguish the cases in which  $T$  is folded or unfolded.

If  $T$  is folded, the hypothesis ensures that  $\text{Ext}(d, T) = \text{Ext}(d, T') = v \neq \perp$  for some  $T' \in \text{unf}(T)$ . For the inductive hypothesis,  $\text{sd}(v) \preceq d$  and  $v : T'$ . For Lemma A.1,  $v : T$  and the thesis is proven.

If  $T$  is unfolded, there are three cases to examine.

(i)  $T = T_1 \wedge T_2$ . For the hypothesis and the definition of  $\wedge_{\perp}$ ,  $\text{Ext}(d, T) = v = v_1 \wedge v_2$ , where  $v_1 = \text{Ext}(d, T_1) \neq \perp$  and  $v_2 = \text{Ext}(d, T_2) \neq \perp$ . For the inductive hypothesis,  $\text{sd}(v_1) \preceq d, v_1 : T_1, \text{sd}(v_2) \preceq d$ , and  $v_2 : T_2$ . For Definition 4.1 and inclusion scheme (DOC),  $\text{sd}(v) = \text{sd}(v_1 \wedge v_2) = \text{sd}(v_1) \wedge \text{sd}(v_2) \preceq d \wedge d = d$ . For typing scheme (REC),  $v : T$  and the thesis is proven.

(ii)  $[l : T']$ , with  $T' = \text{set}(T'')$ . For  $h > 1$ ,  $\text{FV}_1(d) \neq \emptyset$ . For the hypothesis,  $\text{Ext}(d, T) = v = [l = \{v_1, \dots, v_n\}]$ , where  $v_i = \text{Ext}(d_i, T'') \neq \perp, d_i \in \text{FV}_1(d)$ , for each  $i \in [1, n]$  and some  $n \in \mathbf{N}$ . For the inductive hypothesis,  $\text{sd}(v_i) \preceq d$  and  $v_i : T''$ . For Definition 4.1 and scheme (DOC),  $\text{sd}(v) = \text{sd}([l = \{v_1, \dots, v_n\}]) = \bigwedge_{i=1}^n \langle l \rightarrow \text{sd}(v_i) \rangle \preceq \bigwedge_{i=1}^n \langle l \rightarrow d_i \rangle \preceq d$ . For typing schemes (SREC) and (SET),  $v : T$  and the thesis is proven.

(iii)  $[l : T']$ , with  $T' \neq \text{set}(T'')$ . For the hypothesis,  $\text{Ext}(d, T) = v = [l = v']$ , where  $v' = \text{Ext}(d', T') \neq \perp, d' \in \text{FV}_1(d)$ . In particular,  $d = [l : d'] \wedge d''$ , for some  $d'' \in (D)$ . For the inductive hypothesis,  $\text{sd}(v') \preceq d'$  and  $v' : T'$ . For Definition 4.1 and inclusion scheme (DOC),  $\text{sd}(v) = \text{sd}([l = v']) = \langle l \rightarrow \text{sd}(v') \rangle \preceq d$ . For scheme (SREC),  $v : T$  and the thesis is proven. ■

**Proposition A.3** *Let  $d \in \mathbf{D}, T \in \mathbf{T}, v \in \mathbf{V}$ . If  $v$  is extractable from  $d$  according to  $T$  then  $\text{Ext}(d, T) \neq \perp$ .*

**Proof.** By induction on the height  $h$  of the proof tree of  $v : T$ .

The case  $h = 1$  corresponds to one of the case  $v = b_k$  and  $v = \{\}$ . The second case is excluded by the hypothesis and

<sup>13</sup>Reflexivity can be easily proven by structural induction on  $D$ .

Definition 4.1. In the first case,  $T = B_k, \text{Ext}(d, T) = b_k \neq \perp$  and the thesis is proven.

Assume now an execution tree of height  $h > 1$  and that the thesis is proven for all execution trees of height  $h - 1$ . Let us distinguish the cases in which  $T$  is folded or unfolded.

If  $T$  is folded, the hypothesis  $v : T$  and Lemma A.1 ensure that  $v : T'$ , for some  $T' \in \text{unf}(T)$ , and thus that  $v$  is extractable from  $d$  according to  $T'$ . For the inductive hypothesis,  $\text{Ext}(d, T') \neq \perp$  and thus  $\text{Ext}(d, T) \neq \perp$ .

If  $T$  is unfolded, there are three cases to examine.

(i)  $T = T_1 \wedge T_2$ . For typing scheme (REC) there exist  $v_1, v_2 \in \mathbf{V}$  such that  $v_1 : T_1, v_2 : T_2$ , and  $v = v_1 \wedge v_2$ . For the definition of  $\wedge$ , Definition 4.1, inclusion schemes (EMP) and (DOC), and the hypothesis  $\text{sd}(v) \preceq d, \text{sd}(v_1) = \text{sd}(v_1) \wedge \langle \rangle \preceq \text{sd}(v_1) \wedge \text{sd}(v_2) = \text{sd}(v_1 \wedge v_2) = \text{sd}(v) \preceq d$ . Similarly,  $\text{sd}(v_2) \preceq d$ . For the inductive hypothesis,  $\text{Ext}(d, T_1) \neq \perp$  and  $\text{Ext}(d, T_2) \neq \perp$ . For the definition of  $\wedge_{\perp}$ ,  $\text{Ext}(d, T) = \text{Ext}(d, T_1) \wedge_{\perp} \text{Ext}(d, T_2) \neq \perp$ , and the thesis is proven.

(ii)  $[l : T']$ , with  $T' = \text{set}(T'')$ . The thesis follows immediately from  $\text{Ext}(d, T) = [l = v']$ , for some  $v' \in (V)$ .

(iii)  $[l : T']$ , with  $T' \neq \text{set}(T'')$ . For the hypothesis  $v : T$  and scheme (SREC),  $v = [l = v']$ , and  $v' : T'$ . For the hypothesis  $\text{sd}(v) \preceq d$  and Definition 4.1,  $\text{sd}(v) = \text{sd}([l = v']) = \langle l \rightarrow \text{sd}(v') \rangle \preceq d$ . From inclusion scheme (DOC),  $d = [l : d'] \wedge d''$ , for some  $d', d'' \in (D)$ , where  $\text{sd}(v') \preceq d'$ . For the inductive hypothesis,  $\text{Ext}(d', T') = v'' \neq \perp, \text{Ext}(d', T') = [l : v''] \neq \perp$ , and the thesis is proven. ■