

Turning DTDs into specialized tree-automata-based schemata to match a collection of marked-up documents

Rafael C. Carrasco
Dept. de Lenguajes y
Sistemas Informáticos
University of Alicante
E-03071 Alicante, Spain
carrasco@dlsi.ua.es

Alejandro Bia
Miguel de Cervantes Digital
Library / Dept.L.S.I.
University of Alicante
E-03071 Alicante, Spain
alex.bia@ua.es

Mikel L. Forcada
Dept.L.S.I. / interNOSTRUM
University of Alicante
E-03071 Alicante, Spain
mlf@internostrum.com

ABSTRACT

Regular tree automata (RTA) or, equivalently, forest regular grammars (FRG) have been recently proposed for use as XML (extended markup language) schemata. They are more powerful than usual XML DTDs (document-type definitions), make the implementation, optimization and pruning of XML queries easier and allow for the implementation of context-sensitive content models. We describe a method for the automatic generation of a specialized RTA-based schema from a source DTD and a sample of marked-up documents showing context-sensitive behaviour in content models. It creates the smallest RTA-based schema with which all the XML documents in the sample comply and which does not accept any documents not valid according to the original DTD. In this way, new files can be created, parsed, and queried using the specialized schema but still be compliant with the original DTD. The tool is currently being tested at the Miguel de Cervantes digital library at the University of Alacant (<http://cervantesvirtual.com>).

Keywords

automatic learning, feature extraction, grammatical inference, document analysis, document markup, digital libraries, tree automata, document type definitions

1. INTRODUCTION

An extended markup language (XML) document-type definition (DTD) specifies the elements that are allowed in a document of this type. Document types are defined by extended context-free grammars in which the right-hand sides of productions are unambiguous regular expressions [8]. To address some of the problems posed by XML DTDs, regular tree automata (RTA), or, equivalently, forest regular grammars (FRG), have been recently proposed for use as XML schemata [13, 19]; RTA-based schemata are more powerful

than usual XML DTD. For example, XML queries [13] are more easily implemented, optimized and pruned in terms of RTA-based schemata. In particular, RTA-based schemata allow the implementation of context-sensitive content models [19], which allow an element to have different behaviors depending on its placement in the document structure, i.e. depending on its ancestors in the document hierarchy. For instance, it is common sense to expect paragraphs to be allowed to have footnotes in general except when they are included within certain particular structures as could be an abstract or a metadata section, such as the `teiHeader` of the TEI (Text Encoding Initiative) encoding scheme [21, 18, 10], without having to label each kind of element with a different name in the text itself. Regular tree automata and languages also provide a powerful formalism to address document and schema transformations [16].

In the ensuing discussion, we will use the following DTD (inspired in [19]) as an example:

```
<!ELEMENT DOC (HEAD,BODY)>
<!ELEMENT HEAD (TITLE,ABSTRACT)>
<!ELEMENT TITLE (#PCDATA)>
<!ELEMENT ABSTRACT (P)+>
<!ELEMENT BODY (P)+>
<!ELEMENT P (#PCDATA,(FOOTNOTE,#PCDATA)*)>
<!ELEMENT FOOTNOTE (#PCDATA)>
```

1.1 Regular tree automata

A *regular tree automaton* [13] is a 5-tuple $M = (\Sigma, D, Q, \delta, F)$ where Σ is a finite set of element types, D is a finite set of data types, Q is a finite set of states, δ is a function $\delta : \Sigma \times E \rightarrow Q$ where E is either an element of D or a regular expression over $Q \cup D$ (this is a slight variation from the definition in [13]), and $F \subset Q$ is the set of final states. The automaton walks the document tree bottom-up and computes, using function $\delta()$ a state at each node from the label at the node itself and the states of its children and accepts a document tree only if the state eventually computed for the root node is in F . The DTD above may be turned into a regular tree automaton in which $\Sigma = \{\text{DOC, HEAD, TITLE, ABSTRACT, BODY, P, FOOTNOTE}\}$, $D = \{\#\text{PCDATA}\}$, $Q = \{q_{\text{DOC}}, q_{\text{HEAD}}, q_{\text{TITLE}}, q_{\text{ABSTRACT}}, q_{\text{BODY}}, q_{\text{P}}, q_{\text{FOOTNOTE}}\}$, $F = \{q_{\text{DOC}}\}$,

