

# Schema validation via streaming circuits\*

Filip Murlak  
University of Warsaw  
fmurlak@mimuw.edu.pl

Charles Paperman  
University of Warsaw  
paperman@mimuw.edu.pl

Michał Pilipczuk  
University of Warsaw  
michal.pilipczuk@mimuw.edu.pl

## ABSTRACT

XML schema validation can be performed in constant memory in the streaming model if and only if the schema admits only trees of bounded depth—an acceptable assumption from the practical view-point. In this paper we refine this analysis by taking into account that data can be streamed block-by-block, rather than letter-by-letter, which provides opportunities to speed up the computation by parallelizing the processing of each block. For this purpose we introduce the model of streaming circuits, which process words of arbitrary length in blocks of fixed size, passing constant amount of information between blocks. This model allows us to transfer fundamental results about the circuit complexity of regular languages to the setting of streaming schema validation, which leads to effective constructions of streaming circuits of depth logarithmic in the block size, or even constant under certain assumptions on the input schema. For nested-relational DTDs, a practically motivated class of bounded-depth XML schemas, we provide an efficient construction yielding constant-depth streaming circuits with particularly good parameters.

## Categories and Subject Descriptors

H.2.1 [Database Management]: Logical Design; F.4.3 [Mathematical Logic and Formal Languages]: Formal Languages

## Keywords

semi-structured data; XML; streaming; schema validation; Boolean circuits; nested-relational DTDs

\*These results were obtained when Charles Paperman and Michał Pilipczuk held post-doc positions at the Warsaw Centre of Mathematics and Computer Science. Michał Pilipczuk was supported by the Foundation for Polish Science (FNP) via the START stipend program. Filip Murlak was supported by Poland's National Science Centre grant no. 2013/11/D/ST6/03075.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*PODS'16, June 26–July 01, 2016, San Francisco, CA, USA*

© 2016 ACM. ISBN 978-1-4503-4191-2/16/06...\$15.00

DOI: <http://dx.doi.org/10.1145/2902251.2902299>

## 1. INTRODUCTION

Over tree-structured data, like XML documents or JSON files, schemas impose restrictions on the structure of trees modeling the data. Popular schema formalisms, like DTDs and especially XML Schema, are able to express very complex properties, bringing their expressive power close to tree automata, which are often used as theoretical abstractions of schemas. With such expressive power, the task of schema validation—that is, verifying that a given data instance conforms to the schema—is not entirely trivial even if we have direct access to the whole data instance. When the data are streamed, schema validation becomes a major challenge.

In their seminal paper [31], Segoufin and Vianu consider streaming validation in constant memory. An algorithm over streamed data that works in constant memory can be seen as a finite automaton. Whether such an algorithm exists for a given schema depends on whether the set of word representations of the instances of the schema is a regular language. Segoufin and Vianu show that the word representation of a regular tree language (covering all popular schema formalisms) is regular if and only if there exists a uniform bound on the depth of the trees in the language. In this paper we refine this result by looking more closely at the way data are streamed. Due to the result of Segoufin and Vianu, we focus mostly on tree languages of bounded depth.

Our starting point is the observation that data need not be fed to the algorithm letter-by-letter. For instance, if the data stream serves as an abstraction of sequential access to a mass storage device, the algorithm is fed entire blocks of data that are fetched to a moderately-sized cache. This requires evaluating the finite automaton over a word read in block-sized portions. We could process each block sequentially in time linear in its size, but we would like to do better, assuming certain ability to parallelize computation over each block. As a model of parallelism we choose Boolean circuits. The most important reason is that their relation with regular languages is well understood and documented [3, 19, 20, 32], but Boolean circuits have several other advantages. On one hand, they are very close to hardware implementation: from a Boolean circuit of small depth one can directly obtain a hardware description that could be compiled into, say, an FPGA. On the other hand, Boolean circuits are also a commonly accepted theoretical model for higher-level parallelism, providing abstraction for various concrete practical models. For example, on a multi-core machine different cores could be assigned to evaluating different parts of the circuit. We remark that combining the challenges of streaming data access and parallelization has been considered from the practical

perspective [5, 6, 12, 17], in particular in the context of the standard MapReduce approach (see e.g. [24]); however, to the best of our knowledge, hardly any theoretical models have been proposed so far.

In order to reconcile the random-access parallelism of Boolean circuits with the streaming setting, we introduce a model of computation called streaming circuits. Intuitively, a streaming circuit takes a block of the input word together with additional feedback information of constant size (the state of the underlying finite automaton) and outputs updated feedback. This model allows us to talk about the complexity of streaming algorithms in a way that does not abstract away the size of the block, and to transfer the huge body of results on the circuit complexity of regular languages to the streaming setting. It also avoids the inherent flaw of the classical Boolean circuit setting: the nonuniformity. While having a separate circuit for each size of the input data is entirely impractical, in our setting this is not an issue any more, as the circuit is nonuniform in the block size, which can be chosen and fixed in advance. With the right choice of the block size, the time needed to process a block can (potentially) be made close to the time needed to fetch it, thus avoiding bottlenecks.

Any finite automaton can be transformed into a streaming circuit of chosen block size. The challenge is to get the circuit as simple as possible. In the context of schema validation we are interested in how the complexity of tree language is reflected in the streaming circuits recognizing their word encodings. As we shall see, one can always build an  $\text{NC}^1$ -streaming circuit for any bounded-depth regular tree language. That is, in the block-by-block access model, one can efficiently do streaming validation in parallel—by a circuit that has polynomial size, constant fan-in, and logarithmic depth. Can we do better than that? For any class  $\mathcal{C}$  of circuits one can ask about  $\mathcal{C}$  streaming circuits for regular tree languages. A full positive answer to this question should include:

- an algorithm to decide if a given tree language has a  $\mathcal{C}$  streaming circuit;
- an algorithm to construct a recognizing  $\mathcal{C}$  streaming circuit, if it exists; and
- a syntactic fragment (a restricted schema language) that corresponds to languages that can be recognized with a  $\mathcal{C}$  streaming circuit.

From the practical point of view, the crucial part of the answer is a restricted schema language guaranteeing feasibility and an efficient algorithm to build the circuit from the schema definition in this restricted language. Ideally, the schema language should cover all feasible schemas, but this should not be achieved at the expense of its simplicity and usability.

We consider two restricted classes of circuits:  $\text{AC}^0$  and  $\text{WLAC}^0$ . Recall that  $\text{AC}^0$  comprises circuit families with polynomial size and constant depth bounds, and an  $\text{AC}^0$  circuit family is in  $\text{WLAC}^0$  (for *wire-linear*  $\text{AC}^0$ ) if the number of wires is bounded linearly. Unlike for  $\text{NC}^1$ , not all bounded-depth regular tree languages admit  $\text{AC}^0$  streaming circuits, but we can decide effectively if a given bounded-depth regular language admits one. We also show that if one additionally assumes that the tree language is definable in first order logic, then the answer is always affirmative. For

a practically relevant class of languages defined by nested-relational DTDs [1, 2], we provide an efficient construction of  $\text{AC}^0$  streaming circuits with particularly good properties. For  $\text{WLAC}^0$  we observe that one cannot even check the correctness of the usual XML encodings of bounded-depth trees. We propose a new encoding, enriched with the information about ancestors of nodes. Under the new encoding we can validate nested-relational schemas with  $\text{WLAC}^0$  streaming circuits. We also show that this encoding can be computed from the usual encoding by an  $\text{AC}^0$  streaming circuit that does not depend on the schema, only on the depth and the alphabet.

## 2. STREAMING CIRCUITS

In this article we use several well-known classes of circuits. We recall the basics very briefly and refer to the book by Straubing [32] for a more detailed presentation.

*Basics of circuit complexity.* We work with Boolean circuits with AND, OR, and NOT gates, taking as input words over alphabet  $\Sigma = \{a_1, a_2, \dots, a_k\}$ . The letters of the input word are encoded *in unary*; that is, each input gate is modeled with  $k$  binary gates, and letter  $a_i$  is encoded as the binary sequence  $0^{i-1}10^{k-i}$ . A family

$$(C_n)_{n \in \mathbb{N}}$$

of circuits recognizes a language  $L \subseteq \Sigma^*$  if  $C_n$  has  $n$  input gates and a single binary output gate, and returns 1 if and only if the input word is in  $L$ . We shall refer to this model of recognition as the random-access model. Whenever we consider the size of the circuit, we mean the number of gates.

Since we consider mainly regular languages (of words and of trees), we restrict ourselves to languages recognizable with  $\text{NC}^1$  circuit families: Boolean circuit families of polynomial size, logarithmic depth and bounded fan-in. Two other classes of interest are  $\text{AC}^0$  circuit families, which have polynomial size, constant depth and unbounded fan-in, and  $\text{WLAC}^0$  circuit families, which are  $\text{AC}^0$  circuit families with a linear number of wires (hence, also gates). The interaction of these classes with the class of regular languages is well understood, and we will use this knowledge to design adequate devices in the context of streaming schema validation.

We shall say that a language is in class  $\mathcal{C}$  if it is recognized by a  $\mathcal{C}$  circuit family. Some important examples separating the three classes described above:

- the parity language  $(c + ac^*a)^*$  is in  $\text{NC}^1$  (with linear-size circuits) but not in  $\text{AC}^0$  [14];
- $(c + ac^*b)^*$  is in  $\text{AC}^0$  (with quadratic-size circuits) but not in  $\text{WLAC}^0$  [19].

*Streaming circuits.* In the streaming circuit setting, a single circuit is used to recognize words of all lengths, by processing them sequentially in blocks of fixed size with the help of a feedback mechanism.

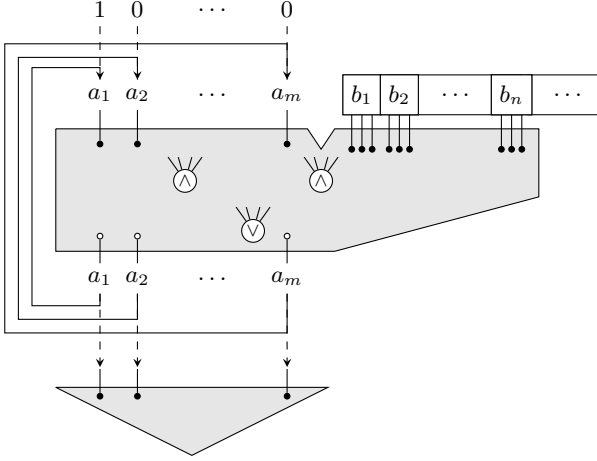
A *streaming circuit* over alphabet  $\Sigma$  with block size  $n$  and feedback size  $m$  is a circuit  $C$  with  $n$  input gates,  $m$  feedback gates and  $m$  output gates, together with an acceptor circuit  $A$  with  $m$  input gates and 1 output gate. The computation of such a circuit on an input word  $w$  is carried out in stages. In each stage the circuit  $C$  is given the output from the

previous stage (initially, the bit sequence  $10^{m-1}$ ) and the next size  $n$  block of the input word (in unary encoding, as explained above). The output of the last stage is fed to the acceptor circuit  $A$  and the word  $w$  is *accepted* if and only if the acceptor circuit  $A$  returns 1. If the last block of the input word is shorter than  $n$  symbols, a designated padding symbol  $\$$  (encoded as a sequence of zeros) is used to fill it up. More formally, let  $u_0 = 10^{m-1}$  and for  $i < \lceil \frac{|w|}{n} \rceil$  let

$$u_{i+1} = C(u_i, w_{ni+1}w_{ni+2} \dots w_{ni+n})$$

where  $w_j = \$$  for  $j > |w|$ ; the word  $w$  is accepted if

$$A(u_{\lceil \frac{|w|}{n} \rceil}) = 1.$$



Such a streaming circuit can be interpreted as a deterministic automaton over the alphabet  $\Gamma = \Sigma^n$ : the state space is  $\{0, 1\}^m$  with the initial state  $10^{m-1}$ , and the transition function and the set of accepting states are given by circuits  $C$  and  $A$ . Consequently, languages recognized by streaming circuits are regular.

