

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

June 6, 2002

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 [BPSMM00] and its predecessor, SGML [Int86], provides two major mechanisms for representing information: *attributes* and *elements*. There has been much debate (e.g., [OAS02]) 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 [BPSMM00], W3C XML Schema [Fal01], and RELAX [Mur01]) 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 [Cla01b], 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 latter 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 [Cla01a], 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 [HP00, FFM⁺01, Toz01, MSV00, Mur97]. (Indeed, we used our boolean algorithms for implementing the typechecker of the XDuce language [HP00].)

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 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 computation 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 that we have seen.

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. Several algorithms that are omitted from the body can be found in Appendixes A and B. We give formalizations and correctness proofs for the boolean algorithms in Appendix C.

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

$$\text{author}^+ \text{title} \text{publisher}^?$$

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 similar notation to regular expressions. For example, the following expression

$$\text{@key @year}^?$$

requires a `key` attribute and optionally allows a `year` attribute. (We prepend an `at-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 (`+`) 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. For example, in the expression

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

we can put any attribute with a name from the `myns` name space, that is, any name of the form `myns:name`. (Name spaces are a prefixing mechanism for

names in XML documents. See [BHLC99] for the details.) In general, attribute repetitions are useful in making a schema “open” so that anyone can put her own pieces of information in unused attributes. Apart from above-mentioned name sets from a specific name space, we can think of the following kinds of name sets: (1) the set of all names, (2) a set of all names except some specific names, and (3) a set of all names except those from some specific name spaces.

2.1 Data model

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

- α is a set of pairs of a name and a value, and
- β is a sequence of pairs of a name and a value.

A pair in α and a pair in β are called *attribute* and *element*, respectively. In the formalization, attributes associate names with values, rather than with strings as in the above examples, for simplicity and generality. We write ϵ for an empty sequence and $\beta_1\beta_2$ for the concatenation of sequences β_1 and β_2 . For convenience, we define several notations for values.

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

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 types.

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 defined by the following

syntax in conjunction with element expressions e .

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

Note that we stratify compound expressions c and element expressions e so that we can ensure attribute expressions not to be enclosed by repetitions. (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.) We define $\text{FV}(c)$ for the set of variables appearing in c and $\text{FV}(G)$ for $\bigcup_{x \in \text{dom}(G)} F(x)$. 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 $\text{att}(c_1) \cap \text{att}(c_2) \neq \emptyset$.¹

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 Appendix B.2.

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

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

A “single-attribute expression” $@N[x]$, which describes a single attribute with a name from N , would also be useful. However, treatment of such an expression is a bit complicated. In the case that N is finite, $@N[x]$ can be rewritten as follows, where $N = \{n_1, \dots, n_k\}$:

$$@N[x] \equiv @\{n_1\}[x]^+ | \dots | @\{n_k\}[x]^+$$

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

On the other hand, our framework has no way to represent the case N is infinite.² 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.

$$\frac{\forall i. (n_i \in N \quad G \vdash v_i \in G(x)) \quad k \geq 1 \quad n_i \neq n_j \text{ for } i \neq j}{G \vdash @n_1[v_1] \dots @n_k[v_k] \in @N[x]^+} \text{ (T-ATTRREP)}$$

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

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

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

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

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

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. We provide here an informal description of our algorithm. Appendix A provides a formalization.

²This restriction has been adopted in RELAX NG.

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 [MLM01]. 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 Σ 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* [HU79]. 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 do not fully present the rest of validation in this paper, but refer to a note [MLM01]. For interested readers, we give a brief sketch of validation algorithms. Understanding of this sketch is not required to read the rest of this paper. Intuitively speaking, we validate a value by assigning variables to all elements (including descendant elements) in this value. Variable assignment must conform to grammars, and we ensure this conformance by executing element automata per element. If we can successfully assign variables to all elements, we report that the value is valid.

3.2 Handling attributes

Attribute-element constraints pose a significant challenge. Since the occurrence order of attributes is insignificant, we cannot use traditional automata for examining attributes in sequence. To illustrate, suppose that a value is constrained by the compound expression

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

As we have observed in Section 2, this compound expression allows $b[\epsilon]@a[\epsilon]$. However, if we assume that the compound expression is a regular expression having $@a[\epsilon]^+$, $@b[\epsilon]^+$, $a[\epsilon]$, and $@b[\epsilon]$ as symbols and create an automaton by some well-known algorithm, the created automaton fails to accept $@b[\epsilon]a[\epsilon]$.

A straightforward 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.

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

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. We say that such compound expressions are in the *guarded normal form*. Any compound expression can be normalized to this normal form, although we do not further consider this normalization in this paper.³ After normalization, we can create an element automaton for each of the element expressions. Given a value $\langle \alpha, \beta \rangle$, we first validate α against each element-free compound expression. Note that we have to recursively validate values of attributes in α , since these values may contain attributes and elements. If α is found valid, we examine β by invoking the element automaton constructed from the corresponding element expression. 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. 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) we introduce special atoms to a compound expression; (2) we create an automaton from the compound expression, (3) given an attribute set, we create an element automaton from this automaton by rewriting some of its transitions. 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 fail to report invalid values invalid. We then present the first step as a remedy to this problem.