and

$$\begin{aligned}
\delta(\text{DOC}, q_{\text{HEAD}} q_{\text{BODY}}) &= q_{\text{DOC}} \\
\delta(\text{HEAD}, q_{\text{TITLE}} q_{\text{ABSTRACT}}) &= q_{\text{HEAD}} \\
\delta(\text{TITLE}, \#PCDATA) &= q_{\text{TITLE}} \\
\delta(\text{ABSTRACT}, q_p^+) &= q_{\text{ABSTRACT}} \\
\delta(\text{BODY}, q_p^+) &= q_{\text{BODY}} \\
\delta(\text{P}, \#PCDATA(q_{\text{FOOTNOTE}} \#PCDATA)^*) &= q_p \\
\delta(\text{FOOTNOTE}, \#PCDATA) &= q_{\text{FOOTNOTE}}
\end{aligned} \tag{1}$$

where it can be seen that there is a one-to-one correspondence between states in Q and elements in Σ ; this is not the general case, but is instead due to the fact that the RTA has been derived from a DTD.

1.2 Forest regular grammars

A forest regular grammar (FRG, [19]) $G = (Q, \Sigma, D, P, S)$ has a finite set of variables Q , a finite set of labels Σ , a finite set of data types D , a set of productions P of the form $V \rightarrow \langle \Sigma \rangle \alpha \langle / \Sigma \rangle$, where α is any regular expression over $Q \cup D$, and a start symbol $S \in Q$. The RTA in eq. (1) may be written as a FRG having variables $Q = \{q_{\text{DOC}}, q_{\text{HEAD}}, q_{\text{TITLE}}, q_{\text{ABSTRACT}}, q_{\text{BODY}}, q_p, q_{\text{FOOTNOTE}}\}$, document labels $\Sigma = \{\text{DOC}, \text{HEAD}, \text{TITLE}, \text{ABSTRACT}, \text{BODY}, \text{P}, \text{FOOTNOTE}\}$, data types $D = \{\#PCDATA\}$, $S = \{\text{DOC}\}$ and the following set of productions P : follows:

$$\begin{aligned}
q_{\text{DOC}} &\rightarrow \langle \text{DOC} \rangle q_{\text{HEAD}} q_{\text{BODY}} \langle / \text{DOC} \rangle \\
q_{\text{HEAD}} &\rightarrow \langle \text{HEAD} \rangle q_{\text{TITLE}} q_{\text{ABSTRACT}} \langle / \text{HEAD} \rangle \\
q_{\text{TITLE}} &\rightarrow \langle \text{TITLE} \rangle \#PCDATA \langle / \text{TITLE} \rangle \\
q_{\text{ABSTRACT}} &\rightarrow \langle \text{ABSTRACT} \rangle q_p^+ \langle / \text{ABSTRACT} \rangle \\
q_{\text{BODY}} &\rightarrow \langle \text{BODY} \rangle q_p^+ \langle / \text{BODY} \rangle \\
q_p &\rightarrow \langle \text{P} \rangle \#PCDATA (q_{\text{FOOTNOTE}} \#PCDATA)^* \langle / \text{P} \rangle \\
q_{\text{FOOTNOTE}} &\rightarrow \langle \text{FOOTNOTE} \rangle \#PCDATA \langle / \text{FOOTNOTE} \rangle
\end{aligned} \tag{2}$$

As in the case of the RTA in eq. (1), the fact that the FRG is derived from a DTD results in a set of states Q which is essentially the same as the set of labels Σ .