In fact, streaming circuits give precise description of the implementation of finite automata over words read by fixed-size blocks. Indeed, if one takes an automaton and makes it read blocks of  $n$  letters instead of one letter, one obtains an automaton with the same state space over the alphabet  $\Sigma^n$ , but with the set of transitions growing exponentially with  $n$  (if  $|\Sigma| \geq 2$ ). Circuits allow us to represent (and carry out) transitions more efficiently.

We can therefore talk about streaming-circuit complexity of regular languages: a regular language  $L$  has *streaming-circuit complexity*  $\mathcal{C}$  if for some  $m$  there exists a  $\mathcal{C}$  circuit family  $(C_n)_{n \in \mathbb{N}}$  with  $m$  feedback gates and  $m$  output gates and an acceptor circuit  $A$  with  $m$  input gates, such that for each  $n$ , the streaming circuit  $(C_n, A)$  recognizes  $L$ . We say that  $(C_n)_{n \in \mathbb{N}}$  is a  $\mathcal{C}$  *streaming circuit family* for  $L$  (with feedback  $m$  and acceptor  $A$ ). In general,  $(C_n)_{n \in \mathbb{N}}$  need not have a finite description, but all families constructed in this paper will have one.

**Streaming vs random access.** Over regular languages of words, random-access recognizability and streaming recognizability coincide for reasonable classes of circuits. The following theorem provides efficient translation from circuits to streaming circuits, and *vice versa*. A class  $\mathcal{C}$  of circuits families is *closed under shifts* if for each circuit family  $(C_n)_{n \in \mathbb{N}}$

from  $\mathcal{C}$ , each family obtained from  $(C_{n+1})_{n \in \mathbb{N}}$  by hard-wiring a chosen input gate is in  $\mathcal{C}$ . By *parallel composition* we mean running two circuits on the same input and concatenating the outputs, and by *sequential composition* we mean wiring the output gates of one circuit to the input gates of another circuit.

**THEOREM 1.** *Let  $\mathcal{C}$  be a class of circuit families closed under shifts, parallel compositions, and sequential compositions with constant-size circuits. Then a regular language  $L \subseteq \Sigma^*$  has streaming-circuit complexity  $\mathcal{C}$  if and only if it is in  $\mathcal{C}$ .*

*More precisely, if the recognizing family of circuits has depth and size bounded by non-decreasing functions  $d(n)$  and  $s(n)$ , the resulting streaming circuit for block size  $n$  has feedback  $k$ , depth  $d(n + k + k^2) + \mathcal{O}(1)$ , and size  $k^3 \cdot s(n + k + k^2) + \mathcal{O}(k^3)$ , where  $k$  is the number of states of the minimal deterministic automaton for  $L$ .*

**PROOF.** The left-to-right implication is almost immediate. Let  $(C_n)_{n \in \mathbb{N}}$  be a  $\mathcal{C}$  streaming circuit family for  $L$  with feedback  $m$  and acceptor  $A$ . The family of circuits recognizing  $L$  is

$$(A(C_n(10^{m-1}, \cdot)))_{n \in \mathbb{N}};$$

that is, we hardwire the initial values in the feedback gates of  $C_n$ , and feed the output of  $C_n$  to  $A$ . Since circuit  $A$  is fixed, by the closure properties of  $\mathcal{C}$ , the resulting circuit family is indeed in  $\mathcal{C}$ , and so is  $L$ .

The right to left implication is more complicated. Let  $\mathcal{A}$  be the minimal deterministic automaton for  $L$  and  $(C_n)_{n \in \mathbb{N}}$  a  $\mathcal{C}$  family of circuits recognizing  $L$ . Let  $p, q$  be states of  $\mathcal{A}$ . We shall construct a  $\mathcal{C}$  family of circuits recognizing the language

$$L_{p,q} = \{w \in \Sigma^* \mid \delta_w(p) = q\},$$

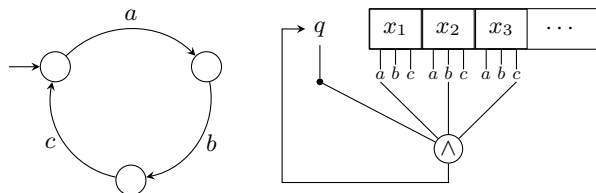
where  $\delta_w(p)$  is the state to which the automaton moves from state  $p$  after reading word  $w$ . Let  $u$  be the shortest word that reaches  $p$  from the initial state; by pumping,  $|u| \leq k$ , where  $k$  is the number of states of  $\mathcal{A}$ . By minimality, for all distinct states  $q, q'$  there is a word  $v$  of size at most  $k^2$  such that  $\delta_v(q) \in F$  if and only if  $\delta_v(q') \notin F$ . Consequently,  $L_{p,q}$  is a Boolean combination of  $k - 1$  *residual* languages of the form

$$u^{-1}L v^{-1} = \{w \mid u w v \in L\},$$

where  $|u| \leq k$  and  $|v| \leq k^2$ . Each of these languages can be recognized with a  $\mathcal{C}$  family of circuits  $(C'_n)_{n \in \mathbb{N}}$  where  $C'_n$  is obtained from  $C_{n+|u|+|v|}$  by hardwiring the first  $|u|$  input gates to  $u$  and the last  $|v|$  input gates to  $v$ . An appropriate Boolean combination of these circuits gives circuits for  $L_{p,q}$ . Assuming that  $i$ -th state is coded as  $0^{i-1}10^{k-i}$  on the feedback gates and the first state is initial, it is easy to obtain a  $\mathcal{C}$  family of streaming circuits for  $L$  from the circuits for languages  $L_{p,q}$ . We simply add on top of these circuits an additional circuit of depth  $\mathcal{O}(1)$  and size  $\mathcal{O}(k^2)$  that computes the state after processing the current block from the previous state passed in the feedback. To verify the required size bound, observe that in total we construct  $k^2$  circuits for languages  $L_{p,q}$ . Every such circuit consists of  $k - 1$  circuits for residual languages, each of size at most  $s(n + k + k^2)$ . Since  $k$  is considered a constant, by the closure properties of  $\mathcal{C}$  we have that the obtained circuit family belongs to  $\mathcal{C}$ .  $\square$

One could easily make the feedback logarithmic in  $k$  by encoding the state in binary. However, this would not improve the parameters of the circuit, as they depend on the number of states, not the size of their representation.

The fact that the circuit has access, thanks to its non-uniformity, to additional numerical information, sometimes allows to simplify drastically the ongoing computation. For instance, the language  $(a_1 a_2 \dots a_n)^*$  requires an automaton with  $n + 1$  states, but assuming block size  $n$  it can be recognized by a streaming circuit with feedback 1, which corresponds to 2 states:



This kind of behavior is difficult to analyze, but always beneficial in the construction of streaming circuits.

### 3. VALIDATION: A GENERAL BOUND

Following Segoufin and Vianu [31], we work with ordered unranked trees, node-labelled with letters from a finite alphabet  $\Sigma$ . We denote by  $\text{Trees}(\Sigma)$  the set of all such trees.

For technical convenience, we model schemas as “previous sibling, last child” tree automata. A *nondeterministic tree automaton*

$$\mathcal{A} = (\Sigma, Q, q_0, F, \delta)$$

consists of a finite input alphabet  $\Sigma$ , a finite set of states  $Q$  with an initial state  $q_0$ , a set of accepting states  $F \subseteq Q$ , and a transition relation

$$\delta \subseteq Q \times Q \times \Sigma \times Q.$$

Being in a node  $v$  of the input tree  $t \in \text{Trees}(\Sigma)$ , the automaton has processed  $t_v$ , the subtree of  $t$  rooted at  $v$ . The state  $q$  for node  $v$  depends on the label  $\sigma$  of  $v$  and the states  $q_1, q_2$  from the previous sibling and the last child of  $v$ , respectively, in the way specified by the transition relation:

$$(q_1, q_2, \sigma, q) \in \delta$$

(in leftmost siblings and leaves we use the initial state  $q_0$  instead of  $q_1$  and  $q_2$ , respectively). The tree  $t$  is *accepted* by  $\mathcal{A}$  if states can be chosen for nodes in such a way that the root gets a state from  $F$ . We write  $L(\mathcal{A})$  for the set of accepted trees. If  $L = L(\mathcal{A})$ , we say that  $L$  is *regular*, and that it is recognized by  $\mathcal{A}$ .

A schema language that is simpler, but often sufficient in practice, is offered by *document type definitions*, or DTDs for short. A DTD

$$\mathcal{D} = (\Sigma, r, P)$$

consists of a finite alphabet  $\Sigma$  with a distinguished root label  $r \in \Sigma$ , and a function  $P$  that assigns to each label  $a \in \Sigma$  a regular expression  $P(a)$  over  $\Sigma$ , called the production for  $a$ , and written as

$$a \rightarrow P(a).$$

A tree  $t \in \text{Trees}(\Sigma)$  conforms to  $\mathcal{D}$  if its root is labelled with  $r$  and for each label  $a \in \Sigma$  and each node  $v$  in  $t$  with label  $a$ ,

the sequence of labels of  $v$ 's children forms a word generated by the regular expression  $P(a)$ .

A practically relevant class of *nested-relational* DTDs, covering a large proportion of real life schemas [4], is obtained by assuming non-recursiveness (that is, no  $a$ -labelled node has an  $a$ -labelled descendant) and allowing only productions of the form

$$a \rightarrow \widehat{a}_1 \widehat{a}_2 \dots \widehat{a}_\ell,$$

where  $a_1, a_2, \dots, a_\ell$  are distinct elements of  $\Sigma$ , and  $\widehat{a}_i$  is equal to  $a_i$ ,  $a_i^? = (\varepsilon + a_i)$ ,  $a_i^*$ , or  $a_i^+ = a_i a_i^*$ .

