

Generating, sampling and counting subclasses of regular tree languages *

Timos Antonopoulos
Hasselt University and
Transnational University of Limburg
timos.antonopoulos@uhasselt.be

Wim Martens
Technische Universität Dortmund
wim.martens@udo.edu

Floris Geerts
University of Edinburgh
fgeerts@inf.ed.ac.uk

Frank Neven
Hasselt University and
Transnational University of Limburg
frank.neven@uhasselt.be

Abstract

To experimentally validate learning and approximation algorithms for XML Schema Definitions (XSDs), we need algorithms to generate uniformly at random a corpus of XSDs as well as a similarity measure to compare how close the generated XSD resembles the target schema. In this paper, we provide the formal foundation for such a testbed. We adopt similarity measures based on counting the number of common and different trees in the two languages, and we develop the necessary machinery for computing them. We use the formalism of extended DTDs (EDTDs) to represent the unranked regular tree languages. In particular, we obtain an efficient algorithm to count the number of trees up to a certain size in an unambiguous EDTD. The latter class of unambiguous EDTDs encompasses the more familiar classes of single-type, restrained competition and bottom-up deterministic EDTDs. The single-type EDTDs correspond precisely to the core of XML Schema, while the others are strictly more expressive. We also show how constraints on the shape of allowed trees can be incorporated. As we make use of a translation into a well-known formalism for combinatorial specifications, we get for free a sampling procedure to draw members of any unambiguous EDTD. When dropping the restriction to unambiguous EDTDs, i.e. taking the full class of EDTDs into account, we show that the counting problem becomes $\#P$ -complete and provide an approximation algorithm. Finally, we discuss uniform generation of single-type EDTDs, i.e., the formal abstraction of XSDs. To this end, we provide an algorithm to generate k -occurrence automata (k -OAs) uniformly at random and show how this leads to uniform generation of single-type EDTDs.

1. Introduction

XML Schema is the accepted industry standard for the specification of schemas for collections of XML documents. At the same time, it is widely recognized that XML Schema

*We acknowledge the financial support of the Future and Emerging Technologies (FET) programme within the Seventh Framework Programme for Research of the European Commission, under the FET-Open grant agreement FOX, number FP7-ICT-233599.

is not a simple language. As it is very unlikely that the World Wide Web Consortium (W3C) will adopt a new schema standard any time soon, several initiatives have been taken to simplify XML Schema. For instance, algorithms have been developed to automatically infer XML Schema Definitions (XSDs) from XML data [8, 9, 10]. We later refer to this setting as the *learning scenario*. Another type of simplification is to let users design a schema in a different, but more user-friendly formalism and then offer the means to automatically convert this schema into an XSD. In general, the latter schema can not be equivalent but, hopefully, constitutes a best approximation in some well-defined way. The latter approach was taken in [18]. We later refer to this setting as the *approximation scenario*. In addition, algorithms to approximate non-deterministic content models by deterministic ones, hereby relieving the user from the Unique Particle Attribution constraint, are studied in [7].

Because it is not always possible to formally prove optimality of the above mentioned types of algorithms, their effectiveness is usually validated by an experimental study using real-world data, for instance using XSDs and corresponding XML corpora found on the web. Unfortunately, as real world data is often only sparsely available, ad-hoc methods are used to generate schemas and corresponding XML corpora. At the same time, a similarity measure is needed that quantifies how closely two unranked regular tree languages resemble each other, and which can be efficiently computed. The aim of this paper is to provide the machinery to efficiently compute the similarity between two tree languages and to provide algorithms to generate a corpus of XSDs uniformly at random. As usual, we use the abstraction of XSDs as single-type unranked regular tree languages [28, 30]. In particular, we consider the following three problems:

- (i) **Counting:** Given a tree language \mathcal{L} and $n \in \mathbb{N}$, compute the number of trees in \mathcal{L} of size n ;
- (ii) **Sampling:** Given a tree language \mathcal{L} , generate uniformly at random a tree $t \in \mathcal{L}$;
- (iii) **Generation:** Given a class of tree languages \mathcal{C} , generate uniformly at random a member $\mathcal{L} \in \mathcal{C}$.

We next provide further motivation and describe our contributions for each of these three problems.

Counting and Sampling. We start by discussing an approach towards a similarity measure for tree languages. To this end, let \mathcal{S} and \mathcal{T} be two tree languages. In the schema learning case described above, \mathcal{T} can be the target language and \mathcal{S} can be the schema inferred by the learning algorithm under consideration. Or, in the second scenario of schema approximation, \mathcal{T} can be the schema designed by the user and \mathcal{S} is an approximation of \mathcal{T} in a certain (simple) subclass of tree languages. This raises the natural question of how closely \mathcal{S} resembles \mathcal{T} . In this paper, we approach this problem by quantifying the number of common and different trees in \mathcal{S} and \mathcal{T} . For instance, one possibility is to define the *similarity* of \mathcal{S} and \mathcal{T} as

$$\text{sim}_{\leq n}(\mathcal{S}, \mathcal{T}) := \frac{\sum_{k=0}^n |(\mathcal{S} \cap \mathcal{T})^{=k}|}{\sum_{k=0}^n |(\mathcal{S} \cup \mathcal{T})^{=k}|},$$

where the set of trees of size k in a language \mathcal{L} is denoted by $\mathcal{L}^{=k}$, and the cardinality of $\mathcal{L}^{=k}$ is denoted by $|\mathcal{L}^{=k}|$. This similarity measure coincides with a measure commonly used when comparing regular *string* languages [7, 8, 9]. Furthermore, this measure has a natural probabilistic interpretation: the similarity between \mathcal{S} and \mathcal{T} is defined as (an approximation) of the expected probability that a tree, chosen uniformly at random from $\mathcal{S} \cup \mathcal{T}$, belongs to $\mathcal{S} \cap \mathcal{T}$. The approximation is realized by restricting attention to trees up to a certain size n . The algorithmic challenge is to efficiently compute $|\mathcal{L}^{=k}|$ for a tree language \mathcal{L} .

For string languages, when \mathcal{L} is represented by a deterministic finite automaton, the counting problem reduces to counting the number of accepting paths in a graph; an easy exercise in dynamic programming. However, when \mathcal{L} is represented by an NFA the problem becomes #P-complete [23]. We establish a similar dichotomy for *tree languages*.

Three classes of unranked regular tree languages are of immediate interest to us: single-type, restrained competition, and bottom-up deterministic EDTDs. Whereas single-type EDTDs correspond to the core of XML Schema [28, 30], restrained competition EDTDs correspond to EDTDs that can be correctly typed in a one-pass pre-order manner [28]. Both of these classes are deterministic in a top-down sense and are strict subclasses of the unranked regular tree languages. Moreover, the single-type EDTDs are known to be a strict subclass of the restrained-competition EDTDs [28]. The class of bottom-up deterministic EDTDs are deterministic in a bottom-up sense and correspond to the full class of unranked regular tree languages. We observe that while every restrained-competition EDTD is equivalent to a bottom-up deterministic EDTD, there is in general no efficient translation. Indeed, in some cases an exponential size increase can not be avoided.

In fact, we consider the class of *unambiguous* EDTDs in which any tree can have at most one valid typing. We observe that every single-type, restrained-competition and bottom-up deterministic EDTD is in effect an unambiguous EDTD. As a consequence, it suffices to develop counting and sampling algorithms for unambiguous EDTDs only. Rather than providing an ad-hoc dynamic programming solution to count the number of trees of a certain size in an unambiguous EDTD, we exhibit a mapping from the class of unambigu-

ous EDTDs into a (recursive) *combinatorial specification*. The latter is a formalism defined by Flajolet, Zimmermann and Van Cutsem [17] and provides an elegant way to derive counting and sampling algorithms. We show that in the case of unambiguous EDTDs, these algorithms are also efficient. In addition, we show how to incorporate *shape constraints* into combinatorial specifications. These are numerical constraints on the depth and width of trees in relation to the total size of the tree. For instance, to avoid string-like trees, we can restrict the depth of a tree to be at most logarithmic in the total number of nodes. In this way, the computation of the similarity of two tree languages can be restricted to trees of a certain shape (which is not necessarily regular).

Finally, when going beyond unambiguous EDTDs, the counting problem becomes intractable. That is, for general EDTDs, we show that computing the number of trees of a certain size is #P-complete. However, we do provide a pseudo-polynomial approximation algorithm based on a similar result for context-free grammars [20].

Generation. To assess the average behaviour of an algorithm, one can test it on a substantial input set drawn uniformly at random. This approach makes sense when no or few real-world data is available and opens up the possibility to quantify the quality of the obtained results in terms of confidence intervals.

In this paper, we consider the problem of generating XSDs uniformly at random. That is, for each n , every non-isomorphic XSD of size n must be generated with the same probability. This definition is the same as for the random generation of deterministic finite automata [2, 5]. Furthermore, since XSDs can be modelled as top-down DFAs that map states to content models [28, 26], we can extend methods for DFA generation to XSDs.

Unfortunately, current DFA generation methods do not constrain the occurrence of alphabet symbols, a constraint important for XSDs. Indeed, it has been noted in [8] that content models in XSDs contain large alphabets but every alphabet symbol occurs only a small number of times. We have referred to such expressions with alphabet symbol occurrence up to k as *k-OREs* (*k*-occurrence regular expressions) and to their automata counterparts as *k-OAs* (*k*-occurrence automata). In this paper, we provide an algorithm to generate uniformly at random deterministic *k-OAs* and show how this leads to uniform XSD generation.

Outline. In Section 2, we introduce the necessary definitions concerning automata, regular expressions and abstractions of XML schema languages. We study the counting and sampling problem for tree languages in Section 3 and 4, respectively. The uniform generation problem for XSDs is discussed in Section 5. Finally, Section 6 contains related work and the paper is concluded in Section 7. Proofs are referred to the Appendix.

2. Preliminaries

We define regular expressions, automata and XML Schema languages. First, we fix some basic notation.

2.1 String languages

Strings. For any two integers $n, m \in \mathbb{N}$ where $n \leq m$, we denote by $[n, m]$ the set of all the integers j such that

$n \leq j \leq m$. A *symbol* is an element of the alphabet Σ and a *string* w is a finite sequence of symbols $\sigma_1 \cdots \sigma_n$ for some $n \in \mathbb{N}$. We assume that the alphabet Σ is finite. We define the length of a string $w = \sigma_1 \cdots \sigma_n$, denoted by $|w|$, as n and we also refer to $|w|$ as the size of w . The empty string is denoted by ε . If w_1 and w_2 are two strings, we denote their concatenation by $w_1 \cdot w_2$ or simply by $w_1 w_2$. The set of all strings is denoted by Σ^* and a *string language* is a subset of Σ^* . If L_1 and L_2 are two string languages, their concatenation is defined as the set $\{w_1 w_2 \mid w_1 \in L_1, w_2 \in L_2\}$, and is denoted by $L_1 \cdot L_2$ or simply by $L_1 L_2$. For a string language L and for any $k \in \mathbb{N}$, we denote by $L^{=k}$ the set of strings in L that have length or size k , namely $L^{=k} = L \cap \Sigma^k$.

Automata. A *non-deterministic finite automaton* (NFA) A is a tuple $(\Sigma, Q, I, F, \delta)$, such that Q is a finite set of states, $I \subseteq Q$ is the set of initial states, F is the set of final states, and δ is the transition function of the automaton, defined as $\delta : Q \times \Sigma \rightarrow 2^Q$, mapping each pair of a state and symbol to a set of states. A run ρ of A on some string $w = a_1 \cdots a_n$ is a sequence of states q_0, \dots, q_n , such that $q_0 \in I$ and for each $i \in [1, n]$, $q_i \in \delta(q_{i-1}, a_i)$. Furthermore, when q_n is a member of F , we say that a run is accepting. The string language accepted by A is denoted by $\mathcal{L}(A)$ and is defined as the set of strings w for which there exists an accepting run of A on w . A non-deterministic finite automaton A is said to be *deterministic* (or A is a DFA) if the transition function maps each state/symbol-pair to a *singleton* set.

Regular expressions. The set of *regular expressions* (REs) over Σ is defined recursively as follows. The empty string ε and every symbol in Σ is a regular expression and if r_1 and r_2 are regular expressions, then so are $r_1 \cdot r_2$, $r_1 + r_2$, r^+ and r^* . The string language defined by a regular expression r , is denoted by $\mathcal{L}(r)$ and is defined as follows. If $r = \varepsilon$ then $\mathcal{L}(r) = \{\varepsilon\}$ and if $r = \sigma$ for some $\sigma \in \Sigma$, then $\mathcal{L}(r) = \{\sigma\}$. If $r = r_1 \cdot r_2$ then $\mathcal{L}(r) = \mathcal{L}(r_1)\mathcal{L}(r_2)$, if $r = r_1 + r_2$ then $\mathcal{L}(r) = \mathcal{L}(r_1) \cup \mathcal{L}(r_2)$, and finally if $r = r_1^+$ then $\mathcal{L}(r) = \{w \mid w = w_1 \cdots w_n \text{ for some } n \geq 1 \text{ and } \forall i \in [1, n], w_i \in \mathcal{L}(r_1)\}$. For any regular expression r , the regular expression r^* is equivalent to the regular expression $r^+ + \varepsilon$ and $r^?$ is used to abbreviate $r + \varepsilon$.

For any regular expression r , we denote by \bar{r} the regular expression obtained from r by replacing, for each i and each $a \in \Sigma$, the i -th occurrence of a by a_i . For example, if $r = ab^*(b + a^+)$, then $\bar{r} = a_1 b_1^+(b_2 + a_2^+)$. The XSD and DTD specifications (to be defined later) restrict regular expressions to be deterministic. A regular expression r is *deterministic* or *1-unambiguous* if there are no strings $w \cdot a_i \cdot v$ and $w \cdot a_j \cdot v'$ in $\mathcal{L}(\bar{r})$ such that $i \neq j$ [14]. We recall that a deterministic regular expression can be translated into an equivalent DFA in quadratic time [12].

2.2 XML schema languages

Trees. A set of strings S is *prefix closed* if for every string $s \in S$ and any prefix s_p of s , s_p is also in S . A *tree* t over an alphabet Σ is a tuple $(\text{Nodes}, \text{lab}, \Sigma)$ where Nodes , the set of nodes of t , is a finite prefix closed set of strings over the natural numbers, such that if $v \cdot i \in \text{Nodes}$ then $v \cdot i' \in \text{Nodes}$ for all $i' < i$, and $\text{lab} : \text{Nodes} \rightarrow \Sigma$ is a labelling function assigning symbols of Σ to each node in Nodes . The size of a tree equals its number of nodes. A node $v \in \text{Nodes}$ is a leaf node if there is no $v' \in \text{Nodes}$ such that v is a prefix of v' .

The root of t is the empty string in Nodes . The children of a node v in t are all nodes $v' \in \text{Nodes}$ such that $v' = v \cdot i$ for $i \in \mathbb{N}$. The subtree of a tree t at a node v of t is the set of nodes with prefix v . For the tree consisting of a single leaf node v labelled with the symbol σ , we write $\sigma(\varepsilon)$, and for any node v labelled with σ and having subtrees t_1, t_2, \dots, t_n rooted at its children, we write $\sigma(t_1, t_2, \dots, t_n)$, denoting the subtree of t rooted at v . The set of all trees over Σ is denoted by Trees_Σ and a tree language \mathcal{T} over Σ is a subset of Trees_Σ . The set of trees over Σ that have exactly k nodes is denoted by $\text{Trees}_\Sigma^{=k}$, for $k \in \mathbb{N}$. For a tree language \mathcal{T} , $\mathcal{T}^{=k}$ denotes the set of trees with m nodes, namely $\mathcal{T}^{=k} = \mathcal{T} \cap \text{Trees}_\Sigma^{=k}$.

