

# Efficient enumeration algorithms for annotated grammars

Antoine Amarilli

LTCI, Télécom Paris, Institut polytechnique de Paris  
antoine.amarilli@telecom-paris.fr

Louis Jachiet

LTCI, Télécom Paris, Institut polytechnique de Paris  
louis.jachiet@telecom-paris.fr

Martín Muñoz

PUC & IMFD

mmunos@uc.cl

Cristian Riveros

PUC & IMFD

cristian.riveros@uc.cl

We introduce annotated grammars, an extension of context-free grammars which allows annotations on terminals. Our model extends the standard notion of regular spanners, and is more expressive than the extraction grammars recently introduced by Peterfreund. We study the enumeration problem for annotated grammars: fixing a grammar, and given a string as input, enumerate all annotations of the string that form a word derivable from the grammar. Our first result is an algorithm for unambiguous annotated grammars, which preprocesses the input string in cubic time and enumerates all annotations with output-linear delay. This improves over Peterfreund's result, which needs quintic time preprocessing to achieve this delay bound. We then study how we can reduce the preprocessing time while keeping the same delay bound, by making additional assumptions on the grammar. Specifically, we present a class of grammars which only have one derivation shape for all outputs, for which we can enumerate with quadratic time preprocessing. We also give classes that generalize regular spanners for which linear time preprocessing suffices.

# 1. Introduction

Arguably the most fundamental problem in database research is query evaluation: given as input a query and data, we must find the results of the query over the data. Database theory research has studied the complexity of such problems for decades. However, in some contexts, for instance over large datasets, the usual complexity measures are not well-suited to this study. Indeed, the number of query results might be so large that it is unreasonable in practice to produce all of them. Further, in theoretical terms, the complexity of an algorithm may be dominated by the cost of writing the full output, hiding the actual complexity of the computation. For this reason, a significant line of research on query evaluation has adopted the perspective of *enumeration algorithms*. Instead of explicitly producing all results, the task is to *enumerate* them, in any order and without repetition. The cost of the algorithm is then measured across two dimensions: the *preprocessing time*, which is the time needed to read the input and prepare an enumeration data structure; and the *delay*, the worst-case time elapsed between any two solutions while enumerating using the data structure.

This study of enumeration algorithms has managed in some cases to achieve a *constant-delay* guarantee. In this case, once the algorithm has preprocessed its input, the delay between any two outputs is *constant*, i.e., it is independent from the input. Of course, the challenge is to achieve this strong guarantee after a preprocessing phase that runs in a limited amount of time – in particular, one that does not explicitly materialize all solutions. Starting with the work of Durand and Grandjean [16], researchers have designed constant-delay algorithms for several query evaluation problems, e.g., the evaluation of some queries over relational databases [10, 23], query evaluation over dynamic data [11, 22], query evaluation over graph data [21, 24], among others [33].

One area where enumeration algorithms have been especially successful is the problem of *information extraction*, studied through the lens of *document spanners* [17]. In this data management task, the data is a textual document (i.e., a string), and the query is a declarative specification of information to extract from the text, formalized as a *spanner*. The spanner describes *mappings*, which are possible choices of how to map variables to substrings of the document (called spans). The enumeration problem is then to enumerate all mappings of a spanner on an input document, i.e., to enumerate efficiently all possible results for the information extraction task. The work by Florenzano et al. [18] showed that the task could be solved with preprocessing linear in the document and polynomial in a finite deterministic automaton describing the spanner, improving on a theoretical result by Bagan [9]; and this was extended in [5] to spanners described using nondeterministic automata or regular expressions.

However, while regular spanners are natural, they do not capture all possible information extraction tasks. More expressiveness is needed for extraction over structured data (e.g., XML, or JSON documents), over the source code of programs, or possibly over natural language texts. We believe that a natural way to address such tasks is to move from finite automata to *context-free grammars* (CFGs). Context-free grammars are a well-known formalism: they extend regular expressions and are commonly used, e.g., in programming language design. Common verification tasks on textual representations of

tree documents can be expressed using CFGs, and so can parsing tasks, e.g., to extract subexpressions from source code data. However, CFGs do not describe *captures*, i.e., they do not specify how to extract the parts of interest of an input document, and thus cannot be used directly for information extraction.

This question of information extraction with grammars was studied by Peterfreund in very recent work [30]. This paper proposed a formalism of *extraction grammars*, which are CFGs extended via special terminals that describe the endpoints of spans. Further, it presents an algorithm to enumerate the mappings captured by *unambiguous* extraction grammars on an input document. However, while the algorithm achieves constant-delay, the preprocessing bound is significantly worse than in the case of regular spanners: it is quintic in the document, and exponential in the number of variables of the grammar. This complexity is also worse than CFG parsing, e.g., the standard CYK parsing algorithm runs in cubic time in the input string.

Our goal in this paper is to study the enumeration problem for CFGs while achieving better complexities. Our algorithms ensure a constant-delay guarantee when outputs have constant size, and more generally ensure *output-linear* delay when this is not the case: the delay is linear in the size of each produced solution. Within this delay bound, the preprocessing time has lower complexity: it is at worst cubic in the input document, and improves to quadratic or even linear time for restricted classes. We achieve these results by proposing a new formalism to extend CFGs, called *annotated grammars*, on which we impose an unambiguity restriction similar to that of [30]. Let us present our specific contributions.

**Contributions.** Our first contribution is to introduce *annotated grammars* (Section 2). They are a natural extension of CFGs, where terminals are optionally annotated by the information that we wish to extract. We then study the problem, given an annotated grammar  $\mathcal{G}$  and document  $s$ , of enumerating all annotations of  $s$  that are derived by  $\mathcal{G}$ . This captures the enumeration problems for regular spanners [18, 7], nested words [27], and even the extraction grammars of [30] (we explain this in Section 6). As we explain, we aim for *output-linear delay*, which is the best possible delay in our setting where the solutions to output may have non-constant size.

Our second contribution is to study the enumeration problem for *unambiguous* annotated grammars (Section 3), that do not produce multiple times the same annotation of an input string. This is a natural restriction to avoid duplicate results, which is also made in [30]. For such grammars, we present an algorithm to enumerate the annotations produced by a grammar  $\mathcal{G}$  on a string  $s$  with output-linear delay (independent from  $\mathcal{G}$  or  $s$ ), after a preprocessing time of  $\mathcal{O}(|\mathcal{G}| \cdot |s|^3)$ , i.e., cubic time in  $s$ , and linear time in  $\mathcal{G}$ . This improves over the result of [30] whose preprocessing is quintic. Our algorithm has a modular design: it follows a standard design of a CFG parsing algorithm, but uses the abstract data structure of [27] to represent the sets of annotations and combine them with operators in a way that allows for output-linear enumeration. We further show a conditional lower bound on the best preprocessing time that can achieve output-linear delay, by reducing from the standard task of checking membership to a CFG, and using the lower bound of [1]. We show that the preprocessing time must be  $\Omega(|s|^{\omega-c})$  for every

$c > 0$ , where  $\omega$  is the Boolean matrix multiplication exponent.

Our third contribution is to improve the preprocessing time by imposing a different requirement on grammars. Thus, we introduce *rigid* annotated grammars (Section 4) where, for every input string, all annotations on the string are intuitively produced by parse trees that have the same shape. In contrast with general annotated grammars, we show that rigid annotated grammars can always be made unambiguous, so that our algorithm applies to them. But we also show that, under this restriction, the data complexity of our algorithm goes down from cubic to quadratic time. Further, achieving sub-quadratic preprocessing time would imply a sub-quadratic algorithm to test membership to an unambiguous CFG, which is an open problem.