In the context of streaming processing we need a string representation of trees. Under the *XML encoding* trees are represented as words over  $\Sigma \cup \overline{\Sigma}$ , the elements of  $\Sigma$  and  $\overline{\Sigma}$  being, respectively, the *opening* and *closing tags*,

$$\text{flat} \left( \begin{array}{c} a \\ / \quad \backslash \\ t_1 \quad \dots \quad t_k \end{array} \right) = a \cdot \text{flat}(t_1) \dots \text{flat}(t_k) \cdot \bar{a}.$$

We call  $\text{flat}(t)$  the *flattening* of  $t$ , and

$$\text{flat}(L) = \{ \text{flat}(t) \mid t \in L \}$$

the *flattening* of  $L$ . Another natural possibility is the *term encoding*, which is similar to the XML encoding except that we only have one closing tag symbol  $\#$ ,

$$\text{flat}_{\#} \left( \begin{array}{c} a \\ / \quad \backslash \\ t_1 \quad \dots \quad t_k \end{array} \right) = a \cdot \text{flat}_{\#}(t_1) \dots \text{flat}_{\#}(t_k) \cdot \#.$$

Intuitively, the XML encoding corresponds to recognition by a *visibly push-down automaton* [22] (or input driven automaton [13]). The term encoding requires only one stack symbol, which corresponds to visibly counter automata [8]. In the sequel we work with the XML encoding for concreteness, but for most of our results, the choice of the encoding does not matter.

As observed by Segoufin and Vianu [31], the flattening of a regular tree language is a regular word language if and only if the tree language has bounded depth (there exists a uniform bound on the depth of all trees in  $L$ ).

**PROPOSITION 1** (SEGOUFIN, VIANU [31]). *For each regular tree language  $L$  the following conditions are equivalent:*

- $L$  has bounded depth;
- $\text{flat}(L)$  is a regular word language.

Thus, to have any chance for streaming-circuit validation, we restrict our attention to bounded-depth trees. For practical purposes this is an acceptable assumption, as real-life schemas tend to be bounded-depth [4].

Translation from bounded-depth tree automata to *deterministic* word automata over encodings involves only single-exponential blow-up.

**PROPOSITION 2.** *Let  $\mathcal{A}$  be a tree automaton with  $k$  states recognizing a bounded-depth language  $L \subseteq \text{Trees}(\Sigma)$ . One can construct a deterministic automaton with  $\mathcal{O}(|\Sigma|^k \cdot 2^{k^2})$  states recognizing  $\text{flat}(L)$ .*

**PROOF.** As  $L(\mathcal{A})$  has bounded depth, each accepting run of  $\mathcal{A}$  uses each state at most once on each branch of the input tree. Indeed, if this was not the case, one could construct an arbitrarily deep tree accepted by  $\mathcal{A}$  by repeating the part of the tree corresponding to the segment of the branch between

two occurrences of the same state. Consequently,  $L(\mathcal{A})$  has depth at most  $k$ .

Let  $\mathcal{B} = (\Sigma, Q, q_0, \delta, F)$  be a deterministic automaton recognizing  $L$  obtained from  $\mathcal{A}$  by the standard power-set construction; we have  $|Q| = 2^k$ . The automaton for  $\text{flat}(L)$  simulates stack of depth at most  $k$  in its states. Its state-space is

$$(\Sigma \times Q)^{\leq k} \cup \{\perp, \top\},$$

where the empty sequence  $\varepsilon$  is the initial state and  $\top$  is the only final state. The transitions are given as follows: upon reading symbol  $\sigma \in \Sigma$  in state  $\alpha \notin \{\perp, \top\}$ ,

- if  $|\alpha| < k$ , move to  $\alpha(\sigma, q_0)$ ,
- if  $|\alpha| = k$ , move to  $\perp$ ;

upon reading symbol  $\bar{\sigma} \in \bar{\Sigma}$  in state  $\alpha \notin \{\perp, \top\}$ ,

- if  $\alpha = \varepsilon$ , move to  $\perp$ ,
- if  $\alpha = (\sigma, q_1)$  and  $\delta(q_0, q_1, \sigma) \in F$ , move to  $\top$ ,
- if  $\alpha = (\sigma, q_1)$  and  $\delta(q_0, q_1, \sigma) \notin F$ , move to  $\perp$ ,
- if  $\alpha = \beta(\sigma', q')(\sigma, q)$ , move to  $\beta(\sigma', \delta(q', q, \sigma))$ ,
- if  $\alpha = \beta(\tau, q)$  with  $\tau \neq \sigma$ , move to  $\perp$ ;

upon reading any symbol in state  $\perp$  or  $\top$ , move to  $\perp$ .  $\square$

We finish this section with a general  $\text{NC}^1$ -upper bound for streaming circuit complexity of bounded-depth regular tree languages. Assuming the usual interpretation of  $\text{NC}^1$  as the class of problems that can be solved efficiently in parallel, this shows that streaming validation parallelizes. The bound combines Theorem 1, Proposition 2 and a folklore fact that all regular languages are in  $\text{NC}^1$  (i.e., random-access validation parallelizes).

**PROPOSITION 3.** *Each regular language  $L$  can be recognized by an  $\text{NC}^1$  streaming circuit family. More precisely, for a given block size  $n$  one can construct a recognizing circuit with feedback  $k$ , depth  $\mathcal{O}(\log n)$ , and size  $\mathcal{O}(k^3 n)$ , where  $k$  is the number of states of the minimal deterministic automaton for  $L$ .*

**PROOF.** The first part of the claim follows by Theorem 1 from the fact that all regular languages are in  $\text{NC}^1$ . To achieve the claimed bounds, we adapt the standard construction to obtain directly a streaming circuit.

Let  $\mathcal{A} = (\Sigma, Q, q_0, \delta, F)$  be the minimal automaton for  $L$ , and let  $|Q| = k$ . Any function  $f: Q \rightarrow Q$  can be represented as a  $k \times k$  binary matrix, which in turn can be seen as a  $k^2$ -bit word. From a single input letter  $\sigma$  one can compute function  $\delta_\sigma$ , represented as a  $k^2$ -bit word, using a circuit that has depth 1 and size  $\mathcal{O}(k^2)$ ; note that this circuit hardwires the transition function of  $\mathcal{A}$ . Given two  $k \times k$  matrices (as  $k^2$ -bit words), one can compute their product (over the Boolean algebra) with a circuit of depth 2 and size  $\mathcal{O}(k^3)$ . Finally, for an input word  $w$  of length  $n = 2^k$  one can compute the function  $\delta_w$  by first computing functions  $\delta_{w_i}$ , and then composing them in pairs to obtain functions for pairs of consecutive letters, quadruples, octuples, etc. The resulting circuit  $C_n$  has depth  $\mathcal{O}(\log n)$  and size  $\mathcal{O}(k^3 n)$ . If  $|w|$  is not a power of two, take  $C_{2^{\lceil \log n \rceil}}$  and hardwire identity function in place of the functions for last  $2^{\lceil \log n \rceil} - |w|$  input letters;

this preserves the bounds. From the family  $(C_n)_{n \in \mathbb{N}}$  one easily constructs a streaming circuit family for  $L$ : one simply computes the value of the function computed by  $C_n$  on the state represented (in unary) by the values in the feedback gates. This can be done easily with a circuit of depth 2 and size  $\mathcal{O}(k^2)$ .  $\square$

Propositions 2 and 3 immediately yield the following.

**THEOREM 2.** *For each regular bounded-depth language  $L \subseteq \text{Trees}(\Sigma)$ ,  $\text{flat}(L)$  can be recognized by an  $\text{NC}^1$  streaming circuit family. More precisely, for block size  $n$ , one can construct a recognizing streaming circuit with feedback  $\mathcal{O}(|\Sigma|^k \cdot 2^{k^2})$ , depth  $\mathcal{O}(\log n)$ , and size  $\mathcal{O}(|\Sigma|^{3k} \cdot 2^{3k^2} \cdot n)$ , where  $k$  is the number of states of the given nondeterministic automaton recognizing  $L$ .  $\square$*

As remarked in [31], if the input tree language is given as a DTD with productions defined by *unambiguous* regular expressions (as required by the DTD specification), one can construct a finite automaton recognizing the flattening, whose number of states is bounded by a polynomial of degree at most  $|\Sigma|$ . This immediately improves the bounds in the theorem above.

As we have seen in Theorem 1, for regular word languages streaming and random-access recognition with circuits is the same for any reasonable class of circuits. For flattenings of regular tree languages, this is not the case. In the streaming model, the flattening must be regular to be recognized by any circuit family; in the random-access model, the flattening of each regular tree language can be recognized by an  $\text{NC}^1$  circuit family [13].

## 4. VALIDATION IN CONSTANT DEPTH

In the last section we saw a generic construction translating a description of a bounded-depth tree language (an automaton or a DTD) into a streaming circuit with relatively good parameters: logarithmic depth, constant fan-in, and polynomial size. However, this construction is largely suboptimal as shown by the following example.

**EXAMPLE 1.** *Let  $L$  be given by the following DTD*

$$r \rightarrow a^*;$$

*that is,  $L$  consists of trees of the form*

$$\begin{array}{c} r \\ \swarrow \quad \searrow \\ a \quad \dots \quad a \end{array}$$

*Then,  $\text{flat}(L) = r(a\bar{a})^* \bar{r}$  can be easily recognized by an  $\text{AC}^0$  streaming circuit family (of linear size).*

On the other hand, any regular language not in  $\text{AC}^0$  gives a depth-2 regular tree language whose flattening is not recognizable by an  $\text{AC}^0$  streaming circuit family.

**EXAMPLE 2.** *From the parity lower bound [14] we immediately get that for the tree language  $L$  given by*

$$r \rightarrow (ab^*a + b)^*,$$

*$\text{flat}(L)$  cannot be recognized by an  $\text{AC}^0$  family of streaming circuits.*

Yet again, a simple modification can turn a hard tree language into an easy one, by adding more structure.

EXAMPLE 3. For the language  $L$  given by the DTD

$$r \rightarrow (c + b)^*, \quad c \rightarrow ab^*a,$$

the flattening  $\text{flat}(L)$ , given by

$$r(c(a\bar{a}(b\bar{b})^*a\bar{a})\bar{c} + b\bar{b})^*\bar{r},$$

is recognized by an  $\text{AC}^0$  streaming circuit family.  $\clubsuit$ <sup>1</sup>

As regular languages in  $\text{AC}^0$  have an exact logical characterization, and membership in  $\text{AC}^0$  is decidable, one could decide whether the flattening of a given regular bounded-depth tree language can be recognized by an  $\text{AC}^0$  family of circuits. To explain it in more detail, we need to recall some classical results in descriptive complexity. Then, we shall look at practical fragments of the DTD formalism that guarantee the existence of  $\text{AC}^0$  streaming circuit families.

**First order logic and constant-depth circuits.** We consider first order logic over words, encoded as relational structures over universe  $\{0, \dots, n-1\}$  where  $n$  is the length of the word. Formulas are generated by the first order logic grammar with two kinds of atomic predicates: the *letter predicates* of the form  $\mathbf{a}(x)$  that are true if and only if the position  $x$  in the word is labelled by  $a$ , and *numerical predicates* which are predicates speaking about the word stripped of labels. For conciseness, we also allow numerical constants  $\text{min}$  and  $\text{max}$  for the first and last positions in the word.