We present details of the second step. 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 element automata by rewriting some transitions of the constructed automaton. Given an attribute set α , we replace each $@N[x]$ -transition by a null transition if (1) more than one attribute in α is valid against

³Recall that we treat $@N[x]^+$ as atomic and do not allow $@N[x]$ to be enclosed by repetitions. If we allowed $(@N[\epsilon] a[\epsilon])^+$, we were unable to normalize it to the guarded normal form.

$@N[x]^+$ and (2) no other attributes in α 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 accepts the following four “sequences”.

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

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 accepts the following sequences.

$$\begin{array}{l} b[\epsilon] \\ a[\epsilon]b[\epsilon] \end{array}$$

The only element in $@a[\epsilon]b[\epsilon]$ is $b[\epsilon]$, which is accepted by the element automaton. We can thus correctly report that $@a[\epsilon]b[\epsilon]$ is valid.

Here is the pitfall. 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 *present* in values are *required* by compound expressions.

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. Intuitively speaking, $!@N$ means that no attributes have names in N .

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 α , we replace each $!@N$ -transition by a null transition if no attributes in α 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:

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

From the saturated expression, we create an automaton. It accepts the following four “sequences”.

$$\begin{aligned}
& !@(\mathcal{N} \setminus \{a, b\}) @a[\epsilon]^+ @b[\epsilon]^+ \\
& !@(\mathcal{N} \setminus \{a, b\}) @a[\epsilon]^+ !@b b[\epsilon] \\
& !@(\mathcal{N} \setminus \{a, b\}) !@a a[\epsilon] @b[\epsilon]^+ \\
& !@(\mathcal{N} \setminus \{a, b\}) !@a a[\epsilon] !@b b[\epsilon]
\end{aligned}$$

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 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.

4 Boolean Operations

In this section, we present our algorithms for computing intersections of and differences between compound expressions.

The key technique in our algorithms is *partitioning*. Consider first the following example.

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

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 such a blow-up 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 on 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 [Vou01]. 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 “normalization.” Next, we give some definitions for partitioning and then present our intersection algorithm. Our difference algorithm and implementation techniques can be found in Appendix B.

4.1 Normalization

Our boolean algorithms take two grammars and first pass these to the *normalization* 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. An actual algorithm for normalization is presented in Appendix B.1.

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 $@a[x]^+ (@b[x]^+ | b[x])$ is equal to $(@b[x]^+ | b[x]) @a[x]^+$. On the other hand, $(@a[x]^+ | a[x]) (@b[x]^+ | b[x])$ is *not* equal to $(@b[x]^+ | b[x]) (@a[x]^+ | a[x])$ 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.

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

In implementation, it is desirable to represent equal expressions (in terms of \equiv) by an identical data structure. For example, we may represent an expression by a pair of a bag of expressions and a list of expressions (e.g., represent $(a[x]((@b[x]^+ c[x]) \epsilon))$ $@d[x]^+$ by the pair of the bag of $@b[x]^+$ and $@d[x]^+$ and the list of $a[x]$ and $c[x]$).

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. 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 ϵ). 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. Roughly, it works as follows. 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) $\text{inter}(c_1 (c_2 | c_3) c_4, d) = \text{inter}(c_1 c_2 c_4, d) | \text{inter}(c_1 c_3 c_4, d)$
- 8) $\text{inter}(c, d_1 (d_2 | d_3) d_4) = \text{inter}(c, d_1 d_2 d_4) | \text{inter}(c, d_1 d_3 d_4)$

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 notice that it must be that either (1) one of the argument expressions contains at least one union (except for those enclosed by repetitions) or (2) none of them does so. The first case is handled by rule 6, 7 and 8. The actual algorithm applies rule 6 as often as possible, but rule 7 and 8 can always serve as fall backs. (Note here that we implicitly use the congruence relation \equiv so that we can view, e.g., an expression $c_1 | c_2$ as $\epsilon (c_1 | c_2) \epsilon$ to fit it in the form that the rules accept.) In the second case, both arguments have the form $c_1 \dots c_n$ where each c_i is either an attribute repetition or an element expression. The case that both arguments have length zero or one is handled by rule 1 through 5. The remaining case is therefore that either argument expression has length two or more. By recalling that attributes are normalized, we can always find a proper partition for such an expression; therefore this case is handled by rule 6. (Note again that we implicitly use here that \equiv has the partial commutativity.)

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. The proof is presented in Appendix C.

Theorem 1 *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 $G \vdash v \in d$.*

5 Related Work

James Clark [Cla02] 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 a notion derivatives of regular expressions [Brz64, BS86] and proceeds by rewriting a compound expression each time it consumes an attribute or element. It seems that our automata approach has an advantage in performance since his approach involves frequent operations on expressions during validation. The actual comparison is left for the future work, however.

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 [HP91].) Buneman and Pierce have investigated record types with the union operator [BP98]. 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 [FCB02]. 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 [Vou01]. 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 [HP00]. 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 B.3 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, and Benjamin Pierce for precious comments and suggestions. Haruo Hosoya has been supported by Japan Society for the Promotion of Science while working on this paper.