1.3 Beyond DTDs

But regular-tree automata and forest regular grammars are more powerful than what these examples show because the elements in Q need not be the same as the labels in Σ . In particular, one could choose elements of Q to be partitions or specializations of the labels in Σ . Let us consider the example considered before, in which paragraphs may have footnotes in general but do not have footnotes when they are inside an abstract. It is not possible to write an XML DTD for this case, but it is quite straightforward to write a RTA or a FRG for it; only the FRG will be shown:

$$\begin{aligned}
q_{\text{DOC}} &\rightarrow \langle \text{DOC} \rangle q_{\text{HEAD}} q_{\text{BODY}} \langle / \text{DOC} \rangle \\
q_{\text{HEAD}} &\rightarrow \langle \text{HEAD} \rangle q_{\text{TITLE}} q_{\text{ABSTRACT}} \langle / \text{HEAD} \rangle \\
q_{\text{TITLE}} &\rightarrow \langle \text{TITLE} \rangle \#PCDATA \langle / \text{TITLE} \rangle \\
q_{\text{ABSTRACT}} &\rightarrow \langle \text{ABSTRACT} \rangle q_{p_NOFN}^+ \langle / \text{ABSTRACT} \rangle \\
q_{\text{BODY}} &\rightarrow \langle \text{BODY} \rangle q_p^+ \langle / \text{BODY} \rangle \\
q_p &\rightarrow \langle \text{P} \rangle \#PCDATA (q_{\text{FOOTNOTE}} \#PCDATA)^* \langle / \text{P} \rangle \\
q_{p_NOFN} &\rightarrow \langle \text{P} \rangle \#PCDATA \langle / \text{P} \rangle \\
q_{\text{FOOTNOTE}} &\rightarrow \langle \text{FOOTNOTE} \rangle \#PCDATA \langle / \text{FOOTNOTE} \rangle
\end{aligned} \tag{3}$$

In this FRG, state q_{p_NOFN} is the special case (no footnotes inside an abstract), whereas q_p is the general case.

1.4 Previous work

Previous work has addressed the task of identifying a DTD from examples of XML documents.¹ A common difficulty in this approach is the need to find a correct degree of generalization. Some practical tools as FRED [20] let the users customize their preferred degree of generalization. Ahonen [2, 4] builds a (k, h) -testable model for the element contents and needs non-trivial further generalization in order to disambiguate the model [3]. Young-Lai and Tompa [23] rely on a stochastic approach to control over-generalization, based in turn on the algorithm by Carrasco and Oncina [12]. Presumably, the stochastic approach needs large collections of hand-tagged documents. Pizza-Chef [9] is a tool to generate TEI-compliant DTDs customized to satisfy specific markup needs by means of an HTML form with optional pizza-restaurant-like choices; although it is an excellent tool, it does not automatically learn a DTD from examples, and is not a general purpose tool, being limited to produce DTDs derived only from the TEI encoding scheme. More recently, Bia, Carrasco and Forcada [6] have proposed a method to automatically select from a source DTD only those features that are used by a set of valid sample documents and eliminate the rest of them, obtaining a narrow-scope DTD which defines a subset of the original markup scheme. This pruned DTD can be used to build new documents of the same markup subclass. The simplified DTD can be updated immediately in the event that new features are added to (or even eliminated from) the sample set of XML files. This process can be repeated as often as needed to generate an updated DTD. In particular, Bia et al.'s [6] technique may be used to build a one-document DTD, i.e. the minimal markup schema derived from the general DTD with which a given XML document complies.

1.5 Objective

This paper describes a method for the automatic generation of a specialized RTA-based schema from a source DTD and a sample of marked-up documents showing context-sensitive behaviour in content-models. The purpose is to create the smallest RTA-based schema with which the whole sample of XML documents complies and which in turn does not accept any documents not valid according to the original DTD. In this way, new files can be created, parsed, and processed for queries [13] using the specialized schema but still being compliant with the original, DTD; for example, some particular context-sensitivities observed a given set of digital library texts, encoded following a general XML DTD such as `teixlite.dtd` may be better described and processed with the new RTA-based schema. The tool is currently being tested at the Miguel de Cervantes digital library². This work is part of a larger project in the field of text markup and derived applications [7, 5, 6].