DTDs and extended DTDs. A DTD over some finite alphabet Σ is a tuple $D = (\Sigma, R, d, S_d)$ where R is a set of deterministic regular expressions over Σ , d is a function that maps symbols in Σ to expressions in R , and $S_d \subseteq \Sigma$ is the set of start symbols. We refer to the regular expressions in R as the *content models* of the DTD. A finite tree t is *valid with respect to a DTD D* or *satisfies D* , if its root is labelled by an element of S_d and, for every node labelled with some $a \in \Sigma$, the sequence $a_1 \cdots a_n$ of labels of its children, is in the language defined by $d(a)$.

A DTD-DFA (Σ, A, d, S_d) over some finite alphabet Σ is a DTD whose content models are represented by the DFAs in the finite set A , instead of regular expressions.

An extended DTD (EDTD) over a finite alphabet Σ is a tuple $(\Sigma, \Delta, R, d, S_d, \mu)$, where Δ is a finite set of types, (Δ, R, d, S_d) is a DTD and μ is a mapping from Δ to Σ . A tree t is *valid with respect to an EDTD D* or *satisfies D* if $t = \mu(t')$ for some tree t' that satisfies the DTD (Δ, R, d, S_d) , where μ is extended to trees. We call t' a witness to t .

An EDTD-DFA $(\Sigma, \Delta, A, d, S_d, \mu)$ over a finite alphabet Σ is an EDTD where (Δ, A, d, S_d) is a DTD-DFA.

The tree language consisting of trees that are valid with respect to a DTD or EDTD D is denoted by $\mathcal{L}(D)$. An EDTD D is *reduced* if, for every type τ , there exists a witness tree $t \in \mathcal{L}(D)$ such that the label τ occurs somewhere in t . Any EDTD can be transformed to an equivalent reduced EDTD in polynomial time [1, 27]. In the following, we assume that all EDTDs are reduced.

Let D be an EDTD $(\Sigma, \Delta, d, S_d, \mu)$. Then, for any $\tau \in \Delta$, we denote by D_τ the EDTD $(\Sigma, \Delta, d, \{\tau\}, \mu)$. In particular the set of start symbols S_d of D is changed to $\{\tau\}$. We use the same notation for EDTD-DFAs.

Subclasses of EDTDs. We recall the following subclasses of EDTDs: *single-type EDTDs*, *restrained competition EDTDs*, and *bottom-up deterministic EDTDs*. Intuitively, these classes have the following significance. Single-type EDTDs are the formal abstraction of XSDs [28] and are therefore central in this paper. The class of restrained competition EDTDs corresponds to the EDTDs that can be correctly typed in a one-pass preorder manner [28]. This means that, when visiting the children of a node from left to right it is clear which type is associated with each node without looking ahead at the nodes to the right. Restrained competition EDTDs form a strict superclass of the single-type EDTDs. Finally, bottom-up deterministic EDTDs are a class of EDTDs that are equally expressive as general EDTDs, i.e., they recognize all regular tree languages. They correspond to bottom-up deterministic tree automata [13].

More formally, let $D = (\Sigma, \Delta, R, d, S_d, \mu)$ be an EDTD.

- D is *single-type* if S_d does not contain two conflicting types and no regular expression in R contains two conflicting types. Here, two types $\tau \neq \tau'$ *conflict* if $\mu(\tau) = \mu(\tau')$.
- D is *restrained competition* if S_d does not contain two conflicting types and all regular expressions in R restrain competition. Here, a regular expression r over Δ *restrains competition* if there are no strings $w\tau v$ and $w\tau'v'$ in $\mathcal{L}(r)$ with $\tau \neq \tau'$ and $\mu(\tau) = \mu(\tau')$.
- D is *bottom-up deterministic*, if for any two distinct types $\tau_1, \tau_2 \in \Delta$, it holds that $\mathcal{L}(d(\tau_1)) \cap \mathcal{L}(d(\tau_2)) = \emptyset$.

These notions are defined analogously for EDTD-DFAs. The class of all single-type (resp. restrained competition, bottom-up deterministic) EDTDs is denoted by EDTD^{st} (resp. EDTD^{rc} , EDTD^{bud}).

We note that translating between EDTDsts and EDTD^{bud}s gives rise to unavoidable exponential blow-ups. The following proposition holds for all formalisms used for representing content models of EDTDs in this paper.

PROPOSITION 2.1. There is a class $(D_n)_{n \in \mathbb{N}}$ of EDTDsts such that each D_n has size $O(n)$ and the smallest EDTD^{bud} for $\mathcal{L}(D_n)$ has size $2^{\Omega(n)}$. Likewise, there is a class $(D_n)_{n \in \mathbb{N}}$ of EDTD^{bud}s such that each D_n has size $O(n)$ and the smallest EDTDst and EDTD^{rc} for $\mathcal{L}(D_n)$ has size $2^{\Omega(n)}$. \square

To conclude this section, we next provide two examples of EDTDs that will also be used in Section 3.

EXAMPLE 2.2. Consider the EDTDst $D_1 = (\Sigma, \Delta, d, S_d, \mu)$ with $\Sigma = \{a\}$, $\Delta = \{\tau_o, \tau_e\}$, $d(\tau_e) = (\tau_o\tau_o)^*$, $d(\tau_o) = \tau_e(\tau_e\tau_e)^*$, $S_d = \{\tau_e\}$ and $\mu(\tau_o) = \mu(\tau_e) = a$. Then D_1 defines trees of even height where each node at even height has an even number of children and each node at odd height has an odd number of children. Here, the root has height 0. Let D_2 be the EDTD $(\Sigma, \Delta, d', S_d, \mu)$, where d' is such that $d'(\tau_o) = \tau_e(\tau_e\tau_e)^*$ and $d'(\tau_e) = (\tau_o\tau_o\tau_o\tau_o)^*$. Then, D_2 defines trees where a node at odd height has an odd number of children, but nodes at even height have 0 (mod 4) number of children. \square

2.3 Unambiguous EDTDs

We next define the class of *unambiguous* EDTDs and show that this class captures the single-type, restrained competition, and bottom-up deterministic EDTDs previously defined.

DEFINITION 2.3. An EDTD or EDTD-DFA is *unambiguous*, denoted by EDTD^{un} , if every tree $t \in \mathcal{L}(D)$ has a unique witness tree t' with $\mu(t') = t$. \square

In the remainder of the paper, we regularly use the following observation. The correctness of this observation immediately follows by contraposition.

OBSERVATION 2.4. If an EDTD or EDTD-DFA is unambiguous then, for all types $\tau \in \Delta$ and any two distinct trees t_1, t_2 over Δ , if t_1 and t_2 are witnesses to trees in $\mathcal{L}(D_\tau)$, then $\mu(t_1) \neq \mu(t_2)$. \square

PROPOSITION 2.5. Let $D = (\Sigma, \Delta, d, S_d, \mu)$ be an EDTD. If D is single-type, restrained competition, or bottom-up deterministic then D is unambiguous. \square

The following result readily follows from the standard product construction of automata (see, e.g., [18]). We add the observation that, if the input EDTDs are EDTD^{un}s, then the product EDTDs for the union and intersection are also EDTD^{un}s.

PROPOSITION 2.6. Let D_1 and D_2 be two EDTD-DFA^{un}s. Then we can construct, in quadratic time, an EDTD-DFA^{un} for $\mathcal{L}(D_1) \cup \mathcal{L}(D_2)$ and an EDTD-DFA^{un} for $\mathcal{L}(D_1) \cap \mathcal{L}(D_2)$. \square

Finally, we recall that deciding whether a given EDTD is in one of the particular classes we use here is in polynomial time.

PROPOSITION 2.7 ([28, 32]). Deciding whether a given EDTD is a EDTDst, EDTD^{rc}, EDTD^{bud}, or EDTD^{un} is in PTIME. \square

3. Counting Tree Languages

In this section, we consider the counting problem for tree languages $\mathcal{L}(D)$, where D is an EDTD. More specifically, we show that the counting problem can be efficiently solved when D is unambiguous, even when shape constraints are provided. These results imply that the similarity between unambiguous EDTDs can be efficiently computed. In contrast, the counting problem is shown to be #P-complete for general EDTDs. It does, however, allow for a randomised approximation scheme.

More formally, the general counting problem for languages can be stated as follows:

DEFINITION 3.1. For a class of languages \mathcal{C} , given a language $C \in \mathcal{C}$ and $m \in \mathbb{N}$, we define $\#C$ as the problem of finding the number of members in C of size m . \square

For instance, #DFA reduces to counting the number of paths in a graph, while #NFA is known to be #P-complete [23]. Here, we consider #EDTD and #EDTD-DFA^{un} and establish our results by exploiting the close relationship between (unambiguous) EDTDs and derivation trees of (unambiguous) context-free grammars. We first make this relationship precise.

3.1 From EDTDs to CFGs

Recall that a *context-free grammar* (CFG) G is a tuple (N, Σ, R, S) such that N is a finite set of non-terminal symbols, Σ is a set of terminal symbols, R is a subset of $N \times (N \cup \Sigma)^*$ and $S \in N$ is the start symbol. We denote the tuples $(V, w) \in R$ by $V \rightarrow w$.

A string $w \in (N \cup \Sigma)^*$ is derived from some non-terminal symbol $V \in N$, denoted by $V \Rightarrow w$, if $V \rightarrow w$. The transitive and reflexive closure of \Rightarrow is denoted by \Rightarrow^* . The language accepted by G is the set of strings $w \in \Sigma^*$ such that $S \Rightarrow^* w$, and is denoted by $\mathcal{L}(G)$. If $T \in N$, then $\mathcal{L}(G_T)$ is the set of strings w such that $T \Rightarrow^* w$. A context-free grammar G is *unambiguous* if, for every string $w \in \mathcal{L}(G)$, w has exactly one derivation tree for G .

It is well-known that regular expressions and (deterministic) finite automata can be translated to equivalent context-free grammars. Here, we use a variation of this result for which the actual translation is referred to the Appendix:

LEMMA 3.2. Let r be a regular expression and let A be a DFA over alphabet Δ . Then we can construct in linear time a CFG G with start symbol S such that $\tau_1 \cdots \tau_n \in \mathcal{L}(r)$ (resp. $\in \mathcal{L}(A)$) if and only if $S \Rightarrow^* \tau_1 \cdots \tau_n$ in G . Furthermore, if r (resp. A) is deterministic, then G is unambiguous. \square

In the remainder of this section, if r is a regular expression, we denote by $\text{CFG}(r, V)$ the set of CFG rules obtained by taking rules of G from Lemma 3.2, replacing the start symbol S by V , and replacing each terminal symbol $\tau \in \Delta$ in the derivation rules by a non-terminal T_τ . That way, we have that $\tau_1 \cdots \tau_n \in \mathcal{L}(r)$ if and only if $V \Rightarrow^* T_{\tau_1} \cdots T_{\tau_n}$ in $\text{CFG}(r, V)$. We define $\text{CFG}(A, V)$ for a DFA A similarly.

Let $D = (\Sigma, \Delta, X, d, S_d, \mu)$ be an EDTD, where X is a set of DFAs or REs. Let $R_{N, \Sigma}$ denote the class of context-free grammar rules over the set of terminal symbols Σ and the set of non-terminal symbols N . Let $\psi_D : \Delta \rightarrow \wp(R_{N, \Sigma'})$, for $\Sigma' = \Sigma \cup \{[,]\}$, be a function mapping types to sets of CFG rules, defined as:

$$\psi_D(\tau) = \{T_\tau \rightarrow \sigma[R_\tau]\} \cup \text{CFG}(d(\tau), R_\tau)$$

where $\sigma = \mu(\tau)$. In the following, we assume that the non-terminals in the rules $\text{CFG}(d(\tau), R_\tau)$ that are not of the form T_τ , for some $\tau \in \Delta$, are not used elsewhere. However, this can always be achieved by renaming non-terminals accordingly. Let Ψ_D be the set $\{\psi_D(\tau) \mid \tau \in \Delta\} \cup \{S \rightarrow T_\tau \mid \tau \in S_d\}$. Notice that, for each type τ there exists exactly one rule in Ψ_D whose left-hand side is T_τ .

The next lemma shows how EDTDs and CFGs are related:

LEMMA 3.3. For every EDTD $D = (\Sigma, \Delta, X, d, S_d, \mu)$ the CFG $G_D = (N, \Sigma \cup \{[,]\}, \Psi_D, S)$ is such that for all $n \in \mathbb{N}$, $|\mathcal{L}(D)^{=n}| = |\mathcal{L}(G_D)^{=3n}|$. Furthermore, if D is an EDTD-DFA^{un}, then G is unambiguous. \square

PROOF SKETCH. We show that $|\mathcal{L}(D)^{=n}| = |\mathcal{L}(G_D)^{=3n}|$ by establishing an isomorphism str between the languages $\mathcal{L}(D)$ and $\mathcal{L}(G_D)$. More precisely, $\text{str} : \text{Trees}_\Sigma \rightarrow (\Sigma \cup \{[,]\})^*$ is inductively defined as follows:

$$\begin{aligned} \text{str}(\sigma(\varepsilon)) &= \sigma \cdot [], \\ \text{str}(\sigma(t_1, \dots, t_m)) &= \sigma \cdot [\text{str}(t_1) \cdots \text{str}(t_m)]. \end{aligned}$$

It then suffices to show that for each n , $\text{str} : \text{Trees}_\Sigma^m \rightarrow (\Sigma \cup \{[,]\})^{3m}$ is injective. \square

3.2 #EDTD

It readily follows from the hardness of #NFA [23] that #EDTD is hard as well.

PROPOSITION 3.4. #EDTD is #P-complete. \square

The relationship between EDTDs and CFGs as specified in Lemma 3.3, however, can be used to provide a randomised approximation scheme for #EDTD.

We recall the notion of randomised approximation schemes

for languages from Gore et al. [20]. A *randomised approximation scheme for languages* is a randomised procedure that takes as input a description for a language $L \subseteq \Sigma^*$ and a tolerance $\varepsilon > 0$, and produces as output a number \hat{L} such that $(1 + \varepsilon)^{-1}|L| \leq \hat{L} \leq (1 + \varepsilon)|L|$ with probability at least $\frac{3}{4}$. For instance, an approximation scheme exists for #CFG:

THEOREM 3.5 ([20]). There is a randomized approximation scheme for #CFG, i.e., finding the number of elements of size m of a language defined by a given CFG G , with running time $\varepsilon^{-2}(m|G|)^{O(\log m)}$. \square

The approximation scheme for #EDTDs then immediately follows from Lemma 3.3 and Theorem 3.5.

COROLLARY 3.6. For an EDTD D , there is a randomized approximation scheme for finding the number of elements of size n of the language $\mathcal{L}(D)$, which runs in time $\varepsilon^{-2}(3n|D|)^{O(\log n)}$. \square

3.3 #EDTD-DFA^{un}

We again rely on the relationship between EDTD-DFA^{un}s and unambiguous CFGs as specified in Lemma 3.3. That is, we provide an efficient algorithm for $\#\mathcal{L}(G_D)$ (and hence for $\#\mathcal{L}(D)$) by formulating the CFG G_D as a so-called combinatorial specification [17], which we recall below. The advantage of such a specification is two-fold: (1) one can easily bound the complexity of computing the number of objects of a certain size; and (2) one obtains a general sampling procedure for objects in the specification (cf. Section 4)

Combinatorial specifications. A *combinatorial class* is a finite or denumerable set on which a size function is defined, satisfying the following two conditions:

- (i) the size of an element is a non-negative integer,
- (ii) the number of elements of any given size is finite.

