

# Validation and Boolean operations for Attribute-Element Constraints

Haruo Hosoya  
Kyoto University  
hahosoya@kurims.kyoto-u.ac.jp

Makoto Murata  
IBM Tokyo Research Laboratory  
mmurata@trl.ibm.co.jp

## Abstract

Algorithms for validation and boolean operations play a crucial role in developing XML processing systems involving schemas. Although much effort has previously been made for treating *elements*, very few studies have paid attention to *attributes*. This paper presents a validation and boolean algorithms for Clark’s attribute-element constraints. Although his mechanism has a prominent expressiveness and generality among other proposals, treating this is algorithmically challenging since naive approaches easily blow up even for typical inputs. To overcome this difficulty, we have developed (1) a two-phase validation algorithm that uses what we call *attribute-element automata* and (2) intersection and difference algorithms that proceed by a “divide-and-conquer” strategy.

## 1 Introduction

XML [3] and its predecessor, SGML [19], provides two major mechanisms for representing information: *attributes* and *elements*. There has been much debate (e.g., [25]) about attributes versus elements for the past ten years. However, this debate has reached no conclusions, and we continue to have both elements and attributes. In fact, many XML-based languages (e.g., XHTML and SVG) use both of them quite heavily.

Schema languages for XML have evolved by increasing expressiveness, allowing finer- and finer-grained controls to the structure of documents. At first, the designers of schema languages (DTD [3], W3C XML Schema [10], DSD [20], and RELAX [23]) mainly paid attention to elements. On the other hand, the treatments of attributes have been rather simplistic in early schema languages. Recently, James Clark, in his schema language TREX [7], proposed a description mechanism for attributes that has a comparable power to existing ones for elements.

The kernel of Clark’s proposal is a uniform and symmetric mechanism for representing constraints on elements and those on attributes. We call this mechanism *attribute-element constraints*. The expressive power yielded by attribute-element constraints is quite substantial. For example, by using this mechanism, we can specify a constraint that allows different subelements, depending on attribute values. Such a situation may arise when we want to specify different permissible subelements for `<div class="chapter">` and `<div class="section">` where the former case allows `section` in the subelements while the lat-

ter case does not. As another example, we can specify that either an attribute or element is used to represent some piece of information. This ability would be needed when we want to allow an attribute `name` to represent a simple string name or an element `name` to represent a composite value of `first` and `last` names.

Although the expressiveness of attribute-element constraints has already been established [6], algorithmic studies needed for the practical uses are yet to be done. In this paper, we present our algorithms for (1) validation and (2) boolean operations that treat Clark’s attribute constraints. The importance of validation has been widely known for a long time, while that of boolean operations have just been recognized in the recent series of papers on typechecking for XML processing languages [16, 11, 26, 21, 22]. For example, intersection is used in almost all of these languages as the core of type inference mechanisms for their pattern matching facilities; difference is used in XDuce type inference for precisely treating data flow through prioritized pattern clauses [17]; and, perhaps most importantly, the inclusion test, which is difference operation followed by the emptiness test, is used for “subtyping” in many of the above-mentioned languages.

Exactly because of the ability to mix constraints for elements and those for attributes, developing above-mentioned algorithms become challenging. Naively, one may attempt to use traditional automata techniques, as usually done in the case of element-only data. This is possible since elements are ordered sequences. However, as XML defines, attributes are unordered sets; therefore the logic behind attribute constraints is quite different from traditional automata. Another attempt might be to transform attribute-element constraints so that attribute constraints are isolated from element ones. However, this turns out to blow up in common examples, as will be discussed in Section 3.2.

We have developed a validation and boolean algorithms that overcome this difficulty. Our validation algorithm uses a “two-phase” strategy. We use a variant of traditional automata that mix both attribute and element constraints. In the first phase, we validate the attributes, rewriting the automaton into one that represents only the element constraints. In the second phase, we validate the elements against the rewritten automaton. By this technique, we can achieve quite a simple validation algorithm that does not incur a blow-up. For the boolean algorithms, we adopt a “divide-and-conquer” strategy. It exploits that it is often the case that we can partition given constraint formulas into orthogonal subparts. In that case, we proceed the compu-

tation with the subparts separately and then combine the results. Although this technique cannot avoid an exponential explosion in the worst case, it appears to work well for practical inputs in our experiences.

The rest of this paper is organized as follows. The next section gives a more motivation for attribute-element constraints and basic definitions including the data model and the syntax and semantics of constraints. Section 3 and Section 4 each present the validation algorithm and the boolean algorithms. Section 5 discusses the relationship with other work and Section 6 concludes the paper. Some formalizations and correctness proofs omitted in this paper can be found in the full version of this paper [15].

## 2 Attribute-element constraints

In our framework, data are XML documents with the restriction that only elements and attributes can appear. That is, we consider trees where each node is given a name and, in addition, associated with a set of name-string pairs. In XML jargon, a node is called element and a name-string pair is called attribute. For example, the following is an element with name `article` that has two attributes `key` and `year` and contains four child elements—two with `authors`, one with `title`, and one with `publisher`.

```
<article key="HosoyaMurata2002" year="2002">
  <author> ... </author>
  <author> ... </author>
  <title> ... </title>
  <publisher> ... </publisher>
</article>
```

The ordering among sibling elements is significant, whereas that among attributes is not. The same name of elements can occur multiple times in the same sequence, whereas this is disallowed for attributes.

Attribute-element constraints describe a pair of an element sequence and an attribute set. Let us illustrate the constraint mechanism. Element expressions describe constraints on elements and are regular expressions on names. For example, we can write the following expression

`author+ title publisher?`

to represent that a permitted sequence of child elements are one or more `author` elements followed by a mandatory `title` element and an optional `publisher` element. (Note that the explanation here is informal: for brevity, we show constraints only on names of attributes and elements. The actual constraint mechanism formalized later can also describe contents of attributes and elements.) Attribute expressions are constraints on attributes and have a notation similar to regular expressions. For example, the following expression

`@key @year?`

requires a `key` attribute and optionally allows a `year` attribute. (We prepend an `@`-sign to each attribute name in order to distinguish attribute names from element names.) A more complex example would be:

`@key ((@year @month?) | @date)`

That is, we may optionally append a `month` to a `year`; or we can replace these two attributes with a `date` attribute.

Attribute expressions are different from usual regular expressions in three ways. First, attribute expressions describe (unordered) sets and therefore concatenation is commutative. Second, since names cannot be duplicated in the same set, we require expressions to permit only data that conform to this restriction (e.g., `(@a | @b) @a?` is forbidden). Third, for the same reason, repetition (<sup>+</sup>) is disallowed in attribute expressions. We provide, however, “wild-card” expressions that allow an arbitrary number of arbitrary attributes from a given set of names (discussed later).

Attribute-element expressions or compound expressions allow one expression to mix both attribute expressions and element expressions. For example, we can write the following.

`@key @year? author+ title publisher?`

This expression requires both that the attributes satisfy `@key @year?` and that the elements satisfy `author+ title publisher?`. The next example is a compound expression allowing either a `key` attribute or a `key` element, not both.

`(@key | key) @year? author+ title publisher?`

In this way, we can express constraints where some attributes are interdependent with some elements. (Note that we can place attribute expressions anywhere—even after element expressions.) In the extreme, the last example could be made more flexible as follows

`(@key | key)`  
`(@year | year)?`  
`(@author | author+)`  
`(@title | title)`  
`(@publisher | publisher)?`

where every piece of information can be put in an attribute or an element.

In addition to the above, we provide attribute repetition expressions, which allow zero or more attributes with arbitrary names chosen from a given set of names. In general, attribute repetitions are useful in making a schema “open” so that users can put their own pieces of information in unused attributes. For example, when we want to require `key` and `year` attributes but optionally permit any other attributes, we can write the following expression (where `(*\key\year)` represents the set of all names except `key` and `year`).