The method proposed in this paper aims at obtaining a close description of a set of XML documents, but at the price of having to use a more powerful formulation, that of RTA-based schemata; it goes one step further from [5] by allowing elements to have content models that depend on the particular context where the element appears.

A further application of this technique is to generate statistics that may help DTD designers improve their markup schemes. Information about the frequency of use of certain

¹DTD inference from XML queries has been addressed by [17] and [15].

²<http://cervantesvirtual.com/>

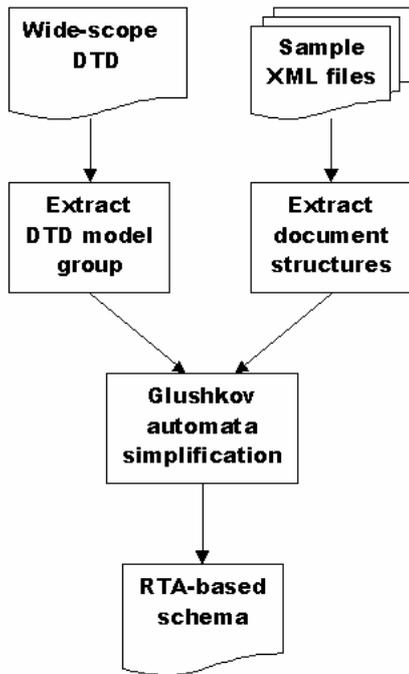


Figure 1: Architecture of the tool to build RTA-based schemata

elements within others helps us to detect unusual structures that could reflect markup mistakes or DTD features that allow for unwanted generalization.

The advantages of following a de-facto standard, or at least a commonly used markup scheme amongst digital libraries (as the TEI scheme is), like text interchangeability or inter-project cooperation —just to mention a few—, the goal of using XML from the very start of the project, as well as the need to spare the cost of developing our own DTD are some of the reasons why the `teixlite.dtd` has been chosen at the Miguel de Cervantes Digital Library. But even such a simple XML DTD, complete enough to structurally markup prose, verse and drama, is too complex for markup beginners, contains unnecessary elements, and allows the insertion of tags in contexts where they are not desirable. On the other hand, it contains only element definitions that are independent of the context, while context-dependencies are common in texts. In the future, as we plan to implement advanced (Boolean) queries, we need to go beyond the DTD formulation [13]. The method proposed in this paper allows us to convert the `teixlite.dtd` into a RTA-based schema which still allows only documents that are TEI-compliant and thus benefit from the advantages of sharing a common DTD with other international digitization projects.

We started by defining what kinds of modifications will be allowed in order to make markup simpler to use but keeping TEI compatibility (except for minor exceptions). In particular, we allowed for the following changes:

- To specify a set of normalized values for some attributes in order to enforce their use instead of free data entry.

- To add new attributes only in a few necessary cases (this is the only exception that may keep our files from being TEI compliant, but they can be easily removed anytime we want full TEI compatibility).
- To impose restrictions in element inclusion rules in order to eliminate the possibility of including certain elements at certain levels of the markup.
- To make some optional elements or attributes mandatory, following our specific markup norms.
- To eliminate optional elements we will not use to simplify the markup task and to avoid possible errors.

It may be easily admitted that simplifying a DTD and turning it into a specialized RTA-based schema by hand is tedious and error-prone. Using our method, the task may however be performed automatically by processing a set of sample documents representative of all kinds of documents we need to markup together for a given application.

2. THE PROCESS

A diagram of the process —a generalization of the one described in [5, 6]— is shown in figure 1. The key notion is that of *element-in-context* (EiC), that is, a specialized version of an element for each possible element enclosing it (for instance, a paragraph (P) inside an **ABSTRACT** may be described by the EiC **P-IN-ABSTRACT**, which is however still marked `<P>...</P>` in the text. As the diagram shows, the initial DTD is processed to extract the structure of the markup models and a Glushkov automaton [11] is built for each one (that is, for each regular expression).