References

- [BHLC99] Tim Bray, Dave Hollander, Andrew Layman, and James Clark. Namespaces in XML. <http://www.w3.org/TR/REC-xml-names>, 1999.
- [BP98] Peter Buneman and Benjamin Pierce. Union types for semistructured data. In *Internet Programming Languages*, volume 1686 of *Lecture Notes in Computer Science*. Springer-Verlag, September 1998. Proceedings of the International Database Programming Languages Workshop.
- [BPSMM00] Tim Bray, Jean Paoli, C. M. Sperberg-McQueen, and Eve Maler. Extensible markup language (XMLTM). <http://www.w3.org/XML/>, 2000.
- [Brz64] Janusz A. Brzozowski. Derivatives of regular expressions. *Journal of the ACM*, 11(4):481–494, October 1964.
- [BS86] G. Berry and R. Sethi. From regular expressions to deterministic automata. *Theoretical Computer Science*, 48(1):117–126, 1986.
- [Cla01a] James Clark. The Design of RELAX NG. <http://www.thaiopensource.com/relaxng/design.html>, 2001.
- [Cla01b] James Clark. TREX: Tree Regular Expressions for XML. <http://www.thaiopensource.com/trex/>, 2001.
- [Cla02] James Clark. <http://www.thaiopensource.com/relaxng/implement.html>, 2002.
- [Fal01] David C. Fallside. XML Schema Part 0: Primer, W3C Recommendation. <http://www.w3.org/TR/xmlschema-0/>, 2001.

- [FCB02] Alain Frisch, Giuseppe Castagna, and Véronique Benzaken. Semantic subtyping. In *Seventeenth Annual IEEE Symposium on Logic In Computer Science*, 2002.
- [FFM⁺01] Peter Fankhauser, Mary Fernández, Ashok Malhotra, Michael Rys, Jérôme Siméon, and Philip Wadler. XQuery 1.0 Formal Semantics. <http://www.w3.org/TR/query-semantics/>, 2001.
- [HP91] Robert Harper and Benjamin Pierce. A recrd calculus based on symmetric concatenation. In *Proceedings of the Eighteenth Annual ACM Symposium on Principles of Programming Languages, Orlando FL*, pages 131–142. ACM, January 1991. Extended version available as Carnegie Mellon Technical Report CMU-CS-90-157.
- [HP00] Haruo Hosoya and Benjamin C. Pierce. XDuce: A typed XML processing language. In *Proceedings of Third International Workshop on the Web and Databases (WebDB2000)*, volume 1997 of *Lecture Notes in Computer Science*, pages 226–244, May 2000.
- [HP01] Haruo Hosoya and Benjamin C. Pierce. Regular expression pattern matching for XML. In *The 25th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 67–80, January 2001.
- [HU79] John E. Hopcroft and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, 1979.
- [HVP00] Haruo Hosoya, Jérôme Vouillon, and Benjamin C. Pierce. Regular expression types for XML. In *Proceedings of the International Conference on Functional Programming (ICFP)*, pages 11–22, September 2000. Full version under submission to TOPLAS.
- [Int86] International Organization for Standardization. *Information Processing – Text and Office Systems – Standard Generalized Markup Language (SGML)*, 1986.
- [MLM01] M. Murata, D. Lee, and M. Mani. Taxonomy of XML schema languages using formal language theory. <http://citeseer.nj.nec.com/murata00taxonomy.html>, 2001.
- [MSV00] Tova Milo, Dan Suciu, and Victor Vianu. Typechecking for XML transformers. In *Proceedings of the Nineteenth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, pages 11–22. ACM, May 2000.
- [Mur97] Makoto Murata. Transformation of documents and schemas by patterns and contextual conditions. In *Principles of Document Processing '96*, volume 1293 of *Lecture Notes in Computer Science*, pages 153–169. Springer-Verlag, 1997.
- [Mur01] Makoto Murata. RELAX (REgular LAnguage description for XML). <http://www.xml.gr.jp/relax/>, 2001.
- [OAS02] OASIS. SGML/XML elements versus attributes. <http://xml.coverpages.org/elementsAndAttrs.html>, 2002.
- [Toz01] Akihiko Tozawa. Towards static type inference for XSLT. In *Proceedings of ACM Symposium on Document Engineering*, 2001.
- [Vou01] Jérôme Vouillon. Interleaving types for XML. Personal communication, 2001.

A Addendum Algorithms for Validation

In this appendix, we introduce two algorithms. One algorithm saturates compound expressions by introducing non-existent attribute declarations. The other algorithm creates element automata by rewriting some transitions of attribute-element automata.

A.1 Saturating compound expressions

We introduce an operator for saturating compound expressions by introducing non-existent attribute declarations. For each subexpression of the form $c_1 | c_2$, this operator introduces two non-existent attribute declarations.

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

Recall that we forbid expressions with overlapping attribute name sets to be concatenated. This restriction ensures that non-existent attribute declarations do not conflict with expressions of the form $@N[x]$.