EXAMPLE 4. The language

$$\{a^n b^n \mid n \in \mathbb{N}\}$$

is defined by the formula

$$\text{max} \equiv 1 \pmod{2} \wedge \forall y \ y < \frac{\text{max}}{2} \leftrightarrow \mathbf{a}(y).$$

It is well known that word languages definable in FO with arbitrary numerical predicates are exactly languages in  $\text{AC}^0$  (see [16] for instance). The simplest way to translate an FO sentence  $\varphi$  into a constant-depth circuit is to introduce a gate for each subformula  $\psi(x_1, \dots, x_k)$  and each choice of positions  $i_1, \dots, i_k$  in the word. The most external logical symbol in  $\psi$  determines the type of the gate:  $\vee, \wedge, \neg$  correspond to OR, AND or NOT gates, and quantifiers are interpreted as disjunctions and conjunctions over all positions of the word. The gate is connected to the gates corresponding to appropriate subformulas of  $\psi(x_1, \dots, x_k)$  with variables valued accordingly. The gates for the letter predicates are simply the binary input gates encoding the input symbols. If  $P$  is a numeric predicate, the gate for  $P(i_1, \dots, i_k)$  is either constantly 0 or constantly 1, depending on  $P$  and  $i_1, \dots, i_k$  (this is where we use non-uniformity). The depth of this circuit is bounded by the depth of the formula, seen as a term. The number of gates is bounded by  $\|\varphi\| \cdot n^k$ , where  $\|\varphi\|$  is the number of different subformulas in  $\varphi$  and  $k$  is the maximal number of free variables in a subformula.

This construction can be optimized for  $\text{FO}^k$ , that is, for formulas using (and reusing) only  $k$  variables (see [21, 26]).

<sup>1</sup>Claims marked with  $\clubsuit$  can be verified using, e.g., an on-line tool available at: [paperman.cadilhac.name/sage](http://paperman.cadilhac.name/sage).

Such a formula can be written in a normal form in which quantification is always of the form

$$\exists x_1 \delta(x_1, \dots, x_k) \wedge \psi_2 \wedge \dots \wedge \psi_k$$

such that  $\delta(x_1, \dots, x_k)$  is a quantifier-free formula using only numerical predicates, and the set of free variables of  $\psi_j$  does not contain  $x_j$ . Then, it essentially suffices to have gates for subformulas with at most  $k-1$  free variables: for each valuation of variables  $x_2, \dots, x_k$ , we have an OR gate connected to the ANDs of the gates for  $\psi_j$  with the variable  $x_1$  valued in all ways that make  $\delta(x_1, \dots, x_k)$  hold. The size of the resulting circuit is bounded by  $\|\varphi\| \cdot n^{k-1}$ .

**Regular languages and logic.** The connection between logic and regular languages is a field of research on its own that takes its root into the celebrated results of McNaughton and Papert [25] and Schützenberger [29], who characterized regular languages of  $\text{FO}[\prec]$ , that is languages definable in first order logic with the linear order over positions. By extending this result to a slightly more complicated fragment, and by using the parity lower bounds for  $\text{AC}^0$ , Barrington et al. [3] proved that regular languages in  $\text{AC}^0$  are exactly those definable in  $\text{FO}[\prec, \text{MOD}]$ ; that is, in first order logic with (strict) order and the unary modulo predicates of the form  $x \equiv r \pmod{q}$  for arbitrary  $r, q \in \mathbb{N}$ . Furthermore, this class of regular languages has decidable membership thanks to its algebraic characterization.

Thus, we get the following corollary from Theorem 1.

COROLLARY 1. Given a bounded-depth regular tree language  $L$ , one can decide if  $\text{flat}(L)$  can be recognized with an  $\text{AC}^0$  streaming circuit family; a recognizing streaming circuit for a given block size can be constructed effectively.

Checking whether a regular word language is definable in  $\text{FO}[\prec, \text{MOD}]$  is PSPACE-complete [10]. The algorithm to construct a circuit from an automaton runs in time linear in the size of the *syntactic monoid* of the recognized language, and is therefore efficient as long as the syntactic monoid is not too large. In general, the syntactic monoid has size at most exponential in the size of the minimal automaton recognizing the language.

A more practical approach to providing  $\text{AC}^0$  streaming circuits is to define a subclass of DTDs that are directly transformable into  $\text{AC}^0$  streaming circuit families. In order to identify such a subclass, we first show that for bounded-depth tree languages, definability in FO is equivalent to  $\text{FO}[\prec]$ -definability of the flattening.

**FO-definable tree languages.** First order logic over trees uses the letter predicates  $\mathbf{a}(x)$  and the navigational predicates  $\text{child}(x, y)$ ,  $\text{descendant}(x, y)$ ,  $\text{nextSibling}(x, y)$ , and  $\text{followingSibling}(x, y)$ , which hold if and only if  $y$  is respectively a child, descendant, the next sibling or a following sibling of  $x$ .

EXAMPLE 5. The language of trees over a single-letter alphabet with only one branch is defined by the formula

$$\forall x, y \ \text{descendant}(x, y) \vee \text{descendant}(y, x).$$

Note that the flattening of this language is exactly the one given in Example 4 (with  $b = \bar{a}$ ).

We begin with a lemma which shows that, assuming bounded depth, the tree structure can be recovered from the

flattening with FO[<] formulas. In the flattening, we think of the positions with the opening tags as the ones representing the nodes of the tree.

LEMMA 1. *For all  $d > 0$  there exist FO[<] formulas*

$$\text{tree}_d(x, y), \quad \text{forest}_d(x, y)$$

*expressing that the segment from  $x$  to  $y$  is, respectively, the flattening of a tree of depth at most  $d$ , and a concatenation of such flattenings.*

*Consequently, for all  $d > 0$  there exist FO[<] formulas*

$$\text{child}_d(x, y), \quad \text{desc}_d(x, y), \quad \text{next}_d(x, y), \quad \text{foll}_d(x, y)$$

*expressing (over flattenings of trees of depth at most  $d$ ) that the node represented by position  $x$  and the node represented by position  $y$  are, respectively, in relation child, descendant, next sibling, and following sibling.*

PROOF. The formulas  $\text{tree}_d$  and  $\text{forest}_d$  are defined by mutual recursion. Let

$$\text{tree}_0(x, y) = \text{forest}_0(x, y) = \text{false}.$$

For  $d > 0$ , the formula  $\text{forest}_d(x, y)$  expresses that each position between  $x$  and  $y$  has a matching position between  $x$  and  $y$  such that the corresponding segment is the flattening of a tree of depth at most  $d$ ,

$$\begin{aligned} \text{forest}_d(x, y) &:= x < y \wedge \\ &\wedge \forall u \in (x, y) \exists v \in (x, y) (\text{tree}_d(u, v) \vee \text{tree}_d(v, u)), \end{aligned}$$

and  $\text{tree}_d(x, y)$  checks that  $x$  and  $y$  are labelled by matching tags and the segment between them is a concatenation of the flattenings of trees of depth at most  $d - 1$ ,

$$\text{tree}_d(x, y) := \left( \bigvee_{a \in \Sigma} \mathbf{a}(x) \wedge \bar{\mathbf{a}}(y) \right) \wedge \text{forest}_{d-1}(x, y).$$

It is not difficult to verify that  $\text{tree}_d$  and  $\text{forest}_d$  indeed define, respectively, flattenings of trees of depth at most  $d$  and concatenations of such flattenings.

The remaining formulas,

$$\text{desc}_d(x, y) := \exists x' \text{tree}_d(x, x') \wedge y \in (x, x') \wedge \left( \bigvee_{a \in \Sigma} \mathbf{a}(y) \right),$$

$$\begin{aligned} \text{child}_d(x, y) &:= \text{desc}_d(x, y) \wedge \\ &\wedge \neg \exists z \in (x, y) \text{desc}_d(x, z) \wedge \text{desc}_d(z, y), \end{aligned}$$

$$\text{foll}_d(x, y) := x < y \wedge \exists z \text{child}_d(z, x) \wedge \text{child}_d(z, y),$$

$$\text{next}_d(x, y) := \text{foll}_d(x, y) \wedge \neg \exists z \in (x, y) \text{foll}_d(x, z),$$

are straightforward.  $\square$

Lemma 1 and Theorem 1 give the following result.

THEOREM 3. *For each bounded-depth tree language  $L$ ,  $L$  is FO-definable if and only if  $\text{flat}(L)$  is FO[<]-definable. In consequence, flattenings of FO-definable tree languages are recognized by  $\text{AC}^0$  streaming circuit families; the recognizing streaming circuit for a given block size can be constructed effectively.*

PROOF. The formula defining the flattening of  $L$  is obtained by taking the conjunction of  $\text{tree}_d(\min, \max)$  and the formula defining  $L$  with each occurrence of child, descendant, next-sibling, and following-sibling replaced with the appropriate formula given by Lemma 1.

For the converse implication, we begin by rewriting each FO[<]-formula over the flattenings so that quantification is in one of the following forms:

$$\exists x^o \left( \bigvee_{a \in \Sigma} \mathbf{a}(x^o) \right) \wedge \varphi, \quad \exists x^c \left( \bigvee_{a \in \Sigma} \bar{\mathbf{a}}(x^c) \right) \wedge \varphi;$$

the resulting formula is at most exponentially larger. Then, each variable has its type, opening or closing. We now rewrite each atomic predicate for all possible types of variables. For an opening variable  $x^o$ ,  $\mathbf{a}(x^o)$  remains unchanged and  $\bar{\mathbf{a}}(x^o)$  is rewritten as *false*; for a closing variable  $x^c$ ,  $\mathbf{a}(x^c)$  is rewritten as *false*, and  $\bar{\mathbf{a}}(x^c)$  as  $\mathbf{a}(x^c)$ . For a closing variable  $x^c$  and an opening variable  $y^o$ , the atomic formula  $x^c < y^o$  is rewritten as  $\text{right}(x^c, y^o)$ , where  $\text{right}(x, y)$  is the formula

$$\begin{aligned} \exists z \exists z' &(\text{descendant}(z, x) \vee z = x) \wedge \\ &\wedge \text{followingSibling}(z, z') \wedge \\ &\wedge (\text{descendant}(z', y) \vee z' = y). \end{aligned}$$

Similarly, the formulas

$$\begin{aligned} &\text{descendant}(x^o, y^o) \vee \text{right}(x^o, y^o), \\ &\text{descendant}(x^o, y^c) \vee \text{right}(x^o, y^c) \vee \\ &\quad \vee \text{descendant}(y^c, x^o) \vee x^o = y^c, \\ &\text{right}(x^c, y^c) \vee \text{descendant}(y^c, x^c) \end{aligned}$$

are used for the remaining three cases. We obtain a formula of FO on trees, easily seen to be equivalent to the original formula of FO[<] on flattenings.  $\square$

Theorem 3 gives an effective sufficient condition for the existence of an  $\text{AC}^0$  streaming circuit family for the flattening: as FO[<] definability is decidable for regular languages, so is FO-definability for bounded-depth regular tree languages. We remark that for regular tree languages of unbounded depth, it is a major open problem whether FO-definability is decidable [7].