If \mathcal{A} is a class, the size of an element $a \in \mathcal{A}$ is denoted by $|a|$. The set of objects in \mathcal{A} of size n is denoted by \mathcal{A}_n . The *counting sequence* of a combinatorial class is the sequence of integers $(A_n)_{n \geq 0}$ where $A_n = |\mathcal{A}_n|$ is the number of objects in class \mathcal{A} that have size n . Two combinatorial classes \mathcal{A} and \mathcal{B} are said to be *combinatorially isomorphic*, written $\mathcal{A} \cong \mathcal{B}$ if and only if their counting sequences are identical. This condition is equivalent to the existence of a bijection from \mathcal{A} to \mathcal{B} that preserves size.

A calculus for combinatorial classes introduced in [17], is presented below. Here, \mathcal{E} and \mathcal{Z} are atoms that denote the classes containing exactly one object of size 0 and size 1 respectively.¹ In the following, we allow different instantiations $\mathcal{Z}_a, \mathcal{Z}_b, \dots$ of the same atom \mathcal{Z} . Let \mathcal{B} and \mathcal{C} be combinatorial classes. Then the combinatorial class $\mathcal{A} = \mathcal{B} + \mathcal{C}$ is the disjoint union of the classes \mathcal{B} and \mathcal{C} . In particular, $\mathcal{Z} + \mathcal{Z}$ contains two objects of size 1. Furthermore, $\mathcal{A} = \mathcal{B} \times \mathcal{C}$ denotes the combinatorial class $\{\alpha = (\beta, \gamma) \mid \beta \in \mathcal{B}, \gamma \in \mathcal{C}\}$ and for each $\alpha = (\beta, \gamma) \in \mathcal{A}$, the size of α is the sum of the sizes of β and γ . For each $\alpha = (\beta_1, \dots, \beta_n) \in \mathcal{A}$, the size of α is the sum of the sizes of β_i for $i \in [1, n]$. For all n the following hold when \mathcal{A} is a combinatorial class:

$$\begin{aligned} \text{if } \mathcal{A} = \mathcal{B} + \mathcal{C} \text{ then } |A_n| &= |B_n| + |C_n|, \\ \text{if } \mathcal{A} = \mathcal{B} \times \mathcal{C} \text{ then } |A_n| &= \sum_{k=0}^n |B_{n-k}| \cdot |C_k| \end{aligned}$$

¹ \mathcal{E} is denoted as $\mathbf{1}$ in [17].

In the following, we assume an infinite set of variables C , C_0, C_1, \dots . Each variable will define a combinatorial class $\mathcal{L}(C)$.

DEFINITION 3.7 ([17]). A *specification* for (C_1, \dots, C_n) is a collection of n equations, with the i -th equation being of the form

$$C_i := \Psi_i(C_1, \dots, C_n)$$

where Ψ_i is a term built from \mathcal{E}, \mathcal{Z} and the C_j , using the constructors $+$ and \times . For each $j \in [1, n]$, let $C_j^0 = \emptyset$ and for each $i \in \mathbb{N}$, let $C_j^{i+1} = \Psi_j(C_1^i, \dots, C_n^i)$. Then $\mathcal{L}(C_j)$ is defined to be $\bigcup_{i \geq 0} \mathcal{L}(C_j^i)$. For $k \in \mathbb{N}$, we denote by $\mathcal{L}(C_j)^{=k}$ the objects in $\mathcal{L}(C_j)$ of size k . \square

We say that a specification is in *normal form* if each equation is either a single atom, or a single operation $C_i := C_j + C_k$ or $C_i := C_j \times C_k$.

THEOREM 3.8 ([17]). Given a specification for (C_1, \dots, C_n) in normal form and an integer k , the counting sequence up to size k can be computed in $O(n \cdot k^2)$ arithmetic operations. \square

From CFGs to combinatorial specifications. We know from Lemma 3.3 that for a given EDTD-DFA^{un} $D = (\Sigma, \Delta, X, d, S_d, \mu)$, the corresponding CFG $G_D = (N, \Sigma \cup \{[,]\}, \Psi_D, S)$ is unambiguous. Furthermore, it is well-known ([17]) that an unambiguous CFG in Chomsky normal form can be translated into a linear-size combinatorial specification (C_1, \dots, C_n) in normal form by simply replacing concatenation (\cdot) by \times , disjunction (\cup) of rules with the same left hand side by $+$, and finally each $\sigma \in \Sigma$ by \mathcal{Z}_σ . Hence, together with Theorem 3.8 we immediately get:

THEOREM 3.9. For an EDTD-DFA^{un} $D = (\Sigma, \Delta, X, d, S_d, \mu)$, the number of trees in $\mathcal{L}(D)$ of size up to k can be computed using $O(|\Delta| |Q_{\max}|^2 k^2)$ arithmetic operations, where Q_{\max} denotes the largest state space of an automaton in X . \square

PROOF. It suffices to show that the Chomsky normal form G'_D of the CFG $G_D = (N, \Sigma \cup \{[,]\}, \Psi_D, S)$ has at most $O(|\Delta| |Q_{\max}|^2)$ non-terminals, which we prove in the Appendix. \square

We stress that the size of the numbers $|\mathcal{L}(D)^{=k}|$ can grow very fast. To implement the algorithm of Theorem 3.9, a Mathematical Software package is needed. Actually, Maple provides an implementation in the `combstruct` module of the combinatorial specifications. We implemented our specification for the EDTD-DFA^{un} D_1 given in Example 2.2. As an illustration, we computed the number of trees of size 1001 valid with respect to D_1 and obtained a number with 314 decimals:

51879502371239317320511752369544517561698193655988840423158521214
8190894888949535843265681593434395020810002443582868233520387650
9254373728438806292876525845302947032070990934669778240958562432
2318852268438965431780372366645013594586870608079034900002010371
2015230396579554922650323287553303269884549851688819208474

The computation remains under the 60 seconds on a 1.8GHz iMac with 1GB of RAM.

3.4 Shape Constraints

Given an EDTD^{un} $D = (\Sigma, \Delta, A, d, S_d, \mu)$, it is often desirable to count the number of trees in $\mathcal{L}(D)$ that satisfy certain *shape constraints*. Here, by shape constraints we mean certain restrictions on the allowed combinations of the size, depth and/or width of trees in the language. More formally, a shape constraint on the depth (δ) (resp. branching width (w)) of trees consists of a function $\phi_\delta(k)$ (resp. $\phi_w(k)$) that assigns to each tree of size k its maximal allowed depth (resp. branching width). For instance, to avoid string-like trees one can take $\phi_\delta(k) = \log k$; to only consider binary trees one simply lets $\phi_w(k) = 2$. As previously described, the counting sequences of trees in $\mathcal{L}(D)$ can be computed using the combinatorial specification corresponding to the CFG G_D . In the presence of shape constraints, we need to augment this specification with parameters corresponding to the depth and width of objects.

We next describe in detail the specification for G_D : Let $D = (\Sigma, \Delta, A, d, S_d, \mu)$ be an EDTD-DFA^{un}. Let $A = \{A_\tau \mid \tau \in \Delta\}$ such that, for each $\tau \in \Delta$, $A_\tau = (\Delta, Q_\tau, \delta_\tau, q_{\tau,0}, F_\tau)$ is the DFA such that $d(\tau) = A_\tau$ and let $Q_\tau = \{q_{\tau,0}, \dots, q_{\tau,m_\tau}\}$. We assume w.l.o.g. that the Q_τ are pairwise disjoint and also disjoint with Δ . Also, let $\text{init} : \{A_\tau \mid \tau \in \Delta\} \rightarrow \{q_{\tau,0} \mid \tau \in \Delta\}$ be the function mapping each automaton to its initial state. Finally, we let $Q = \bigcup_{\tau \in \Delta} Q_\tau$.

Given the maximal tree depth d and width w , the specification is defined over the set of variables $\mathbf{Var}(d, w) = \{R_\tau^{\leq \delta, \leq w}, R_q^{\leq \delta, \leq w}, \mathcal{Z}_{\mu(\tau)} \mid \tau \in \Delta, q \in Q, \delta \in [1, d], w \in [0, w]\}$ and is given by the following set of equations:

$$\text{For } \tau \in \Delta, (\delta, w) \in \mathbb{N}^2, \delta \geq 1 :$$

$$T_\tau^{(\leq \delta, \leq w)} := \mathcal{Z}_{\mu(\tau)} \times \mathcal{Z}_\tau \times R_{\text{init}(d(\tau))}^{(\leq \delta-1, \leq w)} \times \mathcal{Z}_\tau.$$

$$\text{For } q \in Q, (\delta, w) \in \mathbb{N}^2, \delta \geq 1 :$$

$$R_q^{(\leq \delta, \leq w)} := \sum_{\substack{\tau \in \Delta, q' \in Q \\ q' \in \delta_{\tau'}(q, \tau)}} \left(T_\tau^{(\leq \delta, \leq w)} \times R_{q'}^{(\leq \delta, \leq w-1)} \right) \underbrace{+}_{\text{iff } q \in F_{\tau'}} \mathcal{E}.$$

We denote by $\mathcal{L}(D)^{=(k, \leq d, \leq w)}$ the set of trees of size k , maximal depth d and maximal width w . A straightforward generalization of Lemma 3.3 then shows that $|\mathcal{L}(D)^{=(k, \leq d, \leq w)}| = |(\mathcal{L}(\sum_{\tau \in S_d} T_\tau^{(\leq d, \leq w)}))^{=3k}|$. Hence, Theorem 3.9 implies:

COROLLARY 3.10. For an EDTD-DFA^{un} $D = (\Sigma, \Delta, X, d, S_d, \mu)$, size k , depth d and width w , the cardinality of $\mathcal{L}(D)^{=(k, \leq d, \leq w)}$ can be computed using $O(|\Delta| |Q_{\max}|^2 \cdot d \cdot w \cdot k^2)$ arithmetic operations, where Q_{\max} denotes the largest state space of an automaton in X . \square

PROOF. This follows immediately from Theorem 3.9. Indeed, it is easily verified that the normalization of the specification $\sum_{\tau \in S_d} T_\tau^{(\leq d, \leq w)}$ contains $O(|\Delta| |Q_{\max}|^2 \cdot d \cdot w)$ non-terminals. \square

We remark that when shape constraints $\phi_\delta(k)$ and $\phi_w(k)$ are provided, $|\mathcal{L}(D)^{=(k, \phi_\delta(k), \phi_w(k))}|$ is easily obtained from Corollary 3.10. Moreover, when only $\phi_\delta(k)$ or $\phi_w(k)$ is provided one simply removes the w or δ parameter, respectively, from the above specification and the complexity is adjusted correspondingly. Finally, we observe that when no shape

constraint is specified, the specification reduces to the one for G_D .

3.5 Similarity Measure

We return to computing the similarity between two tree languages as defined in the introduction using the machinery obtained above. Specifically, for tree languages S and T , define,

$$\text{sim}_{\leq n}(S, T) := \frac{\sum_{k=0}^n |(S \cap T)^{=k}|}{\sum_{k=0}^n |(S \cup T)^{=k}|},$$

where $\frac{0}{0}$ is taken to be 1.

Then we can prove the following result:

PROPOSITION 3.11. Assume S and T are specified as unambiguous EDTD-DFAs. Then for any n , $\text{sim}_{\leq n}(S, T)$ can be computed using $O(|\Delta| |Q_{\max, S}|^2 |Q_{\max, T}|^2 n^2)$ arithmetic operations, where $Q_{\max, S}$ and $Q_{\max, T}$ denote the largest state space of an automaton in S and T , respectively. \square

PROOF. Since $\text{sim}_{\leq n}(S, T)$ requires both $|(S \cap T)^{=k}|$ and $|(S \cup T)^{=k}|$, for $k \in [0..n]$, it suffices to bound the operations needed to compute these quantities. By Proposition 2.6, EDTD^{un}s can be computed for $S \cap T$ and $S \cup T$. Hence, all $|(S \cap T)^{=k}|$ for $k \in [0..n]$ can be computed from the specification of $S \cap T$ using $O(|\Delta| |Q_{\max, S}|^2 |Q_{\max, T}|^2 n^2)$ operations, where $Q_{\max, S}$ and $Q_{\max, T}$ denote the largest state space of an automaton in S and T , respectively. Indeed, this follows from Theorem 3.9 and the fact that the automata in $S \cap T$ consist of product automata of S and T . Due to trees common to S and T , we cannot use $S \cup T$. Instead, we simply use $|S^{=k}| + |T^{=k}| - |(S \cap T)^{=k}|$ for the counting sequence of the union of S and T . From Theorem 3.9 it follows again that these quantities can be computed up to $k = n$ using $O(|\Delta_S| |Q_{\max, S}|^2 n^2)$, $O(|\Delta_T| |Q_{\max, T}|^2 n^2)$ and $O(|\Delta| |Q_{\max, S}|^2 |Q_{\max, T}|^2 n^2)$ operations, respectively. As a consequence, $\text{sim}_{\leq n}(S, T)$ requires $O(|\Delta| |Q_{\max, S}|^2 |Q_{\max, T}|^2 n^2)$ operations. \square

To illustrate feasibility, we used our implementation in Maple to compute $\text{sim}_{\leq 100}(D_1, D_2) = 2.405906249 \cdot 10^{-7}$ taking D_1 and D_2 as defined in Example 2.2. The score was computed in as little as a few seconds.

The above definition of the function sim is just one possibility. A related but more general approach would be to consider a probability distribution p on the natural numbers, and define $\text{sim}_p(S, T)$ as

$$\sum_{n \geq 0} p(n) \frac{|(S \cap T)^{=n}|}{|(S \cup T)^{=n}|} = \sum_{t \in S \cap T} p(|t|) \frac{1}{|(S \cup T)^{=|t||}}, \quad (1)$$

where $|t|$ denotes the size of t . This means that $\text{sim}_p(S, T)$ is the expected probability that a tree t drawn from $S \cup T$ with probability $p(|t|)$ belongs to $S \cap T$. So, this measure assigns the same probability to trees of equal size.

4. Sampling tree languages

We next turn to the problem of sampling trees of a certain size in a tree language \mathcal{L} uniformly at random. The sampling procedure closely follows the general uniform sampling

methodology for combinatorial classes, as outlined in [17, 19, 25]. The general result for sampling objects of size k in a combinatorial class is as follows:

THEOREM 4.1 ([17]). Any combinatorial specification for (C_1, \dots, C_n) in normal form has a random generation routine that uses precomputed tables of size $O(nk)$ and achieves $O(nk \log k)$ worst case time complexity. The computation of the tables requires $O(nk^2)$ operations. \square

In other words, for any tree language \mathcal{L} that is equivalent to a combinatorial specification, one *automatically* obtains a sampling procedure. This is in particular true for EDTD- DFA^{un} s with and without shape constraints, as we have seen in the previous section.

5. Uniform XSD generation

In this section, we provide an algorithm to generate uniformly at random XSDs of a given size. A first step towards XSD generation is sampling of content models.

5.1 Generating Content Models

Almeida et al. provide a uniform sampler for deterministic connected complete DFAs through a string representation of the automata [2]. We next generalize their approach to deterministic k -OAs, by additionally allowing to parameterize on occurrences of alphabet symbols. While the underlying ideas remain the same, the new parameters introduce a higher level of complexity. Furthermore, we employ the formalism of combinatorial specifications to automatically obtain a generation procedure (cf Section 4).

k -Occurrence Automata As mentioned in the introduction, regular expressions in real-world XSDs can have large alphabets, but each of these alphabet symbols typically occurs only a small number of times.² A *k -occurrence regular expression* or *k -ORE* is a regular expression where every alphabet symbol occurs at most k times. For instance, $a \cdot (a + b)^*$ is a 2-ORE. However, we will not consider these regular expressions as such.