The XML sample files are then preprocessed to extract the EiCs; that is, to label the contexts (content models) in which these elements are used, and to establish their nesting patterns.³ We keep track of the EiCs used in the sample files and mark the visited states of the Glushkov automaton corresponding to each content model for each possible EiC, to simplify the corresponding regular expressions. Finally, we eliminate unused EiCs, merge EiCs having identical content models into single EiCs having more general contexts (which will constitute the states of the RTA or the variables of the FRG), and update the names of EiCs in the right parts of element-in-context definitions to account for these mergings. It has to be stressed that the right parts of all EiC are all simplified (pruned) versions of the regular expressions defining the content of the corresponding context-free element in the original DTD, with simplifications being different in each context.

For the implementation of the the RTA-based-schemata building toolkit we need both an XML parser and a DTD parser. We assume that both the XML sample files and the source DTD are well-formed and already valid, so there is no need to build validating parsers. In particular, regular expressions are parsed against the extended Backus-Naur form (EBNF) grammar described in the following section, although indeed, XML forces stricter parentization patterns.

2.1 An example

Consider a set of XML documents and the following DTD defining them:

³The concept of EiC is closely related to that of *k*-testable tree languages as introduced by Knuutila [14] and Verdú et al. [22].

4. REGULAR EXPRESSION PRUNING

The process by means of which the regular expression describing the content of each EiC is simplified is based on a bottom-up parse of the original regular expression. A syntax-directed definition [1] describing this process is shown in Appendix B. The process replaces any unwitnessed position in the expression E_r by the regular expression corresponding to the empty set (\emptyset); then the expression is projected into the $\text{reg}(\Sigma)$ space; finally, the resulting regular expression is rearranged to avoid using symbols not in Σ .

The following simplification rules, used in the last step, preserve unambiguity as the resulting expression after each replacement defines exactly the same language:

$$\begin{aligned}
 \emptyset, E &= E, \emptyset = \emptyset+ = \emptyset \\
 \emptyset|E &= E|\emptyset = E \\
 \emptyset^* &= \emptyset? = \lambda \\
 \lambda, E &= E, \lambda = E \\
 \lambda|E &= E|\lambda = \begin{cases} E & \text{if empty}(E) \\ E? & \text{otherwise} \end{cases} \\
 \lambda^* &= \lambda+ = \lambda? = \lambda
 \end{aligned} \tag{9}$$

Here, λ is a special symbol denoting the empty string, not allowed in a valid regular expression, and $\text{empty}()$ is a Boolean function determining whether the regular expression accepts the empty string or not (the way to compute it efficiently can be found in the appendix and in [11]). The resulting expression is a simplified one containing neither \emptyset nor λ symbols, except when the result is exactly \emptyset , in which case the element will be removed or an empty content model will be substituted.

5. CONCLUSIONS AND FUTURE WORK

We have developed at the Miguel de Cervantes Digital Library a method which has been used to automatically generate RTA-based XML schemata from general DTD. On this first stage, we have addressed the simplification of element type descriptions based on sample files and the use of context-sensitive content models for elements, called elements in context. On a second stage, we plan to add automatic elimination or addition of attributes. We also plan to collect statistics to detect unusual patterns that may reflect markup mistakes.