The condition given by Theorem 3 is not necessary. As shown in the next example, capturing the entire class of bounded-depth regular tree languages admitting  $\text{AC}^0$  streaming circuit families for the flattenings would require intricate artificial syntactic restrictions over the basic formalism, with an unclear gain in expressivity.

EXAMPLE 6. *Consider the following two DTDs:*

$$r \rightarrow (aa)^*, \quad a \rightarrow (bb)^*; \quad r \rightarrow (aa)^*, \quad a \rightarrow (bbb)^*.$$

*The flattening of the language given by the left DTD is in  $\text{AC}^0$ , whereas for the right DTD it is not.  $\clubsuit$*

The argument in Theorem 3 does not give good complexity bounds: the FO formula for the flattening has size linear in the original formula (and exponential in the depth), but the automaton constructed from the formula, needed to invoke Theorem 1, may have non-elementary size. And even if there was a more efficient way to do it, FO is not a natural schema definition language. A desirable language should be a natural fragment of a known schema definition language. We discuss such a fragment in the following subsection.

We finish this subsection with a remark that without the bounded-depth assumption one can recognize flattenings of FO-definable tree languages in  $\text{TC}^0$  in the random-access model. Recall that  $\text{TC}^0$  is defined like  $\text{AC}^0$ , except that Majority gates can also be used.

PROPOSITION 4. *Let  $L$  be an FO-definable language of trees. Then  $\text{flat}(L)$  is in  $\text{TC}^0$ .*

PROOF SKETCH. We repeat the proof of Lemma 1 and Theorem 3, but this time using circuits rather than formulas. In the course of the structural induction, we need to deal with formulas with free variables. We work with words over the alphabet

$$\{0, 1\}^n \times (\Sigma \cup \bar{\Sigma}),$$

for sufficiently large  $n$ . Over such words we evaluate a formula  $\varphi(x_1, x_2, \dots, x_n)$  by assuming that  $x_i$  is assigned the *unique* position that has 1 in the  $i$ -th coordinate of its label.

It is sufficient to prove that the formula  $\text{tree}(x_i, x_j)$ , saying that the infix from position  $x_i$  to position  $x_j$  is the flattening of a tree, is definable in  $\text{TC}^0$ : to conclude we simply note that  $\text{TC}^0$  is closed under FO quantification (and Boolean connectives), so all predicates from Lemma 1, and all FO formulas using them, can be defined in  $\text{TC}^0$  as well.

To express  $\text{tree}(x_i, x_j)$  with a  $\text{TC}^0$  circuit, we proceed as follows. For every position  $y \in [x_i, x_j]$  which is labelled by an opening tag, we find a position  $z \in [y + 1, x_j]$  such that  $z$  is labelled by the matching closing tag and the number of open tags between  $y$  and  $z$  is exactly the number of closing tags. This last test can be implemented as the AND of two Majority gates, checking that the number of opening tags is at most the number closing tags, and *vice versa*.  $\square$

Even the flattening of the set of all trees is  $\text{TC}^0$ -hard, but not all flattenings of regular languages of unbounded depth are: for instance the language of trees with only one branch over unary alphabet is FO-definable and its flattening is in  $\text{AC}^0$  (see Example 4).

### A practical formalism for validation in constant-depth.

A natural formalism allowing validation with constant-depth streaming circuits can be obtained by restricting the productions in DTDs to FO[<]-definable languages. Over words, being FO[<]-definable is equivalent to being definable with a *star-free* regular expression [25]; that is, an expression built from symbols from the alphabet and the empty set by means of concatenation and all Boolean operations, including complement.

COROLLARY 2. *Each language  $L \subseteq \text{Trees}(\Sigma)$  defined by a non-recursive DTD with star-free productions can be defined in FO on trees, and so can be recognized by an  $\text{AC}^0$  streaming circuit family.*

PROOF. For each  $a \in \Sigma$  there is an FO[<] sentence  $\varphi_a$  defining the word language generated by the production for  $a$ . As the DTD is non-recursive, all generated trees have depth at most  $|\Sigma|$ . We define  $L$  with the formula

$$\exists x \mathbf{r}(x) \wedge \forall y \text{-descendant}(y, x) \wedge \forall z \bigwedge_{a \in \Sigma} \mathbf{a}(z) \rightarrow \widehat{\varphi}_a(z),$$

where  $r$  is the root label of the DTD, and the formula  $\widehat{\varphi}_a(z)$  is obtained from the sentence  $\varphi_a$  by replacing each occurrence of the predicate  $<$  with the predicate  $\text{followingSibling}$  and restricting all quantifiers to the children of  $z$ ; that is, subformulas  $\exists y \psi$  and  $\forall y \psi$  are replaced, respectively, with

$$\exists y \text{child}(z, y) \wedge \psi \quad \text{and} \quad \forall y \text{child}(z, y) \rightarrow \psi$$

(assuming that variable  $z$  is not used in  $\varphi_a$ ).  $\square$

The popular class of nested-relational DTDs is a special case, which admits constant-depth streaming circuits with particularly good parameters.

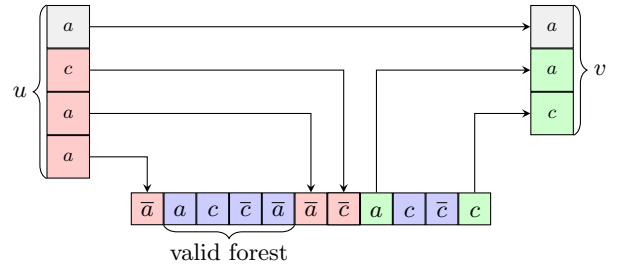
THEOREM 4. *The flattening of each language  $L \subseteq \text{Trees}(\Sigma)$  defined by a nested-relational DTD can be recognized by an  $\text{AC}^0$  streaming circuit family with feedback  $\mathcal{O}(d \cdot |\Sigma|)$ , depth  $\mathcal{O}(d)$  and size  $\mathcal{O}(d^3 \cdot |\Sigma|^2 \cdot n^2)$  for block size  $n$ , where  $d \leq |\Sigma|$  is the maximal depth of trees in  $L$ .*

PROOF. We shall directly construct a streaming circuit, using FO formulas over separate blocks of the input word as an intermediate formalism. Since the language is defined by a nested-relational DTD, its depth is bounded by some  $d \leq |\Sigma|$ . Before we look at the DTD any further, we construct a circuit that for each position  $x$  in the block computes  $\text{open}(x) \in \Sigma^{\leq d}$ , the sequence of unmatched opening tags in the prefix of the entire input word up to (and including) position  $x$ . Note that  $\text{open}(x)$  is equal to the sequence of labels on the path from the root to the node corresponding to  $x$  in the encoded tree, including this node if the tag of  $x$  is opening, and not including it otherwise. As the feedback we shall use  $\text{open}(\text{min} - 1)$ , where  $\text{min} - 1$  is the last position of the previous block. We use the padding symbol  $\$$  (encoded as a sequence of zeros) to fill up  $\text{open}(\text{min} - 1)$  to  $d$  symbols.

First, we compute the values of the formula  $\text{forest}_d(x, y)$  from Lemma 1 for all positions within the block. Note that  $\text{forest}_d(x, y)$  has quantifier rank  $\mathcal{O}(d)$  and uses only 3 variables. Its size is exponential if the recursive definition is unravelled, but it has only  $\mathcal{O}(d + |\Sigma|)$  different subformulas. Thus, the standard translation for formulas with 3 variables gives a streaming circuit of depth  $\mathcal{O}(d)$  and size  $\mathcal{O}((d + |\Sigma|) \cdot n^2)$ ; the circuit has  $n^2$  output gates, one for each pair of values of  $x$  and  $y$ . From now on, we shall treat  $\text{forest}_d(x, y)$  as an atomic formula.

Similarly, we can assume the existence of atomic formulas  $\bar{f}_p(x)$  for  $p = 1, 2, \dots, d$  expressing that position  $x$  has the closing tag corresponding to the  $p$ -th letter stored in the feedback. Their values for all  $x$  can be computed with a circuit of depth  $\mathcal{O}(1)$  and size  $\mathcal{O}(d \cdot |\Sigma| \cdot n)$ .

How does  $v = \text{open}(x)$  depend on  $u = \text{open}(\text{min} - 1)$ ? The situation can be illustrated as follows:



We express this with formulas

$$\varphi_{i,a}(x)$$

for  $i = 1, 2, \dots, d$ , and  $a \in \Sigma$ , saying that the  $i$ -th symbol of  $\text{open}(x)$  is  $a$ , assuming that the feedback stores  $\text{open}(\text{min} - 1)$ . The formula is a disjunction over  $0 \leq j \leq k, \ell \leq d$  with  $\ell \geq i$  of conditions saying that

- the first  $\$$  in the feedback is at position  $k + 1$ ,
- there exist positions  $x \geq y_\ell > y_{\ell-1} > \dots > y_{j+1}$  with opening tags,



- there exist positions  $y_{j+1} > z_{j+1} > z_{j+2} > \dots > z_k$  with closing tags  $\bar{f}_{j+1}, \bar{f}_{j+2}, \dots, \bar{f}_k$ , such that
- all segments between these positions (including  $\text{min}-1$  and  $x$ ) are flattenings of forests, and
- $a = f_i$  if  $i \leq j$ , and  $a(y_i)$  if  $i > j$ ,

where  $f_p$  is the symbol in the  $p$ -th feedback gate.

This resulting formula uses  $\mathcal{O}(d)$  quantifiers and can be written with 2 variables only. It has  $\mathcal{O}(d^3)$  different subformulas (with  $\text{forest}_d(x, y)$  and  $\bar{f}_p(x)$  treated as atomic formulas). Thus, the standard translation gives a streaming circuit of depth  $\mathcal{O}(d)$  and size  $\mathcal{O}(d^3 \cdot n^2)$  for block size  $n$ ; the circuit has  $n$  output gates, one for each value of  $x$ . Note that we cannot use the optimized construction giving linear-size circuit, because the formula  $\text{forest}_d$  does not define a numerical predicate; that is, it depends on the labels of positions.

The  $d \cdot |\Sigma|$  circuits for formulas  $\varphi_{i,a}(x)$  compute  $\text{open}(x)$  for all  $x$ . If for some  $x$  and  $i$  we have 0 for all  $a$ 's, it means that  $\text{open}(x)$  has less than  $i$  symbols. Note that the total size of the constructed circuits is  $\mathcal{O}(d^4 \cdot |\Sigma| \cdot n^2)$ .