The reason is that there is little known on uniform regular expression generation. One approach could be to use for instance the context-free grammar of [24] defining all almost reduced regular expressions. However, this would not exclude different regular expressions defining the same language. Rejection sampling can not be used as there is no notion of minimal regular expression. Therefore, instead of using k -OREs, we turn to the corresponding (but slightly larger) class of k -occurrence automata as defined next. Note that these automata are node labelled.

DEFINITION 5.1. A *k -occurrence automaton* A (k -OA) over Σ is a tuple $(V, E, I, F, \text{lab}, e)$ where V is a finite set of states, $E \subseteq (V \times V)$ is the edge relation, $I \subseteq V$ is the initial set of states, F is the set of final states, $\text{lab} : V \rightarrow \Sigma$ is the labeling function, and e is a Boolean which is true when A accepts the empty string. We have the additional requirement that every Σ -symbol labels at most k states. \square

We say that two k -OAs $A_1 = (V_1, E_1, I_1, F_1, \text{lab}_1, e_1)$ and $A_2 = (V_2, E_2, I_2, F_2, \text{lab}_2, e_2)$ are isomorphic, if there exists a bijective function $\beta : V_1 \rightarrow V_2$ such that for all $v_1, v_2 \in V_1$,

²Actually, most alphabet symbols occur only once.

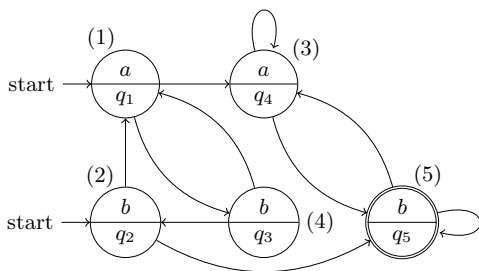


Figure 1: k -OA A

$(v_1, v_2) \in E_1$ if and only if $(\beta(v_1), \beta(v_2)) \in E_2$, $v_1 \in I_1$ if and only if $\beta(v_1) \in I_2$, $v_1 \in F_1$ if and only if $\beta(v_2) \in F_2$, $\text{lab}_1(v_1) = \text{lab}_2(\beta(v_1))$ and $e_1 = e_2$.

A k -OA A is *deterministic* if for every state $s \in V$ and $\sigma \in \Sigma$ there is at most one state $s' \in V$ such that $(s, s') \in E$ and $\text{lab}(s') = \sigma$, and there are no two distinct states $s, s' \in I$ with the same label. A k -OA A is *complete* if for every state $s \in V$ and every $\sigma \in \Sigma$, there exists a state $s' \in V$ such that $\text{lab}(s') = \sigma$ and $(s, s') \in E$, and for every label σ there is a state $s \in I$ with $\text{lab}(s) = \sigma$. Finally, a k -OA A is *connected*, if every state s is reachable from an initial state in I . We call a deterministic, complete and connected k -OA, *admissible*.

String encoding of k -OAs. We next provide a string representation of k -OAs that is inspired by [2]. If A is a deterministic k -OA over Σ with M states, then for every $s \in V$ we denote by $N(s)$ the neighbours of s , which are the states $s' \in V$ such that $(s, s') \in E$. A string encoding, similar to the one introduced in [2], is presented in what follows. This encoding is canonical in the sense that isomorphic automata have the same encoding. For an alphabet Σ of size ℓ , let $o : \Sigma \rightarrow [0, \ell - 1]$ be a total order over the symbols of Σ . Given that ordering o , define a canonical total order $c : V \rightarrow [1, M]$ as follows. First, for each state $s \in I$ let $c(s) = o(\text{lab}(s)) + 1$. Then, traverse the automaton in a breadth-first way, where at each state s , assign to each of the neighbours of s , that are not yet in the domain of c , and in the order induced by o , the smallest number $n \in [1, M]$ that has not been assigned to a state.

Consider the k -OA A shown in Fig. 1, over the alphabet $\Sigma = \{a, b\}$, where $o(a) = 0$ and $o(b) = 1$, and let $\varepsilon \in \mathcal{L}(A)$. Notice that A is both connected and complete. The canonical total order $c : V \rightarrow [1, M]$ is defined as follows. The initial state q_1 with label a is assigned the number 1, and the initial state q_2 is assigned the number 2. Now, traversing the automaton in a breadth-first order, the number 3 is assigned to the state q_4 , which is a neighbour of q_1 and has label a . Similarly, the state q_3 , which is a neighbour of q_1 with label b is assigned the number 4. Finally, proceeding to the neighbours of q_2 , the state q_5 is assigned the number 5 and all states are now ordered. The ordering of the states is annotated between parentheses in Fig. 1.

Given an admissible k -OA A over Σ with M states, the string encoding of A is a string $\text{enc}(A) = S_1 \cdot S_2 \cdot s$ of length $(M + 1) \cdot \ell + M + 1$, where S_1 is the substring encoding the transitions of the automaton and is of length $(M + 1) \cdot \ell$, S_2 is the substring encoding the set of final states and is of length M , and finally s is 1 if A accepts the empty string and 0 otherwise. The substring $S_1 = s_0 \cdots s_{(M+1) \cdot \ell - 1}$ is such that

for each $j \in [0, \ell - 1]$ $s_j = c(i)$ where $i \in I$ and $o(\text{lab}(i)) = j$, and for each $j' \in [1, M]$ and $j \in [0, \ell - 1]$, $s_{j' \cdot \ell + j}$ is equal to $c(v)$ where v is the state of A such that $\text{lab}(v) = o^{-1}(j)$ and $(c^{-1}(j'), v) \in E$. Informally, the latter denotes that $s_{j' \cdot \ell + j}$ is the number corresponding to the state reached from state with number j' reading the alphabet symbol $o^{-1}(j)$. For the substring $S_2 = s_{(M+1) \cdot \ell} \cdots s_{(M+1) \cdot \ell + M - 1}$ encoding the set of final states, for each $j \in [0, M - 1]$, $s_{(M+1) \cdot \ell + j}$ is 1 if the state $c^{-1}(j)$ is final and 0 otherwise.

For example, the string encoding for the k -OA A shown in Fig. 1, is the string

$$\underbrace{12}_{0} \underbrace{34}_{1} \underbrace{15}_{2} \underbrace{35}_{3} \underbrace{12}_{4} \underbrace{35}_{5} \underbrace{00001}_{\text{final}} \underbrace{1}_{\varepsilon}.$$

The first two digits of the string encode that state 1 (q_1) is the initial state with label a and that state 2 (q_2) is the initial state with label b . The next two digits encode the outgoing transitions of state 1. The outgoing a -transition goes to state 3 (q_4) and the outgoing b -transition goes to state 4 (q_3). The outgoing transitions for states 2–5 are encoded similarly.

The lemma below assures the correctness of the encoding:

- LEMMA 5.2. 1. The function c is a bijection for admissible k -OAs.
2. For two admissible k -OAs A_1 and A_2 with n states over an alphabet Σ of size ℓ , with a total order $o : \Sigma \rightarrow [0, \ell - 1]$ defined on this alphabet, if $\text{enc}(A_1) = \text{enc}(A_2)$ then A_1 is isomorphic to A_2 . \square

Characterisation of string encodings. We next provide a characterisation of strings that correspond to encodings of admissible k -OAs. We will leverage upon this characterisation and give a combinatorial specification for these string encodings later on.

For every string $s = s_0 \cdots s_{(n+1) \cdot \ell - 1}$, $n \in \mathbb{N}$, we let $(f_i)_{i \in [1, n]}$ be a sequence of numbers such that for each $i \in [1, n]$, f_i denotes the first position in the string s where the number i appears. Then consider the strings $s_0, \dots, s_{(n+1) \cdot \ell - 1}$ that satisfy the following 4 rules:

- (A1) $\forall i \in [1, n - 1], f_i < f_{i+1}$,
(A2) $\forall i \in [1, n], f_i < i \cdot \ell$,
(A3) $\forall i \in [0, \ell - 1], |\{s_p \mid p = i \pmod{\ell}\}| \leq k$,
(A4) $\forall i, i' \in [0, \ell - 1]$, where $i \neq i'$,
 $\{s_p \mid p = i \pmod{\ell}\} \cap \{s_p \mid p = i' \pmod{\ell}\} = \emptyset$.

Intuitively, the above rules express the following: Rule (A1) expresses that for each state i , the first time i appears in the string is before the first time state $i + 1$ appears in the string. Rule (A2) expresses that for each state i , the first occurrence of i is in the part of the string encoding the transitions of the first $i - 1$ states, which means that state i is reachable from a state $i' < i$. Rule (A3) expresses that for any symbol σ , there are at most k different states that can be reached by reading σ . Finally, rule (A4) expresses that each state has a unique label.

The precise relationship between the set of strings satisfying (A1) – (A4) and k -OAs is given by the following lemma:

LEMMA 5.3. For each $n \in \mathbb{N}$, enc is a bijection from the set of non-isomorphic admissible k -OAs with n states, over an alphabet Σ of size ℓ to the strings of size $(n+1) \cdot \ell + n + 1$ whose prefix $s = s_0, \dots, s_{(n+1) \cdot \ell - 1}$ satisfies the rules (A1) – (A4) above, and whose suffix $s' = s_{(n+1) \cdot \ell} \dots s_{(n+1) \cdot \ell + n}$ uses only 0 and 1. \square

In other words, Lemma 5.3 characterizes the strings corresponding to encodings of k -OAs. We next use this characterization to build a combinatorial specification for the set of strings encoding k -OAs.

Combinatorial specification of k -OAs. Let Σ be an alphabet of size ℓ and let A be an admissible k -OA over Σ with n states. Recall that $\text{enc}(A) = S_1 \cdot S_2 \cdot s$ is a string of length $(n+1)\ell + n + 1$, where S_1 is the substring of length $(n+1)\ell$ encoding the transitions of A , S_2 is the substring of length n encoding the set of final states, and finally s is 1 if A accepts the empty string and 0 otherwise. Furthermore, the prefix of S_1 of size ℓ is equal to $1 \cdot 2 \dots \ell$ and the suffix $S_2 \cdot s$ solely contains 0 or 1. Such fixed strings can easily be combinatorially specified. Indeed, for $i \in [1, n]$, let \mathcal{Z}_i denote the atom corresponding to state i . Then, $\mathcal{Z}_1 \times \dots \times \mathcal{Z}_\ell$ corresponds to the prefix $1 \cdot 2 \dots \ell$, whereas for $\mathcal{B} := \mathcal{Z}_0 + \mathcal{Z}_1$, we have that $\mathcal{B}^n \times \mathcal{B}$ corresponds to the set of all possible suffices of size $n+1$ in encodings of k -OAs.

Given these, it remains to specify the remaining symbols in S_1 . We need the following notation: For $m \in [\ell+1, n]$, let $\overline{W}^m = [W_0, \dots, W_{\ell-1}]$ be a partition of $[1, m]$ in which each part is of cardinality at most k . We denote by $\overline{W}_{m+1, j}^m$ the partition of $[1, m+1]$ obtained from \overline{W}^m by adding $\{m+1\}$ to W_j . Finally, if W is a subset of $[1, n]$ then \mathcal{W} denotes the specification $\sum_{i \in W} \mathcal{Z}_i$.

From Lemma 5.3 it follows that the class of strings corresponding to k -OAs of size n can be combinatorially specified as:

$$\begin{aligned} \mathcal{O}\mathcal{A}_n := & \mathcal{Z}_1 \times \dots \times \mathcal{Z}_\ell \times \\ & \left(\sum_{p=\ell}^{\ell^2 + (\ell-1)p - \ell - 1} \left(\prod_{i=0}^{p-\ell-1} \mathcal{Z}_{i \pmod{\ell}} \right) \times \mathcal{Z}_{\ell+1} \times S_{\ell+1}^{(j', j)} [\overline{W}_{\ell+1, j \pmod{\ell}}^\ell] \right) \\ & \times \mathcal{B}^M \times \mathcal{B}, \end{aligned}$$

where j' and j are such that $p = j'\ell + j$, provided that the class $S_{\ell+1}^{(j', j)} [\overline{W}_{\ell+1, j \pmod{\ell}}^\ell]$ consists of all strings of length $n \cdot \ell$ satisfying rules (A1) – (A4), with a fixed prefix $s_0 \dots s_{j'\ell + j}$, where $\ell+1$ occurs in the string at position $j'\ell + j$ for the first time, and for each $i \in [0, \ell-1]$, and $q \in [0, j'\ell + j]$, it holds that $\{s_q \mid q = i \pmod{\ell}\} = W_i$.

We next claim that the specification given in Figure 2 defines precisely this class of strings. Intuitively, for $\ell < m \leq n$, the equations in the specification reflect that the strings satisfying the rules (A1) – (A4) and in which letter m appears for the first time at position $j'\ell + j$ is equal to the (disjoint) union of possible strings where $m+1$ appears for the first time in one of the positions between $j'\ell + j + 1$ and $m \cdot \ell - 1$. In the specification this union is partitioned as follows: The class Q_1 covers the case where between the positions $j'\ell + j$ and $(j'+1) \cdot \ell + j$ no new state $m+1$ is introduced, whereas the class Q_2 covers the cases where at some position J between the positions $j'\ell + j$ and $(j'+1) \cdot \ell + j$, the state $m+1$ is introduced. Finally, two cases need special attention: First,

in case that $j' = m - 1$, the rule (A2) implies that symbol m must appear directly afterwards. In the specification this is reflected by the class defined by Q_3 . Second, when $m = n$ and $j' = n - 1$, the end of string is nearby and the recursive specification must end. The last equation in the specification deals with this case. Indeed, it simply pads the string with symbols until the desired length of $n\ell$ is reached.

THEOREM 5.4. For $\ell, k, n, j', j, n_0, \dots, n_{\ell-1} \in \mathbb{N}$ and $m \geq \ell$, the class $\mathcal{S}_m^{(j', j)} [W_0, \dots, W_{\ell-1}]$ defined by the specification in Figure 2 corresponds to the class of strings of length $(n+1) \cdot \ell$ satisfying rules (A1) – (A4), with a fixed prefix $s_0 \dots s_{j'\ell + j}$, where m occurs in the string at position $j'\ell + j$ for the first time, and for each $i \in [0, \ell-1]$, and $q \in [0, j'\ell + j]$, it holds that $\{s_q \mid q = i \pmod{\ell}\} = W_i$. \square

As a consequence, the specification $\mathcal{O}\mathcal{A}_n$ consists of all strings that correspond to the encoding of a k -OA of size n .

LEMMA 5.5. Let $n, k, \ell \in \mathbb{N}$. The number of operations required to compute $|\mathcal{O}\mathcal{A}_n|$ is bounded by $O(n^2 \cdot \ell^2 \cdot k^\ell)$. \square

PROOF. This is a direct consequence of Theorem 3.8 and the fact that normal form of the specification $\mathcal{O}\mathcal{A}_n$ has at most $O(n^2 \cdot \ell^2 \cdot k^\ell)$ non-terminals. \square

Uniform generation procedure for k -OAs. Given the specification for k -OAs one can rely on Theorem 4.1 to obtain a uniform random generation procedure. This procedure, referred to as k -OA-uniform-generation(n, ℓ), takes as input the size ℓ of the alphabet and the desired number of states n , and returns uniformly at random (with probability $1/|\mathcal{O}\mathcal{A}_n|$) a string $\text{enc}(A)$ encoding a k -OA A over Σ and n states.

From Theorem 4.1 and Lemma 5.5 it follows immediately that k -OA-uniform-generation(n, ℓ) requires $O(n^2 \cdot \ell^2 \cdot k^\ell n \log(n))$ worst case time complexity.

