

# Foundations of Fast Communication via XML

Welf Löwe, Markus L. Noga, Thilo S. Gaul

Institut für Programmstrukturen und Datenorganisation, Universität Karlsruhe  
Postfach 6980, 76128 Karlsruhe, Germany

E-mail: {loewe | noga | gaul}@ipd.info.uni-karlsruhe.de

## Abstract

*Communication with XML often involves pre-agreed document types. In this paper, we propose an offline parser generation approach to enhance online processing performance for documents conforming to a given DTD. Our examination of DTDs and the languages they define demonstrates the existence of ambiguities. We present an algorithm that maps DTDs to deterministic context-free grammars defining the same languages. We prove the grammars to be LL(1) and LALR(1), making them suitable for standard parser generators. Our experiments show the superior performance of generated optimized parsers. Our results generalize from DTDs to XML Schema specifications with certain restrictions, most notably the absence of namespaces, which exceed the scope of context-free grammars.*

## 1 Introduction

The Extensible Markup Language (XML) [15] provides for logical markup: documents are tagged according to their content structure, not their visual appearance in specific presentation media. Sets of documents can be specified with Document Type Definitions (DTDs). A given document can be validated against a given DTD.

One area of application for XML is web serving. Content originating with diverse sources such as static documents or database queries is represented in XML. These documents are mapped to presentation formats like HTML through Extensible Stylesheet Language Transformation (XSLT) scripts or cascades thereof. Although certain generalizations are possible, a given script usually works only with documents conforming to a specific DTD. Parsing the document may account for half the processing time or more, so performance matters when serving dynamic documents.

Over the past years, XML also became popular for general-purpose platform-independent data exchange. When connecting software components, XML processing layers compete with middleware architectures such as Corba [14] and (D)COM [6]. In this arena, performance is crucial, and the DTDs are always available at compile time.

Most available XML parsers are generic: they read arbitrary well-formed XML, analyze the required DTD and validate the document against it. Along with many others, the parsers provided by the Apache project [2], James Clark [6] and IBM [10] fall into this category. As generic parsers cannot use DTDs to optimize parsing, faster solutions are possible.

We take the approach of specific parsing, where a parser applies to documents conforming to a given DTD only. DTDs may change rapidly in software evolution, so manual implementation of specific parsers is a non-option. Compiler construction practice shows deterministic parsers to be manifestly faster than their ambiguous equivalents. Thus, we focus on the automatic generation of deterministic parsers.

At first glance, DTDs are ill-suited to this technique: while the XML specification [15] states deterministic ones as desirable in appendix E, the entire appendix is non-normative. Many real-world DTDs, including some published in W3C documents, are in fact ambiguous.

For regular languages, the existence of deterministic grammars and corresponding automata is a well-known result from formal language theory. Due to its matching opening and closing tags, which correspond to nested parentheses in programming languages, XML is not regular. It belongs to the larger class of context free languages. In general, the question if there is an equivalent deterministic grammar for a given context free grammar is not decidable.

This paper proves the existence of equivalent deterministic grammars for DTDs, i.e., for a subset of the context free languages including XML. The proof is constructive and defines a mapping algorithm. The generated grammars are both  $LL(1)$  and  $LALR(1)$ . Thus, standard parser generators for  $LL(k)$  grammars or subsets thereof, such as `ell` [9,16] and `javacc` [21] can be employed to generate efficient parsers, as well as tools for  $LALR(1)$  grammars like `yacc` [12], `bison` [8] and `lalr` [9,16]. Our experiments show them to outperform generic parsers by an order of magnitude.

The body of this paper is organized as follows: Section 2 defines notions and cites results basic to our approach. Section 3 demonstrates DTD ambiguities and shows the transformation of a DTD into a context free grammars. The generated grammars are proven to be  $LL(1)$  and  $LALR(1)$ . Section 4 generalizes these results to some extensions of XML, that is, XML Schema specifications and namespaces. Section 5 compares the performance of parsers generated with our techniques to generic ones. We examine related work in section 6. Section 7 concludes the paper with directions for future work. The appendix contains additional basic definitions and a larger example.

## 2 Basic Definitions and Results

The following definitions and theorems can be found in every compiler construction or formal languages textbook. They are included here to introduce notational conventions. For details, we refer to [1,20]. We assume the reader is familiar with the notions of substitution and rewrite systems. Definitions of the notions of a formal language, a grammar, deterministic and ambiguous grammar, regular expressions and languages, and context-free grammars and languages can be found in the appendix.

These two theorems state that it is trivial to find a deterministic grammar for a regular language, but hard for a context free language.

### **Theorem 1 Deterministic Regular Grammars**

*There is an algorithm to find a deterministic grammar for every regular language.*

### **Theorem 2 Deterministic Context Free Grammars**

*It is not decidable whether a context-free language can be defined by a deterministic grammar.*

Unfortunately, XML is not regular. To establish our results, we need the  $SLL(1)$  and  $SLR(1)$  grammar classes, which in turn require a basic understanding of deterministic LL and LR parsing strategies.

Both variants operate on an input stream of terminal symbols and an analysis stack containing terminal and non-terminal symbols. Deterministic LL parsers operate top-down. They derive words from the starting symbol  $Z$ , which is initially on the analysis stack. In each step, an LL parser distinguishes three cases:

- (1) If the top of the analysis stack is a terminal symbol  $t$  and the current input symbol is  $t$ , the top of the analysis stack is removed and the next input symbol is read.
- (2) If the top of the analysis stack is a terminal symbol  $t$  and the current input symbol is  $t' \neq t$ , parsing stops with an error (ambiguous parsers back-track in this case).
- (3) If the top of the analysis stack is a non-terminal symbol  $N$ , it is replaced by the right-hand side of a production with left-hand side  $N$ .

Analysis terminates successfully iff both input stream and analysis stack are empty.

Deterministic LR parsers work bottom-up. They start with an empty analysis stack and reduce the input word stepwise to the starting symbol  $Z$ . In each step, LR parsers either

- (I) shift terminal symbols onto the stack or
- (II) reduce the topmost symbols on the stack to a terminal  $N$  if there is a production  $N \rightarrow s$  whose right-hand side equals the topmost symbols reversed.

Analysis terminates successfully iff the input stream is empty and the analysis stack contains only  $Z$ .

For both types of parser, choosing the correct action is critical. LL parsers must load the correct right-hand side of a production in step (3) in the presence of multiple alternatives. LR parsers must decide whether to shift or to reduce if reduction is possible (shift-reduce conflict) and whether to reduce to  $N$  or  $N'$  if multiple alternatives are applicable (reduce-reduce conflict).

It is desirable to define a language with a grammar that allows the efficient resolution of these conflicts, to allow efficient parsing without backtracking. Two examples of such grammar classes are  $SLL(k)$  and  $LR(0)$ , whose original definitions can be found in [15] and [6].

**Definition 1**  $SLL(k)$  Grammar

Let  $G$  be a context free grammar.  $G$  is  $SLL(k)$  iff the next  $k$  input symbols are always sufficient for an LL parser to choose the correct right-hand side to load in step (3).