Local correctness of the encoding can be checked by verifying that  $\text{open}(x)$  is nonempty throughout the computation; as soon as it becomes empty, the remaining positions in the block should store the padding symbol  $\$,$  and it should be the last block. All this can be tested with a circuit of depth  $\mathcal{O}(1)$  and size  $\mathcal{O}(n^2)$ .

Once we have built the circuit computing  $\text{open}(x)$  for each position  $x$ , we have access to the label of the parent of each node whose closing tag is within the block: it is the last letter of  $\text{open}(x)$ . The labelling restrictions enforced by the DTD can be expressed with the formula

$$\begin{aligned} & \forall x \in [\text{min}-1, \text{max}] \\ & \bigwedge_{a \in \Sigma} (\mathbf{a}(x) \rightarrow \bigvee_{c \in N(a)} \mathbf{c}(x+1)) \wedge \\ & \bigwedge_{a, b \in \Sigma} (\text{open}(x) \in \Sigma^* a \mathcal{S}^* \wedge \bar{\mathbf{b}}(x) \rightarrow \bigvee_{c \in N(a, b)} \mathbf{c}(x+1)). \end{aligned}$$

where  $N(a)$  is the set of labels that can succeed the opening tag  $a$ , and  $N(a, b)$  is the set of labels that can succeed the closing tag  $\bar{b}$  in the scope of tag  $a$ . Both sets are determined by the production for  $a$ . Assume the production is

$$a \rightarrow \hat{a}_1 \hat{a}_2 \dots \hat{a}_k.$$

There are two cases, depending on whether there is  $i$  such that  $\hat{a}_i$  is either  $a_i^+$  or  $a_i$ . If there is such  $i$ , let us take the minimal one. Then,  $N(a) = \{a_1, a_2, \dots, a_i\}$ . If there is no such  $i$ ,  $N(a) = \{a_1, a_2, \dots, a_k\} \cup \{\bar{a}\}$ . The set  $N(a, a_j)$  is characterized analogously: if there is  $i > j$  such that  $\hat{a}_i$  is either  $a_i^+$  or  $a_i$ , then for the minimal such  $i$  we have

$$N(a, a_j) = \begin{cases} \{a_j, a_{j+1}, \dots, a_i\} & \text{for } \hat{a}_j \in \{a_j^*, a_j^+\}, \\ \{a_{j+1}, \dots, a_i\} & \text{for } \hat{a}_j \in \{a_j, a_j^?\}. \end{cases}$$

If no such  $i$  exists,

$$N(a, a_j) = \begin{cases} \{a_j, a_{j+1}, \dots, a_k\} \cup \{\bar{a}\} & \text{for } \hat{a}_j \in \{a_j^*, a_j^+\}, \\ \{a_{j+1}, \dots, a_k\} \cup \{\bar{a}\} & \text{for } \hat{a}_j \in \{a_j, a_j^?\}. \end{cases}$$

Note that to evaluate the formula we need access to  $\text{open}(\text{min}-1)$  and the tag at position  $\text{min}-1$ . We have

already pointed out that  $\text{open}(\text{min}-1)$  is included in the feedback; now we see that also the label at position  $\text{min}-1$  should be a part of the feedback. Assuming unary encoding of letters, the size of the feedback is  $\mathcal{O}(|\Sigma|^2)$ . The standard translation of the formula gives a circuit of depth  $\mathcal{O}(1)$  and size  $\mathcal{O}((|\Sigma|^2 + d \cdot |\Sigma|) \cdot n)$ . The output of the circuit is used to propagate the information about the lack of error so far, which is done by means of a designated error feedback gate.

Combining the two stages we obtain a circuit of depth  $\mathcal{O}(d) = \mathcal{O}(|\Sigma|)$  and size  $\mathcal{O}(d^3 \cdot |\Sigma|^2 \cdot n^2)$ .  $\square$

Let us remark that the construction can be extended to productions with thresholds  $a^{\ell \cdot k}$ ,  $a^{\geq k}$ , at the cost of including more information in the feedback: for each label in  $\text{open}(y)$  we would need the number of its repetitions among its siblings so far (up to threshold  $k$ ).

## 5. WIRE-LINEAR CIRCUITS

While having an  $\text{AC}^0$  streaming circuit family guarantees depth independent of the block size, it is still possible that the number of gates and the number of wires makes implementation for larger block sizes unreasonable. We now turn to *wire-linear circuit families*,  $\text{WLAC}^0$ ; that is, bounded-depth circuit families in which the number of wires (and thus the number of gates) grows linearly with the size of the input (or block in case of streaming circuits). As for  $\text{AC}^0$ ,  $\text{WLAC}^0$  has been studied and regular languages in  $\text{WLAC}^0$  are characterized.

*Regular languages in  $\text{WLAC}^0$ .* The logical characterization of regular languages in  $\text{WLAC}^0$  is given by the following result, extending a similar characterization of regular languages with a *neutral letter* in  $\text{WLAC}^0$  [19].

**THEOREM 5** ([26]). *A regular language is in  $\text{WLAC}^0$  if and only if it is definable in  $\text{FO}^2[+1, <, \text{MOD}]$ .*

Note that the signature includes the successor relation  $+1$ , which cannot be defined from  $<$  with just 2 variables.

Unlike for  $\text{AC}^0$ , both directions are involved. The lower bound relies on an effective algebraic characterization of languages definable in  $\text{FO}^2[+1, <, \text{MOD}]$  from [11] (which also makes definability decidable) and the fact that the language  $(c + ac^*b)^*$  is not in  $\text{WLAC}^0$  [19]. The upper bound, which we care mostly about, uses a clever circuit construction for *prefix functions* [9]. For completeness, we sketch this construction and explain how to use it to construct a  $\text{WLAC}^0$  circuit family from a formula of  $\text{FO}^2[+1, <, \text{MOD}]$ . To get a  $\text{WLAC}^0$  streaming circuit family we use Theorem 1.

Consider the language

$$\Sigma^* a \Sigma^* a \Sigma^*$$

of words with at least two letters  $a$ . It can be defined by the formula

$$\exists x \mathbf{a}(x) \wedge \exists y y < x \wedge \mathbf{a}(y)$$

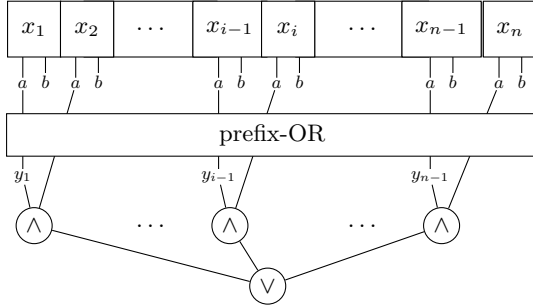
The standard translation from  $\text{FO}$ , which introduces a gate for each subformula with free variables valuated in all possible ways, gives a circuit of quadratic size. The optimized construction for  $\text{FO}^2$  formulas gives a circuit with linearly many gates, but quadratically many wires: for each value of  $x$  we have an OR gate connected to the circuits for  $\mathbf{a}(y)$  for all  $y < x$ .

To obtain a wire-linear circuit, we use *prefix functions*. The prefix-OR is a function  $f: \{0, 1\}^n \rightarrow \{0, 1\}^n$  such that

$$f(u)_i = \bigvee_{j \leq i} u_j;$$

suffix-OR, prefix-AND, and suffix-AND are defined similarly. A  $\text{WLAC}^0$  circuit for prefix-OR is constructed by evaluating prefix-OR naively (with quadratically many wires) over the ORs of size- $\sqrt{n}$  blocks, and then over each block separately with the additional knowledge of the bit computed by the first stage for the previous block. If we use a separate circuit for each block, we get a circuit with  $\mathcal{O}(n\sqrt{n})$  wires. To avoid this we note that we need to compute the prefix-OR only for the single block where 0's switch to 1's in the prefix-OR for block ORs. The remaining prefix and suffix functions can be computed similarly. For more details we refer to the original article [9].

Coming back to our example, a  $\text{WLAC}^0$  circuit for the language  $\Sigma^* a \Sigma^* a \Sigma^*$  can be obtained by computing the prefix-OR of *being letter a* and checking if there exists an input gate which contains  $a$  such that the prefix-OR for the previous position evaluates to 1:



This construction can be nested: for the language of words having at least  $k$  occurrences of  $a$ , one uses  $k$  prefix-OR circuits with  $n$  inputs interleaved with  $k$  layers of AND gates of fan-in 2, with the last layer of AND gates connected to a single OR gate. The size of the resulting circuit is therefore  $\mathcal{O}(kn)$  and is independent from the size of the input alphabet.

To build a circuit for an arbitrary  $\text{FO}^2$  formula we proceed by structural induction over formulas in the classical normal form. The basic cases are unary predicates  $\mathbf{a}(x)$  and  $x \equiv r \pmod q$ , for which the naïve construction gives wire-linear circuits. As  $\text{WLAC}^0$  is closed under Boolean connectives, the only difficulty in the inductive step is the quantification. In the normal form, the quantification is always of the form  $\exists y \delta(x, y) \wedge \varphi(y)$  where  $\delta(x, y)$  only uses predicates  $x < y$  and  $x = y + k$  for  $k \in \mathbb{Z}$ . We deal with it like in the example above, by computing prefix functions for  $\varphi(1), \dots, \varphi(n)$  and then for each  $x$  adding an OR gate wired appropriately to the outputs of  $\varphi(1), \dots, \varphi(n)$  and the prefix functions; details can be adapted from [19] or found in [26].

**Tree languages and  $\text{WLAC}^0$  validation.** Using the described results on regular languages in  $\text{WLAC}^0$ , and Theorem 1, we obtain the following corollary.

**COROLLARY 3.** *Given a regular bounded-depth tree language  $L$ , one can decide if  $\text{flat}(L)$  can be recognized with a  $\text{WLAC}^0$  streaming circuit family; a recognizing streaming circuit for a given block size can be constructed effectively.  $\square$*

XML flattenings of the sets of trees of depth 1 and 2 (over a singleton alphabet), that is, languages  $(ab)^*$  and  $(a(ab)^*b)^*$ , are in  $\text{WLAC}^0$ , but for depth 3 the flattening is the language  $(a(a(ab)^*b)^*b)^*$ , which is not in  $\text{WLAC}^0$ . ( $\clubsuit$ ) Hence, unlike for full FO, the two-variable fragment of FO over trees is not captured by  $\text{WLAC}^0$  (and thus by  $\text{FO}^2$  over words). In fact, it is not easy to imagine a nontrivial schema formalism that guarantees recognizability with a  $\text{WLAC}^0$  streaming circuit family. We shall eventually propose such a formalism, but for now we shift the perspective and ask: what if we change the way trees are encoded as words?