Finally, we observe that an alternative generation procedure for k -OAs, similar to the one presented in [2], can be devised. The difference is that while generating the string encoding of the admissible k -OAs, one needs to additionally keep track of which alphabet symbols have been assigned to which states, so that no symbol is assigned to more than k states.

5.2 XSD generation

In what follows, we present an algorithm that generates uniformly at random an XSD, whose content models are described by k -OAs, for some $k \in \mathbb{N}$. We make use of the notion of DFA-based XSDs (which is equivalent to XSDs and introduced in [26]) with k -OA content models.

DEFINITION 5.6. A DFA-based k -OA XSD is a tuple $D = (\Sigma, A, \lambda)$, where A is a DFA with set of states Δ over the alphabet Σ , and λ is a function mapping each non-initial state q of A to some k -OA over Σ . A tree t is valid with respect to D if for every node v of t , the sequence $a_1 \dots a_n$ of labels of its children is in $\mathcal{L}(\lambda(q))$, where q is the state reached by A when started in its initial state, and by reading the string of labels on the path from the root to v . \square

It therefore suffices to generate DFA-based k -OA XSDs.

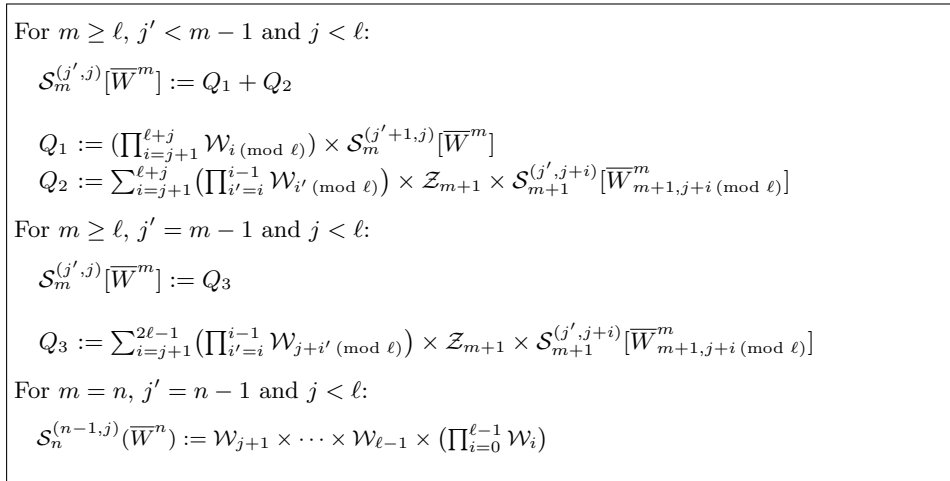
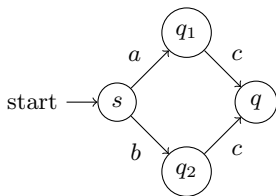


Figure 2: Combinatorial specification of k -OAs.

EXAMPLE 5.7. We provide an example of a DFA-based k -OA XSD. Consider the following DFA without final states:



Furthermore, we define $\lambda(q_1) = k\text{OA}(c)$, $\lambda(q_2) = k\text{OA}(cc)$, and $\lambda(q) = k\text{OA}(\varepsilon)$, where $k\text{OA}(r)$ denotes the minimal k -OA for the regular expression r . This DFA-based k -OA XSD accepts two trees, namely $a(c)$ and $b(c, c)$. \square

We next discuss the algorithm we need for generating the DFA of a DFA-based k -OA XSD. Example 5.7 illustrates that we need to generate possibly incomplete connected DFAs. Indeed, we cannot simply resort to methods generating minimal DFAs, since the DFA in Example 5.7, having no final states, is obviously not minimal. Furthermore, even when we would define the states that can be reached when reading an entire path from root to leaf in a tree of the language to be final, the resulting DFA is not minimal. Indeed, this approach would define state q to be a final state, but then minimization of the DFA in Example 5.7 would merge states q_1 and q_2 , thereby changing the language of the DFA-based XSD. In fact, it can be shown that the DFA-based XSD from Example 5.7 is minimal [29].

This justifies our choice for a method that generates possibly incomplete connected DFAs. More precisely, for $m \in \mathbb{N}$ and any ordered alphabet Σ where $|\Sigma| = \ell$, we consider the algorithm `generate-DFA`(m, ℓ) from Bassino et al. that generates a possibly incomplete connected DFA with m states (and no final states) over the alphabet Σ [4]. Note that the generated DFAs are not necessarily minimal.

We next show to generate the k -OAs. Whereas we already know to generate k -OAs of a given size, we now need to uniformly range over the sizes of the k -OAs. Let A be a k -OA over an alphabet Σ such that $|\Sigma| = \ell$. Then notice that the number of states for A can only be larger than or equal

to ℓ and smaller than or equal to $k \cdot \ell$. For any two values n_1, n_2 such that $n_1, n_2 \in [\ell, k \cdot \ell]$, the number of k -OAs with n_1 states and the number of k -OAs with n_2 states is not necessarily the same, and therefore, to choose uniformly at random the number of states of a k -OA, these numbers have to be taken into account. The algorithm `generate-k-OA-num-states` relies on the specification \mathcal{OA}_i for k -OAs in order to compute the number of k -OAs with i states over Σ . We next use these numbers to produce a uniform distribution over the number of states for the automaton.

Algorithm 1 `generate-k-OA-num-states`

Input: ℓ
 $\max = \sum_{i=1}^{k \cdot \ell} |\mathcal{OA}_i|;$
return $j \in [1, k \cdot \ell]$ with probability $\frac{|\mathcal{OA}_j|}{\max}$

Algorithm 2 `generate-k-OA-XSD`

Input: m, ℓ
Output: DFA-based XSD (A, λ)
 $A = \text{generate-DFA}(m, \ell);$
for every state $j \in [2, m]$ of A **do**
 $\ell_j = |\{\sigma \mid \delta_A(j, \sigma) \text{ is defined in } A\}|$
 $n_j = \text{generate-k-OA-num-states}(\ell_j);$
end for
for every state $j \in [2, m]$ of A **do**
repeat
 $\lambda(j) = k\text{-OA-uniform-generation}(n_j, \ell_j);$
until `is-minimal-complete-k-OA`($\lambda(j)$)
end for

Consider then Algorithm 2. The algorithm `generate-k-OA-XSD`, given m and ℓ , first generates a (possibly incomplete) DFA A with m states over the alphabet $[0, \ell - 1]$ uniformly at random. We assume that A 's states are numbered from 1 to m and its initial state is 1. Then, for each $j \in [2, m]$, i.e., for each state apart from the initial one, we compute the alphabet size ℓ_j of the k -OA that will be associated to state j . This alphabet size must correspond precisely to the number of outgoing transitions of state j (see Example 5.7) in order for the DFA-based k -OA XSD to be well-defined.

Then we choose the number of states n_j of a k -OA for state j uniformly at random using `generate- k -OA-num-states`. Finally, we generate a k -OA for state j of A with n_j number of states and alphabet size ℓ_j using `k -OA-uniform-generation`.

Our aim for Algorithm 2 is to produce, uniformly at random, non-isomorphic *admissible* DFA-based k -OA XSDs, that is, DFA-based XSDs in which the inner DFA is connected and the k -OAs are admissible and minimal.

We note that testing minimality of a deterministic k -OA can be done by a simple adaptation of the standard DFA minimization algorithms. In particular, only states bearing the same alphabet symbol are allowed to be merged.

PROPOSITION 5.8. Testing minimality of a deterministic k -OA can be done in PTIME. \square

Therefore, the test `is-minimal-complete- k -OA($s_{\text{start}}, \text{end}$)` runs in polynomial time.

The above discussion leads to the following result.

THEOREM 5.9. Algorithm `generate- k -OA-XSD` generates, given m and ℓ , uniformly at random an admissible DFA-based k -OA XSD with m types and alphabet size ℓ . \square

Finally, we note that the DFA-based k -OA XSDs generated by `generate- k -OA-XSD` are not necessarily minimal and there can be two such automata A_1 and A_2 generated by the algorithm that generate the same language. We can, therefore, go through a rejection stage that checks if the resulting automaton is minimal, and keeps generating such automata until a minimal one is found. This test whether the generated DFA-based k -OA XSD is minimal can be performed in polynomial time using Proposition 5.8 and the minimization algorithm for XSDs from [29].

As a final remark, we note that generating EDTD^{un}s in a similar manner is computationally harder. Indeed, as an immediate consequence from [11], already testing minimality of a EDTD^{un} is coNP-complete.

6. Related Work

Sampling. Our approach towards counting of tree languages is based on the recursive method, which was initiated by Nijenhuis and Wilf [31], and then formalized by Flajolet, Zimmermann and Van Cutsem [17] in the more general setting of combinatorial specifications. In the current paper we only use a restricted class of specifications (called context-free) in that only atoms, union and product are allowed. General recursive specifications allow many more operators (cf. [17]).

The computational complexity of variants of computing the number of strings of given length in context-free languages is investigated by Bertoni et al. [6]. We choose to employ the method of going through the combinatorial specification as it gives rise to an immediate implementation in Maple and is versatile enough to incorporate shape constraints (which are not context-free definable).

Sampling of trees that adhere to a probabilistic tree model is investigated in [15, 16]. In particular in [16] a sampling procedure for trees that additionally satisfy a bottom-up tree automaton is provided. These methods differ from ours in

that trees are sampled in accordance with their probabilities specified by the probabilistic model rather than uniformly.

XSD generation. There has been substantial work on the uniform generation of regular languages represented by DFAs. To the best of our knowledge, no algorithm exists that uniformly generates *minimal* DFAs. The works [2, 5] and [4] do consider the non-isomorphic uniform generation of admissible and connected DFAs, respectively, but need rejection sampling to sample minimal DFAs. This approach, however, has no proven guarantees on running times.

Our string encoding for k -OAs is inspired by [2]. For the generation procedure, however, we rely on a sampling procedure for combinatorial specifications [17]. Although Héam, Nicaud and Sylvain [21] show that non-isomorphic deterministic tree-walking automata can be generated uniformly at random through an encoding into string transducers, it is not clear how to use their results to generate XSDs.

To the best of our knowledge, the present paper presents the first step uniform XSD generation. The papers [9, 8] only dealt with DTDs which reduce to regular expressions. In [10], the experimental validation used one real-world XSD and 8 hand-crafted XSDs. The XSD generation algorithm presented in this paper could be used to generate a benchmark of XSDs. We did not address generation of XML corpora adhering a given schema as is for instance implemented in ToXGene [3].

7. Conclusion

In this paper, we presented a first step towards the foundation for an experimental testbed for XSD generating algorithms. We addressed uniform XSD generation as well as the machinery to compute similarity measures based on counting of trees of a certain size in tree languages. Finally, we provided a sampling procedure for (unambiguous) tree languages using the formalism of combinatorial specifications.

An initial implementation in Maple shows that the approach through combinatorial specifications is promising. Although the approach to assess similarity through counting of the number of different and common trees is intuitive, in depth experimental validation of efficiency and effectiveness remains needed to obtain a concrete robust similarity measure.

Directions for future work include the following:

The complexity of the generation algorithm for k -OAs reveals an exponential behavior in the size of the alphabet. The main reason is that implicitly for each alphabet symbol it needs to be remembered how many times it has already occurred. Fortunately, in real-world content models the far majority of the symbols occur only once. It would be interesting to see how this constraint can be incorporated into the algorithm.

Furthermore, it would be most interesting to extend the XSD generation algorithm to generate k -OREs rather than k -OAs. This would require a useful canonical representation for regular expressions.

Finally, we want to explore the possibility of using specifications (possibly extended with probabilities) to get a non-uniform sampling procedures of trees in a tree language. This would be particularly useful in the probabilistic XML setting, among others.

8. References

- [1] J. Albert, D. Giammerresi, and D. Wood. Normal form algorithms for extended context free grammars. *Theoretical Computer Science*, 267(1–2):35–47, 2001.
- [2] M. Almeida, N. Moreira, and R. Reis. Enumeration and generation with a string automata representation. *Theoretical Computer Science*, 387(2):93–102, 2007.
- [3] D. Barbosa, A. O. Mendelzon, J. Keenleyside, and K. A. Lyons. ToXgene: a template-based data generator for XML. In *International Symposium on Management of Data (SIGMOD)*, page 616, 2002.
- [4] F. Bassino, J. David, and C. Nicaud. Enumeration and random generation of possibly incomplete deterministic automata. *Pure Mathematics and Applications*, 19(2–3):1–16, 2008.
- [5] F. Bassino and C. Nicaud. Enumeration and random generation of accessible automata. *Theoretical Computer Science*, 381(1–3):86–104, 2007.
- [6] A. Bertoni, M. Goldwurm, and N. Sabadini. The complexity of computing the number of strings of given length in context-free languages. *Theoretical Computer Science*, 86(2):325–342, 1991.
- [7] G. J. Bex, W. Gelade, W. Martens, and F. Neven. Simplifying XML Schema: effortless handling of nondeterministic regular expressions. In *International Symposium on Management of Data (SIGMOD)*, pages 731–744, 2009.
- [8] G. J. Bex, W. Gelade, F. Neven, and S. Vansummeren. Learning deterministic regular expressions for the inference of schemas from XML data. In *International World Wide Web Conference (WWW)*, pages 825–834, 2008.
- [9] G. J. Bex, F. Neven, T. Schwentick, and S. Vansummeren. Inference of concise regular expressions and DTDs. *ACM Transactions on Database Systems*, 2010.
- [10] G. J. Bex, F. Neven, and S. Vansummeren. Inferring XML Schema Definitions from XML data. In *International Conference on Very Large Data Bases (VLDB)*, pages 998–1009, 2007.
- [11] H. Björklund and W. Martens. The tractability frontier for NFA minimization. In *International Colloquium on Automata, Languages and Programming (ICALP)*, pages 27–38, 2008.
- [12] A. Brüggemann-Klein. Regular expressions into finite automata. In *Latin American Symposium on Theoretical Informatics (LATIN)*, pages 87–98, 1992.
- [13] A. Brüggemann-Klein, M. Murata, and D. Wood. Regular tree and regular hedge languages over unranked alphabets: Version 1, april 3, 2001. Technical Report HKUST-TCSC-2001-0, The Hongkong University of Science and Technology, 2001.
- [14] A. Brüggemann-Klein and D. Wood. One-unambiguous regular languages. *Information and Computation*, 142(2):182–206, 1998.
- [15] S. Cohen, B. Kimelfeld, and Y. Sagiv. Incorporating constraints in probabilistic XML. *ACM Transactions on Database Systems*, 34(3):1–45, 2009.
- [16] S. Cohen, B. Kimelfeld, and Y. Sagiv. Running tree automata on probabilistic XML. In *International Symposium on Principles of Database Systems (PODS)*, pages 227–236, 2009.
- [17] P. Flajolet, P. Zimmermann, and B. Van Cutsem. A calculus for the random generation of labelled combinatorial structures. *Theoretical Computer Science*, 132(2):1–35, 1994.
- [18] W. Gelade, T. Idziaszek, W. Martens, and F. Neven. Simplifying XML Schema: Single-type approximations of regular tree languages. In *International Symposium on Principles of Database Systems (PODS)*, 2010.
- [19] M. Goldwurm. Random generation of words in an algebraic language in linear binary space. *Information Processing Letters*, 54:229–233, 1995.
- [20] V. Gore, M. Jerrum, S. Kannan, Z. Sweedyk, and S. R. Mahaney. A quasi-polynomial-time algorithm for sampling words from a context-free language. *Information and Computation*, 134(1):59–74, 1997.
- [21] P.-C. Héam, C. Nicaud, and S. Schmitz. Random generation of deterministic tree (walking) automata. In *International Conference on Implementation and Application of Automata (CIAA)*, pages 115–124, 2009.
- [22] J.E. Hopcroft, R. Motwani, and J.D. Ullman. *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley, 3 edition, 2007.
- [23] S. Kannan, Z. Sweedyk, and S. R. Mahaney. Counting and random generation of strings in regular languages. In *ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 551–557, 1995.
- [24] J. Lee and J. Shallit. Enumerating regular expressions and their languages. In *International Conference on Implementation and Application of Automata (CIAA)*, pages 2–22, 2004.
- [25] H. G. Mairson. Generating words in a context-free language uniformly at random. *Information Processing Letters*, 49(2):95–99, 1994.
- [26] W. Martens, F. Neven, and T. Schwentick. Simple off the shelf abstractions of XML Schema. *Sigmod RECORD*, 36(3):15–22, 2007.
- [27] W. Martens, F. Neven, and T. Schwentick. Complexity of decision problems for XML schemas and chain regular expressions. *SIAM Journal on Computing*, 39(4):1486–1530, 2009.
- [28] W. Martens, F. Neven, T. Schwentick, and G.J. Bex. Expressiveness and complexity of XML Schema. *ACM Transactions on Database Systems*, 31(3):770–813, 2006.
- [29] W. Martens and J. Niehren. On the minimization of XML Schemas and tree automata for unranked trees. *Journal of Computer and System Sciences*, 73(4):550–583, 2007.
- [30] M. Murata, D. Lee, M. Mani, and K. Kawaguchi. Taxonomy of XML schema languages using formal language theory. *ACM Transactions on Internet Technology*, 5(4):660–704, 2005.
- [31] A. Nijenhuis and H. Wilf. *Combinatorial algorithms*. Academic Press Inc., 1979.
- [32] H. Seidl. Deciding equivalence of finite tree automata. *SIAM Journal on Computing*, 19(3):424–437, 1990.