**Definition 2**  $LR(0)$  Grammar

Let  $G$  be a context free grammar.  $G$  is  $LR(0)$  iff regular matches on the analysis stack state are always sufficient for an LR parser to correctly resolve all shift-reduce and reduce-reduce conflicts.

To decide whether a grammar is  $SLL(k)$ , we define the  $k$ -head of a word  $w$ , denoted by  $k:w$ , the set of all  $k$ -heads of terminals of a word  $FIRST_k(w)$  and the set of all  $k$ -heads of terminals following a word  $FOLLOW_k(w)$ . Given a grammar  $G=(T, N, P, Z)$  with terminals  $T$ , non-terminals  $N$ , productions  $P$  and starting symbol  $Z$ :

**Definition 3 The  $k$ -head and First and Follow Sets**

Let  $\#$  be a symbol,  $\# \notin T$ , which marks the end of a word. Let  $w \in (T \cup N)^*$ , let  $\overset{*}{\vdash}$  denote a derivation of arbitrary length. The  $k$ -head of a word  $w$  is defined as

$$\begin{aligned} k:w &= a, & \text{when } w=ay, |a|=k \text{ and} \\ k:w &= w\#, & \text{when } |w|<k. \end{aligned}$$

The first and follow sets, respectively, are defined by:

$$\begin{aligned} FIRST_k(w) &= \{ r \mid \exists v \in T^* \text{ with } w \overset{*}{\vdash} v \text{ and } r=k:v \} \\ FOLLOW_k(w) &= \{ r \mid \exists v \in (T \cup N)^* \text{ with } Z \overset{*}{\vdash} uvv \text{ and } r \in FIRST_k(v) \} \end{aligned}$$

Remark: We omit the subscript  $k$  for  $k = 1$ .

**Definition 4  $SLR(k)$  Grammar**

Let  $G$  be a context free grammar.  $G$  is  $SLR(k)$  iff

- for all pairs of productions  $N_x \rightarrow x$  and  $N_y \rightarrow x$  causing a reduce-reduce conflict in an  $LR(0)$  parser,  $FOLLOW_k(N_x) \cap FOLLOW_k(N_y) = \emptyset$ , i.e., the next  $k$  input symbols determine whether to reduce to  $N_x$  or to  $N_y$ .
- for all pairs of productions  $N_x \rightarrow x$  and  $N_y \rightarrow xy$  causing a shift-reduce conflict in an  $LR(0)$  parser,  $FOLLOW_k(N_x) \cap FIRST_k(y FOLLOW_k(N_y)) = \emptyset$ , i.e., the next  $k$  input symbols determine whether to reduce to  $N_x$  or to shift the symbols belonging to  $y$ .

Remark:  $SLR(k)$  grammars allow to resolve the remaining conflicts of an  $LR(0)$  parser using the next  $k$  input symbols.

Using the first and follow sets, the question whether a given grammar is  $SLL(k)$  can be decided with the following theorems:

**Theorem 3  $SLL(k)$  Property**

A grammar is  $SLL(k)$  iff for all pairs of productions  $N \rightarrow x$  and  $N \rightarrow y$  with  $x \neq y$ :

$$FIRST_k(x FOLLOW_k(N)) \cap FIRST_k(y FOLLOW_k(N)) = \emptyset.$$

From a practical viewpoint, narrow definitions of language classes are less important than the potential to apply existing tools. Although  $SLL(k)$  and  $SLR(1)$  are stricter language classes, the applicability of  $LL(k)$  and  $LALR(1)$  parser generators is the dominant concern. The original definition of look-ahead- $LR(1)$  -  $LALR(1)$  - grammars and parser generation algorithms are not essential to this paper, but covered in [6] and [20]. The following theorems deliver the results of interest:

**Theorem 4  $SLL(k)$  and  $LL(k)$  Grammars**

$SLL(k) \subset LL(k)$ ,  $SLL(1) = LL(1)$ .

**Theorem 5  $LR(k)$ ,  $SLR(k)$  and  $LALR(k)$  Grammars**

$LR(0) \subset SLR(k) \subset LALR(k) \subset LR(k)$ .

### 3 DTDs and Grammars

This section briefly examines well-formed XML and DTDs with their components. Examples for ambiguities in DTDs are given. Algorithms mapping individual DTD components and entire DTDs to equivalent grammars are described.

### 3.1 XML and DTDs

XML documents consist of elements, attributes and text. Text is a character sequence in the specified coding system, which must not contain the '<' character used in the representation of elements. An attribute is a triple of attribute name, '='-symbol and quoted attribute value (e.g. `a="foo"`), separated by arbitrary whitespace. We omit whitespace in our presentation.

Elements are defined recursively. They consist of the following sequence:

- (1) a start tag containing the element name (e.g. `<x`)
- (2) a sequence of attributes, followed by the '>'-symbol (e.g. `a="foo" b="bar">`)
- (3) a sequence of elements or text
- (4) an end tag repeating the element name (e.g. `</x>`).

The sequences in (2) must not contain multiple occurrences of the same attribute name.

In general, (3) may also contain comments, entities, processing instructions and unparsed data sections. We ignore them to simplify this presentation, as they are simple regular tokens. For empty elements the sequence in (3) has no entries. They can be represented using a short-cut that closes the element after the attribute definitions with '/>' (e.g. `<x a="foo" b="bar" />`).

DTDs define the structure of sets of documents. In this paper, we assume that all implicit and external components are inlined. Then, DTDs contain two kinds of definition components, element definitions and attribute definitions. Element definitions have the form:

```
<!ELEMENT x ( content_model_x )>
```

where `x` is the name of the element to be defined, and `content_model_x` is either a regular expression over element names and `#PCDATA` or `#EMPTY`, which denotes a regular expression accepting only the empty string. `#PCDATA` is a placeholder for possibly empty text. Regular expressions may employ the well-known operators iteration ('\*' and '+'), alternative choice ('|'), option('?') and sequence(','), as well as grouping parentheses.

Attribute definitions take the form

```
<!ATTRIBUTE x a t o>
```

where `x` is the element name for which attribute `a` is defined as type `t`, with lower and upper bounds  $l, u \in \{0, 1\}$  on its occurrence given in `o`. W.l.o.g., we ignore that DTDs also allow multiple attributes of an element to be specified in a single definition. In XML documents conforming to the DTD, an element's attributes may appear in any order.

### 3.2 Ambiguous DTDs

Although a DTD defines the structure of documents concisely, the derivation of an XML document conforming to this DTD may be ambiguous.

We give three examples of ambiguities in content models.

**Example 1:** An ambiguous content model

---

```
<!ELEMENT a (#PCDATA | b)*>
```

---

According to the definition, an `a` may contain an arbitrary sequence of `#PCDATA` and `b` children. `#PCDATA` is in turn a possibly empty string. Therefore, this XML fragment cannot be parsed deterministically:

---

```
<a></a>
```

---

It may be interpreted as an empty iteration of `(#PCDATA | b)*`. Alternatively, it can be seen as a single, empty `#PCDATA` element.

This more complex example is an excerpt of a DTD in an XML Schema working draft [19].