`@key @year @(*\key\year)*`

Although our formulation will not include a direct treatment, attribute repetitions can be even more useful if combined with name spaces. (Name spaces are a prefixing mechanism for names in XML documents; see [2] for the details.) For example, when we are designing a schema in the name space `myns`, we can write the following to permit any attributes in difference name spaces (where `myns:*` means “any names in name space `myns`”).

`@myns:key @myns:year @(*\myns:*)*`

Apart from the kinds of name sets described above, the following can be useful: (1) the set of all names, (2) the set of all names in a specific name space, and (3) the set of all names except those from some specific name spaces. (In fact, these are exactly the ones supported by RELAX NG [9].)

## 2.1 Data model

We assume a countably infinite set  $\mathcal{N}$  of *names*, ranged over by  $a, b, \dots$ . We define *values* inductively as follows: a *value*  $v$  is a pair  $\langle \alpha, \beta \rangle$  where

- $\alpha$  is a set of pairs of a name and a value, and
- $\beta$  is a sequence of pairs of a name and a value.

A pair in  $\alpha$  and a pair in  $\beta$  are called *attribute* and *element*, respectively. In the formalization, attributes associate names with values and therefore may contain elements. This may appear odd since XML allows only strings to be contained in attributes (and the above examples followed this). Our treatment is just for avoiding the need to introduce another syntactic category for attribute contents and thereby simplifying the formalism. We write  $\epsilon$  for an empty sequence and  $\beta_1\beta_2$  for the concatenation of sequences  $\beta_1$  and  $\beta_2$ . For convenience, we define several notations for values.

$$\begin{aligned} a[v] &\equiv \langle \emptyset, \langle a, v \rangle \rangle \\ @a[v] &\equiv \langle \{ \langle a, v \rangle \}, \epsilon \rangle \\ \langle \alpha_1, \beta_1 \rangle \langle \alpha_2, \beta_2 \rangle &\equiv \langle \alpha_1 \cup \alpha_2, \beta_1\beta_2 \rangle \\ \epsilon &\equiv \langle \emptyset, \epsilon \rangle \end{aligned}$$

For example,  $@a[v] @b[w] c[u]$  means  $\langle \{ \langle a, v \rangle, \langle b, w \rangle \}, \langle c, u \rangle \rangle$ . We write  $\mathcal{V}$  for the set of all values.

## 2.2 Expressions

Let  $\mathcal{S}$  be a set of sets of names where  $\mathcal{S}$  is closed under boolean operations. In addition, we assume that  $\mathcal{S}$  contains at least the set  $\mathcal{N}$  of all names and the empty set  $\emptyset$ . Each member  $N$  of  $\mathcal{S}$  is called *name set*.

We next define the syntax of expressions for attribute-element constraints. As already mentioned, our expressions describe not only top-level names of elements and attributes but also their contents (which are omitted in the examples in the beginning of this section). Since, moreover, we want expressions to describe arbitrary depths of trees, we introduce recursive definitions of expressions, that is, grammars.

We assume a countably infinite set of *variables*, ranged over by  $x, y, z$ . We use  $X, Y, Z$  for sets of variables. A *grammar*  $G$  on  $X$  is a finite mapping from  $X$  to compound expressions. Compound expressions  $c$  are defined by the following syntax in conjunction with element expressions  $e$ .

$$\begin{aligned} c &::= @N[x]^+ \\ &\quad c|c \\ &\quad cc \\ &\quad e \\ e &::= N[x] \\ &\quad \epsilon \\ &\quad e|e \\ &\quad ee \\ &\quad e^+ \end{aligned}$$

Note that we stratify compound expressions  $c$  and element expressions  $e$  so that expressions for attributes never appear under the repetition operator. (We treat attribute repetitions as atomic since attribute names cannot be used multiple times and therefore attribute repetitions have quite different properties from the Kleene star.) For convenience, we abbreviate  $@\{a\}[x]^+$  by  $@a[x]^+$  and  $\{a\}[x]^+$  by  $a[x]^+$ .

We write  $\text{dom}(G)$  for the domain of a grammar  $G$ . We define  $\text{FV}(c)$  as the set of variables appearing in  $c$  and  $\text{FV}(G)$  as  $\bigcup_{x \in \text{dom}(G)} \text{FV}(x)$ . We require any grammar to be “self-contained,” i.e.,  $\text{FV}(G) \subseteq \text{dom}(G)$ .

We define  $\text{elm}(c)$  as the union of all the element name sets (the  $N$  in the form  $N[x]$ ) appearing in the expression  $c$ . Similarly,  $\text{att}(c)$  is the union of all the attribute name sets (the  $N$  in the form  $@N[x]^+$ ) appearing in the expression  $c$ . We forbid expressions with overlapping attribute name sets to be concatenated. That is, any expression must not contain an expression  $c_1 c_2$  with  $\text{att}(c_1) \cap \text{att}(c_2) \neq \emptyset$ .<sup>1</sup>

In addition, we impose a restriction on the form  $@N[x]^+$ : it must be that either the name set  $N$  is singleton or the content  $x$  is a distinguished variable **any** accepting any values. We assume that any given grammar  $G$  has the following mapping.

$$G(\mathbf{any}) = \mathcal{N}[\mathbf{any}]^* @ \mathcal{N}[\mathbf{any}]^*$$

This restriction is for ensuring closure under complementation. We discuss this in detail in Section 4.4.

By using the above syntax, we can derive the following useful forms.

$$\begin{aligned} c^* &\equiv c^+ | \epsilon \\ c? &\equiv c | \epsilon \\ @N[x]^* &\equiv @N[x]^+ | \epsilon \end{aligned}$$

Our expressions do not include a “single-attribute expression”  $@N[x]$  for describing a single attribute with a name from  $N$ . One may wonder if we can encode it somehow by other supported expressions. The answer is yes in the case that  $N$  is finite.<sup>2</sup> Indeed, if  $N$  is finite, then we can rewrite  $@N[x]$  as follows, where  $N = \{a_1, \dots, a_k\}$ :

$$@N[x] \equiv @\{a_1\}[x]^+ | \dots | @\{a_k\}[x]^+$$

In the case that  $N$  is infinite, our framework has no way to represent single-attributes expressions. We cannot simply add single-attribute expressions with infinite name sets since it breaks closure under complementation. For example, the complementation of  $@N[\mathbf{any}]$  would be zero, two, or more attributes with any name and any content. However, there is no way to express “two or more” in our framework. It may appear disappointing, but such a constraint as “only one attribute is present from a given infinite set” seems hardly useful.

## 2.3 Semantics

The semantics of expressions with respect to a grammar is described by the relation of the form  $G \vdash v \in c$ , which is read “value  $v$  conforms to expression  $c$  under  $G$ .” This relation is inductively defined by the following rules.

$$\frac{\forall i. (a_i \in N \quad G \vdash v_i \in G(x)) \quad k \geq 1 \quad a_i \neq a_j \text{ for } i \neq j}{G \vdash @a_1[v_1] \dots @a_k[v_k] \in @N[x]^+} \quad (\text{T-ATTREP})$$

$$\frac{a \in N \quad G \vdash v \in G(x)}{G \vdash a[v] \in N[x]} \quad (\text{T-ELM})$$

<sup>1</sup>The stratification and the disjointness restriction on attribute names are also adopted in RELAX NG.

<sup>2</sup>RELAX NG has single-attribute expressions and adopts this restriction of finiteness.

$$\frac{G \vdash v \in c_1 \quad \text{or} \quad G \vdash v \in c_2}{G \vdash v \in c_1 \mid c_2} \quad (\text{T-ALT})$$

$$\frac{G \vdash v_1 \in c_1 \quad G \vdash v_2 \in c_2}{G \vdash v_1 v_2 \in c_1 c_2} \quad (\text{T-CAT})$$

$$G \vdash \epsilon \in \epsilon \quad (\text{T-EPS})$$

$$\frac{\forall i. G \vdash v_i \in c \quad k \geq 1}{G \vdash v_1 \dots v_k \in c^+} \quad (\text{T-PLU})$$

Note that rules T-ALT and T-CAT treat alternation and concatenation both in compound expressions and element expressions.