A. Appendix

Proofs for Section 2

PROPOSITION 2.1. *There is a class $(D_n)_{n \in \mathbb{N}}$ of EDTDsts such that each D_n has size $O(n)$ and the smallest EDTD^{bud} for $\mathcal{L}(D_n)$ has size $2^{\Omega(n)}$. Likewise, there is a class $(D_n)_{n \in \mathbb{N}}$ of EDTD^{bud}s such that each D_n has size $O(n)$ and the smallest EDTDst and EDTD^{rc} for $\mathcal{L}(D_n)$ has size $2^{\Omega(n)}$.*

PROOF. The class $(D_n)_{n \in \mathbb{N}}$ of EDTDsts defines the unary trees that, when read as a string from root to leaf, obey the regular expression $(a+b)^n a(a+b)^*$. The fact that the minimal EDTD^{bud} for $\mathcal{L}(D_n)$ has size $2^{\Omega(n)}$ immediately follows from the fact that the smallest DFA for $(a+b)^n a(a+b)^*$ has size $2^{\Omega(n)}$. The direction from EDTD^{bud}s to EDTDsts is analogous. \square

PROPOSITION 2.5. *Let $D = (\Sigma, \Delta, d, S_d, \mu)$ be an EDTD. If D is single-type, restrained competition, or bottom-up deterministic then D is unambiguous.*

PROOF. If D is single-type, then D is also restrained competition. Therefore, let us first assume that D is restrained competition. We prove that D is also unambiguous. Towards a contradiction, assume that D is not unambiguous. Then there is a type $\tau \in \Delta$ and two distinct trees t_1, t_2 over Δ that are witnesses to some tree $t \in \mathcal{L}(D_\tau)$ and are such that $\mu(t_1) = \mu(t_2)$. Let v be a node in t such that the type of v in t_1 , say τ_1 , is different from the type of v in t_2 , say τ_2 . Moreover, let v be such that none of its ancestors or left siblings have this property. Notice that v cannot be the root of t since the type of the root is τ in both t_1 and t_2 . Let v' be the parent of v in t , and let its type (which is the same both in t_1 and t_2) be τ' . Then the child-string of v' in t_1 is of the form $w\tau_1 u_1$ and the child-string of v' in t_2 is of the form $w\tau_2 u_2$ by the assumption that v is the leftmost child with this above property. Therefore, the regular language associated with τ' does not restrain competition, which is a contradiction.

Finally, assume that D is bottom-up deterministic. Towards a contradiction, assume D is not unambiguous. Then there is a type $\tau \in \Delta$ and two distinct trees t_1, t_2 over Δ that are witnesses to some tree $t \in \mathcal{L}(D_\tau)$ and are such that $\mu(t_1) = \mu(t_2)$. Notice that t_1, t_2 and t are the same trees with different type-labeling. Let v be a node in t such that the type τ_1 of v in t_1 is different from the type τ_2 of v in t_2 , and such that none of v 's descendants have this property. Since none of v 's descendants have this property, v 's child string is the same in t_1 and t_2 . Denote this child string by $s \in \Delta^*$. However, this means that $s \in \mathcal{L}(d(\tau_1)) \cap \mathcal{L}(d(\tau_2))$, which contradicts that D is bottom-up deterministic. \square

PROPOSITION 2.7. *Deciding whether a given EDTD is a EDTDst, EDTD^{rc}, EDTD^{bud}, or EDTD^{un} is in PTIME.*

PROOF. The result for EDTDst and EDTD^{rc} is proved in [28]. Testing whether an EDTD is a EDTD^{bud} simply boils down to testing emptiness of finite string automata, which is in PTIME.

Finally, Seidl has shown that testing m -ambiguity of a ranked tree automaton is in PTIME, for each fixed m [32]. When translating an EDTD to a ranked tree automaton using a variant of the standard first-child next-sibling encoding (in polynomial time), the resulting ranked tree automaton

is unambiguous (i.e., 1-ambiguous in Seidl's terminology) if and only if the original EDTD was unambiguous. Therefore, testing whether an EDTD is an EDTD^{un} is also in PTIME. \square

Proofs for Section 3.1

Let \mathcal{REG}_Σ denote the class of regular expressions over the alphabet Σ , and let N be an infinite set of non-terminal symbols. Also, let $R_{N,\Sigma}$ denote the class of context-free grammar rules over the set of terminal symbols Σ and the set of non-terminal symbols N . We define the function $\varphi : N \times \mathcal{REG}_\Sigma \rightarrow \varphi(R_{N,\Sigma})$ to be a function mapping each regular expression to a set of context-free grammar rules. To this end, we slightly adapt a translation given by Hopcroft et al. [22]. The function is defined inductively as follows:

$$\begin{aligned} \varphi(R, r) &= \{R \rightarrow \varepsilon\} && \text{if } r = \varepsilon, \\ \varphi(R, r) &= \{R \rightarrow \sigma\} && \text{if } r = \sigma, \\ \varphi(R, r) &= \{R \rightarrow R_1, R \rightarrow R_2\} \cup \\ &\quad \varphi(R_1, r_1) \cup \varphi(R_2, r_2) && \text{if } r = r_1 + r_2, \\ \varphi(R, r) &= \{R \rightarrow R_1 R_2\} \cup \\ &\quad \varphi(R_1, r_1) \cup \varphi(R_2, r_2) && \text{if } r = r_1 \cdot r_2, \\ \varphi(R, r) &= \{R \rightarrow RR_1, R \rightarrow \varepsilon\} \cup \\ &\quad \varphi(R_1, r_1) && \text{if } r = r_1^*, \end{aligned}$$

where the non-terminal symbols introduced in the rules above, are not used elsewhere.

We recall the translation from finite automata to context-free grammars. To this end, let $A = (\Sigma, Q, I, F, \delta)$ be an NFA. Then, a context-free grammar G^A such that $L(G^A) = L(A)$ can be constructed as follows. The grammar $G^A = (N_A, \Sigma, R_A, S_A)$ is defined with $N_A = S_A \uplus \{R_q \mid q \in Q\}$ and the set of rules R_A which is defined as the union of all rules of the form

- $S_A \rightarrow R_q$, for every $q \in I$,
- $R_q \rightarrow \varepsilon$ for every $q \in F$, and
- $R_{q_1} \rightarrow aR_{q_2}$ for every $a \in \Sigma$ and $q_2 \in Q$ such that $q_2 \in \delta(q_1, a)$.

Observe that if A is a DFA then G^A is an unambiguous CFG. In particular, there is only one rule of the form $S_A \rightarrow R_q$ since I is a singleton.

These translations are standard from the literature. It is well-known that they obey the properties listed in the following lemma.

LEMMA 3.2. *Let r be a regular expression and let A be a DFA over alphabet Δ . Then we can construct in linear time a CFG G with start symbol S such that $\tau_1 \dots \tau_n \in \mathcal{L}(r)$ (resp. $\mathcal{L}(A)$) if and only if $S \Rightarrow^* \tau_1 \dots \tau_n$ in G . Furthermore, if r (resp. A) is deterministic, then G is unambiguous.*

LEMMA A.1. *Let $D = (\Sigma, \Delta, X, d, \{s_d\}, \mu)$ be an EDTD or EDTD-DFA and $G = (N, \Sigma \cup \{[,]\}, \Psi_D, T_{s_d})$ be its associated context-free grammar. Then, for any $R \in N$, and for any $w \in (\Sigma \cup \{[,]\})^*$ such that $R \Rightarrow^* w$, it holds that $|w| = 3n$ for some $n \in \mathbb{N}$.*

PROOF. We first prove the lemma for the EDTD case and then for the EDTD-DFA case. Let D be an EDTD as defined in the lemma, and G the corresponding context-free grammar, and let R be any non-terminal symbol in N . Let $R \Rightarrow^* w$. We proceed by induction on the number of derivation steps n . Suppose that $n = 1$ and hence that $R \Rightarrow w$. By definition, the only rule in Ψ_D that produces a string of only terminal symbols, is $R \Rightarrow \varepsilon$, and $|\varepsilon| =$

$0 = 3 \cdot 0$. Suppose then that the statement holds for all $k < K$ derivation steps, for some $K \in \mathbb{N}$, and consider the case where the number of derivation steps is K . Suppose that $R \Rightarrow R_1 R_2$ and therefore there exist w_1, w_2 such that $w = w_1 w_2$ and $R_1 \Rightarrow w_1$ and $R_2 \Rightarrow w_2$. Then by the inductive hypothesis, $|w_1| = 3n_1$ and $|w_2| = 3n_2$ for $n_1, n_2 \in \mathbb{N}$. Therefore, $|w| = 3(n_1 + n_2)$. The argument is similar for the other cases that are in the image of φ given some regular expression r .

The only remaining case is where $R \Rightarrow \sigma \cdot [R']$, where $R' \in N$. Therefore, if $R \Rightarrow^K w$, there exists w' such that $w = \sigma \cdot [w']$ and $R' \Rightarrow^{(K-1)} w'$. By the inductive hypothesis, $|w'| = 3n$ for some $n \in \mathbb{N}$ and therefore $|w| = 3n + 3 = 3(n + 1)$.

The case where D is an EDTD-DFA is similar to the previous case, with the difference that, in the inductive step we need to consider a production of the form $R \Rightarrow T_\tau R'$. Again, by induction, we have that there exist w_1, w_2 such that $w = w_1 w_2$ and $T_\tau \Rightarrow w_1$ and $R' \Rightarrow w_2$ and by the inductive hypothesis, $|w_1| = 3n_1$ and $|w_2| = 3n_2$ for $n_1, n_2 \in \mathbb{N}$. \square

We prove that, using the translation from EDTD-DFA to CFGs, EDTD-DFA^{un} are translated to unambiguous CFGs.

LEMMA A.2. *If $D = (\Sigma, \Delta, A, d, S_d, \mu)$ is an EDTD-DFA^{un}, then its associated context-free grammar $G = (N, \Sigma \cup \{[,]\}, \Psi_D, S)$ is unambiguous.*

PROOF. Let w be a string in $\mathcal{L}(G)$. Let $\Sigma' = \Sigma \uplus \{[,]\}$. We want to show that there exists a unique derivation tree for w . We show the equivalent statement that each w has a unique *left derivation*, i.e., a derivation in which at each step the leftmost non-terminal is replaced.

We show by induction on n that for any $B \in N \cup \Sigma'$ and any string $w \in (\Sigma')^*$, $B \Rightarrow^n w$ implies that there is a unique left derivation for w from B . Let $n = 0$. Then $w = \sigma$ for some $\sigma \in \Sigma$ or $w = \varepsilon$. In both cases, $B = w$ and therefore there is a unique left derivation for w . Suppose then that the statement is true for all $n < K$ for some $K \in \mathbb{N}$, and let $B \Rightarrow^K w$.

Suppose first that B is equal to $R_{q_\tau} \in R_{A_\tau}$ for some τ and $q_\tau \in A_\tau$. There are two cases: either $w = \varepsilon$ or $w = a \cdot w'$ where $a \in \Sigma$. In the first case, the only possible left derivation for w from B is the one with the single production R_{q_τ} , if q_τ is a final state from A_τ . In the second case, let $w = a \cdot w'$ where $a \in \Sigma$. Then there exists a unique $q'_\tau \in A_\tau$ such that $R_{q_\tau} \Rightarrow a R_{q'_\tau}$ and $R_{q'_\tau} \Rightarrow^{(K-1)} w'$ because A_τ is deterministic. By the inductive hypothesis, w' has a unique left derivation from $R_{q'_\tau}$, and therefore w also has a unique left derivation from R_{q_τ} .

Suppose then that B is equal to T_τ for some τ . For each τ there is a unique rule $T_\tau \rightarrow \sigma[S_{A_\tau}]$ with T_τ at the left hand side, and therefore w is of the form $\sigma \cdot w'$, where $S_{A_\tau} \Rightarrow^{(K-1)} w'$. By the inductive hypothesis, w' has a unique left derivation from S_{A_τ} , and therefore so does w from $B = T_\tau$.

Let B be equal to S_{A_τ} for some $\tau \in \Delta$. Since D is a DFA, for each $\tau \in \Delta$, there exists a unique $q \in A_\tau$ such that $S_{A_\tau} \rightarrow R_q$. Therefore, $R_q \Rightarrow^{(K-1)} w$, and by the inductive hypothesis w has a unique left derivation from R_q and therefore also from S_{A_τ} .

Finally let B be equal to the start symbol S . Since D is an EDTD, there are no two distinct types $\tau_1, \tau_2 \in \Delta$ such that $\mu(\tau_1) = \mu(\tau_2)$ and $\tau_1, \tau_2 \in S_d$. So suppose that $w = \sigma \cdot w'$.

Then by the assumption above, there exists a unique type $\tau \in S_d$ such that $T_\tau \Rightarrow \sigma[S_{A_\tau}]$. In particular, $T_\tau \Rightarrow^{(K-1)} w$ and therefore there is a unique left derivation of w from T_τ by the inductive hypothesis. Hence w has a unique left derivation from S . \square

LEMMA 3.3. *For every EDTD $D = (\Sigma, \Delta, X, d, S_d, \mu)$ the CFG $G = (N, \Sigma \cup \{[,]\}, \Psi_D, S)$ is such that for all $n \in \mathbb{N}$, $|\mathcal{L}(D)^{=n}| = |\mathcal{L}(G)^{=3n}|$. Furthermore, if D is an EDTD-DFA^{un}, then G is unambiguous.*

PROOF. We will show by induction on n that, for any $n \in \mathbb{N}$ and any type τ , \mathbf{str} is a bijection between $\mathcal{L}(D_\tau)^{=n}$ and $\mathcal{L}(G_{T_\tau})^{=3n}$.³ Namely the following statement holds. For any tree t , if $t \in \mathcal{L}(D_\tau)^{=n}$ then $\mathbf{str}(t) \in \mathcal{L}(G_{T_\tau})^{=3n}$ and, for any string w , if $w \in \mathcal{L}(G_{T_\tau})^{=3n}$ then $w = \mathbf{str}(t)$ for some tree t such that $t \in \mathcal{L}(D_\tau)^{=n}$. Since the function \mathbf{str} is injective and well-defined, the statement of the Lemma follows.