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

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

A.2 Rewriting attribute-element automata

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 introduce an operator for rewriting attribute-element automata. Given a set α of attributes, the rewriting operator creates an element automaton from an attribute-element automaton $(Q, q^{\text{init}}, Q^{\text{fin}}, \delta)$. We borrow Q, q^{init} , and Q^{fin} as the state set, start states, and final states of the element automaton, respectively. However, we create the transition relation δ' of the element automaton by rewriting δ . For convenience, we use α_N to denote a subset of α such that α_N contains all attributes in α having names in N . We do not rewrite transitions having labels of the form $N[x]$. The set of these transitions is denoted by δ'_1 . For each transition in δ having $@N[x]$ as a label, we introduce a null transition if α_N is non-empty and the value of any attribute in α_N is

valid against x with respect to G . The set of these transitions is denoted by δ'_2 . For each transition in δ having a label of the form $!@N$, we introduce a null transition if α_N is empty. The set of these transitions is denoted by δ'_3 . We define δ' as the union of δ'_1 , δ'_2 , and δ'_3 . That is,

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

where

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

Note that $G \vdash v \in x$ in the definition of δ'_2 requires recursive validation.

B Addendum Algorithms for Boolean Operations

B.1 Normalization algorithm

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 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 set $N' \cap N$ and the set $N' \setminus N$. We separately shred all the name sets of the first form and those of the second form. We then combine the results with 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.

We can show the following expected properties.

Lemma 1 *Let $S = \text{shred}(\{N_1, \dots, N_k\})$.*

1. *For all $N, N' \in S$, either $N = N'$ or $N \cap N' = \emptyset$.*
2. *For all N_i , there is uniquely $S' \subseteq S$ such that $N_i = \bigcup S'$.*

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.

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

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

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

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 $@N[x]^*$ instead of $@N[x]^+$. 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.

Lemma 2 *Let G be a grammar on X and G' be the normalization of G . Then, $G \vdash v \in G(x)$ iff $G' \vdash v \in G'(x)$ for all v and $x \in X$.*

B.2 Difference

The difference algorithm is basically similar to the intersection algorithm except that we use product construction and subset construction at the same time. As before, let grammars F on X and G on Y be given and have been normalized

simultaneously. 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 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. This function is defined as follows. (Here, D ranges over sets of compound expressions and \uplus is the disjoint union.)

- 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) $\text{diff}(@\{n\}[x]^+, \{@\{n\}[y_1]^+, \dots, \@\{n\}[y_k]^+\}) = \@\{n\}[\langle x, \{y_1, \dots, y_k\} \rangle]^+$
- 6) $\text{diff}(@N[\mathbf{any}]^+, \{@N[\mathbf{any}]^+\}) = \emptyset$
- 7) $\text{diff}(c, \{d_1, \dots, d_k\}) = \big|_{I \subseteq \{1, \dots, k\}} \text{diff}(c', \{d'_i \mid i \in I\}) \text{diff}(c'', \{d''_i \mid i \in \{1, \dots, k\} \setminus I\})$
if $(c', c''), (d'_1, d''_1), \dots, (d'_k, d''_k)$ 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) $\text{diff}(c, \{d_1 (d_2 \mid d_3) d_4\} \uplus D) = \text{diff}(c, \{d_1 d_2 d_4, d_1 d_3 d_4\} \uplus D)$

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 diff^{reg} (explained below). Rules 2, 3, and 4 removes, from the set in the second argument, an expression that is disjoint with the first argument. The remaining base cases are when all the expressions in the arguments are attribute repetitions 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 right hand side is somewhat complicated and may need some explanations, in particular, why we need the “for all subsets $I \subseteq \{1, \dots, k\}$ ” formula. Let us take a value v that is in c and *not* in any d_i . The subsetting formula is relevant to the “not in d_i ” part. Since all d_i 's can be partitioned to $d'_i d''_i$, the value v can also be partitioned to $v' v''$ such that v' contains at most the names in d'_i and v'' contains those in d''_i . Thus, it must be that either (a) $v' \not\subseteq d'_i$ or (b) $v'' \not\subseteq d''_i$. Now, note that the statement “(a) or (b)” must hold for *all* i . Collect all i 's satisfying (a). Then I appearing in the complex formula is the set of such i 's. Likewise, $\{1, \dots, k\} \setminus I$ is the set of i 's satisfying (b). This technique of subsetting has repeatedly been used in the literature. Interested readers are referred to [HVP00, HP01, FCB02].

The function diff^{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:

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

Traditional formal language theories typically study complementation rather than difference. One might wonder why we do not do the same. The reason is that, even if we consider complementation, its subcomputation needs to calculate differences anyway. For example, suppose that we have an expression $c\bar{d}$ where $N = \text{att}(c)$ and $N' = \text{elm}(c)$ with $N' \cap \text{elm}(d) = \emptyset$. Then, we would compute the complementation $\overline{c\bar{d}}$ by

$$((@N[\mathbf{any}]^* N'[\mathbf{any}]^*) \setminus c) d \mid c((@N[\mathbf{any}]^* \overline{N'}[\mathbf{any}]^*) \setminus d)$$

which uses differences. (Note that it is wrong to answer $\overline{c\bar{d}} \mid c\bar{d}$ since the names contained in \bar{c} may overlap with those 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 $@N[\bar{x}]^+$ 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.