### 3 Validation

In this section, we present a validation technique for values containing both attributes and elements.

#### 3.1 Attribute-free Validation

We first consider validation for attribute-free grammars. To validate values against element expressions, we introduce element automata, which are variations of usual string automata. Element automata play critical roles in most validation algorithms [24]. Later in this section, we introduce attribute-element automata by extending element automata.

An attribute-free grammar is a finite mapping from variables to element expressions. Recall that an element expression is a regular expression whose atoms are of the form  $N[x]$ . From each element expression, we can construct an automaton whose transition labels are also of the form  $N[x]$ . Formally, we first define automata in the usual way: an *automaton*  $M$  on an alphabet  $\Sigma$  is a tuple  $(Q, q^{\text{init}}, Q^{\text{fin}}, \delta)$  where  $Q$  is a finite set of *states*,  $q^{\text{init}} \in Q$  is an *initial state*,  $Q^{\text{fin}} \subseteq Q$  is a set of *final states*, and  $\delta \subseteq (Q \times \Sigma \times Q) \cup (Q \times Q)$  is a *transition relation* [14]. Note that we have allowed null transitions by incorporating  $Q \times Q$ . Then, an *element automaton* is an automaton over  $\{N[x] \mid N \in S, x \in X\}$ , where  $S$  is a set of name sets and  $X$  is a set of variables. Well-known algorithms for creating automata from regular expressions can be used for creating element automata from element expressions by assuming  $N[x]$  as symbols.

We give a brief sketch of a validation algorithm for attribute-free grammars, although understanding of this sketch is not required to read the rest of this paper. We validate a value by assigning variables to every element (including every descendant element) in this value. The variable assignment must “conform” to the grammar. That is, for each element, let  $x$  be the variable assigned to this element and  $M$  be the element automaton constructed from the element expression defined as  $x$  in the grammar. We execute the automaton for the children of the element in such a way that each child element  $a[v]$  matches a transition labeled  $N[x]$  when  $N$  contains  $a$  and  $x$  is assigned to this element  $a[v]$ . If we can successfully assign variables to all elements, we report that the value is valid. Although this algorithm may be too naive in that it tests all possible assignments, there is a known algorithm that takes only linear time in the size of the input value. Interested readers are referred to a note [24].

#### 3.2 Validation by separating attributes and elements

Attribute-element constraints pose a significant challenge. We cannot directly use element automata described above since our values are set-sequence pairs whereas element automata accept just sequences.

A naive solution is to separate constraints on attribute sets and those on element sequences. For example, we can handle the above compound expression by first converting it to the following expression.

$$\begin{array}{l} (((@a[\epsilon]^+ @b[\epsilon]^+) \epsilon) \\ | (@a[\epsilon]^+ b[\epsilon]) \\ | (@b[\epsilon]^+ a[\epsilon]) \\ | (\epsilon a[\epsilon] b[\epsilon])) \end{array}$$

Observe that this compound expression is a choice of  $c$ - $e$  pairs, where  $c$  is an element-free compound expression and  $e$  is an element expression. Here  $c$  describes *sets* of attributes, while  $e$  describes *sequences* of elements. Any compound expression can be normalized to this normal form, although we do not further describe any specific normalization procedure in this paper.

After normalization, we can easily compare values against normalized compound expressions. Given a value  $\langle \alpha, \beta \rangle$ , we examine if, for some  $c$ - $e$  pair,  $\alpha$  is valid against  $c$  and  $\beta$  is accepted by an element automaton constructed from  $e$ . The rest of validation is the same as in attribute-free grammars.

However, this approach causes combinatory explosion easily. For example, the normalization shown above created  $4 (= 2^2)$  pairs, since one is constructed for each subset of  $\{ @a[\epsilon]^+, @b[\epsilon]^+ \}$ . If we normalize

$$((@a[\epsilon]^+ | a[\epsilon]) (@b[\epsilon]^+ | b[\epsilon]) \dots (@y[\epsilon]^+ | y[\epsilon]) (@z[\epsilon]^+ | z[\epsilon])),$$

we will have  $2^{26}$  pairs.

#### 3.3 Attribute-element automata

We extend our attribute-free validation technique for handling attributes as well as elements. This extension introduces three steps: (1) given a compound expression, we insert special atoms (called “non-existent attribute declaration”) to it; (2) we create an automaton (called “attribute-element automaton”) from the compound expression, (3) given an attribute set, we create an element automaton from this attribute-element automaton by rewriting some of its transitions. In an actual implementation, we can perform the first and second steps statically, i.e., before receiving values. On the other hand, the third step can only be done dynamically. The rest of validation is the same as in attribute-free validation.

Below, for the ease of explanation, we first present the second and third steps. If we perform these two steps only, we sometimes report invalid values valid. We then present the first step as a remedy to this problem.

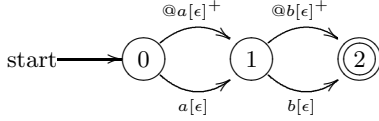
The second step creates an attribute-element automaton. Recall that atoms of compound expressions are of the form  $@N[x]^+$  or  $N[x]$ . We can assume that a compound expression is a regular expression whose atoms are  $@N[x]^+$  or  $N[x]$ , and construct an automaton from this compound expression. Thus, some transitions of this automaton have  $@N[x]^+$  as labels and others have  $N[x]$  as labels.

The third step creates an element automaton by rewriting some transitions of the constructed automaton. Given an attribute set  $\alpha$ , we replace each  $@N[x]$ -transition by a null transition if (1) more than one attribute in  $\alpha$  is valid against  $@N[x]^+$  and (2) no other attributes in  $\alpha$  have names in  $N$ . Otherwise, we remove this transition. We preserve those transitions having  $N[x]$  as labels.

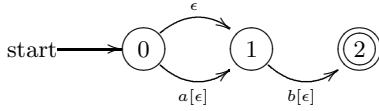
To illustrate, we use the aforementioned expression

$$(@a[\epsilon]^+ | a[\epsilon]) (@b[\epsilon]^+ | b[\epsilon]).$$

An automaton constructed from this compound expression is shown below.



Suppose that we would like to validate  $@a[\epsilon]b[\epsilon]$ . Since  $@a[\epsilon]$  is the only attribute in this value, we replace the transition labeled  $@a[\epsilon]^+$  by a null transition, and remove the transition labeled  $@b[\epsilon]^+$ . The element automaton obtained by this rewriting is shown below. The only element in  $@a[\epsilon]b[\epsilon]$  is  $b[\epsilon]$ , which is accepted by this element automaton. We can thus correctly report that  $@a[\epsilon]b[\epsilon]$  is valid.



However, the element automaton also accepts  $a[\epsilon]b[\epsilon]$  and we thus mistakenly report that  $@a[\epsilon]a[\epsilon]b[\epsilon]$  is valid. Furthermore,  $@f[\epsilon]a[\epsilon]b[\epsilon]$  (where  $@f[\epsilon]$  is undeclared) is also mistakenly reported valid: although we remove the transitions for  $@a[\epsilon]$  and  $@b[\epsilon]$ , the element automaton still accepts  $a[\epsilon]b[\epsilon]$ . What we missed is that, although we ensured that attributes *required* by compound expressions are *present* in values, we did not ensure that attributes *not required* by compound expressions are *absent* in values.

To overcome this problem, we introduce the first step. This step saturates a compound expression by introducing *non-existent attribute declarations*. A non-existent attribute declaration is  $!@N$ , where  $N$  is a name set. Informally,  $!@N$  means that no attributes have names in  $N$ . When a grammar maps a variable  $x$  to some compound expression  $c$ , we prepend a non-existent attribute declaration to  $c$ , which ensures that attributes not required by  $c$  are absent. When an alternative  $c_1 | c_2$  of two compound expressions occurs in  $c$ , we prepend a non-existent attribute declaration to  $c_1$ , which ensures that attributes required by  $c_2$  only are absent. Likewise, we prepend another non-existent attribute declaration to  $c_2$ .