**Example 2:** An ambiguous content model for `complexType`

---

```
<!ELEMENT %complexType; ( (%annotation;)?,
    ( (%facet;)* |
      ( (%element; | %mgs; | %group; | %any;)*,
        (%attribute; | %attributeGroup;)*,
        (%anyAttribute;)?
      )
    )
  )>
```

---

The ambiguity is due to possibly empty iterations composed by alternative and sequence operators. A simplification plainly shows the problem:

---

```
<!ELEMENT a ( x? , ( y* | z* ) )>
```

---

The empty element

---

```
<a></a>
```

---

may be interpreted as the missing `x` and an empty iteration of `y` or, alternatively, missing `x` and an empty iteration of `z`. This content model can be transformed into a deterministic definition:

---

```
<!ELEMENT a ( (x, (y+|z+)?) | y+ | z+ ) >
```

---

Deterministic definitions tend to become large compared to the ambiguous ones. The human reader has no problem to understand the structure defined by ambiguous DTDs. As deterministic DTDs are allowed, they will occur in practice. At first glance, the situation appears even worse: some ambiguous content models have no deterministic equivalent.

**Example 3:** The following ambiguous DTD cannot be defined by a deterministic DTD:

---

```
<!ELEMENT a ( (x|y)* , x , (x|y) ) >
```

---

Moreover, individual deterministic content models do not imply a deterministic overall grammar. We will show in the next subsections that for DTDs, they do.

### 3.3 Grammars for DTD Components

Ambiguities in DTDs are not a property of the languages they define. There is a deterministic context-free language for each DTD defining the same language, as the following lemma shows.

**Lemma 1 Content models and the *SLL(1)* property.**

*For every content model, there is an *SLL(1)* grammar defining the same language.*

**Proof:** We define an algorithm mapping a content model to a context free grammar. We show the grammar to define the same language and prove membership in *SLL(1)*:

Let  $S$  be the set of all element names in a DTD and  $S' = S \cup \{string\}$ , where string must not contain '<'. Replace each occurrence of **#PCDATA** by *string?* in the element definition expressions. Obviously, each of the modified expressions is a regular expression over  $S'$  per Definition 8 in the appendix. We construct a minimal deterministic acceptor from the regular expression with standard formal language techniques:

- (1) Construct a ambiguous acceptor from the regular expression.
- (2) Make the ambiguous acceptor deterministic using the standard power set construction.
- (3) Minimize the deterministic acceptor using the standard equivalence class construction.

Then, we generate the grammar directly from this acceptor:

- (1) For each state  $t$  of the minimal deterministic acceptor, we generate a unique non-terminal  $N(t)$ .
- (2) If  $t$  is an accepting state, we generate the production  $N(t) \rightarrow \epsilon$ .
- (3) If there is a transition from  $t$  to  $t'$  accepting symbol  $s \in S'$ , we generate the production  $N(t) \rightarrow s N(t')$ .

All above transformations are standard (see [20]) and the constructed grammar is proven to define the same language as the original regular expression. It remains to show the resulting grammar is *SLL(1)*. We use Theorem 3. Suppose there are two productions  $N(t) \rightarrow x$  and  $N(t) \rightarrow y$  with  $x \neq y$ . Let  $N(t)$  correspond to state  $t$ . According to our construction, we distinguish two cases:

- (1)  $x = aN(t')$  and  $y = bN(t'')$ , with  $a, b \in S^*$  and non-terminals  $N(t')$  and  $N(t'')$ . As the acceptor is deterministic, it follows  $a \neq b$  and

$$FIRST(x FOLLOW(N(t))) \cap FIRST(y FOLLOW(N(t))) = \{a\} \cap \{b\} = \emptyset$$

- (2)  $x = aN(t')$  and  $y = \epsilon$ . By construction, the follow set of each non-terminal is the virtual end-of-sentence marker '#'. Thus

$$FIRST(x FOLLOW(N(t))) \cap FIRST(y FOLLOW(N(t))) = \{a\} \cap \{\#\} = \emptyset$$

Using this result, a corresponding lemma for attribute definitions can be established.

**Lemma 2 Attribute definitions and the SLL(1) property.**

*For every set of attribute definitions of an element, there is an SLL(1) grammar defining the same language.*

**Proof:** Let  $R_e$  be the finite set of attribute definitions for element  $e$  in the DTD. In conforming XML documents, instances of  $e$  may contain at most one instance of every attribute in  $R_e$ , so every legal attribute instance sequence is finite. As there are only a finite number of permutations, the set of legal attribute instance sequences is also finite. Individual attribute instances are regular as they consist of three regular parts: the regular attribute name, the '=' character and the regular attribute value string.

Finite sets of finite sequences of regular components are regular, so we obtain an SLL(1) grammar by applying the algorithm in Lemma 1 to an appropriate regular expression.

Remark: Note that not all regular grammars are SLL(1). The regular productions  $S \rightarrow a A$  and  $S \rightarrow a B$ , e.g., have the same FIRST set. The construction from the *deterministic* regular acceptor guarantees that such ambiguities do not occur.

The generated grammars capture all aspects of attribute occurrences, but due to their exponential size, they are mostly theoretical in nature. Practical implementations are advised to use the Kleene closure of  $R_e$  instead and limit occurrences through non-grammatical means. However, these results pave the way for the synthesis of grammars for DTD languages in the following section.

According to the transformations defined above, we obtained deterministic grammars for content models and attributes. Their productions have the form

- (1)  $N \rightarrow \epsilon$   
(2)  $N' \rightarrow s N''$

where  $s \in S^*$  are terminal symbols and  $N, N'$  and  $N''$  are non-terminal symbols.

**Lemma 3 Content models, attribute definitions and the SLR(1) property.**

*For every content model as well as every set of attribute definitions of an element, there is an SLR(1) grammar defining the same language*

**Proof:** Since the content model as well as the attribute definitions can be transformed into a grammar of the above form, we can prove the lemma by showing that deterministic grammars of this form are SLR(1) in general. We sketch an LR(0) parsing algorithm, which is an even stronger restriction.

First, we shift the entire word into the analysis stack. Additionally, we perform the state transitions of the deterministic acceptor the grammar is generated from (according to the proof of Lemma 1). If this does not lead to a final state, the sentence does not belong to the language.

- (1) W.l.o.g., let  $t$  be the final state of the acceptor. We reduce with production  $N(t) \rightarrow \epsilon$ .  
(2) Perform (3) until the start symbol is the only symbol of the analysis stack.



- (3) W.l.o.g., let  $wsN(t)$ ,  $w \in T^*$ ,  $s \in T$  and  $N(t) \in N^*$ , be the content of the analysis stack. Then, there must be a production  $N(t) \rightarrow sN(t) \in P$ . Furthermore, accepting the word  $w$  with the deterministic acceptor ends in the state  $t'$ . We reduce with this production.

As no look-ahead is performed, the parser is  $LR(0)$  and thus  $SLR(1)$ .

Remark: It is known that every regular language allows for  $LR(0)$  parsing. We included a constructive proof to prepare our main results in the next section as a generalization of this algorithm.