Our last contribution shows how we can, in certain cases, achieve linear-time preprocessing complexity and output-linear delay (Section 5). This is the complexity of enumeration for regular spanners, and is by definition the best possible. We show that the same complexity can be achieved, beyond regular spanners, for a subclass of rigid grammars, intuitively defined by a determinism requirement. We define it via the formalism of *pushdown annotators* (PDAnn for short), which are the analogue of pushdown automata for CFGs, or the extraction pushdown automata of [30]. We show that PDAnn are equally expressive to annotated grammars, and that rigid CFGs correspond to a natural class of PDAnns where all runs have the same sequence of stack heights. Moreover, we show that we can enumerate with linear-time preprocessing and output-linear delay in the case of *profiled-deterministic* PDAnn, where the sequence of stack heights can be computed deterministically over the run: this generalizes regular spanners and visibly-pushdown automata.

This paper is the extended version of the work published at PODS'22 [8]. It includes complete proofs of all results in the appendix.

**Related work.** We have explained how our work is set in the context of document spanners [17], and in particular of enumeration results for regular spanners [18, 5]. A recent survey of much of this literature can be found in Peterfreund's PhD thesis [29]. The most related work to ours is the more recent introduction of extraction grammars by Peterfreund [30], which we already discussed. Another related work is [27], by some authors of the present paper. This paper studies enumeration for spanners over nested documents, defined as visibly pushdown transducers. In the present paper, we re-use the enumeration data structure of [27], and we consider a transducer model in Section 5 that recaptures some of the results of that work. However, in our problem, we do not require visibility guarantees. This poses new technical challenges: the underlying tree structure (i.e., the parse tree) is not known in advance and generally not unique.

There are also some other extensions of regular spanners that are reminiscent of CFGs, e.g., core spanners (featuring equality) or generalized core spanners (with difference) already introduced in [17], or Datalog evaluated over regular spanners as in [31]. However, to our knowledge, there are no known constant-delay enumeration algorithms in these contexts.

Our study of enumeration for annotated grammars is also reminiscent of enumeration results for queries over trees expressed as tree automata. An algorithm for this was given

by Bagan [9] with linear-time preprocessing and constant-delay in data complexity, for deterministic tree automata, and this was extended in [6] to nondeterministic automata. However, this is again more restricted: evaluating a tree automaton on a tree amounts to evaluating a *visibly pushdown* automaton over a string representation of the tree, which is again more restrictive than general context-free grammars.

## 2. Grammars and Annotators

**Strings and annotations.** Let  $\Sigma$  be a finite alphabet. We write  $\Sigma^*$  for the set of strings over  $\Sigma$ . The *length* of a string  $w = w_1 \cdots w_n \in \Sigma^*$  is  $|w| := n$ . The string of length 0 is written  $\varepsilon$ . We write  $u \cdot v$  or  $uv$  for the *concatenation* of  $u, v \in \Sigma^*$ .

Let  $\Omega$  be a finite set of annotations. An *annotated string* is a string  $\hat{w} \in (\Sigma \cup \Sigma \times \Omega)^*$ . We denote strings by  $w$  and annotated strings by  $\hat{w}$  when this avoids confusion. Intuitively, if  $\hat{w} = \hat{w}_1 \cdots \hat{w}_n$ , then  $\hat{w}_i = (a, \sigma) \in \Sigma \times \Omega$  means that the letter  $a$  at position  $i$  is annotated with  $\sigma$  (called an *annotated letter*) and  $\hat{w}_i \in \Sigma$  means that there is no annotation at position  $i$ . Given an annotated string  $\hat{w} = \hat{w}_1 \cdots \hat{w}_n$ , we denote by  $\text{str}(\hat{w}) = \text{str}(\hat{w}_1) \cdots \text{str}(\hat{w}_n)$  the *unannotated string* of  $\hat{w}$ , i.e.,  $\text{str}((a, \sigma)) := a$  and  $\text{str}(a) := a$ , and we denote by  $\text{ann}(\hat{w}) = \text{ann}(\hat{w}_1, 1) \cdots \text{ann}(\hat{w}_n, n)$  the *annotations* of  $\hat{w}$ , i.e.,  $\text{ann}((a, \sigma), i) := (\sigma, i)$  and  $\text{ann}(a, i) := \varepsilon$ . Note that  $|\text{str}(\hat{w})| = |w|$ , but the length  $|\text{ann}(\hat{w})|$  of  $\text{ann}(\hat{w})$  can be much less than  $|w|$ .

**Annotated grammars.** A *context-free grammar* (CFG) over  $\Sigma$  is a tuple  $G = (V, \Sigma, P, S)$ , where  $V$  is a set of *nonterminals*,  $\Sigma$  is the alphabet (whose letters are called *terminals*),  $S \in V$  is the *start symbol*, and  $P$  is a finite set of *rules* of the form  $X \rightarrow \alpha$  where  $X \in V$  and  $\alpha \in (V \cup \Sigma)^*$ . We assume that  $V$  and  $\Sigma$  are disjoint. In this paper, we extend this definition to an *annotated (context-free) grammar*  $\mathcal{G} = (V, \Sigma, \Omega, P, S)$ , which is simply the CFG  $(V, \Sigma \cup \Sigma \times \Omega, P, S)$ . We use  $G$  to denote a CFG and  $\mathcal{G}$  to denote an annotated grammar. The *terminals* of  $\mathcal{G}$  are letters  $a \in \Sigma$  and annotated letters  $(a, \sigma) \in \Sigma \times \Omega$ .

We recall the semantics of a CFG  $G = (V, \Sigma, P, S)$ . Given a string  $u \in \Sigma^*$ , two strings  $\gamma, \delta \in (V \cup \Sigma)^*$ , and  $X \in V$ , we say that  $uX\delta$  *produces*  $u\gamma\delta$ , denoted by  $uX\delta \Rightarrow_G u\gamma\delta$ , if  $P$  contains the rule  $X \rightarrow \gamma$ . We then say that  $\alpha \in (V \cup \Sigma)^*$  *derives*  $\beta \in (V \cup \Sigma)^*$ , denoted by  $\alpha \Rightarrow_G^* \beta$  or just  $\alpha \Rightarrow^* \beta$ , if there is a sequence of strings  $\alpha_1, \dots, \alpha_m$  with  $m \geq 1$  such that  $\alpha = \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_m = \beta$ . We say that  $G$  *derives*  $\alpha \in (V \cup \Sigma)^*$  if  $S \Rightarrow^* \alpha$ , and define the *language*  $L(G)$  of  $G$  as the set of strings  $\{w \in \Sigma^* \mid S \Rightarrow^* w\}$ . Note that our derivations are leftmost derivations, which is standard for the unambiguity notions that we introduce afterwards. The *language of an annotated grammar*  $\mathcal{G}$  is that of the underlying CFG on the alphabet of terminals  $\Sigma \cup \Sigma \times \Omega$ . In particular,  $L(\mathcal{G})$  is a set of annotated strings.

The purpose of annotated grammars is to consider all possible annotations of an input unannotated string  $w \in \Sigma^*$ . Specifically, the *semantics* of an annotated grammar  $\mathcal{G}$  is the function  $\llbracket \mathcal{G} \rrbracket$  mapping each string  $w \in \Sigma^*$  to the following (possibly empty) set of annotations:  $\llbracket \mathcal{G} \rrbracket(w) := \{\text{ann}(\hat{w}) \mid \hat{w} \in L(\mathcal{G}) \wedge \text{str}(\hat{w}) = w\}$ .

An *output* of evaluating  $\mathcal{G}$  over  $w$  is just an element  $\mu \in \llbracket \mathcal{G} \rrbracket(w)$ . Note that, in the case when  $\Omega = \emptyset$ , for all  $w \in \Sigma^*$  we have  $\llbracket \mathcal{G} \rrbracket(w) = \emptyset$  if  $w \notin L(\mathcal{G})$  and  $\llbracket \mathcal{G} \rrbracket(w) = \{\varepsilon\}$  if

$w \in L(\mathcal{G})$ . So, annotated grammars subsume CFGs. In Section 6, we show that they also subsume the extraction grammars of [30], which implies that annotated grammars are more expressive than regular spanners [18, 7], or even visibly pushdown transducers from [27].