It may be worth seeking for an alternative restriction. Here is one possibility. The above argument suggests that the expression we want should look like

$$@N[\bar{x}]^+ @N[x]^*.$$

This expression breaks another restriction: for the disjointness of the name sets in concatenated expressions. Our idea is to relax this. We redefine $\text{att}(c)$ so that it takes contents of attributes into account:

$$\begin{aligned} \text{att}(@N[x]^+) &= N \times \{v \mid G \vdash v \in G(x)\} \\ \text{att}(c_1 c_2) &= \text{att}(c_1) \cup \text{att}(c_2) \\ \text{att}(c_1 \mid c_2) &= \text{att}(c_1) \cup \text{att}(c_2) \\ \text{att}(e) &= \emptyset \end{aligned}$$

We then require, as before, that $\text{att}(c_1) \cap \text{att}(c_2) = \emptyset$ for each concatenation expression $c_1 c_2$. This alternative seems plausible because, given a value v satisfying $c_1 c_2$, we can uniquely determine which of c_1 and c_2 each attribute in v belongs to. However, the correctness remains to be an open question. (We would need an addition restriction that requires expressions for attribute contents to be simple enough so that the disjointness becomes trivial. Also, some external mechanism would be necessary for enforcing attribute names to be disjoint since this alternative restriction does not require this and real XML documents must satisfy this.)

B.3 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 the techniques that we used in our implementation. More discussions can be found in [HVP00, HP01].

B.3.1 Top-down strategy

The intersection algorithm presented above takes grammars F (on X) and G (on Y) and calculates the product of F and G for the whole domain $X \times Y$. However, the actual algorithm takes “start variables” x_0 and y_0 in addition and the useful part of the output grammar is the one that defines the reachable variables from the output start variable $\langle x_0, y_0 \rangle$. Our observation is that such reachable variables are typically much fewer than the whole $X \times Y$. Exploiting this, we construct the output grammar *lazily* from the start variable $\langle x_0, y_0 \rangle$ in a top-down manner. That is, we initialize H to be a grammar just containing the mapping $\langle x_0, y_0 \rangle \mapsto \text{inter}(F(x_0), G(y_0))$. We then take any pair $\langle x, y \rangle \in \text{FV}(H)$ that is not yet in the domain of H , and add the mapping $\langle x, y \rangle \mapsto \text{inter}(F(x), G(y))$. We repeat this until H becomes self-contained. The same technique can be used in the algorithms for product automata, difference grammars, and difference automata.

B.3.2 Sharing

It is quite common that the intersection algorithm encounters, during its computation, a pair of the same expression $c \cap c$. Since the result is trivially c itself whatever c is, we would like to somehow exploit this fact. However, since expressions can be nested and refer to other variables, sameness is a bit complicated. For example, in the following grammars F and G

$$\begin{array}{ll} F(x_1) = a[x_2] & F(x_2) = \epsilon \\ G(y_1) = a[y_2] & G(y_2) = \epsilon \end{array}$$

x_1 and y_1 are the “same,” but detecting this requires a bit of work. In addition, even if we can detect it, we need to copy the whole structure reachable from x_1 to the output grammar.

To deal with this issue, we employ one global grammar that includes all input and output grammars. We modify our intersection algorithm so that it takes only two start variables, proceeds the computation using the global grammar, and adds a new mapping to it whenever necessary. Now, the intersection of x and x can be done by returning x itself. Similarly, we can modify our difference algorithm in a way that the difference between x and Y with $x \in Y$ results in a variable x_\emptyset that is assigned to \emptyset . The same technique can also be used in the product and difference automata algorithms.

B.3.3 Other techniques

Since the function `inter` (and similarly for `diff`) often computes the intersection of the same pair of expressions repeatedly, memoizing the previous computations and reusing the results later are quite helpful. Another technique is to use rules that are specialized to some fixed but frequently used expressions, such as `any` and \emptyset . For example, we can use the special rules $x \cap \mathbf{any} = x$ and $x \setminus Y = x_\emptyset$ when $\mathbf{any} \in Y$.

C Correctness of the Boolean Algorithms

C.1 Element Automata

From the standard automata theory, the following is well-known.

Proposition 1 ([HU79])

1. *There is an algorithm `compile` that constructs an automaton M from a given regular expression e such that $L(M) = L(e)$.*
2. *There is an algorithm `decompile` that constructs a regular expression e from a given automaton M such that $L(e) = L(M)$.*

Given a grammar G on X , an *element automaton* is an automaton on $S \times X$, where S is the set of name sets and. We present the semantics of an element automaton M w.r.t. G by the relation $G \vdash v \in M$ defined as follows:

$$\frac{n_i \in N_i \quad G \vdash v_i \in G(x_i) \quad N_1[x_1] \dots N_k[x_k] \in L(M)}{G \vdash n_1[v_1] \dots n_k[v_k] \in M}$$

Then, we can easily show that this semantics satisfies the following properties.

Corollary 1

1. $G \vdash v \in e$ iff $G \vdash v \in \text{compile}(e)$.
2. $G \vdash v \in M$ iff $G \vdash v \in \text{decompile}(M)$.