### 3.4 Grammars for entire DTDs

A DTD  $D$  and a root element  $Z$  together define a language  $L_{D,Z}$ . It defines algorithms mapping a DTD  $D$  and a root element  $Z$  to a grammar  $G$  with  $L(G)=L_{D,Z}$ . Furthermore, we prove that  $G$  is  $LL(1)$  and  $LALR(1)$ . Due to Theorem 4 and Theorem 5, it is sufficient to show that  $G$  is  $SLL(1)$  and  $SLR(1)$ , respectively.

Proceeding from the grammars for DTD components derived in the previous section, we are now ready to generate useful grammars for a language  $L_{D,Z}$  defined by a DTD  $D$  and a starting element  $Z$ .

In the following, we denote starting and closing tags of an element with name “**element**” by  $element_l$  and  $element_r$ , respectively. We define symbols  $element_l = \langle \mathbf{element} \rangle$  and  $element_r = \langle / \mathbf{element} \rangle$ .

#### Lemma 4 Context free grammars and DTDs.

There is a context-free grammar for every pair of DTD  $D$  and starting element  $Z$  defining the same language  $L_{D,Z}$ .

**Proof:** We give an algorithm to construct a such a context free grammar  $G=(T, N, P, S)$  with  $L(G)=L_{D,Z}$ .

Let  $T$  be a set of terminal symbols, including *string*, *quotes*, ‘=’, ‘>’, and ‘/>’. Additionally,  $T$  contains pairs of symbols  $element_l$  and  $element_r$ , for every element name “**element**” in the DTD, and a symbol for every attribute name “**attribute**” in the DTD. They are assumed to originate with a standard longest-match regular scanner.

Let  $C = \{C_x, C_y, \dots\}$  be the set of grammar productions for the content models of elements named “**x**”, “**y**”, ... resulting from the construction from Lemma 1. Let  $A = \{A_x, A_y, \dots\}$  be the corresponding set of grammar productions for attributes of elements named “**x**”, “**y**”, ... from Lemma 2. W.l.o.g., we assume all non-terminals of productions in  $C$  and  $A$  to be disjoint; the set  $N_C$  and  $N_A$  denote the union of the element and the attribute grammars, respectively. Let  $N_E = \{N_x, N_y, \dots\}$  and  $B_E = \{B_x, B_y, \dots\}$  be sets of non-terminal symbols, pairwise disjoint and disjoint from those in  $N_C$  and  $N_A$  and corresponding to elements named “**x**”, “**y**”, ... We then define  $N = N_E \cup B_E \cup N_C \cup N_A$ .

Except for the symbol *string*, we replace all terminals  $x$  in productions in  $C$  by non-terminals  $N_x$  and denote the set of transformed productions by  $K = \{K_x, K_y, \dots\}$  (Remember: the terminals from productions in  $C$  are element names.) We introduce additional productions  $E_x$  for each element name  $x$ :

- (1)  $N_x \rightarrow x_l A_x B_x$
- (2)  $B_x \rightarrow \text{'>'} K_x x_r$

For all content model productions  $C_x$  with  $\epsilon \in L(C_x)$ , we additionally include this production:

- (3)  $B_x \rightarrow \text{'/>'}$

Let  $E = \{E_x, E_y, \dots\}$  be the set of grammar productions as defined above for elements named “ $\mathbf{x}$ ”, “ $\mathbf{y}$ ”, ... We set  $P = K \cup A \cup E$ . Finally, we set  $S = N_Z$  for the starting element  $Z$ , which completes the grammar construction.

By construction,  $L(G) = L_{D,Z}$ . As the left sides of all productions  $p \in P$  consist of but one non-terminal,  $G$  is context free.

In the appendix. the above transformations are demonstrated on an example. Now, we are ready to prove our main results as additional properties of the constructed grammar:

**Theorem 6 SLL(1) grammars and DTDs.**

*There is an SLL(1) grammar for every pair of DTD  $D$  and starting element  $Z$  defining the same language  $L_{D,Z}$ .*

**Proof:** We prove that the theorem holds for the grammars  $G$  generated by the algorithm in Lemma 4. We show the SLL(1)-property for all pairs of productions  $N \rightarrow v$  and  $N \rightarrow w$ ,  $v \neq w$ . By construction, the non-terminals from the productions in  $C$  (generating the content model) and in  $A$  (generating the attribute definitions) are pair wise disjoint, disjoint from each other from all other non-terminals. Hence, Lemma 1 and Lemma 2 continue to hold in the context of the grammar  $G$ . There is only one more case: productions  $B_x \rightarrow '>' K_x$  and  $x_1 B_x \rightarrow '/>'$  have the same left-hand sides. However, the right-hand sides start with different terminals ('>' vs. '/>'), so the SLL(1) property holds for all productions.

Theorem 6 shows how to generate deterministic SLL(1) grammars for XML documents conforming to a specific DTD. The same problem is not computable for context free languages in general. The class of recursive descent parsers corresponding to SLL(1) grammars is one of the fastest known for context-free languages. As the requirement for deterministic element definitions in the XML 1.0 Specification is non-normative, this result is practically relevant – highly performant parsers may be generated even from ambiguous DTDs.

Top down parsing is quite intuitive and allows direct implementations as well as generation from an LL grammar using tools like `javacc` or `ell`. However, `yacc`, `bison`, `lalr` and many other parser generators accept the more powerful LR grammars, actually LALR(1) grammars. We therefore establish

**Theorem 7 LALR(1) grammars and DTDs.**

*There is an LALR(1) grammar for every pair of DTD  $D$  and starting element  $Z$  defining the same language  $L_{D,Z}$ .*

**Proof:** We prove that the grammar generated by the construction algorithm given in the proof of Lemma 4 is LALR(1). Actually, we prove the stronger SLR(1) property for these grammars. Let  $F_x$  be the deterministic acceptor for the content model of element  $x$  as generated in the proof of Lemma 1.

Central to parsing is the following procedure analyzing one complete element. This procedure is initially called for the top element  $Z$  and, recursively, whenever an element start tag appears in the input:

- (1) Shift symbol  $x_1$ . Set the current state to the initial state of  $F_x$ .
- (2) Analyze attributes of element  $x$  using productions  $A_x$  according to the SLR(1) algorithm given in Lemma 3. Instead of shifting the whole input string, shift only until the next input symbol is either '>' or '/>'. Reduce the attributes to  $A_x$ . Return with error if reduction is not possible or neither '>' nor '/>' is detected in the input.
- (3) Shift the next symbol.