Towards ensuring tractability, we call a CFG  $G$  *unambiguous* if for every  $w \in L(G)$  there is a unique derivation of  $w$  by  $G$ . We call an annotated grammar  $\mathcal{G}$  *unambiguous* if the underlying CFG over  $\Sigma \cup \Sigma \times \Omega$  is unambiguous. Intuitively, this means that each output  $\mu \in \llbracket \mathcal{G} \rrbracket(w)$  can be produced in only one way. Remember that there are CFGs  $G$  with no unambiguous CFG  $G'$  *equivalent* to  $G$  (i.e., such that  $L(G') = L(G)$ ), and it is undecidable to check whether an input CFG is unambiguous, or has an equivalent unambiguous CFG. The same is immediately true for annotated grammars.

**Problem statement.** The goal of this paper is to study how to efficiently enumerate the annotations of an annotated grammar:

**Input:** An annotated grammar  $\mathcal{G}$  and a string  $s \in \Sigma^*$   
**Output:** Enumerate the outputs of  $\llbracket \mathcal{G} \rrbracket(s)$

We work in the standard computational model of Random Access Machines (RAM) with logarithmic word size and uniform cost measure, having addition and subtraction as basic operations [2]. The size of  $\mathcal{G}$  is measured as the sum of rule lengths.

As the set of outputs  $\llbracket \mathcal{G} \rrbracket(w)$  can be large, we work in the framework of *enumeration algorithms*. Such algorithms consist of two *phases*. First, in the *preprocessing phase*, the algorithm receives the input annotated grammar  $\mathcal{G}$  and string  $w$ , and produces some *index* structure  $D$ . The *preprocessing time* is the worst-case running time of this preprocessing phase, measured as a function of the input, i.e., in terms of  $w$  and  $\mathcal{G}$  when studying *combined complexity*, and in terms of  $w$  only when studying *data complexity*.

Second, in the *enumeration phase*, the algorithm can use  $\mathcal{G}$ ,  $w$ , and  $D$ , and must produce all outputs of  $\llbracket \mathcal{G} \rrbracket(w)$  one after the other and without repetitions. The *delay* of this phase is the worst-case time to produce any of the outputs, i.e., for  $N$  the number of outputs, if we call  $\text{time}_0$  the moment the preprocessing ends,  $\text{time}_i$  the moment the algorithm finishes producing the  $i$ -th output with  $1 \leq i \leq N$ , and  $\text{time}_{N+1}$  the moment when the algorithm terminates, then the delay is the maximum of the values  $(\text{time}_i - \text{time}_{i-1})$  for any  $0 < i \leq N + 1$ . We aim for *output-linear* delay [19] (also called linear delay [14], or constant delay [33] for constant-sized outputs), where the delay is linear in the size of each produced output, and is independent from the input (i.e., from  $w$  and  $\mathcal{G}$ ). The *memory usage* of the algorithm is the maximum memory used across both phases, including the size of  $D$ .

The ultimate goal of this paper is to find enumeration algorithms to enumerate the outputs of annotated grammars with linear preprocessing and output-linear delay. However, as we will see, this goal is not always realistic, so we will initially settle for a higher processing time, i.e., quadratic or cubic, before presenting classes with linear preprocessing in data complexity. We present our first results towards this goal in this next section.

### 3. Unambiguous Grammars

In this section we start presenting our results and show a first algorithm to enumerate the outputs of an annotated grammar on an input string. The algorithm applies to any unambiguous annotated grammar, and ensures cubic-time preprocessing and output-linear delay in data complexity; in terms of combined complexity, the preprocessing is linear on the grammar. This improves the result by Peterfreund [30], which had quintic-time preprocessing in data complexity.

The section is structured as follows. We present a general-purpose enumeration data structure called *enumerable sets*, which is the basis of our enumeration algorithms. We then introduce the *arity-two normal form* for annotated grammars, designed to ensure efficient enumeration, and which can be enforced in linear time. After this, we present our algorithm and state our main result (Theorem 3). Last, we state a conditional data complexity lower bound.

**Enumerable sets.** The preprocessing phase of our enumeration algorithm builds data structures representing the set of outputs to enumerate. For this, we essentially re-use the  $\varepsilon$ -ECS data structure of [27], but for convenience we present them in a self-contained way for our context, and name them *enumerable sets*. We now define them and state that we can enumerate their contents with output-linear delay (Theorem 1). The enumeration phase of our algorithm simply enumerates the outputs of an enumerable set using this delay guarantee.

An *enumerable set* is a representation of a set of strings over some alphabet  $\mathcal{O}$ . For our case, we want strings of  $\mathcal{O}^*$  to describe outputs, so  $\mathcal{O}$  consists of pairs of annotations with positions of the input string  $w$ , i.e.,  $\mathcal{O} := \Omega \times \{1, \dots, |w|\}$ .

The basic enumerable sets are:

- *empty*, the empty set;
- *singleton( $\varepsilon$ )*, the singleton set containing the empty string;
- *singleton( $x$ )* for  $x \in \mathcal{O}$ , the singleton set with the single-character string  $x$ .

Enumerable sets can be combined using operators to form more complex enumerable sets. The operators that we consider all take constant-time and are *fully-persistent* [15]. Specifically, given enumerable sets  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , combining them creates an enumerable set without modifying  $\mathcal{D}_1$  and  $\mathcal{D}_2$  (i.e., they can still be used in other operator applications). To make this possible, enumerable sets can share some components, e.g., some parts of the arguments  $\mathcal{D}_1$  and  $\mathcal{D}_2$  can be shared in memory, and the result can also have some parts that are shared with  $\mathcal{D}_1$  and  $\mathcal{D}_2$ . This is similar, e.g., to persistent lists, where we can only extend a list by adding an element to its head: this does not modify the original list, and returns a new list sharing some memory with the original list.

The two operators to combine enumerable sets are:

- The *union* operator *union( $\mathcal{D}_1, \mathcal{D}_2$ )* can be applied if the sets represented by  $\mathcal{D}_1$  and  $\mathcal{D}_2$  are disjoint, intuitively to avoid duplicates. It returns an enumerable set representing the union of these sets,

- The *product* operator  $\text{prod}(\mathcal{D}_1, \mathcal{D}_2)$  can be applied if there are no common letters in the strings of the sets represented by  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , i.e., if  $\mathcal{D}_1$  (resp.  $\mathcal{D}_2$ ) represents  $S_1$  (resp.  $S_2$ ) then the sets  $\{x_1 \in \mathcal{O} \mid x_1 \text{ occurs in some } w_1 \in S_1\}$  and  $\{x_2 \in \mathcal{O} \mid x_2 \text{ occurs in some } w_2 \in S_2\}$  are disjoint. Then, the operation returns an enumerable set  $\mathcal{D}$  which represents the concatenations of the strings in  $\mathcal{D}_1$  and in  $\mathcal{D}_2$ : formally, it represents  $S_1 \cdot S_2 = \{w_1 \cdot w_2 \mid w_1 \in S_1, w_2 \in S_2\}$ .

It is known that enumerable sets can be enumerated efficiently:

**Theorem 1** ([4, 27]). *We can implement enumerable sets such that:*

- *The enumerable sets empty, singleton( $\varepsilon$ ), and singleton( $x$ ) for  $x \in \mathcal{O}$  can be built in constant time;*
- *The union and product operations can be implemented in constant time and in a fully persistent way;*
- *Given an enumerable set, we can enumerate the strings it represents with output-linear delay and memory usage linear in the number of instructions used to build it.*

This was shown in [27, Theorem 9] with the data structure of *Enumerable Compact Sets* ( $\varepsilon$ -ECS). It also follows from the work in [4] on *zero-suppressed d-DNNF circuits*. For the reader's convenience, and to derive the bound on memory usage, we give a self-contained proof in Appendix A.1.