For the base case, let $n = 1$. For the first direction, suppose that for some tree t , $t \in \mathcal{L}(D_\tau)^{=1}$. Then $t = \sigma(\varepsilon)$ for $\sigma = \mu(\tau)$. By definition of $\psi_D(\tau)$, the string $\sigma[]$ is in $\mathcal{L}(G_{T_\tau})^{=3}$ if and only if $\varepsilon \in \mathcal{L}(G_{R_\tau})$. We know that $\varepsilon \in \mathcal{L}(d(\tau))$ and therefore by induction on the structure of $d(\tau)$, as shown in Lemma 3.2, it holds that $\varepsilon \in \mathcal{L}(G_{R_\tau})$.

For the second direction suppose that for some type τ , a string w is in $\mathcal{L}(G_{T_\tau})^{=3}$. Notice that each rule in G that produces terminal symbols, produces 3 terminal symbols and is of the form $T_\tau \rightarrow \sigma[R]$, for some non-terminal symbol R and terminal symbol σ . So suppose that $\sigma[] \in \mathcal{L}(G_{T_\tau})^{=3}$. Then, since the set of rules of G_{T_τ} is equal to $\{S \rightarrow T_\tau \mid \tau' \in S_d\} \cup \bigcup_{\tau' \in \Delta} \psi(\tau')$ with $\psi(\tau') = \{T_\tau \rightarrow \sigma[R]\} \cup \text{CFG}(d(\tau'), R)$, where $\sigma = \mu(\tau')$, it holds by induction on the structure of $d(\tau')$, as shown in Lemma 3.2, that $\varepsilon \in \mathcal{L}(d(\tau'))$ and therefore $\sigma(\varepsilon) \in \mathcal{L}(D_\tau)$.

So suppose that the statement holds for all $n < K$ for some $K \in \mathbb{N}$. We want to show that the statement holds for K .

For the first direction, suppose that, for some tree t and type τ , $t \in \mathcal{L}(D_\tau)^{=K}$. Then $t = \sigma_0(t_1, \dots, t_m)$ for some m such that for each $i \in [1, m]$, $t_i \in \mathcal{L}(D_{\tau_i})^{=n_i}$, where $n_i < K$, $\sum_{i=1}^m n_i = K - 1$, and where $\tau_1 \dots \tau_m \in d(\tau)$. By the inductive hypothesis, $\mathbf{str}(t_i) \in \mathcal{L}(G_{T_{\tau_i}})^{=3n_i}$ for each $i \in [1, m]$. Notice that $K = 1 + \sum_{i=0}^m n_i$ and therefore $3K = 3 + \sum_{i=0}^m 3n_i$. By the definition of $\psi_D(\tau) = \{T_\tau \rightarrow \sigma[R]\} \cup \text{CFG}(d(\tau), R)$, it suffices to show that $\mathbf{str}(t_1) \dots \mathbf{str}(t_m) \in \mathcal{L}(G_R)$. From the assumption above, that for all $i \in [1, m]$, $\mathbf{str}(t_i) \in \mathcal{L}(G_{T_{\tau_i}})^{=3n_i}$, it is enough to show that $R \Rightarrow^* T_{\tau_1} \dots T_{\tau_m}$ by the rules in G . However, this is true by construction of G and by Lemma 3.2.

For the second direction, suppose that for some type τ and string w , $w \in \mathcal{L}(G_{T_\tau})^{=3K}$. Since the set of rules of G_{T_τ} is equal to $\bigcup_{\tau' \in \Delta} \psi_D(\tau')$ and $\psi_D(\tau) = \{T_\tau \rightarrow \sigma[R]\} \cup \text{CFG}(d(\tau), R)$ where $\sigma = \mu(\tau)$, it holds that $w \in \mathcal{L}(G_{T_\tau})^{=3K}$ if and only if there exists string w' such that $w = \sigma \cdot [w']$ and $w' \in \mathcal{L}(G_R)^{=3(K-1)}$. Furthermore, $w' \in \mathcal{L}(G_R)^{=3(K-1)}$ if and only if $R \Rightarrow^* T_{\tau_1} \dots T_{\tau_m}$ for some m and $w' = w_1 \dots w_m$ such that for each $i \in [1, m]$, $w_i \in \mathcal{L}(G_{T_{\tau_i}})^{=n_i}$, where $\sum_{i=1}^m n_i = 3(K-1)$. Notice that for each $i \in [1, m]$, $n_i = 3n'_i$ for some n'_i by Lemma A.1, and therefore $\sum_{i=1}^m n'_i = K - 1$. By the inductive hypothesis, there exist trees t_1, \dots, t_m such that for all $i \in [1, m]$,

³Recall that D_τ and G_{T_τ} denote the EDTD D and the grammar G with start symbols τ and T_τ , respectively.

$w_i = \text{str}(t_i)$ and $t_i \in \mathcal{L}(D_{\tau_i})^{=n_i}$. Furthermore, by induction on the structure of $d(\tau)$, as shown in Lemma 3.2, $R \Rightarrow^* T_{\tau_1} \dots T_{\tau_m}$ if and only if $\tau_1 \dots \tau_m \in \mathcal{L}(d(\tau))$. Therefore, $t = \sigma(t_1, \dots, t_m) \in \mathcal{L}(D_\tau)$ and in addition $t \in \mathcal{L}(D_\tau)^K$.

Finally, if D is an EDTD-DFA^{un}, we have shown in Lemma A.2 that G is an unambiguous grammar. \square

Proofs for Section 3.3

THEOREM 3.9. *For an EDTD-DFA^{un} $D = (\Sigma, \Delta, X, d, S_d, \mu)$, the cardinality $\mathcal{L}(D)^{=k}$ can be computed using $O(|\Delta||Q_{\max}|^2k)$ arithmetic operations, where Q_{\max} denotes the largest state space of an automaton in X .*

PROOF. It suffices to show that the Chomsky normal form G_D of the CFG $G_D = (N, \Sigma \cup \{[,]\}, \Psi_D, S)$ has at most $O(|\Delta||Q_{\max}|^2)$ non-terminals. For this we observe that the number of non-terminals $|N|$ of G_D is $O(|\Delta||Q_{\max}|)$, where Q_{\max} the largest set of states in any automaton in X . Furthermore, a close inspection of the rules in G_D reveals that only two rules are not in normal form: (1) the initial rule corresponding to $S \rightarrow T_{s_d^1} \cup \dots \cup T_{s_d^p}$ where $S_d = \{s_d^1, \dots, s_d^p\}$; and (2) the rules corresponding to the outgoing transitions in each state: $R_q \rightarrow \bigcup_{(a_i, q_j)} a_i \cdot R_{q_j}$, where the (a_i, q_j) range over all $a_i \in \Sigma$ and $q_j \in Q$ such that $q_j \in \delta(q, a_i)$. Hence, after normalization (1) results in $|S_d|$ new non-terminals and (2) in at most $|Q|^2$ new non-terminals. As a consequence, G' has $O(|\Delta||Q_{\max}|^2)$ non-terminals and hence Theorem 3.8 and Lemma 3.3 imply that $\mathcal{L}(D)^{=k}$ can be computed using $O(|\Delta||Q_{\max}|^2 3k) = O(|\Delta||Q_{\max}|^2 k)$ arithmetic operations. \square

Proofs for Section 3.2

PROPOSITION 3.4. *#EDTD is #P-complete.*

PROOF. It is shown in [23] that the problem #NFA is #P-complete, where, given $m \in \mathbb{N}$ and an NFA N , #NFA is the problem of finding $|\mathcal{L}(N) \cap \Sigma^m|$. Let $A = (\Sigma, Q, I, F, \delta)$ be any NFA. We define D to be the EDTD $(\Sigma \uplus \{\sigma_0\}, \Delta = Q \times \Sigma, R, d, S_d, \mu)$, where the following holds. The symbol σ_0 is a distinguished symbol not in Σ . Furthermore, for any $\tau \in Q \times \Sigma$, if $\tau = (q, a)$ then $\mu(\tau) = a$, d is defined as the function mapping any $\tau = (q, a)$ to the disjunction of the elements of the set $\{(q', a') \mid q' \in \delta(q, a')\} \cup \{\varepsilon \mid q \in F\}$, and $S_d = \{(q, \sigma_0) \mid q \in I\}$. Since each $d((q, a))$ is a finite language, it can be represented by a polynomial-size deterministic regular expression. We want to show that for any string $w \in \Sigma^*$, $w \in \mathcal{L}(A) \Leftrightarrow \sigma_0 \cdot w \in \mathcal{L}(D)$.

For the *only if* direction, suppose that for some string $w = a_1 \dots a_n$, $w \in \mathcal{L}(A)$. Then there is a successful run $\rho = q_0, \dots, q_n$ of A on w . To show that $\sigma_0 a_1 \dots a_n$ is tree that satisfies D , it suffices to show that there is a tree $t' = \tau_{\tau_1} \dots \tau_n$ that is a witness to t with $\tau \in S_d$ and such that for all $i \in [1..n]$, $\mu(\tau_i) = a_i$. But $t' = (q_0, \sigma_0) \cdot (q_1, a_1) \cdot \dots \cdot (q_n, a_n)$ is such a tree by definition.

For the *if* direction, suppose that a tree $t = \sigma_0 a_1 \dots a_n$ is in $\mathcal{L}(D)$. Notice that the letter of the root of t is always labelled by σ_0 by definition of D . We want to show that there is an accepting run ρ of A on $w = a_1 \dots a_n$. Let $t' = (q, \sigma_0)(q_1, a_1) \dots (q_n, a_n)$ be the witness of t . Then the run $\rho = q, q_1, \dots, q_n$ is an accepting run of A on w . The translation of NFAs to EDTDs above can be performed in PTIME and therefore #EDTD is #P-hard.

Finally, deciding whether there is a tree of a given size k in the language of an EDTD is in NP. Therefore, the problem

is also in #P \square

Proofs for Section 5

LEMMA 5.2.

1. *The function c is a bijection for admissible k -OAs.*
2. *For two admissible k -OAs A_1 and A_2 with n states over an alphabet Σ of size ℓ , with a total order $o : \Sigma \rightarrow [0, \ell - 1]$ defined on this alphabet, if $\text{enc}(A_1) = \text{enc}(A_2)$ then A_1 is isomorphic to A_2 .*

PROOF. For (1), let A be an admissible k -OA with M states. Since $c(I) = [1, \ell]$, every state s is mapped to some $n \in [1, M]$. To show that every $n \in [1, M]$ is in the image of c , we need to show that c is injective, since $|V| = M$. But notice that every time a number is assigned to a state, it is the smallest number that has not yet been assigned to a state.

For (2), let $A_i = (V_i, E_i, I_i, F_i, \text{lab}_i, e_i)$ for $i \in [1, 2]$, let $s = \text{enc}(A_1)$ and let $A = (V, E, I, F, \text{lab}, e)$ be the following k -OA. The set of states V is equal to $[1, n]$, where n is the largest number appearing in s , for any state $v \in V$, $\text{lab}(v) = o^{-1}(j \pmod{\ell})$, where j is the first position in s where v appears, for any two states $v_1, v_2 \in V$, let $(v_1, v_2) \in E$ if $v_2 = s_{v_1 \cdot \ell + o(\text{lab}(v_2))}$ and let $I = [1, \ell]$. Note that, $\text{enc}(A) = \text{enc}(A_1)$.

We want to show that A_1 is isomorphic to A . By symmetry it will follow that A_2 is also isomorphic to A and therefore isomorphic to A_1 . Let c be the ordering of the states of A_1 and notice that the ordering of the states of A is the identity function. We define the function $\alpha : V_1 \rightarrow V$ that maps a state $s \in V_1$ to the state $c(s) \in V$. Suppose that α is not an isomorphism for the two automata A and A_1 . Then one of the following holds.

- There exists $k \in [1, n]$ such that $\text{lab}(k) \neq \text{lab}_1(c^{-1}(k))$. Then $o(\text{lab}(k)) \neq o(\text{lab}_1(c^{-1}(k)))$ and the first position in $\text{enc}(A)$ that k occurs is different from the first position that k occurs in $\text{enc}(A_1)$.
- The final states of A are different than the final states of A_1 . Then $\text{enc}(A) \neq \text{enc}(A_1)$.
- It's not the case that $\varepsilon \in A \Leftrightarrow \varepsilon \in A_1$. But then the last letter in $\text{enc}(A)$ is different from the last letter in $\text{enc}(A_1)$.
- There exist $k, k' \in [1, n]$ such that k has a transition to k' but the state $s_1 = c^{-1}(k)$ does not have a transition to the state $s_2 = c^{-1}(k')$. Then the $(k \cdot \ell + o(\text{lab}(k)))$ -th position in the string $\text{enc}(A)$ is k' and the same position in $\text{enc}(A_1)$ is not k' , which is a contradiction.

\square

LEMMA 5.3. *For each $n \in \mathbb{N}$, enc is a bijection from the set of non-isomorphic admissible k -OAs with n states, over an alphabet Σ of size ℓ to the strings of size $(n+1) \cdot \ell + n + 1$ whose prefix $s = s_0, \dots, s_{(n+1) \cdot \ell - 1}$ satisfies the rules (A1) – (A4) above, and whose suffix $s' = s_{(n+1) \cdot \ell} \dots s_{(n+1) \cdot \ell + n}$ uses only 0 and 1.*

PROOF. From Lemma 5.2, it follows that the function enc is injective, so it remains to be shown that it is also surjective. Consider any string of size $(n+1) \cdot \ell$ satisfying (A1) – (A4), and let A be the automaton obtained by letting the set of states be $[1, n]$, and for each $i \in [0, \ell - 1]$, let S_i be the set of states $\{s_j \mid j < (n+1) \cdot \ell, j = i \pmod{\ell}\}$.

Then for each $s \in S_i$ let $\text{lab}(s) = i$. According to rule (A3) for each $i \in [0, \ell - 1]$, $|S_i| \leq k$. Furthermore, for each $i \in [\ell, (n+1) \cdot \ell - 1]$, if $i = j' \cdot \ell + j$, with $j < \ell$, let (j', s_i) be a member of E . Finally, let $I = [1, \ell]$ and let e and the final states of A be determined by the suffix of length $n+1$ of the string. This automaton is complete, and it remains to show that it is initially connected. But from rule (A2) for each $i \in [1, n]$, there is $j < i \cdot \ell$ such that $s_j = i$, and therefore state i in A is reachable from state $\lfloor \frac{i}{\ell} \rfloor$. By induction, every state is reachable from one of the states in $[1, \ell]$, the initial states. \square

By Lemma 5.3, sampling of admissible k -OAs reduces to uniform generation of strings satisfying the rules (A1)–(A4). We next show that the specification given in Section 5 correctly generates this set of strings. From Lemma 5.3 it follows that it suffices to show that the class $\mathcal{S}_m^{(j', j)}[\overline{W}]$ is combinatorially equivalent to the set of strings of length $(n+1) \cdot \ell$ satisfying rules (A1)–(A4), with a fixed prefix $s_0 \dots s_{j' \cdot \ell + j}$, where m occurs in the string at position $j' \cdot \ell + j$ for the first time, and for each $i \in [0, \ell - 1]$, and $p \in [0, j' \cdot \ell + j]$, it holds that $|\{s_p \mid p = i \pmod{\ell}\}| = n_i$, where $n_i = |W_i|$.

To this end, we first count the number of valid extensions of such a given prefix and then show that this number coincides with the number of objects in the corresponding class.