We propose an encoding giving even more information than the XML encoding: the path-from-the-root encoding. Trees of depth at most  $d$  over alphabet  $\Sigma$  are encoded as words over alphabet

$$\Sigma \cup \bar{\Sigma} \cup \{\$, \},$$

where  $\$$  is a padding symbol, used to simplify the circuits. For  $0 < i \leq d$ , a word  $u \in \Sigma^{i-1}$ , a tree  $t$  with root  $a \in \Sigma$ , and children that are trees  $t_1, \dots, t_k$  we set

$$\Delta_u^d(t) = ua\$^{d-i} \cdot \Delta_{ua}^d(t_1) \cdots \Delta_{ua}^d(t_k) \cdot u\bar{a}\$^{d-i}$$

and let the *path-from-the-root* encoding of tree  $t$  be

$$\Delta^d(t) = \Delta_\varepsilon^d(t).$$

For instance, for a tree  $t$  consisting of an  $r$ -root with one  $a$ -child and one  $c$ -child,  $\Delta^2(t) = r\$rar\bar{a}rcr\bar{c}\bar{r}\$$ .

For this new encoding, we still have that the flattening

$$\Delta^d(L) = \{\Delta^d(t) \mid t \in L\}$$

of a regular bounded-depth tree language  $L$  is a regular language of words. Moreover, the correctness of the encoding can be checked by a  $\text{WLAC}^0$  streaming circuit family. We shall write  $\Delta(L)$  for the encoding of  $L$ , assuming that the parameter  $d$  is clear from the context.

**LEMMA 2.** *Let  $\text{Trees}_d(\Sigma)$  be the set of trees over  $\Sigma$  of depth at most  $d$ . The language  $\Delta(\text{Trees}_d(\Sigma))$  is recognized by a  $\text{WLAC}^0$  streaming circuit family with feedback  $\mathcal{O}(d \cdot |\Sigma|)$ , depth  $\mathcal{O}(1)$ , and  $\mathcal{O}(|\Sigma| \cdot (n + d))$  wires.*

**PROOF.** To encode consecutive tags in the flattening, the path-from-the-root encoding uses blocks of  $d$  consecutive symbols: at most  $d$  symbols on the path to the root, padded to length  $d$  with symbol  $\$$ . These blocks will be called *d-slabs*. For simplicity, in the following we assume that  $n$  is divisible by  $d$ , so that  $d$ -slabs fit exactly into blocks processed by the circuit. If this is not the case, we can adapt the construction by passing the previous incomplete block in the feedback, together with the length  $r < d$  of the passed fragment, encoded in unary. This increases the number of feedback gates by  $\mathcal{O}(d)$ . All the congruences used in the following constructions can be easily adjusted to take in the account the fact that the  $d$ -slabs in the block are shifted by  $r$ .

The feedback consists of the last  $d$ -slab of the previously processed block, and a special error gate that passes the information whether an error has been encountered so far. In the first block, we assume that the initial feedback encodes a  $d$ -slab consisting only of symbols  $\$$ . As such  $d$ -slabs will never appear again throughout the computation, from this feedback we can recognize that the first block is being processed. Henceforth we assume that the circuit has the access to symbols at positions between  $\text{min} - d$  and  $\text{min} - 1$ , or to the information that these positions do not exist.

Let us introduce an auxiliary formula  $\text{last}(x)$  expressing that  $x$  is the last position with a non-padding symbol within its  $d$ -slab:

$$\text{last}(x) := \neg\$(x) \wedge (\$(x+1) \vee x \equiv d-1 \pmod{d}).$$

We can compute this formula for all positions  $x$  using  $\mathcal{O}(n)$  gates and depth  $\mathcal{O}(1)$ .

Let us now see how to verify that the block's description is correct. For this, we check that certain conditions hold for every position  $x \in [\min -d, \max -d]$  (respectively,  $x \in [0, \max -d]$  for the first block) using a formula whose encoding as a circuit will be straightforward. We first express that the padding symbols behave as expected:

$$\begin{aligned} x \equiv 0 \pmod{d} &\rightarrow \neg\$(x), \\ \$(x) \wedge \neg\$(x+1) &\rightarrow x \equiv d-1 \pmod{d}, \\ \$(x) \wedge \neg\$(x+d) &\rightarrow \text{last}(x-1). \end{aligned}$$

The first implication asserts that no  $d$ -slab contains the  $\$$  symbol at its front. The second one asserts that after any  $\$$  symbol, we necessarily have  $\$$  symbols up to the end of the  $d$ -slab. The third one asserts that the only  $\$$  symbol that can be replaced by a letter in the next  $d$ -slab is the first one. Finally, we introduce the following implications for all  $a \in \Sigma$ :

$$\begin{aligned} \mathbf{a}(x) &\rightarrow \bar{\mathbf{a}}(x+d) \vee (\mathbf{a}(x+d) \wedge \neg\text{last}(x+d)), \\ \bar{\mathbf{a}}(x+d) &\rightarrow \text{last}(x+d) \wedge \mathbf{a}(x). \end{aligned}$$

The first implication asserts that an opening tag either changes to the matching closing tag or remains the same and another symbol is added after it. In particular, an opening tag is never replaced by  $\$$ . The second implication asserts that a closing tag has to be produced from the last position of the previous  $d$ -slab. In particular, in the next  $d$ -slab it can be replaced by an opening tag or a  $\$$  symbol, but not by a closing tag. The described formula can be turned directly into a wire-linear streaming circuit with a single gate of unbounded fan-in, connected to  $|\Sigma| \cdot n$  subcircuits of size  $\mathcal{O}(1)$  and constant fan-in.

The circuit described above verifies that the description of the block is correct. The error gate passed as the feedback to the next iteration is its conjunction with the error gate received in the feedback from the previous iteration. The acceptor circuit just checks that no error has occurred.  $\square$

Even if some structural properties of the input trees are accessible under the new encoding,  $\text{WLAC}^0$  circuits still fail to capture  $\text{FO}^2$  definable tree languages, as shown in the next example.

EXAMPLE 7. Let  $L$  be the language given by the DTD

$$a \rightarrow b^*, \quad b \rightarrow (c+d)^*$$

restricted to trees in which one of the root's children has only  $c$ -children. It is clearly definable in  $\text{FO}^2$ . The flattening  $\Delta(L)$  is (up to relabelling) the language

$$(a(b(\bar{c}\bar{c} + \bar{d}\bar{d})^*\bar{b})^*\bar{a} \cap \Sigma^*b(\bar{c}\bar{c})^*\bar{b}\Sigma^*,$$

which is not in  $\text{WLAC}^0$ .  $\clubsuit$

This latter example works mainly because the language is not defined by a DTD. Of course, it is hopeless to believe that we can have a good streaming circuit for any bounded-depth DTD, but what if we restrict the productions

to  $\text{FO}^2[<]$ -definable languages? As it turns out we can still get a tree language whose path-from-the-root encoding is not recognizable with a  $\text{WLAC}^0$  streaming circuit family.

EXAMPLE 8. For the language  $L$  given by the DTD

$$r \rightarrow (a)^*, \quad a \rightarrow b^*c^*b^*,$$

the flattening  $\Delta(L)$  is (up to relabelling) the language

$$r(a(\bar{b}\bar{b})^*(\bar{c}\bar{c})^*(\bar{b}\bar{b})^*\bar{a})^*\bar{r},$$

which is not in  $\text{WLAC}^0$ .  $\clubsuit$

What we can tackle are nested-relational DTDs.

PROPOSITION 5. Let  $L$  be a tree language defined by a nested-relational DTD. Then  $\Delta(L)$  is definable in  $\text{FO}^2[<, +1, \text{MOD}]$  and is recognized by a  $\text{WLAC}^0$  streaming circuit family with feedback  $\mathcal{O}(|\Sigma|^2)$ , depth  $\mathcal{O}(1)$ , and  $\mathcal{O}(|\Sigma|^2 \cdot n)$  wires for block size  $n$ .  $\square$

PROOF. It suffices to combine the construction from Lemma 2, which guarantees correctness of the encoding, with the second stage of the construction in Theorem 4.  $\square$

Unlike for Theorem 4, we cannot directly extend the construction to DTDs with more general thresholds.

EXAMPLE 9. For the language  $L$  defined by the DTD

$$r \rightarrow a^*, \quad a \rightarrow bb, \quad b \rightarrow c^*$$

the flattening  $\Delta(L)$  is (up to relabelling) the language

$$r(ab(\bar{c}\bar{c})^*\bar{b}\bar{b}(\bar{c}\bar{c})^*\bar{b}\bar{a})^*\bar{r},$$

which is not in  $\text{WLAC}^0$ .  $\clubsuit$

A new encoding may seem an easy way out, since we add exactly the information needed to validate nested-relational DTDs with circuits having good parameters. However, it is possible to benefit from this solution even without using the path-from-the-root encoding: if we have control over the design of the schema, by adjusting the tags we can assume that each label uniquely determines the label of the father. This ensures  $\text{WLAC}^0$  validation at the cost of a mild restriction of the practical scope of nested-relational DTDs. For instance, when a relational database is exported to XML (a major usage of nested-relational DTDs), the XML schema is obtained from a *covering forest* of the ER diagram [23]. As long as each entity in the ER diagram is covered by one branch, our condition is satisfied. If some entity is covered by more branches, we can use different element names for each branch and re-unify them into a single entity later. We could also use a *spanning forest* instead and represent the remaining relationships with foreign keys, but verifying such constraints is beyond the scope of our setting.

If we cannot or would not modify the schema nor change the encoding into path-from-the-root, we can enrich the given encoding before feeding it to the validating streaming circuit. The enriching can be done by a fixed transducer (dependent only on the alphabet, not the schema itself), implemented with an  $\text{AC}^0$  streaming circuit family with additional outputs. This should be viewed as a (rather complicated) fixed device that could be optimized and implemented in hardware once for all. Meanwhile, the proper validation stage can be realized with a reprogrammable hardware device of adapted size, that can be readjusted as the schema changes.

A *streaming circuit with output* with input alphabet  $\Sigma$ , output alphabet  $\Gamma$ , block size  $n$ , and feedback size  $m$  is like a streaming circuit, except that it has two kinds of output gates: immediate and pass-on. The output word over a given input word is obtained by concatenating the values on the immediate output gates for subsequent blocks of the input word; the values on the pass-on output gates are sent to the feedback gates when the next block is processed. The immediate output gates encode letters from  $\Gamma$  in unary, just as the input letters are encoded:  $i$ -th letter is encoded as  $0^{i-1}10^{|\Gamma|-i}$ . The output value is returned only if the acceptor circuit returns 1; if it returns 0, the output is undefined.

PROPOSITION 6. *Let  $\Sigma$  be a finite alphabet and let  $d \in \mathbb{N}$ . For trees over  $\Sigma$  of depth at most  $d$ , one can compute the path-from-the-root flattening from the XML flattening with a streaming circuit with output that has feedback  $\mathcal{O}(d \cdot |\Sigma|)$ , depth  $\mathcal{O}(d)$ , and size  $\mathcal{O}(d^4 \cdot |\Sigma| \cdot n^2)$  for block size  $n$ .*

PROOF. The first stage of the construction in the proof of Theorem 4 gives almost the circuit we need: it remains to append  $\bar{a}$  to *open*( $x$ ) for positions  $x$  labelled with  $\bar{a}$ ; this does not influence the bounds.  $\square$

## 6. CONCLUSIONS

We have introduced streaming circuits, which model parallel processing of streamed data in a way compatible with the schema validation task. We have shown that general results on the circuit complexity of regular languages can be used directly to reason about the existence of good streaming circuits for bounded-depth regular tree languages, giving effective but inefficient criteria. For a restricted, but practically crucial, class of languages defined by nested-relational DTDs we have provided a direct construction of circuits with excellent parameters: compositions of a quadratic-size  $\text{AC}^0$  circuits dependent only on the depth of trees and the alphabet, and wire-linear  $\text{AC}^0$  circuits dependent on the DTD. This construction can be extended easily to schemas with more general thresholds. Extending it further would be very relevant practically.