6. REFERENCES

- [1] A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley, Reading, Massachusetts, USA, 1986.
- [2] H. Ahonen. Automatic Generation of SGML Content Models. *Electronic Publishing Origination, Dissemination, and Design*, 8(2/3):195–206, June/Sept. 1995.
- [3] H. Ahonen. Disambiguation of SGML content models. *Lecture Notes in Computer Science*, 1293:27, 1997.
- [4] H. Ahonen, H. Mannila, and E. Nikunen. Generating grammars for SGML tagged texts lacking DTD. *Mathematical and Computer Modelling*, 26(1):1–13, 1997.
- [5] A. Bia and R. C. Carrasco. Automatic DTD simplification by examples. In *ACH/ALLC 2001. The Association for Computers and the Humanities, The Association for Literary and Linguistic Computing, The 2001 Joint International Conference*, pages 7–9, New York University, New York City, June 2001.
- [6] A. Bia, R. C. Carrasco, and M. L. Forcada. Identifying a reduced DTD from marked up documents. In J. S. Sánchez and F. Pla, editors, *Pattern Recognition and Image Analysis (Proceedings of the IX Spanish Symposium on Pattern Recognition and Image Analysis – SNRFAI’2001 conference)*, volume II of *Treballs d’informàtica i tecnologia, num.7*, pages 385–390, Castellón de la Plana, Spain, May 2001. Publicacions de la Universitat Jaume I, 2001.
- [7] A. Bia and A. Pedreño. The Miguel de Cervantes Digital Library: the Hispanic voice on the Web. *LLC (Literary and Linguistic Computing) journal, Oxford University Press*, 16(2):161–177, 2001. Presented at ALLC/ACH 2000, The Joint International Conference of the Association for Literary and Linguistic Computing and the Association for Computers and the Humanities, 21/25 July 2000, University of Glasgow.
- [8] A. Brüggemann-Klein and D. Wood. One-unambiguous regular languages. *Information and Computation*, 142(2):182–206, 1 May 1998.
- [9] L. Burnard. The Pizza Chef: a TEI Tag Set Selector. <http://www.hcu.ox.ac.uk/TEI/pizza.html>, September 1997. (Original version 13 September 1997, updated July 1998; Version 2 released 8 Oct 1999).
- [10] L. Burnard and C. M. Sperberg-McQueen. Tei lite: An introduction to text encoding for interchange (document no: Tei u 5). <http://www.uic.edu/orgs/tei/intros/teiu5.html>, 6 1995.
- [11] P. Caron and D. Ziadi. Characterization of Glushkov automata. *TCS: Theoretical Computer Science*, 233:75–90, 2000.
- [12] R. C. Carrasco and J. Oncina. Learning deterministic regular grammars from stochastic samples in polynomial time. *RAIRO (Theoretical Informatics and Applications)*, 33(1):1–20, 1999.
- [13] B. Chidlovskii. Using Regular Tree Automata as XML Schemas. In J. Hoppenbrouwers, T. de Souza Lima, M. Papazoglou, and A. Sheth, editors, *Proc. IEEE Advances on Digital Libraries Conference 2000*, pages 89–98, 2000.
- [14] T. Knuutila. Inference of k-testable tree languages. In H. Bunke, editor, *Advances in Structural and Syntactic Pattern Recognition (Proc. Intl. Workshop on Structural and Syntactic Pattern Recognition, Bern, Switzerland)*. World Scientific, August 1993.
- [15] B. Ludaescher, Y. Papakonstantinou, P. Velikhov, and V. Vianu. View Definition and DTD Inference for XML. In *Proc. Post-ICDT Workshop on Query Processing for Semistructured Data and Non-Standard Data Formats*, 1999.
- [16] M. Murata. Transformation of documents and schemas by patterns and contextual conditions. In D. W. Charles K. Nicholas, editor, *Proceedings of the Third International Workshop on Principles of Document Processing (PODP’96), Palo Alto, California, USA*, volume 1293 of *Lecture Notes in Computer Science*, pages 153–169. Springer-Verlag, Heidelberg, September 23 1996.

- [17] Y. Papakonstantinou and V. Vianu. DTD Inference for Views of XML Data. In *Principles Of Database Systems (PODS) 2000*, 2000.
- [18] W. Plotkin and C. Sperberg-McQueen. Text Encoding Initiative, homepage (tei@uic.edu). <http://www.uic.edu/orgs/tei/index.html>.
- [19] P. Prescod. Formalizing XML and SGML Instances with Forest Automata Theory. Technical report, University of Waterloo, Department of Computer Science, Waterloo, Ontario, May 1998. Draft Technical Paper.
- [20] K. E. Shafer. Creating DTDs via the GB-Engine and Fred. Technical report, OCLC Online Computer Library Center, Inc., 6565 Frantz Road, Dublin, Ohio 43017-3395, 1995.
- [21] C. M. Sperberg-McQueen and L. Burnard, editors. *Guidelines for Electronic Text Encoding and Interchange (Text Encoding Initiative P3), Revised Reprint, Oxford, May 1999*. TEI P3 Text Encoding Initiative, Chicago - Oxford, May 1994.
- [22] J. Verdú-Mas, J. Calera-Rubio, and R. Carrasco. A comparison of PCFG models. In C. Cardie, W. Daelemans, C. Nédellec, and E. Tjong-Kim-Sang, editors, *Proceedings of CoNLL-2000 and LLL-2000*, pages 123–125, New Brunswick, NJ (USA), 2000. Association for Computational Linguistics.
- [23] M. Young-Lai and F. W. M. Tompa. Stochastic Grammatical Inference of Text Database Structure. *Machine Learning*, 40(2):1, 2000.