Let Σ be an alphabet of size ℓ . For ℓ and k as above, and $m, j', j, n, n_0, \dots, n_{\ell-1} \in \mathbb{N}$, let $N(m, j' \cdot \ell + j, n, n_0, \dots, n_{\ell-1})$ be defined inductively as follows:

$$\begin{aligned} N(m, j' \cdot \ell + j, n, n_0, \dots, n_{\ell-1}) &= 0 \text{ if } \exists i \in [0, \ell - 1] \text{ s.t. } n_i > k, \\ N(m, j' \cdot \ell + j, n, n_0, \dots, n_{\ell-1}) &= 0 \text{ if } j' \geq m, \\ N(m, j' \cdot \ell + j, n, n_0, \dots, n_{\ell-1}) &= 0 \text{ if } m > n \end{aligned} \quad (2)$$

$$N(n, (n-1)\ell + j, n, n_0, \dots, n_{\ell-1}) = \prod_{i=0}^{\ell-1} n_i \cdot \prod_{i=j+1}^{\ell-1} n_i, \quad (3)$$

and for $j' < m - 1$ and $j < \ell$,

$$N(m, j' \cdot \ell + j, n, n_0, \dots, n_{\ell-1}) = N_1 + N_2, \quad (4)$$

and for $j' = m - 1$

$$N(m, j' \cdot \ell + j, n, n_0, \dots, n_{\ell-1}) = N_3, \quad (5)$$

where:

$$\begin{aligned} N_1 &= \left(n_0 \dots n_{\ell-1} \right) \cdot N(m, (j'+1) \cdot \ell + j, n, n_0, \dots, n_{\ell-1}), \\ N_2 &= \sum_{i=1}^{\ell} \left(\left(\prod_{i'=1}^{i-1} (n_{j+i'}) \right) \cdot \right. \\ &\quad \left. N(m+1, j' \cdot \ell + j + i, n, n_0, \dots, n_{j+i \pmod{\ell} + 1}, \dots, n_{\ell-1}) \right), \\ N_3 &= \sum_{i=1}^{2\ell-j-1} \left(\left(\prod_{i'=1}^{i-1} (n_{j+i'}) \right) \cdot \right. \\ &\quad \left. N(m+1, j' \cdot \ell + j + i, n, n_0, \dots, n_{j+i \pmod{\ell} + 1}, \dots, n_{\ell-1}) \right). \end{aligned}$$

Informally, the number of strings satisfying the rules (A1)–(A4) where letter m appears for the first time at position $j' \cdot \ell + j$ is equal to the sum of the number of possible strings where $m+1$ appears for the first time in one of the positions between $j' \cdot \ell + j + 1$ and $m \cdot \ell - 1$. In equation (4), N_1 covers the case where between the positions $j' \cdot \ell + j$ and $(j'+1) \cdot \ell + j$ no new state $m+1$ is introduced, and N_2 covers the cases where at some position

J between the positions $j' \cdot \ell + j$ and $(j'+1) \cdot \ell + j$, the state $m+1$ is introduced. In the case where $j' = m - 1$, since $N(m, (j'+1) \cdot \ell + j, n, n_0, \dots, n_{\ell-1}) = 0$ by the definition of the base case that also complies with the rule (A2), the equation is different in order to take into account the remaining positions until position $m \cdot \ell - 1$.

LEMMA A.3. For $\ell, k, n, j', j, n_0, \dots, n_{\ell-1} \in \mathbb{N}$ and $m \geq 1$, the number $N(m, j' \cdot \ell + j, n, n_0, \dots, n_{\ell-1})$ is the number of strings of length $(n+1) \cdot \ell$ satisfying rules (A1)–(A4), with a fixed prefix $s_0 \dots s_{j' \cdot \ell + j}$, where m occurs in the string at position $j' \cdot \ell + j$ for the first time, and for each $i \in [0, \ell - 1]$, and $p \in [0, j' \cdot \ell + j]$, it holds that $|\{s_p \mid p = i \pmod{\ell}\}| = n_i$.

PROOF. We proceed by inverse induction on m to show that the statement above holds. For the base case, let $m = n$. Notice that if $m = n+1$, $N(n+1, J, n, n_0, \dots, n_{\ell-1}) = 0$ for all values of the other parameters. We show that $N(n, J, n, n_0, \dots, n_{\ell-1})$ is the number of strings described by the Lemma, by inverse induction on J .

Let $J \in [n \cdot \ell, (n+1) \cdot \ell - 1]$. Then $J = n \cdot \ell + j$ for some $j < \ell$, and by equation (2), $N(n, n \cdot \ell + j, n, n_0, \dots, n_{\ell-1}) = 0$ which is the correct number of strings according to the rule (A2).

If $J = (n-1) \cdot \ell + j$ for some $j \in [0, \ell - 1]$, then

$$N(n, (n-1) \cdot \ell + j, n, n_0, \dots, n_{\ell-1}) = \prod_{i=0}^{\ell-1} n_i \cdot \prod_{i=j+1}^{\ell-1} n_i,$$

by equation (3), which is the correct number according to the rules (A1)–(A2).

For the inductive case, suppose that for all $r > J'$ for some $J' \in [0, (n+1) \cdot \ell - 1]$, $N(n, r, n, n_0, \dots, n_{\ell-1})$ is the correct number of strings. Consider the number $N(n, J', n, n_0, \dots, n_{\ell-1})$. If for any $i \in [0, \ell - 1]$, $n_i > k$ then this number is equal to 0 which complies with the rule (A3) of strings. Suppose then that for all $i \in [0, \ell - 1]$, $n_i \leq k$ and let $J' = J'_0 \cdot \ell + J'_1$, for J'_0 maximal. Notice that, since the last state n occurs at position J' for the first time, no state number can appear in the string at a position to the right of the position J' that has not appeared to the left or exactly at the position J' . Therefore, from the rules (A3) and (A4), each position to the right of position J' can be occupied by states that have already appeared before that position, and furthermore, have appeared at the appropriate σ -position. Consider therefore, the next ℓ positions, starting with $J'_0 \cdot \ell + J'_1 + 1$. For this position there are $n_{J'_1+1}$ possible values for the string to comply with the rules (A1)–(A4). Similarly, for the position $J'_0 \cdot \ell + J'_1 + 2$ there are $n_{J'_1+2}$ possible symbols, and so on, until position $J'_0 \cdot \ell + J'_1 + \ell = (J'_0 + 1) \cdot \ell + J'_1$ for which there are $n_{J'_1}$ possible values. By the inductive hypothesis, $N(n, J'_0 \cdot \ell + J'_1 + \ell, n, n_0, \dots, n_{\ell-1})$ is the number of possible strings with a fixed prefix $s_0, \dots, s_{(J'_0 \cdot \ell + J'_1 + \ell)}$, that satisfy the rules (A1)–(A4) and state n appears for the first time at position $J'_0 \cdot \ell + J'_1 + \ell$. Therefore, $N(n, J'_0 \cdot \ell + J'_1, n, n_0, \dots, n_{\ell-1}) = n_0 \dots n_{\ell-1} \cdot N(n, (J'_0 + 1) \cdot \ell + J'_1, n, n_0, \dots, n_{\ell-1})$, which is also what N_1 is defined to be according to equation (4). Notice that $N_2 = 0$, according to equation (2), which complies with the fact that no new state can appear to the right of J' .

Suppose then that for some $M < n$ and all $m > M$, the number $N(m, J, n, n_0, \dots, n_{\ell-1})$ is the correct number for all $j \in [0, (n+1) \cdot \ell - 1]$, and consider the value of

$N(M, J', n, n_0, \dots, n_{\ell-1})$. We show that this is the correct number of strings by inverse induction on J' . For all values of J' larger than $M \cdot \ell$, the number $N(M, J', n, n_0, \dots, n_{\ell-1})$ is equal to 0 which complies with rule (A2). Suppose $J' = (M-1) \cdot \ell + J'_1$. Then, $N(M, (M-1) \cdot \ell + J'_1, n, n_0, \dots, n_{\ell-1})$ is determined by equation (5). The number of strings where state M first appears at position $(M-1) \cdot \ell + J'_1$ is equal to the sum of the number of strings where state $M+1$ appears in some position in the next $2 \cdot \ell - J'_1$ positions, and this is the number given by equation (5).

For the inductive hypothesis, suppose that $N(M, r, n, n_0, \dots, n_{\ell-1})$ is the correct number of strings for all $r > J'$ for some $J' \in [0, (n+1) \cdot \ell - 1]$. Consider $N(M, J', n, n_0, \dots, n_{\ell-1})$. If for any $i \in [0, \ell - 1]$, $n_i > k$ then this number is equal to 0 which complies with the rule (A3) of strings.

Otherwise, the possible strings with n states that satisfy rules (A1) – (A4) and with prefix $s_0 \dots s_{J'}$, where M occurs for the first time at position J' , are the following. Either, no new symbol appears in the following ℓ positions to the right of J' , or at least one does. For the first case, the number of possible strings is then $n_0 \dots n_{\ell-1} \cdot N(M, J' + \ell, n, n_0, \dots, n_{\ell-1})$, where by the inductive hypothesis, $N(M, J' + \ell, n, n_0, \dots, n_{\ell-1})$ is the correct number of the appropriate strings. This is the number given by the term N_1 of equation (4). For the second case, let $J' = J_0 \cdot \ell + J_1$ and let us consider all possible ℓ positions to the right of position J' where the new symbol appears. Suppose that this position is $J' + i$ for $i \leq \ell$. Then, the number of possible strings complying with rules (A1) – (A4) is the following. For position $J' + 1$ there are n_{J_1+1} possible values, for position $J' + 2$ there are n_{J_1+2} possible values, and so on until the position $J' + i$, where the number of allowed strings with prefix $s_0 \dots s_{J'+i}$, and where state $M+1$ first appears at position $J' + i$, is given by $N(M+1, J' + i, n, n_0, \dots, n_{J_1+i \pmod{\ell}}, \dots, n_{\ell-1})$, by the inductive hypothesis. Considering all possible values where the new symbol can appear, we get a sum equal to the term N_2 of equation (4). Notice that any string counted in one of the terms of this sum, is not counted in any other term of this sum. Therefore, $N(M, J', n, n_0, \dots, n_{\ell-1})$ is equal to $N_1 + N_2$, which is what is described by equation (4). \square

We next show that the number of objects in $\mathcal{S}_m^{(p,j)}$ coincides with $N(m, p \cdot \ell + j, n, n_0, \dots, n_{\ell-1})$, from which the desired combinatorial equivalence follows.

LEMMA A.4. *Let $\ell, k, m, n, p, j, n_0, \dots, n_{\ell-1} \in \mathbb{N}$, $\ell < m$. The number of objects in $\mathcal{S}_m^{(p,j)}[\overline{W}]$ is equal to $N(m, p \cdot \ell + j, n, n_0, \dots, n_{\ell-1})$, where for each $i \in [0, \ell - 1]$, $n_i = |W_i|$.*

PROOF. We show by induction on m that the statement holds. Suppose first that $m = n$. We then show that the statement holds by induction on p . For the base case, suppose that $p = n - 1$. Then $\mathcal{S}_n^{(n-1,j)}[\overline{W}] := \prod_{i=j+1}^{\ell-1} \mathcal{W}_i \times \prod_{i=0}^{\ell-1} \mathcal{W}_i$, and $N(n, (n-1)\ell + j, n, n_0, \dots, n_{\ell-1}) = \prod_{i=j+1}^{\ell-1} n_i \cdot \prod_{i=0}^{\ell-1} n_i$ and hence the statement holds for $p = n - 1$. Next, assume that the statement holds for $m = n$ and $p > P$ for some $P \in \mathbb{N}$ and consider the case where $p = P$. The class $\mathcal{S}_m^{(P,j)}[\overline{W}]$ is given by $Q_1 + Q_2$. Consider first Q_1 , defined as $Q_1 = (\prod_{i=j+1}^{\ell+j} \mathcal{W}_i \pmod{\ell}) \times \mathcal{S}_{m,(P+1,j)}[\overline{W}]$. Since for $m-1 = n$, $N(m-1, P \cdot \ell + j, n, n_0, \dots, n_{\ell-1}) = (n_0 \cdot \dots \cdot n_{\ell-1}) \cdot$

$N(m-1, (P+1) \cdot \ell + j, n, n_0, \dots, n_{\ell-1})$. From the inductive hypothesis, we may conclude that the statement holds $m-1 = n$ and $p = P$. Similarly, the class Q_2 is defined by the equation $Q_2 := \sum_{i=j+1}^{\ell+j} (\prod_{i'=i}^{i-1} \mathcal{W}_{i' \pmod{\ell}}) \times \mathcal{Z}_{m+1} \times \mathcal{S}_{m+1,(P,j+i)}[\overline{W}_{m+1,j+i \pmod{\ell}}]$ and N_2 is defined to be equal to $\sum_{i=1}^{\ell} \left(\left(\prod_{i'=1}^{i-1} (n_{P+i'}) \right) \cdot N(m+1, P \cdot \ell + j + i, n, n_0, \dots, n_{P+i \pmod{\ell}} + 1, \dots, n_{\ell-1}) \right)$. We have that Q_2 is undefined whereas N called for $m+1$ is equal to 0 by definition, since $m+1 > n$. Hence, the statement holds for $m = n$ and $p = P$.

Suppose next that the statement holds for all $m > m_0$ for some $m_0 \in \mathbb{N}$, and consider the case where $m = m_0$. For the base case, let $p = m_0 - 1$, we have that the class $\mathcal{S}_{m_0}^{(p,j)}[\overline{W}]$ is defined by $Q_3 := \sum_{i=j+1}^{2\ell-1} (\prod_{i'=i}^{i-1} \mathcal{W}_{j+i' \pmod{\ell}}) \times \mathcal{Z}_{m_0+1} \times \mathcal{S}_{m_0+1}^{(p,j+i)}[\overline{W}_{m_0+1,j+i \pmod{\ell}}]$, whose number of elements of size $(n+1) \cdot \ell$ is equal to $N_3 = \sum_{i=1}^{2\ell-j-1} \left(\left(\prod_{i'=1}^{i-1} (n_{j+i'}) \right) \cdot$

$N(m_0+1, p \cdot \ell + j + i, n, n_0, \dots, n_{j+i \pmod{\ell}} + 1, \dots, n_{\ell-1}) \right)$ by the inductive hypothesis. Finally, assume that the statement holds for all $p > P$ for some $P \in \mathbb{N}$, and consider the case where $p = P$. The class $\mathcal{S}_{m_0}^{(p,j)}[\overline{W}]$ is defined by the equation $Q_1 + Q_2$. Since $Q_1 := (\prod_{i=j+1}^{\ell+j} \mathcal{W}_i \pmod{\ell}) \times \mathcal{S}_{m_0}^{(P+1,j)}[\overline{W}]$, its number of elements is equal to $N_1 = (n_0 \cdot \dots \cdot n_{\ell-1}) \cdot N(m_0, (P+1) \cdot \ell + j, n, n_0, \dots, n_{\ell-1})$, by the inductive hypothesis. Similarly, the number of elements in Q_2 defined by $\sum_{i=j+1}^{\ell+j} (\prod_{i'=i}^{i-1} \mathcal{W}_{i' \pmod{\ell}}) \times \mathcal{Z}_{m_0+1} \times \mathcal{S}_{m_0+1}^{(P,j+i)}[\overline{W}_{m_0+1,j+i \pmod{\ell}}]$ is equal to $N_2 = \sum_{i=1}^{\ell} \left(\left(\prod_{i'=1}^{i-1} (n_{j+i'}) \right) \cdot N(m_0+1, P \cdot \ell + j + i, n, n_0, \dots, n_{j+i \pmod{\ell}} + 1, \dots, n_{\ell-1}) \right)$, by the inductive hypothesis. \square