- a. If it is ' $/>$ ', reduce it to  $B_x$  if such a production exists (otherwise return with error). The top of the stack is now  $x_1 A_x B_x$ . Reduce this to  $N_x$ . Return  $N_x$ .
  - b. If it is ' $>$ ', proceed with (4)
- (4) Test the next symbol.
- a. If it is an opening tag, recursively call this procedure. Afterwards, proceed with (5).
  - b. If it is *string*, shift it and proceed with (5).
  - c. If it is  $x_j$ , proceed with (6).
  - d. Otherwise, return with error.
- (5) Perform a transition in the deterministic acceptor  $F_x$  from the current state. According to step (4), the top of the analysis stack must be one of:
- a. if it is  $N_y$  perform the transition for  $y$  (step (4 a) parsed element  $y$  before).
  - b. if it is *string*, perform the transition for *string* (step (4 b) parsed a *string* before).
- Proceed with (4), or return with error if no transition is applicable.
- (6) If the current state of  $F_x$  is
- a. an accepting state, the top of the stack is a sequence of non-terminal symbols  $N_a N_b N_c \dots$ , where  $a b c \dots$  is a word accepted by  $F_x$ . In other words, it is a word of the content model language  $L(C_x)$  of element  $x$ . In analogy to step (3) in the algorithm from Lemma 3, reduce the analysis stack top to  $K_x$ .
  - b. Otherwise, an error is detected.
- (7) Finally, shift  $x_j$  and reduce  $K_x x_j$  to  $B_x$ . The top of the stack is now  $x_1 A_x B_x$ . Reduce this to  $N_x$  and return  $N_x$ .

The procedure is an *SLR(1)* parsing schema as only the next input symbol decides whether a shift or a reduce is executed, cf. in steps (2) and (4). Reduce-reduce conflicts do not occur.

Remark: The above algorithm schema differs from the algorithms of LR parser generated by standard tools. It is not chosen for its efficiency. Instead, we designed the schema to simplify the proof of the *SLR(1)* property.

## 4 Extensions to XML 1.0

A number of extensions to XML 1.0 have been proposed by the W3C. One of them, XML Schema [19], is a new language to specify document types. Another one, Namespaces in XML [18], introduces a mechanism to prevent name collisions independent of a specification language. Both require changes to the grammars developed in the previous sections.

### 4.1 XML Schema

In DTDs, attribute names are bound to types in an element context only, but element names are bound to a single type for the entire document. XML Schema weakens the latter link by introducing the notion of types, which are defined separately from elements.

There are two varieties of types, simple and complex ones. Simple types pertain to attribute values and text fragments. They supersede the attribute types and `#PCDATA` sections of DTDs. Still, attributes remain delimited by quotes and text sections by elements, so simple types do not threaten their regularity. As they do not pertain to parsing, we ignore them in the remainder of this paper.

Complex types apply to elements. They consist of a list of applicable attributes and a regular content model. Due to the latter, the ambiguities found in DTDs also apply to schemas. Element names remain unique within a complex type context, but the corresponding types are determined by pairs of name and context, as demonstrated in Example 4:

**Example 4:** An XML Schema fragment

---

```

<element name="x" type="A"/>

<complexType name="A">
  <sequence minOccurs="0" maxOccurs="unbounded">
    <element name="x" type="B"/>
  </sequence>
</complexType>

<complexType name="B">
  <choice>
    <element name="x" type="A"/>
    <element name="y" type="B"/>
  </choice>
</complexType>

```

---

The above definitions show that the same element name `x` is bound to different types depending on its context: a toplevel element `x` is of type `A`, i.e. it consists of a sequence of `x` elements of arbitrary length. Those child elements conform to a type `B` different from `A`. They in turn contain either an element `x` of type `A`, or an element `y` of type `B`.

XML Schema also provides for element and attribute wildcards, type overrides from instance documents and namespaces. They will be discussed separately below. For now, we define the subset of restricted XML Schema to be the set of all schemas not employing or admitting wildcards, overrides and namespaces. Using this notion, we generalize our main result.

**Theorem 8**     **Restricted XML Schema and  $SLL(1)$  and  $LALR(1)$  grammars.**

*There is an algorithm to find an  $SLL(1)$  grammar  $G$  defining the same language for every pair of restricted XML Schema and starting element. The same grammar  $G$  is also  $LALR(1)$ .*

**Proof:** We modify the algorithm in Lemma 4 as follows:

First, we generate grammar fragments for complex types. Since their content models are regular expressions differing from those in DTDs only in notation, we can use the construction in Lemma 1. In the resulting grammar for the content of complex type  $t$ , we substitute all terminals  $s$  representing element names with new, unique non-terminals  $N_{(t', s)}$ , where  $t'$  is the type of  $s$  in this context. The grammars for complex types' attributes are generated as in Lemma 2.

For every pair of complex type  $t$  and element name  $s$  in the schema, we then generate productions per Lemma 4, using  $N_{(t, s)}$  as the toplevel non-terminal and the attribute and content model grammars of  $t$ .

The proof of Theorem 6 applies by analogy, as element names remain unique in a complex type.

**Example 5:** Applying Theorem 8 to the fragment in Example 4:

$$P = \left\{ \begin{array}{lll}
 N_{(A,x)} \rightarrow x[A_{(A,x)} B_{(A,x)}] & N_{(B,x)} \rightarrow x[A_{(B,x)} B_{(B,x)}] & N_{(B,y)} \rightarrow y[A_{(B,y)} B_{(B,y)}] \\
 A_{(A,x)} \rightarrow \varepsilon & A_{(B,x)} \rightarrow \varepsilon & A_{(B,y)} \rightarrow \varepsilon \\
 B_{(A,x)} \rightarrow '/>' & & \\
 B_{(A,x)} \rightarrow '>' C_A x_1 & B_{(A,x)} \rightarrow '>' C_B x_1 & B_{(B,y)} \rightarrow '>' C_B y_1 \\
 C_A \rightarrow N_{(B,x)} C_A & C_B \rightarrow N_{(A,x)} C'_B & \\
 C_A \rightarrow \varepsilon & C_B \rightarrow N_{(B,y)} C'_B & \\
 & C'_B \rightarrow \varepsilon & \\
 \end{array} \right\}$$

Each of the above columns contains the productions of a pair  $(t, s)$ , where  $t$  is the type assigned to element name  $s$  within a complex type: the first defines elements  $x$  of type  $A$ , the second elements  $x$  of type  $B$  and the last elements  $y$  of type  $B$ .

## 4.2 Namespaces

Namespaces were introduced in XML to prevent name collisions and increase the interoperability of document types from different sources. They extend the XML 1.0 naming mechanism, which provides for element and attribute names, to a pair-based mechanism: in XML Namespaces, a qualified name consists of a namespace and a local name.

Namespaces are identified by globally unique strings. To ensure international interoperability, URIs usually serve as namespace identifiers in practice. This is only a convention - processors do not treat URIs any different from arbitrary character sequences. Usage of an URI as a namespace does not imply the corresponding schema is in fact available under that address.

Because globally unique strings are unwieldy, XML Namespaces mandates the use of an abbreviation mechanism to conserve space. Under that convention, an instance document defines arbitrary identifiers called prefixes to represent a namespace. Prefixes precede local names, using a colon for separator. Their definitions take the form of virtual attributes with the reserved prefix `xmlns`.