APPENDIX

A. BUILDING THE GLUSHKOV AUTOMATA OF AN EXPRESSION

Given a regular expression r we describe the procedure [11] to build a Glushkov automaton. In the following $n \in \mathbb{N}$ is any position according to a marking E_r of r ; $F, G \in \text{reg}(\mathbb{N})$ denote subexpressions of E_r .

The Boolean function $\text{empty} : \text{reg}(\mathbb{N}) \rightarrow \{\text{TRUE}, \text{FALSE}\}$ is true if the subexpression contains the empty string and can be computed as

$$\begin{aligned}
 \text{empty}(n) &= \text{FALSE} \\
 \text{empty}(F|G) &= \text{empty}(F) \vee \text{empty}(G) \\
 \text{empty}(F, G) &= \text{empty}(F) \wedge \text{empty}(G) \\
 \text{empty}(F^*) &= \text{TRUE} \\
 \text{empty}(F^+) &= \text{empty}(F) \\
 \text{empty}(F?) &= \text{TRUE}
 \end{aligned} \tag{10}$$

The subset $\text{first} : \text{reg}(\mathbb{N}) \rightarrow 2^{\mathbb{N}}$ gives the positions that can be found as the first symbol in a string and is given by

$$\begin{aligned}
 \text{first}(n) &= \{n\} \\
 \text{first}(F|G) &= \text{first}(F) \cup \text{first}(G) \\
 \text{first}(F, G) &= \begin{cases} \text{first}(F) \cup \text{first}(G) & \text{if empty}(F) \\ \text{first}(F) & \text{otherwise} \end{cases} \\
 \text{first}(F^*) &= \text{first}(F) \\
 \text{first}(F^+) &= \text{first}(F) \\
 \text{first}(F?) &= \text{first}(F)
 \end{aligned} \tag{11}$$

The subset $\text{last} : \text{reg}(\mathbb{N}) \rightarrow 2^{\mathbb{N}}$ gives the positions that can

be found as the last symbol in a string and is given by

$$\begin{aligned}
 \text{last}(n) &= \{n\} \\
 \text{last}(F|G) &= \text{last}(F) \cup \text{last}(G) \\
 \text{last}(F, G) &= \begin{cases} \text{last}(F) \cup \text{last}(G) & \text{if empty}(G) \\ \text{last}(G) & \text{otherwise} \end{cases} \\
 \text{last}(F^*) &= \text{last}(F) \\
 \text{last}(F^+) &= \text{last}(F) \\
 \text{last}(F?) &= \text{last}(F)
 \end{aligned} \tag{12}$$

Finally, the subset $\text{follow} : \text{reg}(\mathbb{N}) \rightarrow 2^{\mathbb{N} \times \mathbb{N}}$ gives the pairs of positions that can be found consecutive in a string and is given by

$$\begin{aligned}
 \text{follow}(n) &= \emptyset \\
 \text{follow}(F|G) &= \text{follow}(F) \cup \text{follow}(G) \\
 \text{follow}(F, G) &= \text{follow}(F) \cup \text{follow}(G) \cup \text{last}(F) \times \text{first}(G) \\
 \text{follow}(F^*) &= \text{follow}(F) \cup \text{last}(F) \times \text{first}(F) \\
 \text{follow}(F^+) &= \text{follow}(F) \cup \text{last}(F) \times \text{first}(F) \\
 \text{follow}(F?) &= \text{follow}(F)
 \end{aligned} \tag{13}$$