**Arity-two normal form.** Having presented enumerable sets, we now present the normal form to enforce on annotated grammars. Our results could be shown using the commonly known Chomsky normal form (CNF), but we cannot always obtain an equivalent CNF of linear size from a grammar. For this reason, we use a variant of CNF, the *arity-two normal form* (2NF) [25], which is intuitively like CNF but without disallowing rules of the form  $X \rightarrow Y$  or  $X \rightarrow \varepsilon$ . Formally, we say that an annotated grammar  $(V, \Sigma, \Omega, P, S)$  is in *arity-two normal form* (2NF) if the following hold:

- Every nonterminal  $X$  can derive some string, i.e., there exists  $\hat{w} \in (\Sigma \cup \Sigma \times \Omega)^*$  such that  $X \Rightarrow_{\mathcal{G}}^* \hat{w}$ .
- Every nonterminal  $X$  can be reached from the start symbol  $S$ , i.e., there exists  $\alpha, \beta \in (V \cup \Sigma \cup \Sigma \times \Omega)^*$  such that  $S \Rightarrow_{\mathcal{G}}^* \alpha X \beta$ .
- For every rule  $X \rightarrow \alpha$  in  $P$ , we have  $|\alpha| \leq 2$ , and if  $|\alpha| = 2$  then it consists of two nonterminals.

We can easily translate annotated grammars to 2NF, as in [25]:

**Proposition 2** ([25]). *Given any annotated grammar  $\mathcal{G}$ , we can compute in linear time an annotated grammar  $\mathcal{G}'$  in 2NF such that  $\mathcal{G}$  and  $\mathcal{G}'$  are equivalent. Furthermore, if  $\mathcal{G}$  is unambiguous then  $\mathcal{G}'$  is unambiguous as well.*

The proof is straightforward and shown in Appendix A.2. By Proposition 2, we assume that the input grammar  $\mathcal{G}$  is in 2NF.

We also compute in linear time some more information about  $\mathcal{G}$ . First, we precompute which nonterminals are *nullable*, i.e., are such that  $X \Rightarrow_{\mathcal{G}}^* \varepsilon$ : if we have a rule  $X \rightarrow \varepsilon$  then  $X$  is nullable, and if we have a rule  $X \rightarrow Y$  where  $Y$  is nullable or  $X \rightarrow YZ$  where  $Y$  and  $Z$  are nullable then  $X$  is also nullable. From this information, we further compute for each nonterminal  $Y$  a set  $D[Y]$  of all nonterminals  $X$  such that one of the following rules exist: a rule  $X \rightarrow Y$ , a rule  $X \rightarrow YZ$  where  $Z$  is nullable, or a rule  $X \rightarrow ZY$  where  $Z$  is nullable. We can clearly compute all of this in linear time (note that each rule contributes at most two entries to  $D$ ).

Second, as the grammar is assumed to be unambiguous, it also contains no *cycles*, i.e., there are no sequence of nonterminals  $X_1 \dots X_n$  such that for  $1 \leq i < n$ ,  $X_i \in D[X_{i+1}]$  and  $X_1 \in D[X_n]$ . Indeed, otherwise there would be infinitely many possible derivations of some string starting at  $X_1$ , contradicting the unambiguity of  $\mathcal{G}$ . Thus, we can sort the nonterminals of  $\mathcal{G}$  in *topological order*, by which we mean that when  $Y \in D[X]$  then  $Y$  is enumerated after  $X$ . Intuitively, when we consider a nonterminal  $X$ , we want to be done with processing the nonterminals  $Y$  such that  $X \rightarrow Y$  or  $X \rightarrow YZ$  or  $X \rightarrow ZY$  with  $Z$  nullable. This order can also be computed in linear time.

**Enumeration algorithm.** We now present the preprocessing phase of the enumeration algorithm, formalized as Algorithm 1 where the input string  $w = a_1 \dots a_n$  is assumed nonempty.

The principle of the algorithm is the following:

**Principle 1.** *For every triple of the form  $(i, j, X)$  with  $1 \leq i < j \leq n + 1$  and  $X \in N$ , the table cell  $\mathbb{I}[i][j][X]$  will contain an enumerable set representing the annotations of the string  $a_i \dots a_{j-1}$  that can be derived from symbol  $X$  in the grammar.*

These sets are initialized to be empty. In lines 5–11 of the algorithm, the cells  $\mathbb{I}[i][j][X]$  with  $j - i = 1$  are initialized to consider derivations via “simple rules” of the form  $X \rightarrow a$  or  $X \rightarrow (a, \sigma)$ . (For now, ignore the role of the `endl` table.) Note that the rules of the form  $X \rightarrow \varepsilon$  are considered when defining  $D$  and not further examined by the algorithm. At the end, line 24 returns the enumerable set for the annotations of the entire string derivable from the start symbol, i.e., the outputs of  $\mathcal{G}$  on  $w$ .

The main part of the algorithm consists in satisfying Principle 1 by adding the annotations corresponding to “complex” rules (i.e., of the form  $X \rightarrow YZ$  or  $X \rightarrow Y$ ). At the beginning of the algorithm the cells of the table  $\mathbb{I}$  might lack some annotations corresponding to complex rules, but each cell will be considered complete at some point during the execution, at which point it will satisfy Principle 1 and will not be modified anymore. We define the order in which the cells are considered complete as follows:  $(i, j, X) < (i', j', X')$  when  $j < j'$  or  $(j = j' \wedge i > i')$  or  $(j = j' \wedge i = i' \wedge X < X')$  where we order nonterminals  $X$  and  $X'$  following the topological order from  $D$ .

Consider the *complex derivations* starting from  $X$  of the string  $a_i \dots a_{j-1}$ , i.e., those that begin with a complex rule. We will see here how to reflect them in  $\mathbb{I}[i][j][X]$ . There are two kinds of complex derivations. The first kind is the derivations where we

first rewrite  $X$  to another nonterminal  $Z$  with a rule  $X \rightarrow Z$ , or by rewriting  $X$  to  $YZ$  or  $ZY$  but where  $Y$  is nullable and will be rewritten to  $\varepsilon$ . In these three cases, we have  $X \in D[Z]$ . Thus, we fill the index  $\mathbb{I}[i][j][X]$  with the contents of  $\mathbb{I}[i][j][Z]$ , which is already complete, for  $X \in D[Z]$  (lines 15-16).

The second kind of complex derivation begins with a complex rule  $X \rightarrow YZ$  where neither  $Y$  nor  $Z$  will be rewritten to  $\varepsilon$ . In this case, the set of annotations to add into  $\mathbb{I}[i][j][X]$  using this rule is the union of products of all the  $\mathbb{I}[i][k][Y]$  and  $\mathbb{I}[k][j][Z]$  where  $i < k < j$ . We have  $(i, k, Y) < (k, j, Z) < (i, j, X)$ , so we can fill  $\mathbb{I}[i][j][X]$  with the product of the contents of  $\mathbb{I}[i][k][Y]$  and  $\mathbb{I}[k][j][Z]$ , at the moment where  $\mathbb{I}[k][j][Z]$  is considered complete.

To summarize, from line 12 onwards, the algorithm considers the positions  $j$  in ascending order, and populates all cells  $\mathbb{I}[i][j][X]$  so that they are complete. To do so, we consider the triples  $(k, j, Z)$  by increasing order in our sorting criterion, i.e., by decreasing  $k$ , then increasing  $Z$  in the order of the topological sort. Whenever we consider a cell, it is complete, and we consider its contributions to cells of the form  $\mathbb{I}[i][j][X]$  with  $i = k$  using complex rules of the first kind (lines 15-16), and if it is non-empty we consider how to combine it with a neighboring cell (which is also complete and non-empty) as we explained previously, adding the results to a cell  $\mathbb{I}[i][j][X]$  with  $i < k$  which is not yet complete (lines 17-23).

We now explain the optimization involving the set `endl`. It is not necessary to achieve the cubic running time of this section, but is required for the quadratic bound in Section 4. The optimization is that, when processing the triple  $(k, j, Z)$  and the rule  $X \rightarrow YZ$ , we do not test all the possible cells  $\mathbb{I}[i][k][Y]$ , but only those that are non-empty. Indeed, if  $\mathbb{I}[i][k][Y]$  is empty, then the concatenation of  $\mathbb{I}[i][k][Y]$  with  $\mathbb{I}[k][j][Z]$  is also empty. Thus, we maintain the list `endl[k][Y]` of all the  $i$ 's to consider with  $i < k$ , i.e., those such that  $\mathbb{I}[i][k][Y]$  is non-empty. We initialize this list to be empty, add  $i$  to `endl[k][Y]` whenever  $\mathbb{I}[i][k][Y]$  becomes non-empty (at line 11 in the base case, or at line 21 before adding to an empty cell for the first time). Then, we only consider the indices  $i$  of this list to combine  $\mathbb{I}[i][k][Y]$  with another cell.