The extent of prefix definitions resembles the extent of variable definitions in block-structured languages: a prefix represents the closest like-named definition in an ancestor element or the element itself. Thus, a prefix can be used before its definition, but forward references are limited to the scope of an opening tag. These rules are illustrated by Example 6:

**Example 6:** Namespaces and prefixes

---

```

1 <a xmlns:p="http://www.noga.de/namespaces/XYZ">
2   <p:b/>
3   <p:b xmlns:p="http://www.noga.de/namespaces/ABC" />
4   <p:b/>
5 </a>

```

---

In this example, the prefix `p` represents two different namespaces. The definition in line 1 is in scope for lines 1, 2, 4 and 5 and used in lines 2 and 4. The definition in line 3 is in scope in line 3 and used

there in the element prefix - an instance of a forward reference. Substituting a different identifier for `p`, e.g., `abracadabra`, leaves the content invariant.

The concept of scope exceeds the modeling power of context-free grammars. In block-structured languages, this problem is commonly resolved by describing a superset of the language with a context-free grammar. Any additional constraints are realized in a separate semantic analysis phase operating on the parse tree.

As XML Namespaces limit forward references, prefix resolution only requires a single pass of a stack automaton. Using a suitable representation for qualified names, the grammars derived from DTDs or XML Schemas may be retained for namespace-aware processing. Single-pass treatment remains feasible if superset parsing, prefix resolution and qualified name matching are suitably interleaved.

## 5 Performance

Due to just-in-time compilers (JITs) and garbage collection overheads, the Java platform cannot be fairly compared to native code. We compare the faster implementations in the C/C++ language. Apache's `Xerces-C` and James Clark's `expat` enjoy widespread deployment and are generally considered highly performant.

To establish firm performance figures, we generated a series of test cases from a suitably complex real-world document. Using our `aXMLerat` toolkit, we generated a validating parser from the corresponding DTD. The generated parser builds a tree representation for the entire document in memory. `Xerces-C` was tested in DOM and SAX mode – while both are validating, only the former actually builds a tree representation. `expat` neither validates nor builds a tree representation.

Measurements were taken on an Athlon 850 MHz with 256 MBytes RAM. Because the distributors provide precompiled binaries for the Windows platform, `Xerces-C` 1.3.0 and `expat` 1.2 were tested with `DOMCount` and `SAXCount` resp. `expat -t` under NT 4.0 SP 5. The parser generated by `aXMLerat` was executed under SuSE Linux 7.1 on the identical machine.

Despite validating the XML documents for DTD conformance and building a full tree representation in memory, our generated parser is up to 40% faster than `expat`, which does neither. `Xerces-C`, considered a high-performance implementation by its authors, is not remotely comparable. In DOM mode, it even aborts processing the largest input file due to memory restrictions.

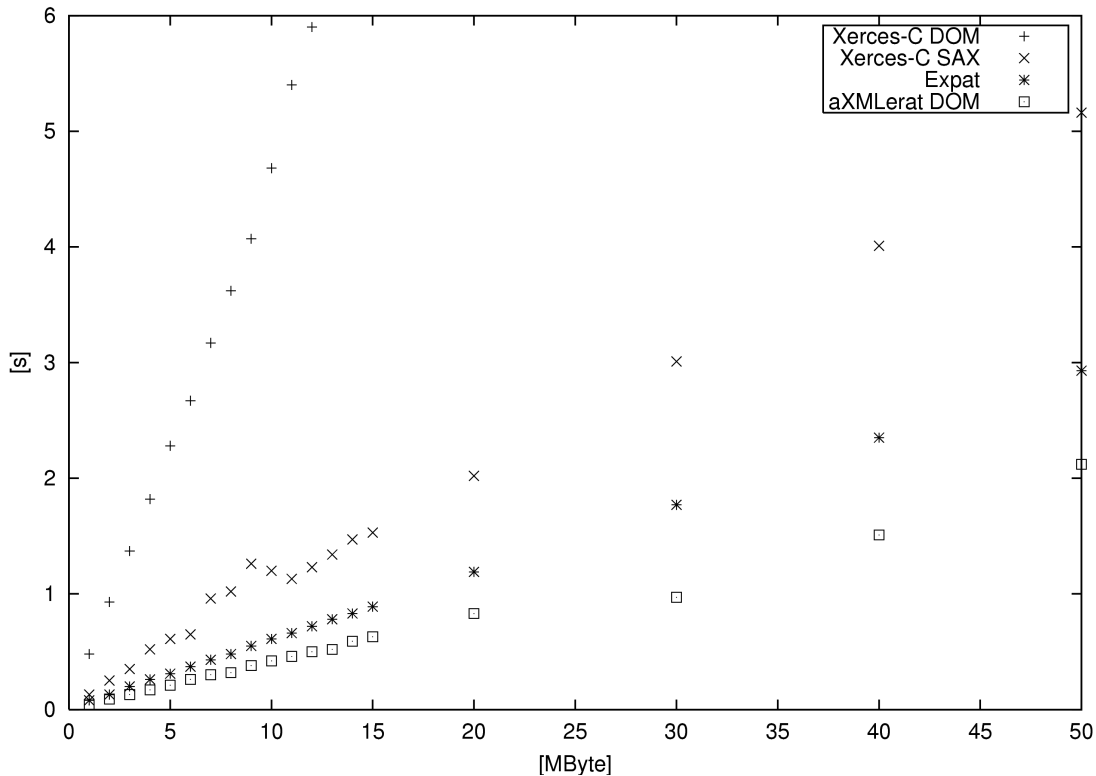


Figure 1: Parsing times over a range of XML file sizes.

## 6 Related Work

The SGML [11] designers solved the ambiguity problem by design restriction: SGML requires content models to be deterministic regular expressions. If this restriction continued to hold for XML, deterministic parsing would be no problem.

Such a restriction would limit the class of languages that DTDs can express. As shown in [5], there are ambiguous regular expressions for which no equivalent deterministic regular expressions exist. As a consequence, in general, deterministic grammars constructed from DTDs and XML Schemas cannot be rewritten into deterministic DTDs or XML Schemas.

A theoretical view on XML languages is given in [4]. The authors ignore attributes, namespaces etc. In their paper, they take a language theoretic point of view. They show that some properties which are not decidable for general context-free languages become decidable for XML grammars. The result most closely related to our work is that DTDs (called XML-grammars) have a unique normal form. However, this normal form can be ambiguous.

There are working parser generators for XML documents. Our own generators in the aXMLerat [3] toolkit exploit the transformations described in the present paper. They generate `lalr` and `javacc` specifications from DTDs and directly generate an LL parser for XML Schema definitions. The XMLBooster [21] uses its own specification language. Thereby, they avoid the problems of ambigu-

ous DTDs and schemas. Using the theoretical results above, it is possible to transform their proprietary specification to DTDs. However, it remains unclear if transformations in the opposite directions are always possible.

## 7 Conclusion

This paper enables the generation of deterministic parsers for XML documents from arbitrary, even ambiguous DTDs.

We investigated the context-free DTD languages. As our examples show, there exist ambiguous DTDs, some of which have no deterministic equivalents. Despite this, all DTDs can be transformed into deterministic grammars with the algorithms in this paper. The generated grammars are both *LL(1)* and *LALR(1)*. Respective transformations are non-computable for context-free grammars in general.