With the above elements, the Glushkov automaton $[8]$ $(Q, \mathbb{N}, \delta, I, F)$ is given by

- $Q = \{I\} \cup \text{sym}(E_r)$
- $\delta(I, a) = \{n \in \text{first}(E_r) : \Phi_r(n) = a\}$
- $\delta(n, a) = \{m \in Q : (n, m) \in \text{follow}(E_r) \wedge \Phi_r(m) = a\}$
- $F = \begin{cases} \{I\} \cup \text{last}(E_r) & \text{if empty}(E_r) \\ \text{last}(E_r) & \text{otherwise} \end{cases}$

where I is any position not in $\text{sym}(E_r)$.

B. SYNTAX-DIRECTED DEFINITION

The process of simplification is described here as a syntax-directed definition [1] which associates translation actions to each rule in grammar 7, actions which compute attributes at each node of the parse tree of the regular expressions from attributes at child nodes. The simplified regular expression is stored in attribute $.s$, whereas attributes $.n$ and $.e$ are auxiliary Boolean attributes which store, respectively, whether the subexpression is simple and whether it accepts the empty string. Function $\text{HasWitnesses}()$ returns TRUE if the symbol in the regular expression has been visited during the parsing of a document in the sample.

PRODUCTION	TRANSLATION
$E_0 \rightarrow T E_1$	<pre> E₀.n := T.n and E₁.n; E₀.e := T.e or E₁.e; if T.w = ∅ then E₀.w := E₁.w else if T.w = λ then if E₁.w = ∅ or T.w = λ or E₁.e then E₀.w := E₁.w else E₀.w := E₁.w “?” endif else if E₁.w = λ then if T.e then E₀.w := T.w else E₀.w := T.w “?” endif else if E₁.w ≠ ∅ then begin E₀.w := T.w “ ” E₁.w; E₁.n := FALSE end endif endif endif endif endif </pre>
$E \rightarrow T$	<pre> E.w := T.w; E.n := T.n; E.e := T.e </pre>
$T_0 \rightarrow F, T_1$	<pre> T₀.n := F.n and T₁.n; T₀.e := F.e and T₁.e; if F.w = ∅ or F.w = λ then T₀.w := T₁.w else if T₁.w = ∅ or T₁.w = λ then T₀.w := F.w else begin T₀.w := F.w “,” T₁.w; T₀.n := FALSE end endif endif endif </pre>
$T \rightarrow F$	<pre> T.w := F.w; T.n := F.n; T.e := F.e </pre>

Figure 2: Syntax-directed definition for the regular expression simplification process (part 1). All attributes are synthetic.

PRODUCTION	TRANSLATION
$F \rightarrow W$	<pre> F.w := W.w; F.n := W.n; F.e := W.e </pre>
$F \rightarrow W^*$	<pre> F.n := TRUE; F.e := TRUE; if W.w = ∅ or W.w = λ then F.w := λ else F.w := W.w “*” endif </pre>
$F \rightarrow W^+$	<pre> F.n := TRUE; F.e := W.e; if W.w = ∅ or W.w = λ then F.w := W.w else F.w := W.w “+” endif </pre>
$F \rightarrow W?$	<pre> F.n := TRUE; F.e := TRUE; if W.w = ∅ or W.w = λ then F.w := λ else F.w := W.w “?” endif </pre>
$W \rightarrow \text{name}$	<pre> W.n := TRUE; W.e := FALSE; if HasWitnesses(name.w) then W.w := name.w else W.w := ∅ endif </pre>
$W \rightarrow (E)$	<pre> W.n := E.n; W.e := E.e; if E.w = ∅ or E.w = λ then W.w := E.w else if E.n then W.w := E.w else begin W.w := “(” E.w “)”; W.n := TRUE end endif endif </pre>

Figure 3: Syntax-directed definition for the regular expression simplification process (part 2). All attributes are synthetic.