We accordingly modify the second and third steps. An element-attribute automaton constructed by the second step has transitions having  $!@N$  as labels. The third step has to rewrite these transitions. Given an attribute set  $\alpha$ , we replace each  $!@N$ -transition by a null transition if no attributes in  $\alpha$  have names in  $N$ . Otherwise, we remove this transition.

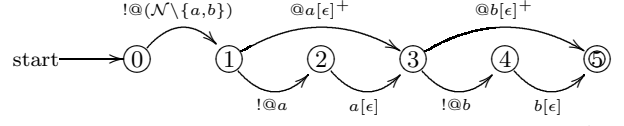
As an example, we again consider

$$(@a[\epsilon]^+ | a[\epsilon]) (@b[\epsilon]^+ | b[\epsilon]).$$

By introducing non-existent attribute declarations, we saturate it as follows:

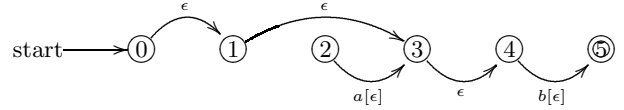
$$!@(\mathcal{N} \setminus \{a, b\}) (@a[\epsilon]^+ | (!@a a[\epsilon])) (@b[\epsilon]^+ | (!@b b[\epsilon])),$$

where  $\mathcal{N} \setminus \{a, b\}$  denotes the set of names except  $a$  and  $b$ . From the saturated expression, we create an automaton shown below.



Now, this automaton has transitions labeled  $@a[\epsilon]^+$  and  $@b[\epsilon]^+$  as well as transitions labeled  $!@a$ ,  $!@b$ , and  $!@(\mathcal{N} \setminus \{a, b\})$ . We call such automata *attribute-element automata*.

Suppose that we want to validate  $@a[\epsilon]a[\epsilon]b[\epsilon]$ . As previously, we replace the transition labeled  $@a[\epsilon]^+$  by a null transition and remove the transition labeled  $@b[\epsilon]^+$ . We further replace the transition labeled  $!@b$  by a null transition, since there are no attributes named  $b$ . We also replace the transition labeled  $!@(\mathcal{N} \setminus \{a, b\})$  by a null transition, since no attributes have names other than  $a$  or  $b$ . But we remove the transition labeled  $!@a$ , since there is an attribute named  $a$ . The element automaton obtained by this rewriting (shown below) accepts  $b[\epsilon]$  only. Thus, we do not mistakenly report that  $@a[\epsilon]a[\epsilon]b[\epsilon]$  is valid.



Let us make sure that undeclared attributes lead to invalidity. Suppose that a given value has an attribute  $@f[\epsilon]$ . Then, the transition labeled  $!@(\mathcal{N} \setminus \{a, b\})$  is removed. Since this transition is the only transition from the start state, the created element automaton accepts no element sequences.

We now formally describe the three steps.

The first step is effected by a saturating operator. This operator introduces non-existent attribute declarations for each alternative of compound expressions.

$$\begin{aligned} \text{sat}(@N[x]^+) &= @N[x]^+ \\ \text{sat}(c_1 | c_2) &= (!@(\text{att}(c_2) \setminus \text{att}(c_1)) \text{sat}(c_1)) \\ &\quad | (!@(\text{att}(c_1) \setminus \text{att}(c_2)) \text{sat}(c_2)) \\ \text{sat}(c_1 c_2) &= \text{sat}(c_1) \text{sat}(c_2) \\ \text{sat}(N[x]) &= N[x] \\ \text{sat}(\epsilon) &= \epsilon \\ \text{sat}(e^+) &= e^+ \end{aligned}$$

Next, we extend the saturating operator for grammars. Recall that a grammar  $G$  is a finite mapping from variables to compound expressions. We define  $\text{sat}(G)$  as a mapping from the same variables to saturated expressions such that

$$\text{sat}(G)(x) = !@(\mathcal{N} \setminus \text{att}(G(x))) \text{sat}(G(x)).$$

The second step creates attribute-element automata from saturated compound expressions. Formally, an *attribute-element automaton*  $M$  is an automaton over  $\{N[x], @N[x], !@N \mid N \in S, x \in X\}$ , where  $S$  is a set of name sets and  $X$  is a set of variables. We assume that a saturated compound expression is a regular expression whose atoms are  $N[x]$ ,  $@N[x]$ , or  $!@N$ . By applying a standard algorithm, we can construct an attribute-element automaton from  $\text{sat}(G)(x)$  for every  $x$ .

Finally, we present the third step as an operator for rewriting attribute-element automata. Given a set  $\alpha$  of attributes, the rewriting operator creates an element automaton from an attribute-element automaton  $(Q, q^{\text{init}}, Q^{\text{fin}}, \delta)$ . The element automaton borrows  $Q, q^{\text{init}}$ , and  $Q^{\text{fin}}$  as the

state set, start states, and final states, respectively. However, the transition relation of the element automaton is created by rewriting  $\delta$ .

First, we preserve transitions for elements. These transitions have labels of the form  $N[x]$  and they are captured by  $\delta'_1$  (shown below). Second, we rewrite transitions for attributes. For convenience, we use  $\alpha_N$  to denote a subset of  $\alpha$  such that  $\alpha_N$  contains all attributes in  $\alpha$  having names in  $N$ . For each transition in  $\delta$  having  $@N[x]$  as a label, we introduce a null transition if  $\alpha_N$  is non-empty and the value of any attribute in  $\alpha_N$  is valid against  $x$  with respect to  $G$ . These transitions are captured by  $\delta'_2$  (shown below). Third, we rewrite transitions for non-existent attribute declarations. For each transition in  $\delta$  having a label of the form  $!@N$ , we introduce a null transition if  $\alpha_N$  is empty. These null transitions are captured by  $\delta'_3$  (shown below). The transition function  $\delta'$  of the element automaton is defined as the union of  $\delta'_1$ ,  $\delta'_2$ , and  $\delta'_3$ . That is,

$$\text{rewrite}_\alpha(Q, q^{\text{init}}, Q^{\text{fin}}, \delta) = (Q, q^{\text{init}}, Q^{\text{fin}}, \delta'),$$

where

$$\begin{aligned} \alpha_N &= \{ @a[v] \in \alpha \mid a \in N \}, \\ \delta' &= \delta'_1 \cup \delta'_2 \cup \delta'_3, \\ \delta'_1 &= (\delta \cap (Q \times \{ N[x] \mid N \in S, x \in X \} \times Q) \\ &\quad \cup (\delta \cap (Q \times Q))), \\ \delta'_2 &= \{ (q_1, q_2) \mid (q_1, @N[x], q_2) \in \delta, \alpha_N \neq \emptyset, \\ &\quad \forall @a[v] \in \alpha_N. (G \vdash v \in x) \}, \\ \delta'_3 &= \{ (q_1, q_2) \mid (q_1, !@N, q_2) \in \delta, \alpha_N = \emptyset \}. \end{aligned}$$

We can prove the following theorem, which justifies our validation algorithm.

**Theorem 1** *Let  $\alpha = \{ @a_1[u_1], @a_2[u_2], \dots, @a_m[u_m] \}$  and  $\beta = b_1[v_1]b_2[v_2] \dots b_n[v_n]$ . Then,  $G \vdash \langle \alpha, \beta \rangle \in G(x)$  if and only if for some  $y_j \in X$  ( $1 \leq j \leq n$ ),*

1.  $G \vdash v_j \in G(y_j)$ ,
2.  $a_k \neq a_l$  ( $k \neq l$ ), and
3.  $b_1[y_1]b_2[y_2] \dots b_n[y_n]$  is accepted by an element automaton  $\text{rewrite}_\alpha(M)$ , where  $M$  is an attribute-element automaton constructed from  $G(x)$ .

## 4 Boolean Operations

In this section, we present our algorithms for computing intersections of and differences between attribute-element grammars.