We now argue that our algorithm is correct, and in particular that (i) we satisfy Principle 1; that (ii) all the unions are disjoint, and that (iii) all the products involve enumerable sets on disjoint alphabets. One can establish (i) by showing by induction over cells that the invariant is correct when each cell is considered complete by our algorithm (and the cell is not changed afterwards). Knowing (i), the first violation of (ii) would witness that the same annotation of some factor  $a_i \cdots a_{j-1}$  can be derived in two different ways from a nonterminal  $X$ , contradicting unambiguity, so there are no violations of (ii). For (iii), we simply observe that, by (i),  $\mathbb{I}[i][j][X]$  only contains pairs of the form  $(\sigma, k)$  for some  $i \leq k < j$ , so we can indeed perform the product of  $\mathbb{I}[i][k][Y]$  and  $\mathbb{I}[k][j][Z]$ .

This establishes that the algorithm is correct. Now, the running time of the preprocessing phase of the algorithm is clearly in  $\mathcal{O}(n^3|\mathcal{G}|)$ , because (1) the `endl` lists are of size  $\mathcal{O}(n)$  at most, and (2) the consideration of all  $Z \in N$  and  $X \in D[Z]$  is in  $\mathcal{O}(|\mathcal{G}|)$ : every  $X \in D[Z]$  corresponds to a rule, so the consideration of all  $Z \in N$  and rules in  $\text{CRule}[Z]$  is in  $\mathcal{O}(|\mathcal{G}|)$ . The enumeration phase is then simply that of Theorem 1. Hence,

---

**Algorithm 1** Preprocessing phase: given a 2NF unambiguous annotated grammar  $\mathcal{G} = (N, \Sigma, \Omega, P, S)$  and a non-empty string  $w = a_1 \cdots a_n$ , compute an enumerable set representing  $\llbracket \mathcal{A} \rrbracket(w)$ .

---

```

1: I  $\leftarrow$  an array  $(n+1) \times (n+1) \times N$  initialized with empty
2: endl  $\leftarrow$  an array  $(n+1) \times N$  initialized with empty lists
3: CRule  $\leftarrow$  an array such that CRule[Z] =  $\{X \rightarrow YZ \in P\}$ 
4: D  $\leftarrow$  an array as described in the presentation of 2NF
5: for  $1 \leq i \leq n$  do
6:   if rule  $(X \rightarrow a_i)$  in  $P$  then
7:      $\mathbb{I}[i][i+1][X] \leftarrow \text{union}(\mathbb{I}[i][i+1][X], \text{singleton}(\varepsilon))$ 
8:   for rule  $(X \rightarrow (a_i, \sigma))$  in  $P$  do
9:      $\mathbb{I}[i][i+1][X] \leftarrow \text{union}(\mathbb{I}[i][i+1][X], \text{singleton}((\sigma, i)))$ 
10:  if  $\mathbb{I}[i][i+1][X] \neq \text{empty}$  then
11:    endl[ $i+1$ ][X].append( $i$ )
12: for  $j = 1$  to  $n+1$  do
13:   for  $k = j-1$  downto  $1$  do
14:    for nonterminal  $Z \in N$  in topological order do
15:     for nonterminal  $X \in D[Z]$  do
16:       $\mathbb{I}[k][j][X] \leftarrow \text{union}(\mathbb{I}[k][j][X], \mathbb{I}[k][j][Z])$ 
17:     if  $\mathbb{I}[k][j][Z] \neq \text{empty}$  then
18:       for rule  $(X \rightarrow YZ)$  in CRule[Z] do
19:        for  $i \in \text{endl}[k][Y]$  do
20:         if  $\mathbb{I}[i][j][X] = \text{empty}$  then
21:           endl[ $j$ ][X].append( $i$ )
22:          $\mathbb{I}[i][j][X] \leftarrow \text{union}(\mathbb{I}[i][j][X],$ 
23:            $\text{prod}(\mathbb{I}[i][k][Y], \mathbb{I}[k][j][Z]))$ 
24: return  $\mathbb{I}[1][n+1][S]$ 

```

---

we have shown that enumeration for unambiguous annotated grammars can be achieved with cubic time preprocessing and output-linear delay:

**Theorem 3.** *Given an unambiguous annotated grammar  $\mathcal{G}$  and an input string  $w$ , we can enumerate  $\llbracket \mathcal{G} \rrbracket(w)$  with preprocessing in  $\mathcal{O}(|w|^3 \cdot |\mathcal{G}|)$  (hence cubic in data complexity), and output-linear delay (independent from  $w$  or  $\mathcal{G}$ ). The memory usage is in  $\mathcal{O}(|w|^3 \cdot |\mathcal{G}|)$ .*

**Lower bounds.** We cannot show a lower bound that matches the complexity of our algorithm, but we can prove that we cannot achieve a preprocessing time better than the time to test whether a string is accepted by a CFG, which is essentially the complexity of Boolean matrix multiplication [1]:

**Proposition 4.** *Let  $\omega$  be the smallest value such that we can multiply two Boolean  $n \times n$  matrices in time  $\mathcal{O}(n^{\omega+o(1)})$ . Then for any  $c > 0$  there is an unambiguous annotated*

grammar  $\mathcal{G}$  such that, given an input string  $w$ , enumerating  $[[\mathcal{G}]](w)$  with output-linear delay requires a preprocessing time of  $\Omega(|w|^{\omega-c})$ .

## 4. Rigid Grammars

In the previous section, we have shown how to enumerate the output of unambiguous annotated grammars on strings, with output-linear delay and cubic preprocessing in the input string. This algorithm has two drawbacks: it requires us to impose that the grammar is unambiguous, and the cubic preprocessing may be expensive.

In this section, we introduce a new class of annotated grammars, called *rigid grammars*. Rigid grammars do not need to be unambiguous, but as we will show a rigid grammar can always be converted to additionally impose unambiguity. The point of rigid grammars is that we can show a quadratic bound on Algorithm 1 for them.

We first define rigid grammars in this section. We then state that we can impose unambiguity for rigid grammars, and derive some consequences about their expressiveness and the complexity of recognizing them. Last, we show a quadratic bound on the preprocessing time for output-linear enumeration for such grammars, and explain why a better bound would be challenging to achieve.

**Rigid grammars.** We first define the restricted notion of grammars that we study. Consider an annotated grammar  $\mathcal{G} = (N, \Sigma, \Omega, P, S)$ , and a string  $\gamma \in (\Sigma \cup (\Sigma \times \Omega) \cup N)^*$  of nonterminals and of terminals which may carry an annotation in  $\Omega$ . We will be interested in the *shape* of  $\gamma$ , written  $\text{shape}(\gamma)$ : it is the string over  $\{0, 1\}$  obtained by replacing every nonterminal of  $N$  in  $\gamma$  by 1 and replacing all terminals (annotated or not) by 0: note that  $|\text{shape}(\gamma)| = |\gamma|$ .

We then say that an annotated grammar  $\mathcal{G}$  is *rigid* if for every string  $w \in \Sigma^*$ , all derivations from the start symbol  $S$  of  $\mathcal{G}$  to an annotated string  $\hat{w}$  of  $w$  have the same sequence of shapes (remember that we only consider leftmost derivations). Formally, there exists a sequence  $s_1, \dots, s_k \in \{0, 1\}^*$  depending only on  $w$  such that for any derivation  $S = \alpha_1 \Rightarrow_{\mathcal{G}} \alpha_2 \Rightarrow_{\mathcal{G}} \dots \Rightarrow_{\mathcal{G}} \alpha_m = \hat{w}$  with  $\text{str}(\hat{w}) = w$ , we have  $m = k$  and  $\text{shape}(\alpha_i) = s_i$  for all  $1 \leq i \leq k$ .