C.2 Intersection

Given two element automata $M_i = (Q_i, q_i^{\text{init}}, Q_i^{\text{fin}}, \delta_i)$ on $S \times X_i$ ($i = 1, 2$), 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 $S \times (X_1 \times X_2)$ where

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

By using the standard proof technique, we can show the following expected property of product element automata.

Lemma 3 *Let M be the product automaton of M_1 and M_2 as defined above. Let v be any value. Suppose that, for any x, y, w with $|w| < |v|$,*

$$H \vdash w \in H(\langle x, y \rangle) \text{ iff } F \vdash w \in F(x) \text{ and } G \vdash w \in G(y).$$

Then, we have

$$H \vdash v \in M \text{ iff } F \vdash v \in M_1 \text{ and } G \vdash v \in M_2.$$

Corollary 2 *Let $\text{inter}^{\text{reg}}(e_1, e_2) = e$. With the same assumption as Lemma 3, we have*

$$H \vdash v \in e \text{ iff } F \vdash v \in e_1 \text{ and } G \vdash v \in e_2.$$

The following technical lemma shows that, if a value can be partitioned by disjoint sets of attributes and elements, then such a partition is unique.

Lemma 4 *Given a value v , if $v = u w = u' w'$ where $(\text{elm}(u) \cup \text{elm}(u')) \cap (\text{elm}(w) \cup \text{elm}(w')) = \emptyset$ and $(\text{att}(u) \cup \text{att}(u')) \cap (\text{att}(w) \cup \text{att}(w')) = \emptyset$, then $u = u'$ and $w = w'$.*

PROOF: Let $v = \langle \alpha, \beta \rangle$. Also, let $u = \langle \alpha_1, \beta_1 \rangle$ and $w = \langle \alpha_2, \beta_2 \rangle$; similarly $u' = \langle \alpha'_1, \beta'_1 \rangle$ and $w' = \langle \alpha'_2, \beta'_2 \rangle$. Suppose $u \neq u'$. Then, either $\alpha_1 \neq \alpha'_1$ or $\beta_1 \neq \beta'_1$.

- When $\alpha_1 \neq \alpha'_1$, either $\alpha_1 \setminus \alpha'_1 \neq \emptyset$ or $\alpha'_1 \setminus \alpha_1 \neq \emptyset$. Therefore $\alpha_1 \cap \alpha'_2 \neq \emptyset$ or $\alpha_2 \cap \alpha'_1 \neq \emptyset$. Either case contradicts the condition $(\text{att}(u) \cup \text{att}(u')) \cap (\text{att}(w) \cup \text{att}(w')) = \emptyset$.
- When $\beta_1 \neq \beta'_1$, either (1) $\beta_1 = \beta'_1 \gamma$ and $\beta_2 = \gamma \beta_2$ for some non-empty γ or (2) $\beta'_1 = \beta_1 \gamma$ and $\beta_2 = \gamma \beta'_2$ for some non-empty γ . Either case contradicts the condition $(\text{elm}(u) \cup \text{elm}(u')) \cap (\text{elm}(w) \cup \text{elm}(w')) = \emptyset$. \square

We use the following function weight from compound expressions to integers for induction in the subsequent proofs.

$$\begin{aligned} \text{weight}(c_1 c_2) &= \text{weight}(c_1) \text{weight}(c_2) \\ \text{weight}(c_1 | c_2) &= 1 + \text{weight}(c_1) + \text{weight}(c_2) \\ \text{weight}(c) &= 1 \quad \text{if } c \neq c_1 c_2 \text{ and } c \neq c_1 | c_2 \end{aligned}$$

The following lemma shows the correctness of the function inter .

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 $G \vdash v \in d$.

PROOF: Let $\text{measure}(v, c, d) = (|v|, \text{weight}(c) + \text{weight}(d), \text{width}(c) + \text{width}(d))$. By induction on the lexicographic order of $\text{measure}(v, c, d)$.

Case: $c = e$ (i.e., $\text{att}(c) = \emptyset$) $d = f$ $b = \text{inter}^{\text{reg}}(e, f)$

The induction hypothesis allows us to use Corollary 2, from which the result follows.

Case: $c = e$ $d = @N[y]^+$ $b = \emptyset$ or
 $c = @N[x]^+$ $d = f$ $b = \emptyset$ or
 $c = @N[x]^+$ $d = @N'[y]^+$ $N \neq N'$ $b = \emptyset$

Trivial.

Case: $c = @N[x]^+$ $d = @N[y]^+$ $b = @N[\langle x, y \rangle]^+$

Clearly, v has the form $@n_1[v_1] \dots @n_k[v_k]$. By definition, $F \vdash v \in c$ and $G \vdash v \in d$ iff $k \geq 1$ and $n_i \in N$ with $F \vdash v_i \in F(x)$ and $G \vdash v_i \in G(y)$. By induction hypothesis, the last two conditions equal to $H \vdash v_i \in H(\langle x, y \rangle)$. Thus, the result follows.