The key technique in our algorithms is *partitioning*. Consider first the following intersection of compound expressions.

$$(@a[x]^+ | a[x]) (@b[x]^+ | b[x]) \cap @a[y]^+ (@b[y]^+ | b[y]) \quad (1)$$

How can we calculate this intersection? A naive algorithm would separate constraints on attribute sets and those on element sequences in the same way as Section 3.2

$$\begin{aligned} & (@a[x]^+ @b[x]^+ | @a[x]^+ b[x] | a[x] @b[x]^+ | a[x] b[x]) \\ & \cap (@a[y]^+ @b[y]^+ | @a[y]^+ b[y]) \end{aligned}$$

and compute the intersection of every pair of clauses on both sides. As before, such use of “distributive laws” makes the

algorithm easily blow up. Fortunately, we can avoid it in typical cases. Note that each expression in the formula (1) is the concatenation of two subexpressions, where the left subexpressions on both sides contain the names  $@a$  and  $a$  and the right subexpressions contain the different names  $@b$  and  $b$ . In such a case, we can compute intersections of the left subexpressions and of the right subexpressions separately, and concatenate the results:

$$((@a[x]^+ | a[x]) \cap @a[y]^+) ((@b[x]^+ | b[x]) \cap (@b[y]^+ | b[y]))$$

The intuition behind why this works is that each “partitioned” expression can be regarded as cross products, and therefore the intersection of the whole expressions can be done by intersecting each corresponding pair of subexpressions. Note also that no subexpression is duplicated by this partitioning process. Therefore the algorithm proceeds linearly in the size of the inputs as long as partitioning can be applied. This idea of splitting expressions into orthogonal parts was inspired by Vouillon’s unpublished work on shuffle expressions [27]. We will discuss the difference of our work from his in Section 5.

In the rest of this section, we first describe a preprocessing procedure called “name set normalization.” After this, we present definitions for partitioning, our intersection and difference algorithms, and several implementation techniques.

### 4.1 Name set normalization

Our boolean algorithms take two grammars and first pass these to the *name set normalization* (we simply say “normalization” from now on) phase before the main algorithm. The goal of normalization is to transform the given grammars to equivalent ones such that all name sets appearing in them are pair-wise either equal or disjoint. The reason for using normalization is to simplify our boolean algorithms. For example, consider the following

$$@N_1[x]^+ @N_2[x]^+ \cap @N_3[y]^+ @N_4[y]^+.$$

If  $N_1$  and  $N_2$  are respectively equal to  $N_3$  and  $N_4$ , computing this intersection is obvious. However, if these are overlapping in a non-trivial way (e.g.,  $@\{a, b\}[x]^+ @\{c, d\}[x]^+ \cap @\{a, c\}[y]^+ @\{b, d\}[y]^+$ ), it will require more work. Normalization releases us from this tedium.

Let  $\{N_1, \dots, N_k\}$  be the set of name sets appearing in the given grammars. (When we are given two grammars, this set includes all the names both in the grammars.) From this, we generate a set of disjoint name sets by the following shred function. (Here,  $\uplus$  is the disjoint union.)

$$\begin{aligned} \text{shred}(\{N\} \uplus S) &= \begin{cases} \{N\} \cup \text{shred}(S) \setminus \{\emptyset\} & \text{if } N \cap (\bigcup S) = \emptyset \\ \{N \setminus (\bigcup S)\} \cup & \\ \quad \text{shred}(\{N' \cap N \mid N' \in S\}) \cup & \\ \quad \text{shred}(\{N' \setminus N \mid N' \in S\}) \setminus \{\emptyset\} & \text{otherwise} \end{cases} \\ \text{shred}(\emptyset) &= \emptyset \end{aligned}$$

That is, we pick one member  $N$  from the input set. If this member is already disjoint with any other member, we continue shredding for the remaining set  $S$ . Otherwise, we first divide each name set  $N'$  in  $S$  into the disjoint sets  $N' \cap N$

and  $N' \setminus N$ . We separately shred all the name sets of the first form and those of the second form. We then combine the results from these and the name set obtained by subtracting all the members in  $S$  from  $N$ . As an example, this function shreds the set  $\{\{a, b\}, \{b, c\}, \{a, c\}\}$  in the following way.

$$\begin{aligned} \text{shred}(\{\{a, b\}, \{b, c\}, \{a, c\}\}) \\ &= \{\{a\}\} \cup \text{shred}(\{\{b\}, \{a\}\}) \cup \text{shred}(\{\{c\}\}) \\ &= \{\{a\}\} \cup \{\{b\}\} \cup \{\{a\}\} \cup \{\{c\}\} \\ &= \{\{a\}, \{b\}, \{c\}\} \end{aligned}$$

This function may blow up since the “otherwise” case above uses two recursive calls to shred where the size of the arguments do not necessarily decrease by half. However, this does not seem to happen in practice since the initial name sets are usually mostly disjoint and therefore the case  $N \cap (\bigcup S) = \emptyset$  is taken in most of the time.

Having computed a shredded set  $S = \text{shred}(\{N_1, \dots, N_k\})$ , we next replace every occurrence of the form  $N[x]$  or  $@N[x]^+$  with an expression that contains only name sets in  $S$ . We obtain such an expression by using the following norm function.

$$\begin{aligned} \text{norm}_S(\emptyset[x]) &= \emptyset \\ \text{norm}_S((N \uplus N')[x]) &= N[x] \mid \text{norm}_S(N'[x]) \quad (N \in S) \\ \text{norm}_S(@\emptyset[x]^+) &= \emptyset \\ \text{norm}_S(@N[x]^+) &= @N[x]^+ \mid ((\epsilon \mid @N[x]^+) \text{norm}_S(@N'[x]^+)) \quad (N \in S) \end{aligned}$$

In the special case that a given  $@N[x]^+$  is unioned with  $\epsilon$ , we can transform it to a somewhat simpler form as follows.

$$\begin{aligned} \text{norm}'_S(@\emptyset[x]^+ \mid \epsilon) &= \epsilon \\ \text{norm}'_S(@N[x]^+ \mid \epsilon) &= (@N[x]^+ \mid \epsilon) \text{norm}'_S(@N'[x]^+ \mid \epsilon) \quad (N \in S) \end{aligned}$$

This specialized rule is important in practice. In our observation,  $@N[x]^+$  almost always appears in the form  $@N[x]^+ \mid \epsilon$  since the user typically writes zero or more repetitions rather than one or more. Moreover, the straightforward form of concatenations yielded by the specialized rule gives more opportunities to the partitioning technique, compared to the complex form yielded by the general rule, where unions and concatenations are nested each other.

## 4.2 Partitioning

In our formalization, it is often convenient to view a nested concatenation of expressions as a flat concatenation and ignore empty sequences in such an expression (e.g., view  $(c_1 (c_2 \epsilon)) c_3$  as  $c_1 c_2 c_3$ ). In addition, we would like to treat expressions to be “partially commutative,” that is, concatenated  $c_1$  and  $c_2$  can be exchanged if one of them is element-free. For example, the expression  $@a[x]^+ (@b[x]^+ \mid b[x])$  is equal to  $(@b[x]^+ \mid b[x]) @a[x]^+$ . On the other hand,  $(@a[x]^+ \mid a[x]) (@b[x]^+ \mid b[x])$  is *not* equal to  $(@b[x]^+ \mid b[x]) (@a[x]^+ \mid a[x])$  since, this time,  $a[x]$  prevents such an exchange. Formally, we identify expressions up to the relation  $\equiv$  defined as follows:  $\equiv$  is the smallest congruence relation including the following.

$$\begin{aligned} c_1 c_2 &\equiv c_2 c_1 && \text{if } \text{elm}(c_1) = \emptyset \\ c_1 (c_2 c_3) &\equiv (c_1 c_2) c_3 \\ c \epsilon &\equiv c \end{aligned}$$