We have seen how to get a constant-depth polynomial-size streaming circuit from a bounded-depth tree language definable in first order logic. This relies on the fact that FO-definability of word languages is decidable and effective. There is a certain trade-off between the depth of the circuits and the degree of the polynomial bounding their size. We have seen that all FO-definable languages have quadratic circuits, but achieving this requires increasing the depth. It is even possible to achieve near-linear upper bound for these languages by increasing the depth sufficiently (see [20]). However, one may want to optimize the depth at the cost of increasing the degree of the polynomial. This question is related to the famous dot-depth conjecture, which is equivalent to deciding levels of the alternation hierarchy of first-order logic. Indeed, if a language belongs to the  $k$ -th level of this hierarchy then it is a finite Boolean combination of regular languages recognized by depth- $k$  circuits. Straubing conjectures that the languages in  $k$ -th level of the hierarchy (enriched with the modulo predicates) are exactly the Boolean combinations of regular languages recognized by depth- $k$  circuits [32].

Another related problem is the question of the circuit complexity of word encodings of regular tree languages (of

unbounded depth). While there are  $\text{TC}^0$ -complete and  $\text{NC}^1$ -complete examples, no good characterizations are known. We conjecture that each regular tree language is either  $\text{NC}^1$ -complete or in  $\text{TC}^0$ . It would be very interesting to have an effective characterization of  $\text{NC}^1$ -complete regular tree languages such that languages that do not satisfy it are in  $\text{TC}^0$ . Such a characterization exists in the classical setting of word languages and relies on an algebraic decomposition of finite automata [32]. Since such a decomposition for tree automata is unknown and related to the open question of deciding FO-definability of regular tree languages [7], we believe this question might be very hard.

Checking correctness of the encoding has a huge impact on validation. A way to isolate it is to consider *weak validation*, where the input is assumed to be a correct encoding of a tree (a *well-formed* document, under XML encoding). While no tree language of unbounded depth can be validated in constant memory, some can be weakly validated [30, 31]. For instance, for the set of trees whose each  $a$  node's leftmost child has label  $b$ , we only need to check that each opening  $a$  tag is followed by an opening  $b$  tag. However, it remains an open question to decide whether a given language can be weakly validated in constant memory. When restricted to bounded-depth tree languages, this question can be seen as a special case of the *separation* problem for regular word languages, which has rich bibliography of its own. The problem of *separation of languages from class  $\mathcal{K}$  by languages from class  $\mathcal{S}$*  is to decide for given languages  $K, M \in \mathcal{K}$  if there exists a language  $S \in \mathcal{S}$  such that  $K \subseteq S$  and  $M \cap S = \emptyset$ . Weak validation for a tree language  $L$  amounts to separating the flattening of  $L$  from the flattening of the complement of  $L$ . Separation of regular word languages by FO[ $\langle \cdot \rangle$ ]-definable languages is known to be decidable [15, 27], and similarly for  $\text{FO}^2[+1, \langle \cdot \rangle]$  [28]. One can use these results as a black box to find, respectively,  $\text{AC}^0$  and  $\text{WLAC}^0$  streaming circuits for the weak validation. Unfortunately, the separation abstraction is too powerful for tree languages of *unbounded* depth: as showed recently by Kopczyński, separation of flattenings of regular tree languages with regular word languages is undecidable under both XML and term encoding [18]. In contrast, the weak validation problem under term encoding is decidable [8], and still open under XML encoding.

## Acknowledgments

We thank the anonymous reviewers of PODS 2016 for their useful and inspiring comments.

## 7. REFERENCES

- [1] Serge Abiteboul, Luc Segoufin, and Victor Vianu. Representing and querying XML with incomplete information. *ACM Trans. Database Syst.*, 31(1):208–254, 2006.
- [2] Marcelo Arenas and Leonid Libkin. XML data exchange: Consistency and query answering. *J. ACM*, 55(2), 2008.
- [3] David A. Mix Barrington, Kevin Compton, Howard Straubing, and Denis Thérien. Regular languages in  $\text{NC}^1$ . *J. Computer and System Sciences*, 44(3):478–499, 1992.
- [4] G. J. Bex, F. Neven, and J. Van den Bussche. DTDs versus XML Schema: a practical study. In *WEBDB 2004*, pages 79–84, 2004.

- [5] David Black-Schaffer. *Block-Parallel Programming for Real-Time Applications on Multi-Core Processors*. PhD thesis, Humboldt-Universität zu Berlin, 2008. Available at [http://cva.stanford.edu/people/davidbbs/black-schaffer\\_thesis\\_final.pdf](http://cva.stanford.edu/people/davidbbs/black-schaffer_thesis_final.pdf).
- [6] David Black-Schaffer and William J. Dally. Block-parallel programming for real-time embedded applications. In *ICPP 2010*, pages 297–306. IEEE Computer Society, 2010.
- [7] Mikołaj Bojańczyk, Howard Straubing, and Igor Walukiewicz. Wreath products of forest algebras, with applications to tree logics. *Logical Methods in Computer Science*, 8(3), 2012.
- [8] Vince Bárány, Christof Löding, and Olivier Serre. Regularity problems for visibly pushdown languages. In Bruno Durand and Wolfgang Thomas, editors, *STACS 2006*, volume 3884 of *Lecture Notes in Computer Science*, pages 420–431. Springer Berlin Heidelberg, 2006.
- [9] Ashok K. Chandra, Steven Fortune, and Richard J. Lipton. Lower bounds for constant depth circuits for prefix problems. In Josep Díaz, editor, *ICALP 1983*, volume 154 of *Lecture Notes in Computer Science*, pages 109–117. Springer, 1983.
- [10] Sang Cho and Dung T. Huynh. Finite-automaton aperiodicity is PSPACE-complete. *Theoretical Computer Science*, 88(1):99 – 116, 1991.
- [11] Luc Dartois and Charles Paperman. Alternation hierarchies of first order logic with regular predicates. In Adrian Kosowski and Igor Walukiewicz, editors, *FCT 2015*, volume 9210 of *Lecture Notes in Computer Science*, pages 160–172. Springer, 2015.
- [12] Tathagata Das, Yuan Zhong, Ion Stoica, and Scott Shenker. Adaptive stream processing using dynamic batch sizing. In Ed Lazowska, Doug Terry, Remzi H. Arpaci-Dusseau, and Johannes Gehrke, editors, *SoCC 2014*, pages 16:1–16:13. ACM, 2014.
- [13] Patrick Dymond. Input-driven languages are in log n depth. *Inf. Process. Lett.*, 26(5):247–250, January 1988.
- [14] Merrick Furst, James B. Saxe, and Michael Sipser. Parity, circuits, and the polynomial-time hierarchy. *Theory of Computing Systems*, 17:13–27, 1984.
- [15] Karsten Henckell. Pointlike sets: the finest aperiodic cover of a finite semigroup. *Journal of Pure and Applied Algebra*, 55(1):85 – 126, 1988.
- [16] Neil Immerman. Languages that capture complexity classes. *SIAM J. Computing*, 16(4):760–778, 1987.
- [17] Rashid Khogali, Olivia Das, and Kaamran Raahemifar. Mobile parallel computing algorithms for single-buffered, speed-scalable processors. In *TrustCom 2013 / ISPA-13 / IUCC-2013*, pages 1832–1839. IEEE Computer Society, 2013.
- [18] Eryk Kopczyński. Invisible pushdown languages. In *LICS 2016 (to appear)*, 2016. Available at <http://arxiv.org/abs/1511.00289>.
- [19] Michal Koucký, Pavel Pudlák, and Denis Thérien. Bounded-depth circuits: separating wires from gates. In Harold N. Gabow and Ronald Fagin, editors, *STOC 2005*, pages 257–265. ACM, 2005.
- [20] Michal Koucký. Circuit complexity of regular languages. *Theory of Computing Systems*, 45(4):865–879, 2009.
- [21] Michal Koucký, Clemens Lautemann, Sebastian Poloczek, and Denis Thérien. Circuit lower bounds via Ehrenfeucht-Fraïssé games. In *CCC 2006*, pages 190–201, 2006.
- [22] Viraj Kumar, P. Madhusudan, and Mahesh Viswanathan. Visibly pushdown automata for streaming XML. In *WWW 2007*, pages 1053–1062. ACM, 2007.
- [23] Dongwon Lee, Murali Mani, Frank Chiu, and Wesley W. Chu. Net & cot: translating relational schemas to XML schemas using semantic constraints. In *Proceedings of the 2002 ACM CIKM International Conference on Information and Knowledge Management, McLean, VA, USA, November 4-9, 2002*, pages 282–291. ACM, 2002.
- [24] Kyong-Ha Lee, Yoon-Joon Lee, Hyunsik Choi, Yon Dohn Chung, and Bongki Moon. Parallel data processing with MapReduce: a survey. *SIGMOD Record*, 40(4):11–20, 2011.
- [25] Robert McNaughton and Seymour Papert. *Counter-Free Automata*. The MIT Press, Cambridge, Mass., 1971.
- [26] Charles Paperman. *Circuits booléens, prédicats modulaires et langages réguliers*. PhD dissertation, Université Paris Diderot, 2014. In French.
- [27] Thomas Place and Marc Zeitoun. Separating regular languages with first-order logic. In Thomas A. Henzinger and Dale Miller, editors, *CSL-LICS 2014*, pages 75:1–75:10. ACM, 2014.
- [28] Thomas Place and Marc Zeitoun. Separation and the successor relation. In Ernst W. Mayr and Nicolas Ollinger, editors, *STACS 2015*, volume 30 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 662–675, Dagstuhl, Germany, 2015. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- [29] Marcel-Paul Schützenberger. On finite monoids having only trivial subgroups. *Information and Control*, 8(2):190–194, 1965.
- [30] Luc Segoufin and Cristina Sirangelo. Constant-memory validation of streaming XML documents against DTDs. In Thomas Schwentick and Dan Suciu, editors, *ICDT 2007*, volume 4353 of *Lecture Notes in Computer Science*, pages 299–313. Springer, 2007.
- [31] Luc Segoufin and Victor Vianu. Validating streaming XML documents. In Lucian Popa, Serge Abiteboul, and Phokion G. Kolaitis, editors, *PODS 2002*, pages 53–64. ACM, 2002.
- [32] Howard Straubing. *Finite Automata, Formal Logic, and Circuit Complexity*. Birkhäuser, Boston, 1994.