Intuitively, the sequence of shapes of a derivation describes the skeleton of the corresponding derivation tree. Thus, a rigid annotated grammar is one where, for each unannotated string, all derivation trees for all annotations of the string are isomorphic (ignoring the labels of nonterminals and the annotation of terminals).

**Rigidity vs unambiguity.** Unambiguity and rigidity for annotated grammars seem incomparable: unambiguity imposes that every annotation is produced by only one derivation, whereas rigidity imposes that all derivations across all annotations have the same shape (but the same annotation may be obtained multiple times).

However, it turns out that, on rigid grammars, we can impose unambiguity without loss of generality: all rigid grammars can be converted to equivalent rigid and unambiguous grammars.

**Theorem 5.** *For any rigid grammar  $\mathcal{G}$  we can build an equivalent rigid and unambiguous grammar  $\mathcal{G}'$ . The transformation runs in exponential time, i.e., time  $\mathcal{O}(2^{|\mathcal{G}|^c})$  for some  $c > 0$ .*

The transformation to impose unambiguity goes via a notion of annotated pushdown automata (introduced in the next section), and is inspired by the determinization procedure for visibly pushdown automata [3], even though rigid grammars generally do not define visibly pushdown languages. The transformation comes at a cost, as it will generally blow up the size of the grammar exponentially.

**Expressiveness of rigid grammars.** Armed with Theorem 5, we study what is the expressive power of rigid grammars. For this, let us first go back to the setting without annotations. Theorem 5 tells us that for (unannotated) CFGs the rigidity requirement is equivalent to the usual unambiguity requirement: each accepted word has a unique derivation. Now, for the case of an annotated grammar  $\mathcal{G}$ , rigidity additionally imposes the requirement that all annotations of an input string have the same parse tree. In particular, the language of the strings where  $\mathcal{G}$  accepts some annotation must be recognizable by a rigid (unannotated) CFG, hence an unambiguous CFG (by Theorem 5). Formally:

**Proposition 6.** *For a rigid grammar  $\mathcal{G}$ , let  $L'$  be the set of strings with nonempty output, i.e.,  $L' = \{w \mid \llbracket \mathcal{G} \rrbracket(w) \neq \emptyset\}$ . Then  $L'$  is recognized by an unambiguous CFG.*

This yields concrete examples of languages (on the empty annotation alphabet) that cannot be recognized by a rigid annotated grammar, e.g., inherently ambiguous context-free languages such as  $L_a = \{a^i b^j c^k \mid i, j, k \geq 1 \wedge (i = j \vee j = k)\}$  on  $\{a, b, c\}^*$  [26]. Proposition 6 also implies that we cannot decide if the language of an annotated grammar can be expressed instead by a rigid grammar, or if an annotated grammar is rigid:

**Proposition 7.** *Given an unannotated grammar  $\mathcal{G}$ , it is undecidable to determine whether  $\mathcal{G}$  is rigid, and it is undecidable to determine whether there is some equivalent rigid grammar  $\mathcal{G}'$ .*

These undecidability results make rigid grammars less appealing, but note that our enumeration algorithm for such grammars applies in particular to decidable grammar classes which are designed to ensure rigidity. For instance, this would be the case of grammars arising from visibly pushdown automata, which we discuss in more detail in the next section.

**Enumeration algorithm.** We now give our algorithm with quadratic preprocessing time for rigid grammars. Given a rigid grammar, we first make it unambiguous if necessary, using Theorem 5, in exponential time in the input grammar. The result is a rigid and unambiguous annotated grammar. Now, we transform it in 2NF like in Section 4: this takes linear time, preserves unambiguity, and one can check that it also preserves rigidity.

Armed with our rigid and unambiguous grammar  $\mathcal{G}$  in 2NF, we can simply use Algorithm 1 to construct a data structure allowing us to enumerate the outputs with output-linear delay. But we now claim that Algorithm 1 runs in time  $\mathcal{O}(|\mathcal{G}| \cdot |w|^2)$  because  $\mathcal{G}$  is rigid.

For this, we study for every nonterminal  $X$  and pair  $1 \leq i \leq j \leq n + 1$  how many times we can consider the cell  $\mathbb{I}[i][j][X]$  in lines 20–23. Whenever we consider it, we witness the existence of a complex rule  $X \rightarrow YZ$  and a value  $k$  such that  $\mathbb{I}[i][k][Y]$  and  $\mathbb{I}[k][j][Z]$  are nonempty (the first is because  $i \in \text{endIn}[k][Z]$ ). Thus, we witness a derivation from  $X$  of some annotation of the string  $a_i \cdots a_{j-1}$  that starts with a rule  $X \rightarrow YZ$  where  $Y$  derives some annotation of the string  $a_i \cdots a_{k-1}$  and  $Z$  derives some annotation of the string  $a_k \cdots a_{j-1}$ . We now claim that, for  $(i, j, X)$ , the rigidity of the grammar ensures that there is only one such value  $k$ . Indeed, assume by contradiction that we have two rules  $X \rightarrow YZ$  and  $X \rightarrow Y'Z'$  and two values  $i \leq k < k' \leq j$  such that  $Y$  and  $Y'$  respectively derive some annotation of the strings  $a_i \cdots a_{k-1}$  and  $a_i \cdots a_{k'-1}$ , and  $Z$  and  $Z'$  respectively derive some annotation of the strings  $a_k \cdots a_{j-1}$  and  $a_{k'} \cdots a_{j-1}$ . Then once we are done rewriting  $Y$  and all the nonterminals that it generates in the first derivation, we obtain a different shape from what we obtain after rewriting  $Y'$  and all the nonterminals it generates in the second derivation, contradicting the rigidity of the grammar.

Thus, whenever we consider the cell  $\mathbb{I}[i][j][X]$  in lines 21–23, it is for one value of  $k$  which is unique for  $(i, j, X)$ , and we thus consider the cell once at most for every complex rule of the grammar with  $X$  as left-hand-side. Thus, we consider the cells of  $\mathbb{I}[i][j]$  at most  $|\mathcal{G}|$  times in total. As there are  $\mathcal{O}(n^2)$  pairs  $(i, j)$ , this ensures that the total running time of the innermost for loop (lines 19–23), and that of the entire algorithm, is indeed in  $\mathcal{O}(|\mathcal{G}| \cdot |w|^2)$ :

**Theorem 8.** *Given a rigid annotated grammar  $\mathcal{G}$  and an input string  $w$ , we can enumerate  $\llbracket \mathcal{G} \rrbracket(w)$  with preprocessing in  $\mathcal{O}(|w|^2)$  data complexity and output-linear delay (independent from  $w$  or  $\mathcal{G}$ ). The combined complexity of the preprocessing is  $\mathcal{O}(2^{|\mathcal{G}|^c} \cdot |w|^2)$  for some  $c > 0$ , or  $\mathcal{O}(|\mathcal{G}| \cdot |w|^2)$  if  $\mathcal{G}$  is additionally assumed to be unambiguous.*

**Optimality.** We now turn to the question of whether the quadratic preprocessing time for rigid grammars is optimal. For this, we notice that the parsing of (unannotated) unambiguous grammars can be performed in quadratic time, but the question of finding a better algorithm was open as of 2012 [32]. Now, this is a special case of our problem, because an unannotated unambiguous grammar is in particular a rigid and unambiguous annotated grammar, and enumerating the outputs of an unannotated grammar just means deciding in constant time after the preprocessing whether the input unannotated string is accepted or not. Thus:

**Proposition 9.** *Any algorithm to enumerate the accepted outputs of a rigid annotated grammar can be used to test if an input string is accepted by an unambiguous unannotated grammar, with same complexity as that of the preprocessing phase.*

For this reason, we leave open the question of whether a better than quadratic preprocessing time can be achieved in this case.

## 5. Pushdown Annotators

We have presented an enumeration algorithm for annotated grammars that achieves quadratic-time preprocessing and output linear delay on rigid annotated grammars. We now study whether the bound can be improved even further to achieve linear-time preprocessing and output-linear delay, which is the best possible data complexity bound in our model.