Now,  $(c'_1, c''_1), \dots, (c'_k, c''_k)$  is a *partition* of  $c_1, \dots, c_k$  if

$$\begin{aligned} c_i &= c'_i c''_i && \text{for all } i \\ (\bigcup_i \text{att}(c'_i)) \cap (\bigcup_i \text{att}(c''_i)) &= \emptyset \\ (\bigcup_i \text{elm}(c'_i)) \cap (\bigcup_i \text{elm}(c''_i)) &= \emptyset. \end{aligned}$$

That is, each  $c_i$  can be split into two subexpressions such that the names contained in all the first subexpressions are disjoint with those contained in all the second subexpressions. We will use partition of 2 expressions ( $k = 2$ ) in the intersection algorithm and that of an arbitrary number of expressions in the difference. The partition is said *proper* when  $0 < \text{width}(c'_i) < \text{width}(c_i)$  for some  $i$ . Here, *width* counts the number of expressions that are concatenated at the top level (except  $\epsilon$ ). That is,

$$\begin{aligned} \text{width}(\epsilon) &= 0 \\ \text{width}(c_1 c_2) &= \text{width}(c_1) + \text{width}(c_2) \\ \text{width}(c) &= 1 && \text{if } c \neq \epsilon \text{ and } c \neq c_1 c_2 \end{aligned}$$

(Note that  $\equiv$  preserves *width*.) This properness will be used for ensuring the boolean algorithms to make a progress.

## 4.3 Intersection

Our intersection algorithm is based on product construction. Let grammars  $F$  on  $X$  and  $G$  on  $Y$  be given. Assume also that normalization has already been performed *simultaneously* on  $F$  and  $G$ . (That is, no pair of a name in  $F$  and another in  $G$  is overlapping.) Now, the intersection of  $F$  and  $G$  is to compute a new grammar  $H$  on  $X \times Y$  that satisfies

$$H(\langle x, y \rangle) = \text{inter}(F(x), G(y))$$

for all  $x \in X$  and  $y \in Y$ .

The function *inter* computes an intersection of compound expressions. It works roughly in the following way. We proceed the computation by progressively decomposing the given compound expressions. At some point, they become attribute-free. Then, we convert the expressions to element automata (defined in Section 3.1), compute an intersection by using a variant of the standard automata-based algorithm, and convert back the result to an expression. Formally, *inter* is defined as follows.

- 1)  $\text{inter}(e, f) = \text{inter}^{\text{reg}}(e, f)$
- 2)  $\text{inter}(e, @N[y]^+) = \emptyset$
- 3)  $\text{inter}(@N[x]^+, f) = \emptyset$
- 4)  $\text{inter}(@N[x]^+, @N'[y]^+) = \emptyset \quad (N \neq N')$
- 5)  $\text{inter}(@N[x]^+, @N[y]^+) = @N[\langle x, y \rangle]^+$
- 6)  $\text{inter}(c, d) = \text{inter}(c_1, d_1) \text{inter}(c_2, d_2)$   
if  $(c_1, c_2), (d_1, d_2)$  is a proper partition of  $c, d$
- 7)  $\text{inter}(c_1 (c_2 \mid c_3) c_4, d) = \text{inter}(c_1 c_2 c_4, d) \mid \text{inter}(c_1 c_3 c_4, d)$
- 8)  $\text{inter}(c, d_1 (d_2 \mid d_3) d_4) = \text{inter}(c, d_1 d_2 d_4) \mid \text{inter}(c, d_1 d_3 d_4)$

The base cases are handled by rules 1 through 5, where each of the arguments is either an element expression (as indicated by the metavariables  $e$  or  $f$ ) or an attribute repetition of the form  $@N[x]^+$ . In rule 1, where both arguments are element expressions, we pass them to another intersection function  $\text{inter}^{\text{reg}}$  specialized to element expressions. This function will be explained below. Rules 2, 3, and 4 return  $\emptyset$  since the argument expressions obviously denote disjoint

sets. (Note that, in rule 4, normalization ensures that  $N$  and  $N'$  are disjoint.) When both arguments are attribute repetitions with the same name set  $N$ , rule 5 yields the same form where the content is the intersection of their contents  $x$  and  $y$ . The inductive cases are handled by rules 6 through 8. Rule 6 applies the partitioning technique already explained. Rules 7 and 8 simply expand one union form appearing in the argument expressions.

Although it may not be obvious at first sight, the presented rules cover all the cases. To see this, first let us view each given expression as a sequence of the form  $c_1 \dots c_k$  where each  $c_i$  is either an attribute repetition  $@N[x]^+$ , a union  $c_1 | c_2$ , or an element expression  $e$ . Notice that it must be that either (1) one of the two given expressions contains at least one union or (2) none of them has a union. The first case is handled by rules 6, 7 and 8. The actual algorithm applies rule 6 as often as possible, but rules 7 and 8 can always serve as fall backs. (One might think that cases like  $\epsilon(c_1 | c_2)\epsilon$  are not handled. However, since, as stated in the previous subsection, we identify expressions up to associativity, partial commutativity, and neutrality of  $\epsilon$ , such cases are already handled by the aforesaid rules.) In the second case, both of the given expressions are sequences consisting of attribute repetitions and element expressions. The case that both sequences have length zero or one is handled by rule 1 through 5. The remaining case is therefore that either sequence has length two or more. By recalling that attributes are normalized, we can always find a proper partition for such expressions; therefore this case is handled by rule 6.

The intersection function  $\text{inter}^{\text{reg}}$  for element expressions performs the following: (1) construct element automata  $M_1$  and  $M_2$  from element expressions  $e_1$  and  $e_2$ , (2) compute the “product automaton”  $M$  from  $M_1$  and  $M_2$ , and (3) convert  $M$  back to an element expression  $e$ . Since well-known conversion algorithms between automata and regular expressions can directly be used for the case of element automata and element expressions (as described in Section 3.1), we use them for (1) and (3) parts of the  $\text{inter}^{\text{reg}}$  function.

The product construction for element automata (used for the (2) part of  $\text{inter}^{\text{reg}}$ ) is slightly different from the standard one. Usually, product construction generates, from two transitions with the same label in the input automata, a new transition with that label in the output automaton. In our case, we generate, from a transition with label  $N[x]$  and another with label  $N[y]$ , a new transition with label  $N[\langle x, y \rangle]$ . Formally, given two element automata  $M_i = (Q_i, q_i^{\text{init}}, Q_i^{\text{fin}}, \delta_i)$  on  $\{N[x] \mid N \in S, x \in X_i\}$  ( $i = 1, 2$ ), the product of  $M_1$  and  $M_2$  is an automaton  $(Q_1 \times Q_2, \langle q_1^{\text{init}}, q_2^{\text{init}} \rangle, Q_1^{\text{fin}} \times Q_2^{\text{fin}}, \delta)$  on  $\{N[\langle x_1, x_2 \rangle] \mid N \in S, x_1 \in X_1, x_2 \in X_2\}$  where

$$\delta = \{ (\langle q_1, q_2 \rangle, N[\langle x_1, x_2 \rangle], \langle q'_1, q'_2 \rangle) \mid (q_i, N[x_i], q'_i) \in \delta_i \text{ for } i = 1, 2 \}.$$

(Note that we use here the assumption that the name sets of elements in the given grammars have been normalized.)

We can prove the following expected property for our intersection algorithm.