Case: $b = \text{inter}(c_1, d_1) \text{inter}(c_2, d_2)$
 $(c_1, c_2), (d_1, d_2)$ is a proper partition of c, d

By the definition of partition, the result suffices to show

$$F \vdash v \in c_1 c_2 \text{ and } G \vdash v \in d_1 d_2 \text{ iff } H \vdash v \in \text{inter}(c_1, d_1) \text{inter}(c_2, d_2) \quad (2)$$

with $(\text{att}(c_1) \cup \text{att}(d_1)) \cap (\text{att}(c_2) \cup \text{att}(d_2)) = \emptyset$ and $(\text{elm}(c_1) \cup \text{elm}(d_1)) \cap (\text{elm}(c_2) \cup \text{elm}(d_2)) = \emptyset$. Since the partition is proper, $\text{measure}(v, c_1, d_1) < \text{measure}(v, c, d)$ and $\text{measure}(v, c_2, d_2) < \text{measure}(v, c, d)$. This allows us to apply the induction hypothesis:

$$F \vdash v \in c_1 \text{ and } G \vdash v \in d_1 \text{ iff } H \vdash v \in \text{inter}(c_1, d_1)$$

and similarly for c_2, d_2 . From this, the “if” direction of the statement (2) is obvious.

To show the other direction, let v satisfy the left hand side of (2). Then, there are v_1 and v_2 where $v = v_1 v_2$ with $F \vdash v_1 \in c_1$ and $F \vdash v_2 \in c_2$; also, there are v'_1 and v'_2 where $v = v'_1 v'_2$ with $G \vdash v'_1 \in d_1$ and $G \vdash v'_2 \in d_2$. From Lemma 4, $v_1 = v'_1$ and $v_2 = v'_2$. By the induction hypothesis shown above, we obtain $H \vdash v_1 \in \text{inter}(c_1, d_1)$ and similarly for v_2 , from which the result follows.

Case: $c = c_1 (c_2 | c_3) c_4$ $b = \text{inter}(c_1 c_2 c_4, d) | \text{inter}(c_1 c_3 c_4, d)$

Clearly $F \vdash v \in c$ if and only if either $F \vdash v \in c_1 c_2 c_4$ or $F \vdash v \in c_1 c_3 c_4$. Note $\text{measure}(v, c_1 c_2 c_4, d) < \text{measure}(v, c, d)$ and similarly for $c_1 c_3 c_4$. The result follows from the induction hypothesis.

Case: $d = d_1 (d_2 | d_3) d_4$ $b = \text{inter}(c, d_1 d_2 d_4) | \text{inter}(c, d_1 d_3 d_4)$

Similar to the previous case. \square

C.3 Difference

Given two element automata $M_i = (Q_i, q_i^{\text{init}}, Q_i^{\text{fin}}, \delta_i)$ on $S \times X_i$ ($i = 1, 2$), the *difference* of M_1 and M_2 is an automaton $(Q, q^{\text{init}}, Q^{\text{fin}}, \delta)$ on $S \times (X_1 \times \mathcal{P}(X_2))$ where:

$$\begin{aligned} Q &= Q_1 \times \mathcal{P}(Q_2) \\ q^{\text{init}} &= \langle q_1^{\text{init}}, \{q_2^{\text{init}}\} \rangle \\ Q^{\text{fin}} &= \{ \langle q_1, P \rangle \mid q_1 \in Q_1^{\text{fin}} \text{ and } P \subseteq Q_2 \text{ and } P \cap Q_2^{\text{fin}} \neq \emptyset \} \\ \delta &= \{ (\langle q_1, \{p_1, \dots, p_k\} \rangle, N[\langle x_1, \{y_1, \dots, y_k\} \rangle], \langle q'_1, \{p'_1, \dots, p'_k\} \rangle) \mid \\ &\quad (q_1, N[x_1, q'_1]) \in \delta_1 \text{ and } \{(p_1, N[y_1, q'_1]), \dots, (p_k, N[y_k, q'_k])\} \subseteq \delta_2 \} \end{aligned}$$

By combining the standard proof technique for product construction and subset construction, we can show the following.

Lemma 5 *Let M be the difference automaton of M_1 and M_2 as defined above. Let v be any value. Suppose that, for any x, y, w with $|w| < |v|$,*

$$H \vdash w \in H(\langle x, Y \rangle) \text{ iff } F \vdash w \in F(x) \text{ and } G \vdash w \notin G(y) \text{ for all } y \in Y.$$

Then, we have

$$H \vdash v \in M \text{ iff } F \vdash v \in M_1 \text{ and } G \vdash v \notin M_2.$$

Corollary 3 *Let $\text{diff}^{\text{reg}}(e_1, e_2) = e$. With the same assumption as Lemma 5, we have*

$$H \vdash v \in e \text{ iff } F \vdash v \in e_1 \text{ and } G \vdash v \notin e_2.$$

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 $d \in D$.*

PROOF: Let $\text{measure}(v, c, D)$ be

$$\left(|v|, \text{weight}(c) + \sum_{d \in D} \text{weight}(d), \text{width}(c) + \sum_{d \in D} \text{width}(d) \right).$$

By induction on the lexicographic order of $\text{measure}(v, c, D)$.

$$\text{Case: } c = e \quad D = \{f_1, \dots, f_k\} \quad b = \text{diff}^{\text{reg}}(e, f_1 \mid \dots \mid f_k)$$