The generated grammars are suitable for common parser generators. With those tools, we automatically generated validating parsers that build full tree representations of documents in memory. Generation eliminates the need to analyze the DTD at runtime. Our specific parsers are manifestly faster than even no-op generic parsers, and beat comparable, fully functioned ones by an order of magnitude.

We developed a generalization to a subset of the XML Schema languages. Possible solutions for the problems arising from namespaces were outlined. As the XML Schema standard is still evolving, this is clearly ongoing work. Experiments on static high-speed XML parsing also continue.

## 8 References

1. A. Aho and J. Ullman: *The Theory of Parsing, Translation, and Compiling*. Prentice-Hall, 1972.
2. Apache Xerces XML Parser for C.  
<http://xml.apache.org/xerces-c>.
3. The aXMLerat tool-box.  
<http://i44pc29.info.uni-karlsruhe.de/B2BWeb>.
4. J. Berstel and L. Boasson: *XML-Grammars*. In proc: Mathematical Foundations of Computer Science 2000, 25<sup>th</sup> International Symposium, LNCS 1893, Springer 2000.
5. A. Brüggemann-Klein: *Regular Expressions into Finite Automata*. Extended abstract in H. Simon (Ed.) LATIN 92, pp 97-98. Springer, 1992. Full version in Theoretical Computer Science 120:197-213, 1993.
6. James Clark: *expat - XML Parser Toolkit*, 2000, <http://www.jclark.com/xml/expat.html>.
7. F. L. DeRemer: *Simple LR(k) Grammars*. Communication of the ACM 14(7), pp. 453-460, 1971.
8. GNU Project. Donnelly and Stallmann: *Bison - Manual Page*, Public Domain Software, 1988, <http://www.gnu.org/manual/bison>.
9. J. Grosch: *Generators for High-Speed Front-Ends*, LNCS 371, pp. 81-92, Springer, Oct. 1988.



10. IBM: *XML Parser for Java*.  
<http://www.alphaworks.ibm.com/tech/xml4j>.
11. Information Processing - Text and Office Systems - Standard Generalized Markup Language (SGML), ISO 8879:1986.
12. S.C. Johnson: *Yacc - Yet Another Compiler-Compiler*, Computer Science Technical Report 32, Bell Telephone Laboratories, Murray Hill, NJ, July 1975.
13. Microsoft: *Common Object Model (COM)*.  
<http://www.microsoft.com/com>.
14. OMG: *Common Object Request Broker Architecture (CORBA)*.  
<http://www.omg.org/cgi-bin/doc?formal/01-02-01>.
15. D. J. Rosenkrantz and R. E. Stearns: *Properties of deterministic Top-Down-Grammars*. Information and Control 17, pp. 226-256, 1970.
16. B. Vielsack: *The parser generators lalr and ell*. Technical report, GMD, Gesellschaft für Mathematik und Datenverarbeitung, Forschungsstelle Karlsruhe, 1988.
17. W3C: *Extensible Markup Language (XML) 1.0 (Second Edition)*.  
<http://www.w3.org/TR/2000/REC-xml-20001006>.
18. W3C: Namespaces in XML.  
<http://www.w3.org/TR/1999/REC-xml-names-19990114>.
19. W3C: *XML Schema Part 1: Structures*.  
<http://www.w3.org/TR/2001/PR-xmlschema-1-20010330>.
20. W. Waite and G. Goos: *Compiler Construction*. Springer, 1984.
21. WebGain: *JavaCC - The Java Parser Generator*.  
[http://www.webgain.com/products/metamata/java\\_doc.html](http://www.webgain.com/products/metamata/java_doc.html).
22. The XML-Booster.  
<http://www.xmlbooster.com>.

## Appendix

### A Basic Definitions

#### Definition 5 Formal Language

An alphabet  $T$  is a finite, non-empty set of symbols. The set of finite strings formed by concatenating symbols from  $T$  is denoted by  $T^+$ .  $T^*$  denotes  $T^+$  augmented by the empty string  $\epsilon$ . Each subset of  $T^*$  is a formal language over the alphabet  $T$ . Its elements are called words.

Formal languages are usually defined by grammars:

#### Definition 6 Grammar and Ambiguous Grammar

A grammar is a quadruple  $G=(T, N, P, Z)$  with an alphabet  $T$ ,  $T$  and  $N$  disjoint,  $Z \in N$ , and  $(T \cup N, P)$  a general rewrite system.  $L(G)$  denotes the formal language defined by  $G$ .  $N$  is called the set of non-terminal symbols of the grammar.

A grammar  $G$  is ambiguous if a sentence of  $L(G)$  may be derived from  $Z$  using at least two different sequences of substitutions from  $P$ . Otherwise it is called deterministic.

Grammars may be classified according to the form of their productions:

#### Definition 7 Regular and Context Free Grammars and Languages

A grammar is regular if each production in  $P$  has the form  $B \rightarrow b$  or the form  $B \rightarrow bC$  with  $B, C \in N$ ,  $b \in T \cup \{\epsilon\}$ .

A grammar is context free if each production in  $P$  has the form  $B \rightarrow x$  with  $B \in N$ ,  $x \in (N \cup T)^*$ .

A language is regular (context free) iff it can be defined by a regular (context free) grammar.

Regular languages may also be defined by regular expressions:

#### Definition 8 Regular Expressions and Regular Languages

Let  $T$  be an alphabet. Then  $\{a \mid a \in T \cup \{\epsilon\}\}$  are regular expressions  $r$  defining the regular languages  $L(r)=\{a\}$ . Let  $r_1$  and  $r_2$  be regular expressions. Then

$r_1 \mid r_2$  defining the regular language  $L(r_1) \cup L(r_2)$ ,

$r_1, r_2$  defining the regular language  $L = \{ ab \mid a \in L(r_1) \wedge b \in L(r_2) \}$

$r_1^*$  defining the regular language  $L = \{ a \dots a \mid a \in L(r_1) \} \cup \{\epsilon\}$

are regular expressions.

Remark: Some definitions of regular expressions include:

$r_1^+$  defining the regular language  $L = \{ a \dots a \mid a \in L(r_1) \}$

$r_1?$  defining the regular language  $L(r_1) \cup \{\epsilon\}$

These expressions can be derived, as  $r_1^+ = r_1^* r_1$  and  $r_1? = r_1 \mid \epsilon$ .

## B An Example based on a DTD

We expand on Example 2:

---

```
<!DOCTYPE a [
  <!ELEMENT a ( x? , ( y* | z* ) )>
  <!ELEMENT x ( #PCDATA )>
  <!ELEMENT y ( #PCDATA )>
  <!ELEMENT z ( #PCDATA )>
]>
```

---

A corresponding ambiguous acceptor for the content model of **a** is given in Figure 2 (left). The resulting deterministic minimal acceptor is displayed in Figure 2 (right).