**Theorem 2** *Let  $H(\langle x, y \rangle) = \text{inter}(F(x), G(y))$ . Then,  $\text{inter}(c, d) = b$  implies that  $H \vdash v \in b$  iff  $F \vdash v \in c$  and  $G \vdash v \in d$ .*

#### 4.4 Difference

The difference algorithm is similar to the intersection algorithm since it uses partitioning, but it is somewhat more complicated because of the need to apply subset construction at the same time. To see this, consider the following difference.

$$@a[x]^+ @b[x]^+ \setminus (@a[y]^+ @b[y]^+ | @a[z]^+ @b[z]^+)$$

First of all, note that each of the left expression and the right expressions under the union can be partitioned to the one with  $@a$  and the one with  $@b$ ; these components can be regarded as orthogonal. We proceed the difference by first subtracting  $@a[y]^+ @b[y]^+$  from  $@a[x]^+ @b[x]^+$ . This yields

$$(@a[x]^+ \setminus @a[y]^+) @b[x]^+ | @a[x]^+ (@b[x]^+ \setminus @b[y]^+).$$

That is, we obtain the union of two expressions, one resulting from subtracting the first component and the other resulting from subtracting the second. This can be understood by observing that “a value  $@a[v]@b[w]$  being not in  $@a[y]^+ @b[y]^+$ ” means “either  $@a[v]$  being not in  $@a[y]^+$  or  $@b[w]$  being not in  $@b[y]^+$ .” Now, back to the original difference calculation, we next subtract the second clause  $@a[z]^+ @b[z]^+$  from the above result. Performing a similar subtraction, we obtain:

$$\begin{aligned} & (@a[x]^+ \setminus (@a[y]^+ | @a[z]^+)) @b[x]^+ \\ & | (@a[x]^+ \setminus @a[y]^+) (@b[x]^+ \setminus @b[z]^+) \\ & | (@a[x]^+ \setminus @a[z]^+) (@b[x]^+ \setminus @b[y]^+) \\ & | @a[x]^+ (@b[x]^+ \setminus (@b[y]^+ | @b[z]^+)) \end{aligned}$$

Thus, the original goal of difference has reduced to the combination of the subgoals of difference. Let us consider only the first difference  $(@a[x]^+ \setminus (@a[y]^+ | @a[z]^+))$  since the other are similar. Since all expressions appearing here are attribute repetitions with  $@a$ , the result is obviously an attribute repetition with  $@a$ . What is less obvious is that its content is the difference between the variable  $x$  and the “union” of the variables  $y$  and  $z$ . In general, given two grammars, we have to compute the difference between a variable from one grammar and a set of variables from the other grammar. This is why the algorithm involves subset construction.

Formally, let grammars  $F$  on  $X$  and  $G$  on  $Y$  be given and have been normalized simultaneously, as before. The difference between  $F$  and  $G$  is to compute a new grammar  $H$  on  $X \times \mathcal{P}(Y)$  that satisfies

$$H(\langle x, Z \rangle) = \text{diff}(F(x), \{G(y) \mid y \in Z\})$$

for all  $x \in X$  and  $Z \subseteq Y$ . The function  $\text{diff}$  takes a compound expression  $c$  and a set of compound expressions  $d_i$  and returns a difference between  $c$  and the union of  $d_i$ ’s. The definition of this function is presented in Figure 1. (Here,  $D$  ranges over sets of compound expressions and  $\uplus$  is the disjoint union.) The base cases are handled by rules 1 through 6. As before, when all the arguments are element expressions, rule 1 passes them to the difference function  $\text{diff}^{\text{reg}}$  (explained below). Rules 2, 3, and 4 remove, from the set in the second argument, an expression that is disjoint with the first argument. The remaining base cases are therefore that all the expressions in the arguments are attribute repetitions



$$\begin{array}{ll}
1) \text{ diff}(e, \{f_1, \dots, f_k\}) & = \text{diff}^{\text{reg}}(e, f_1 \mid \dots \mid f_k) \\
2) \text{ diff}(e, \{@N[y]^+\} \uplus D) & = \text{diff}(e, D) \\
3) \text{ diff}(@N[x]^+, \{f\} \uplus D) & = \text{diff}(@N[x]^+, D) \\
4) \text{ diff}(@N[x]^+, \{@N'[y]^+\} \uplus D) & = \text{diff}(@N[x]^+, D) \quad (N \neq N') \\
5) \text{ diff}(@a[x]^+, \{@a[y_1]^+, \dots, @a[y_k]^+\}) & = @a[\langle x, \{y_1, \dots, y_k\} \rangle]^+ \\
6) \text{ diff}(@N[\mathbf{any}]^+, \{@N[\mathbf{any}]^+\}) & = \emptyset \\
7) \text{ diff}(c, \{d_1, \dots, d_k\}) & = \bigcup_{I \subseteq \{1, \dots, k\}} \begin{cases} \text{diff}(c', \{d'_i \mid i \in I\}) \\ \text{diff}(c'', \{d''_i \mid i \in \{1, \dots, k\} \setminus I\}) \end{cases} \\
& \quad \text{if } (c', c''), (d'_1, d''_1), \dots, (d'_k, d''_k) \text{ is a proper partition of } c, d_1, \dots, d_k \\
8) \text{ diff}(c_1 (c_2 \mid c_3) c_4, D) & = \text{diff}(c_1 c_2 c_4, D) \mid \text{diff}(c_1 c_3 c_4, D) \\
9) \text{ diff}(c, \{d_1 (d_2 \mid d_3) d_4\} \uplus D) & = \text{diff}(c, \{d_1 d_2 d_4, d_1 d_3 d_4\} \uplus D)
\end{array}$$

Figure 1: Difference algorithm

with the same name set. Recall the restriction on attribute repetitions described in Section 2.2: any expression of the form  $@N[x]^+$  must satisfy that either (1)  $N$  is singleton or (2)  $x$  is **any**. Rule 5 handles the first case, returning an attribute repetition with that singleton name set and an appropriate difference form between variables. Rule 6 handles the second case, returning the empty set expression.

The inductive cases are handled by rule 7 through 9. Rule 8 and 9 expand one union form in the argument expressions. Rule 7 is applied when all the arguments can be partitioned altogether. The complex formula involving “for all subsets  $I \subseteq \{1, \dots, k\}$ ” on the right hand side is a generalization of the discussion made in the beginning of this section. This “subsetting” technique has repeatedly been used in the literature. Interested readers are referred to [18, 17, 12].

The function  $\text{diff}^{\text{reg}}$  is analogous to the function  $\text{inter}^{\text{reg}}$  already shown: it constructs element automata  $M_1$  and  $M_2$  from element expressions  $e_1$  and  $e_2$ , then computes the “difference automaton”  $M$  from  $M_1$  and  $M_2$ , and finally converts  $M$  back to an element expression  $e$ . The construction of difference automata uses both product and subset construction. Given two element automata  $M_i = (Q_i, q_i^{\text{init}}, Q_i^{\text{fin}}, \delta_i)$  on  $\{N[x] \mid N \in S, x \in X_i\}$  ( $i = 1, 2$ ), the difference of  $M_1$  and  $M_2$  is an automaton  $(Q, q^{\text{init}}, Q^{\text{fin}}, \delta)$  on  $\{N[\langle x_1, Y_2 \rangle] \mid N \in S, x_1 \in X_1, Y_2 \subseteq X_2\}$  where:

$$\begin{aligned}
Q &= Q_1 \times \mathcal{P}(Q_2) \\
q^{\text{init}} &= \langle I_1, \{I_2\} \rangle \\
Q^{\text{fin}} &= \{\langle q_1, P \rangle \mid q_1 \in F_1 \text{ and } P \subseteq Q_2 \text{ and } P \cap F_2 \neq \emptyset\} \\
\delta &= \{ (\langle q_1, \{p_1, \dots, p_k\} \rangle, N[\langle x_1, \{y_1, \dots, y_k\} \rangle], \langle q'_1, \{p'_1, \dots, p'_k\} \rangle) \mid \\
&\quad (q_1, N[x_1], q'_1) \in \delta_1 \text{ and} \\
&\quad \{(p_1, N[y_1], q'_1), \dots, (p_k, N[y_k], q'_k)\} \subseteq \delta_2 \}
\end{aligned}$$

We can prove the following property expected for the difference algorithm.