The induction hypothesis allows us to use Corollary 3, from which the result follows.

$$\text{Case: } c = e \quad D = \{ @N[y]^+ \} \uplus D' \quad b = \text{diff}(e, D')$$

Since e and $@N[y]^+$ are obviously disjoint, the result follows by the induction hypothesis.

$$\begin{aligned} \text{Case: } c = @N[x]^+ \quad D = \{f\} \uplus D' \text{ or } D = \{ @N'[y]^+ \} \uplus D' \\ b = \text{diff}(@N[y]^+, D') \end{aligned}$$

Similar to the previous case.

$$\text{Case: } c = @N[\mathbf{any}]^+ \quad D = \{ @N[\mathbf{any}]^+ \} \quad b = \emptyset$$

Trivial.

$$\begin{aligned} \text{Case: } c &= @\{n\}[x]^+ & D &= \{ @\{n\}[y_1]^+, \dots, @\{n\}[y_k]^+ \} \\ b &= @\{n\}[\langle x, \{y_1, \dots, y_k\} \rangle]^+ \end{aligned}$$

By definition, $F \vdash v \in c$ and $G \vdash v \notin @\{n\}[y_i]^+$ for all $i = 1, \dots, k$ if and only if $l \geq 1$ and $v = @n[v_1] \dots @n[v_l]$ with $F \vdash v_j \in F(x)$ and $G \vdash v_j \notin G(y_i)$ for all $i = 1, \dots, k$ and $j = 1, \dots, l$. By the induction hypothesis, this is equivalent to saying $H \vdash v \in b$.

$$\begin{aligned} \text{Case: } D &= \{ d_1, \dots, d_k \} \\ b &= \bigsqcup_{I \subseteq \{1, \dots, k\}} \text{diff}(c', \{d'_i \mid i \in I\}) \text{diff}(c'', \{d''_i \mid i \in \{1, \dots, k\} \setminus I\}) \\ & \quad (c', c''), (d'_1, d''_1), \dots, (d'_k, d''_k) \text{ is a proper partition of } c, d_1, \dots, d_k \end{aligned}$$

By the definition of partition, we have $c = c' c''$ and $d_i = d'_i d''_i$ with $(\text{att}(c') \cup \bigcup_i \text{att}(d'_i)) \cap (\text{att}(c'') \cup \bigcup_i \text{att}(d''_i)) = \emptyset$ and $(\text{elm}(c') \cup \bigcup_i \text{elm}(d'_i)) \cap (\text{elm}(c'') \cup \bigcup_i \text{elm}(d''_i)) = \emptyset$. We first show the “only if” direction. Let $H \vdash v \in b$. Then, there are v' and v'' such that $H \vdash v' \in \text{diff}(c', \{d'_i \mid i \in I\})$ and $H \vdash v'' \in \text{diff}(c'', \{d''_i \mid i \in \{1, \dots, k\} \setminus I\})$. By the induction hypothesis, we have

$$F \vdash v' \in c' \text{ and } G \vdash v' \notin d'_i \text{ for all } i \in I$$

and

$$F \vdash v'' \in c'' \text{ and } G \vdash v'' \notin d''_i \text{ for all } i \in \{1, \dots, k\} \setminus I.$$

These imply that $F \vdash v \in c$ and $G \vdash v \notin d_i$ for all $i \in \{1, \dots, k\}$.

We next show the “if” direction. Let $F \vdash v \in c$ and $G \vdash v \notin d_i$ for all $i \in \{1, \dots, k\}$. Then, there are v' and v'' such that $v = v' v''$ with $F \vdash v' \in c'$ and $G \vdash v'' \in c''$. From Lemma 4, v' and v'' are uniquely determined. Therefore either $G \vdash v' \notin d'_i$ or $G \vdash v'' \notin d''_i$. Since this holds for all i , there is $I \subseteq \{1, \dots, k\}$ such that $G \vdash v' \notin d'_i$ for all $i \in I$ and $G \vdash v'' \notin d''_i$ for all $i \in \{1, \dots, k\} \setminus I$. By the induction hypothesis,

$$H \vdash v' \in \text{diff}(c', \{d'_i \mid i \in I\})$$

and

$$H \vdash v'' \in \text{diff}(c'', \{d''_i \mid i \in \{1, \dots, k\} \setminus I\}).$$

The desired result $H \vdash v \in b$ follows from these.

$$\text{Case: } c = c_1 (c_2 | c_3) c_4 \quad b = \text{diff}(c_1 c_2 c_4, D) | \text{diff}(c_1 c_3 c_4, D)$$

The result can be shown straightforwardly by using that

$$F \vdash v \in c_1 (c_2 | c_3) c_4 \text{ iff } F \vdash v \in c_1 c_2 c_4 \text{ or } F \vdash v \in c_1 c_3 c_4$$

plus the induction hypothesis.

$$\begin{aligned} \text{Case: } D &= \{ d_1 (d_2 | d_3) d_4 \} \uplus D' \\ b &= \text{diff}(c, \{d_1 d_2 d_4, d_1 d_3 d_4\} \uplus D') \end{aligned}$$

Similar to the previous case. □