According to the construction in Lemma 1, this acceptor leads to the following *SLL*(1) and *SLR*(1) grammar for the content model of **a**:

- |                                      |                                 |
|--------------------------------------|---------------------------------|
| 1. $N_a(1) \rightarrow \varepsilon$  |                                 |
| 2. $N_a(1) \rightarrow x N_a(2)$     | $N_a(1) \rightarrow N_x N_a(2)$ |
| 3. $N_a(1) \rightarrow y N_a(3)$     | $N_a(1) \rightarrow N_y N_a(3)$ |
| 4. $N_a(1) \rightarrow z N_a(4)$     | $N_a(1) \rightarrow N_z N_a(4)$ |
| 5. $N_a(2) \rightarrow \varepsilon$  |                                 |
| 6. $N_a(2) \rightarrow y N_a(3)$     | $N_a(2) \rightarrow N_y N_a(3)$ |
| 7. $N_a(2) \rightarrow z N_a(4)$     | $N_a(2) \rightarrow N_z N_a(4)$ |
| 8. $N_a(3) \rightarrow \varepsilon$  |                                 |
| 9. $N_a(3) \rightarrow y N_a(3)$     | $N_a(3) \rightarrow N_y N_a(3)$ |
| 10. $N_a(4) \rightarrow \varepsilon$ |                                 |
| 11. $N_a(4) \rightarrow z N_a(4)$    | $N_a(4) \rightarrow N_z N_a(4)$ |

Where applicable, we printed the productions transformed by the construction in Lemma 4 in the right column. The grammars for the content models of **x**, **y** and **z** are obtained accordingly. As they are identical up to renaming non-terminals, we display the one for **x**:

12.  $N_x(1) \rightarrow \varepsilon$
13.  $N_x(1) \rightarrow \textit{string}$

As the example does not define attributes, all grammars for element attributes accept  $\varepsilon$  only. The productions are identical up to renaming of non-terminals, so we display the ones for **a** (14) and **x** (15):

14.  $A_a \rightarrow \varepsilon$
15.  $A_x \rightarrow \varepsilon$

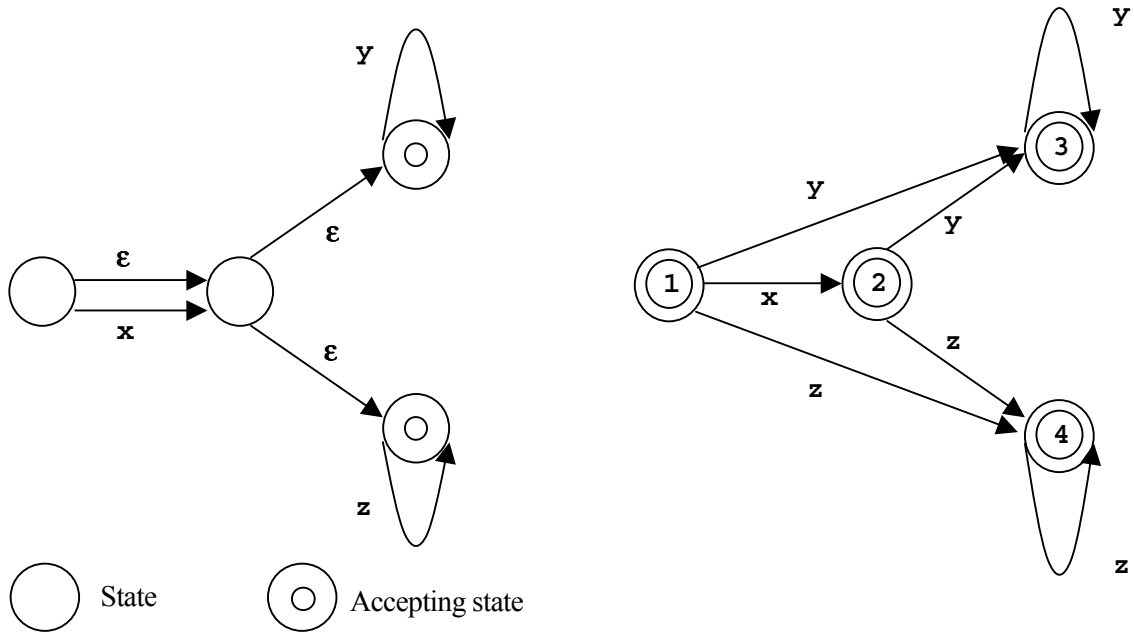


Figure 2: Ambiguous acceptor directly derived from the content model (left). Deterministic acceptor resulting from the construction of (right).

The final construction algorithm of Lemma 4 adds the productions which are identical up to renaming of non-terminals and the opening and closing tag strings. We display those for **a** (16-18) and **x** (19-21) :

- 16.  $N_a \rightarrow \langle \mathbf{a} \rangle A_a B_a$
- 17.  $B_a \rightarrow \rangle N_a(1) \langle / \mathbf{a} \rangle$
- 18.  $B_a \rightarrow \rangle$
- 19.  $N_x \rightarrow \langle \mathbf{x} \rangle A_x B_x$
- 20.  $B_x \rightarrow \rangle N_x(1) \langle / \mathbf{x} \rangle$
- 21.  $B_x \rightarrow \rangle$

We reconsider the example sentence that lead to two derivations for the original grammar.

---

$\langle \mathbf{a} \rangle \langle \mathbf{x} \rangle \dots \langle / \mathbf{x} \rangle \langle / \mathbf{a} \rangle$

---

Its derivation is now unique:

- $N_a \rightarrow \langle \mathbf{a} \rangle A_a B_a$  (16)
- $\rightarrow \langle \mathbf{a} \rangle \varepsilon B_a$  (14)
- $\rightarrow \langle \mathbf{a} \rangle \varepsilon \rangle N_a(1) \langle / \mathbf{a} \rangle$  (17)
- $\rightarrow \langle \mathbf{a} \rangle \varepsilon \rangle N_x N_a(2) \langle / \mathbf{a} \rangle$  ( 2)

$$\rightarrow \langle \mathbf{a} \rangle \varepsilon \langle \mathbf{x} \rangle A_x B_x N_a(2) \langle \mathbf{a} \rangle \quad (19)$$

$$\rightarrow \langle \mathbf{a} \rangle \varepsilon \langle \mathbf{x} \rangle \varepsilon B_x N_a(2) \langle \mathbf{a} \rangle \quad (15)$$

$$\rightarrow \langle \mathbf{a} \rangle \varepsilon \langle \mathbf{x} \rangle \varepsilon \langle \mathbf{x} \rangle N_x(1) \langle \mathbf{x} \rangle N_a(2) \langle \mathbf{a} \rangle \quad (20)$$

$$\rightarrow \langle \mathbf{a} \rangle \varepsilon \langle \mathbf{x} \rangle \varepsilon \langle \mathbf{x} \rangle \text{string} \langle \mathbf{x} \rangle N_a(2) \langle \mathbf{a} \rangle \quad (13)$$

$$\rightarrow \langle \mathbf{a} \rangle \varepsilon \langle \mathbf{x} \rangle \varepsilon \langle \mathbf{x} \rangle \text{string} \langle \mathbf{x} \rangle \varepsilon \langle \mathbf{a} \rangle \quad (5)$$

$$\equiv \langle \mathbf{a} \rangle \langle \mathbf{x} \rangle \dots \langle \mathbf{x} \rangle \langle \mathbf{a} \rangle$$