To achieve this, it is natural to look for a class of grammars having some “deterministic” behavior. Unfortunately, grammars are not convenient for this purpose, and so we move to the equivalent model of pushdown automata. We thus introduce *pushdown annotators* and show that they are equally expressive to annotated grammars. We present syntactic restrictions on pushdown annotators that ensure quadratic-time preprocessing, similarly to rigid grammars. Then, we propose additional deterministic conditions on pushdown annotators that allow for linear-time preprocessing.

**Pushdown annotators.** A *pushdown annotator* (PDAnn) is a tuple  $\mathcal{P} = (Q, \Sigma, \Omega, \Gamma, \Delta, q_0, F)$  where  $Q$  is a finite set of *states*,  $\Sigma$  is the alphabet,  $\Omega$  is a finite set of annotations,  $\Gamma$  is a finite set of *stack symbols*,  $q_0 \in Q$  is the *initial state*, and  $F \subseteq Q$  are the *final states*. We assume that the set  $\Gamma$  of stack symbols is disjoint from  $(\Sigma \cup \Sigma \times \Omega)$ . Finally,  $\Delta$  is a finite set of *transitions* that are of the following kinds:

- *Read-write transitions* of the form  $(p, (a, \sigma), q) \in Q \times (\Sigma \times \Omega) \times Q$ , meaning that, if the next letter of the string is  $a$ , the annotator can go from states  $p$  to  $q$  while reading that letter and writing the annotation  $\sigma$ ;
- *Read-only transitions* of the form  $(p, a, q) \in Q \times \Sigma \times Q$ , meaning that the annotator can go from  $p$  to  $q$  while reading  $a$ ;
- *Push transitions* of the form  $(p, q, \gamma) \in Q \times (Q \times \Gamma)$ , meaning that the annotator can go from  $p$  to  $q$  while pushing the symbol  $\gamma$  on the stack;
- *Pop transitions* of the form  $(p, \gamma, q) \in (Q \times \Gamma) \times Q$ , meaning that, if the topmost symbol of the stack is  $\gamma$ , the annotator can go from  $p$  to  $q$  while removing this topmost symbol  $\gamma$ .

We now give the semantics of PDAnns. Fix a string  $w = w_1 \cdots w_n \in \Sigma^*$ . A *configuration* of  $\mathcal{P}$  over  $w$  is a pair  $C = (q, i) \in Q \times [0, n]$  of the current state and position in  $w$ . An *instantaneous description* of  $\mathcal{P}$  is a pair  $(C, \alpha)$  where  $C$  is a configuration and  $\alpha \in \Gamma^*$  describes the stack. A *run* of  $\mathcal{P}$  over  $w$  is a sequence:

$$\rho := (C_0, \alpha_0) \xrightarrow{t_1} (C_1, \alpha_1) \xrightarrow{t_2} \dots \xrightarrow{t_m} (C_m, \alpha_m) \quad (\dagger)$$

such that  $C_0 = (q_0, 0)$  and  $\alpha_0 = \varepsilon$ , each  $t_k$  is a transition in  $\Delta$ , and for each  $k \in [1, m]$  the following hold:

- if  $t_k$  is a read-write transition  $(p, (a, \sigma), q)$  or a read-only transition  $(p, a, q)$ , then  $\alpha_k = \alpha_{k-1}$ ,  $C_{k-1} = (p, i-1)$ ,  $C_k = (q, i)$  and  $a = a_i$  for some  $i \in [1, n]$ ;

- if  $t_k$  is a push transition  $(p, q, \gamma)$ , then  $\alpha_k = \alpha_{k-1}\gamma$  and for some  $i \in [1, n]$ ,  $C_{k-1} = (p, i)$ ,  $C_k = (q, i)$ ; and
- if  $t_k$  is a pop transition  $(p, \gamma, q)$ , then  $\alpha_{k-1} = \alpha_k\gamma$  and for some  $i \in [1, n]$ ,  $C_{k-1} = (p, i)$ ,  $C_k = (q, i)$ .

We say that  $\rho$  is *accepting* if  $(C_m, \alpha_m) = ((q_f, n), \varepsilon)$  for some  $q_f \in F$ . We define the *annotation* of  $\rho$  as  $\text{ann}(\rho) = \text{ann}(t_1, C_1) \cdots \text{ann}(t_m, C_m)$  such that  $\text{ann}(t, C) = \varepsilon$  if  $t$  is a push, pop, or a read-only transition  $(p, a, q)$ , and  $\text{ann}(t, C) = (i, \sigma)$  if  $t$  is a read-write transition  $(p, (a, \sigma), q)$  and  $C = (q, i)$ . Finally, we define the function  $\llbracket \mathcal{P} \rrbracket$  that maps any  $w \in \Sigma^*$  to its set of outputs:

$$\llbracket \mathcal{P} \rrbracket(w) = \{\text{ann}(\rho) \mid \rho \text{ is an accepting run of } \mathcal{P} \text{ over } w\}.$$

Similarly to annotated grammars, we say that  $\mathcal{P}$  is *unambiguous* if for every  $w \in \Sigma^*$  and output  $\mu$ , there exists at most one accepting run  $\rho$  of  $\mathcal{P}$  over  $w$  such that  $\text{ann}(\rho) = \mu$ .

One can alternatively see a PDAnn as a *pushdown transducer* [12], which is the standard way to extend automata to have an output. However, an important difference is that a PDAnn concisely represents outputs by only writing the annotations and their positions: this can be much smaller than the input string, and cannot easily be encoded as a transducer on a finite alphabet. For instance, where a PDAnn can produce an output such as  $(2, \sigma), (5, \sigma')$ , a transducer would either write  $\sigma\sigma'$  (losing the position information) or  $\_ \sigma \_ \_ \sigma'$  (whose length is always linear in the input) for a special symbol  $\_$ .

**Profiled PDAnns and annotated grammars.** To define the analogue of rigid annotated grammars on PDAnn, we will study the *stack profile* (or simply *profile*) of PDAnn runs, which is informally the sequence of all stack heights. Formally, let  $\mathcal{P}$  be a PDAnn,  $w$  be a string, and consider a run  $\rho$  of  $\mathcal{P}$  over  $w$  like in (†). The *profile*  $\pi$  of  $\rho$  is the sequence  $\pi := |\alpha_0|, \dots, |\alpha_m|$ . We then introduce *profiled* PDAnns by requiring that all accepting runs of the PDAnn on an input string have the same profile (no matter their output). Formally, we say that a PDAnn  $\mathcal{P}$  is *profiled* if, for every string  $w$ , all accepting runs of  $\mathcal{P}$  over  $w$  have the same profile.

As usual for context-free grammars and pushdown automata, the formalisms of annotated grammars and PDAnn have the same expressive power. We call two annotated grammars  $\mathcal{G}$  and  $\mathcal{G}'$  *equivalent* if they define the same functions, i.e.,  $\llbracket \mathcal{G} \rrbracket = \llbracket \mathcal{G}' \rrbracket$ , and extend this notion to PDAnn in the expected way. We then have:

**Proposition 10.** *Annotated grammars and PDAnn are equally expressive. Specifically, for any annotated grammar  $\mathcal{G}$ , we can build an equivalent PDAnn  $\mathcal{P}$  in polynomial time, and vice versa. Further,  $\mathcal{G}$  is unambiguous (resp., rigid) iff  $\mathcal{P}$  is unambiguous (resp., profiled).*

Let us now study the enumeration for PDAnns. We know that the problem for unambiguous PDAnns can be solved via Proposition 10 with cubic-time preprocessing in data complexity and output-linear delay (with Theorem 3). We know that profiled PDAnns can be made unambiguous (via Proposition 10 and Theorem 5) and so that we

can solve enumeration for them in quadratic-time preprocessing in data complexity and output-linear delay (using Theorem 8). We now show that, if we are given a profile of an unambiguous PDAnn  $\mathcal{P}$  on an input string  $w$ , we can use it as a guide to enumerate with linear preprocessing and output-linear delay the set  $\llbracket \mathcal{P} \rrbracket_\pi(w)$  of annotations for that profile, i.e., all  $\text{ann}(\rho)$  such that  $\rho$  is an accepting run with profile  $\pi$  of  $\mathcal{P}$  over  $w$ . Formally:

**Lemma 1.** *Given an unambiguous PDAnn  $\mathcal{P}$ , there exists an enumeration algorithm that receives as input a string  $w$  and a profile  $\pi$  of  $\mathcal{P}$  over  $w$ , and enumerates  $\llbracket \mathcal{P} \rrbracket_\pi(w)$  with output-linear delay after linear-time preprocessing in data complexity.*

This result implies that we could achieve linear-time enumeration over profiled PDAnn if we could easily discover their (unique) profile. We achieve this with *profiled-deterministic PDAnns*.