**Theorem 3** *Let  $H(\langle x, Z \rangle) = \text{diff}(F(x), \{G(y) \mid y \in Z\})$ . Then  $\text{diff}(c, D) = b$  implies that  $H \vdash v \in b$  iff  $F \vdash v \in c$  and  $G \vdash v \notin d$  for all  $d \in D$ .*

So far, we have used the restriction that any attribute repetition has either a singleton name set or the **any** content. What if we remove this restriction? Unfortunately, this seems to break closure under difference. For example, consider computing the difference  $@N[\mathbf{any}]^+ \setminus @N[x]^+$ . The

resulting expression should satisfy the following. Each value in it has a set of attributes all with names from  $N$ . But *at least* one of them has a content *not* satisfying  $x$ . The expression  $@N[\overline{x}]^+$  is not a right answer because it requires *all* attributes to have contents not satisfying  $x$ . Although this argument is not a proof, it gives a strong feeling that there is no answer.

#### 4.5 Implementation techniques

A naive implementation of the algorithms presented above would be prohibitively inefficient. For example, the difference algorithm uses a subset construction, which obviously takes exponential time. Although, in the worst case, we cannot avoid this blow up, we can apply known optimization techniques for making the algorithms much quicker for practical inputs. Below, we briefly explain some of the techniques that we used in our implementation. (More discussions can be found in [18, 17].)

**Top-down strategy** In this, we calculate an intersection (or difference) of two grammars in such a way that the resulting grammar contains only reachable parts from the “start” variable. (Although start variables are not used in the present algorithms, they usually come with input grammars in actual uses.) This strategy is advantageous when reachable variables are much fewer than the whole product or powerset space, which seems to be the case in the domain of XML.

**Sharing** In this, we employ one global grammar. All input and output grammars are parts of it, and, furthermore, these grammars can share some variables. The advantage is that, in trivial cases like intersecting  $x$  and  $x$  or subtracting  $x$  from  $x$ , the algorithms can immediately return with the obvious answer.

#### 5 Related Work

James Clark [8] has designed a different validation algorithm for attribute-element constraints and implemented it in his validators for TREX and RELAX NG. His algorithm is based on derivatives of regular expressions [4, 1] and proceeds by rewriting a compound expression each time it consumes an attribute or element. Meanwhile, our algorithm allows preprocessors to generate element-attribute automata from a given schema. To validate instance documents, validators can use these element-attribute automata

rather than the original schema. Thus, our algorithm is expected to make validators more compact and efficient. To confirm this expectation, we are actively implementing such a preprocessor and validator for RELAX NG.

Our study on attribute constraints has a strong relationship to type theories for record values (i.e., finite mappings from labels to values). Early papers presenting type systems for record types do not consider the union operator and therefore no such complication arises as in our case. (A comprehensive survey of classical records can be found in [13].) Buneman and Pierce have investigated record types with the union operator [5]. Their system does not, however, have any mechanism similar to our attribute repetitions or recursion. The work that is closest to ours may be the recent one by Frisch, Castagna, and Benzaken [12]. In their language called CDuce, they have a typing mechanism for XML attributes that has a similar expressiveness to ours. However, they have not dealt with the algorithmic problems that we did in this paper.

In his unpublished work, Vouillon has considered an algorithm for checking the subset relation between shuffle expressions [27]. His strategy of progressively decomposing given expressions to two orthogonal parts made much influence on our boolean algorithms. The difference is that his algorithm is for a subset checking, which answers just yes or no, and therefore does not incur the complication of switching back and forth between the expression representation and the automata representation, which is needed in our boolean algorithms since they have to reconstruct expressions.

## 6 Future work

In this paper, we have presented our validation and boolean algorithms. We have already implemented them in the XDuce language [16]. For the examples that we have tried, the performance seems quite reasonable. We plan to collect and analyze data obtained from the experiment on the algorithms in the near future. We also feel that we need some theoretical characterization of the algorithms. In particular, our boolean algorithms contain potentials of blow up in many places. Although our implementation techniques presented in Section 4.5 have been sufficient for our examples, one would like to have some more confidence.

## Acknowledgments

We would like to express our warmest thanks to James Clark, Kohsuke Kawaguchi, Benjamin Pierce, and anonymous PLAN-X referees for precious comments and suggestions. Haruo Hosoya has been supported by Japan Society for the Promotion of Science while working on this paper.

## References

- [1] G. Berry and R. Sethi. From regular expressions to deterministic automata. *Theoretical Computer Science*, 48(1):117–126, 1986.
- [2] T. Bray, D. Hollander, A. Layman, and J. Clark. Namespaces in XML. <http://www.w3.org/TR/REC-xml-names>, 1999.
- [3] T. Bray, J. Paoli, C. M. Sperberg-McQueen, and E. Maler. Extensible markup language (XML<sup>TM</sup>). <http://www.w3.org/XML/>, 2000.
- [4] J. A. Brzozowski. Derivatives of regular expressions. *Journal of the ACM*, 11(4):481–494, Oct. 1964.
- [5] P. Buneman and B. Pierce. Union types for semistructured data. In *Internet Programming Languages*, volume 1686 of LNCS. Springer-Verlag, Sept. 1998. Proceedings of the International Database Programming Languages Workshop.
- [6] J. Clark. The Design of RELAX NG. <http://www.thaiopensource.com/relaxng/design.html>, 2001.
- [7] J. Clark. TREX: Tree Regular Expressions for XML. <http://www.thaiopensource.com/trex/>, 2001.
- [8] J. Clark. <http://www.thaiopensource.com/relaxng/implement.html>, 2002.
- [9] J. Clark and M. Murata. RELAX NG. <http://www.relaxng.org>, 2001.
- [10] D. C. Fallside. XML Schema Part 0: Primer, W3C Recommendation. <http://www.w3.org/TR/xmlschema-0/>, 2001.
- [11] P. Fankhauser, M. Fernández, A. Malhotra, M. Rys, J. Siméon, and P. Wadler. XQuery 1.0 Formal Semantics. <http://www.w3.org/TR/query-semantics/>, 2001.
- [12] A. Frisch, G. Castagna, and V. Benzaken. Semantic subtyping. In *Seventeenth Annual IEEE Symposium on Logic In Computer Science*, 2002.
- [13] R. Harper and B. Pierce. A recrd calculus based on symmetric concatenation. In *Proceedings of POPL*, 1991.
- [14] J. E. Hopcroft and J. D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, 1979.
- [15] H. Hosoya and M. Murata. Validation and boolean operations for attribute-element constraints. <http://www.kurims.kyoto-u.ac.jp/~hahosoya/papers/attelm.ps>, 2002.
- [16] H. Hosoya and B. C. Pierce. XDuce: A typed XML processing language (preliminary report). In *Proceedings of Third International Workshop on the Web and Databases (WebDB2000)*, volume 1997 of LNCS, pages 226–244, May 2000.
- [17] H. Hosoya and B. C. Pierce. Regular expression pattern matching for XML. In *Proceedings of POPL*, 2001.
- [18] H. Hosoya, J. Vouillon, and B. C. Pierce. Regular expression types for XML. In *Proceedings of ICFP*, 2000. Full version under submission to TOPLAS.
- [19] ISO. *Information Processing – Text and Office Systems – Standard Generalized Markup Language (SGML)*, 1986.
- [20] N. Klarlund, A. Moller, and M. I. Schwartzbach. “DSD: A Schema Language for XML”. In *ACM SIGSOFT Workshop on Formal Methods in Software Practice*, 2000.
- [21] T. Milo, D. Suciu, and V. Vianu. Typechecking for XML transformers. In *Proceedings of PODS*, 2000.
- [22] M. Murata. Transformation of documents and schemas by patterns and contextual conditions. In *Principles of Document Processing '96*, volume 1293 of LNCS, pages 153–169. Springer-Verlag, 1997.
- [23] M. Murata. RELAX (REgular LAnguage description for XML). <http://www.xml.gr.jp/relax/>, 2001.
- [24] M. Murata, D. Lee, and M. Mani. Taxonomy of xml schema languages using formal language theory, 2001.
- [25] OASIS. SGML/XML elements versus attributes. <http://xml.coverpages.org/elementsAndAttrs.html>, 2002.
- [26] A. Tozawa. Towards static type inference for XSLT. In *Proceedings of ACM Symposium on Document Engineering*, 2001.
- [27] J. Vouillon. Interleaving types for XML. Personal communication, 2001.