**Profiled-deterministic PDAnn.** Let  $\mathcal{P} = (Q, \Sigma, \Omega, \Gamma, \Delta, q_0, F)$  be a PDAnn. We say that a PDAnn  $\mathcal{P}$  is *profiled-deterministic* if, for any string  $w \in \Sigma^*$ , for any two partial runs  $\rho$  and  $\rho'$  of  $\mathcal{P}$  over  $w$  with the same length,  $\rho$  and  $\rho'$  have the same profile.

The relationship between profiled-deterministic PDAnns and deterministic pushdown automata (formally defined in Appendix C) is similar to the relationship between profiled PDAnns and unambiguous pushdown automata (the latter relationship was stated as Proposition 6 in the context of grammars). Specifically:

**Proposition 11.** *For a profiled-deterministic PDAnn  $\mathcal{P}$ , let  $L'$  be the set of strings with nonempty output, i.e.,  $L' = \{w \mid \llbracket \mathcal{P} \rrbracket(w) \neq \emptyset\}$ . Then  $L'$  is recognized by a deterministic pushdown automaton.*

This result gives a concrete picture of the expressive power of profiled-deterministic PDAnn  $\mathcal{A}$ , i.e., as acceptors they are more powerful than the class of *visibly pushdown automata* [3], where each alphabet letter must have a specific effect on the profile. Profiled-deterministic PDAnn are also reminiscent of the *height-determinism* notion introduced for pushdown automata [28], but extend this with the support of annotations.

Profiled-deterministic PDAnn are designed to ensure that they have only one profile (i.e., they are profiled), and further that their unique profile can be constructed in linear time:

**Proposition 12.** *A profiled-deterministic PDAnn  $\mathcal{P}$  is always profiled, and given a string  $w$ , the unique profile of accepting runs of  $\mathcal{P}$  over  $w$  can be computed in linear time in  $w$ .*

Together with Lemma 1, this yields:

**Corollary 1.** *Let  $\mathcal{P}$  be a profiled-deterministic PDAnn. Then for every string  $w$  the set  $\llbracket \mathcal{P} \rrbracket(w)$  can be enumerated with output-linear delay after linear-time preprocessing in data complexity.*

## 6. Application: Document Spanners

We have presented our enumeration results for annotated grammars and pushdown annotators. We conclude the paper by applying them to the standard context of *document spanners* [17] and to the *extraction grammars* recently introduced in [30].

**Mappings, spanners, extraction grammars.** Our paper studies strings  $s \in \Sigma^*$  with letters *annotated* by an annotation from a finite set, but document spanners work with *mappings* selecting so-called *spans* of  $s$ . Formally, a *span* of  $s$  is a pair  $[i, j)$  with  $1 \leq i \leq j \leq |s| + 1$  describing a possibly empty factor of  $s$ . For a finite set of *variables*  $\mathcal{X}$ , a *document spanner* (or just *spanner*) is a function that maps every string  $s \in \Sigma^*$  to a set of *mappings*, where a mapping (intuitively denoting an extraction result) assigns each variable of  $\mathcal{X}$  to a span of  $s$ .

Peterfreund [30] specifies spanners via *extraction grammars*. An *extraction grammar*  $\mathcal{H}$  is simply a CFG over the alphabet  $\Sigma$  extended with *variable operations* of the form  $\vdash_x$  and  $\dashv_x$  for  $x \in \mathcal{X}$ , intuitively denoting the beginning and end of a span for  $x$ . Any such grammar  $\mathcal{H}$  denotes a *language*  $L(\mathcal{H})$  of words with variable operations, called *ref-words*, and gives a function  $\llbracket \mathcal{H} \rrbracket$  associating every string  $w \in \Sigma^*$  to the set of mappings  $\llbracket \mathcal{H} \rrbracket(w)$  defined as follows. We consider every ref-word  $\hat{w}$  of  $L(\mathcal{H})$  whose restriction to  $\Sigma$  is  $w$  and which is *valid*, i.e., for every  $x \in \mathcal{X}$  there is one occurrence of  $\vdash_x$  followed by one occurrence of  $\dashv_x$  in  $\hat{w}$ . Every valid ref-word  $\hat{w}$  defines a mapping that intuitively associates each variable  $x \in \mathcal{X}$  to the span of the characters of  $w$  between  $\vdash_x$  and  $\dashv_x$  in  $\hat{w}$ .

Thus, extraction grammars are like annotated grammars but with variable operations that describe span endpoints (whereas our annotations are arbitrary), and that are expressed as separate variable operation characters (not annotations of existing letters).

**Comparing both formalisms.** Given an extraction grammar  $\mathcal{H}$  on alphabet  $\Sigma$  and with variables  $\mathcal{X}$ , an annotated grammar  $\mathcal{G}$  is *equivalent* to  $\mathcal{H}$  if it is over the alphabet  $\Sigma \cup \{\#\}$  (with a special end-of-word character  $\#$  to account for indexing differences), if its annotation set  $\Omega$  is the powerset of the set  $\bigcup_{x \in \mathcal{X}} \{\vdash_x, \dashv_x\}$  of markers symbols, and if the following holds: for any  $w \in \Sigma^*$ , the outputs  $\llbracket \mathcal{G} \rrbracket(w\#)$  are in one-to-one correspondence with the mappings of  $\llbracket \mathcal{H} \rrbracket(w)$  in the expected way (the formal definition is in Appendix D.1). We show that every extraction grammar has an equivalent annotation grammar in this sense, and the translation further preserves unambiguity:

**Proposition 13.** *Given any extraction grammar  $\mathcal{H}$  with  $k$  variables, we can build an equivalent annotated grammar  $\mathcal{G}$  in time  $\mathcal{O}(9^{3k} \cdot |\mathcal{H}|^2)$ . Moreover, if  $\mathcal{H}$  is unambiguous then so is  $\mathcal{G}$ .*

Hence, our formalism of annotated grammars captures that of extraction grammars. Unfortunately, the translation is exponential, intuitively because  $\Omega$  must cover all possible sets of variable operations: we explain in Appendix D.3 why we believe this to be unavoidable. We note that, in exchange for this, annotated grammars are strictly more *expressive*: each output can annotate an arbitrary number of positions in the string (e.g., every other character), unlike extraction grammars whose mappings have a fixed number of variables.

**Enumeration for extraction grammars.** As extraction grammars can be rewritten to annotated grammars in an unambiguity-preserving way (Proposition 13), we can derive from Theorem 3 an enumeration result for unambiguous extraction grammars with cubic preprocessing time in data complexity.

**Theorem 14.** *Given an unambiguous extraction grammar  $\mathcal{H}$  with  $k$  variables and a string  $s$ , we can enumerate the mappings of  $\llbracket \mathcal{H} \rrbracket(s)$  with preprocessing time  $\mathcal{O}(9^{3k} \cdot |\mathcal{H}|^2 \cdot |s|^3)$  (hence, cubic in  $|s|$ ), and with output-linear delay (independent from  $s$ ,  $k$ , or  $\mathcal{H}$ ).*

In data complexity, this improves over the result of [30] for unambiguous extraction grammars, whose preprocessing time is  $\mathcal{O}(9^{2k} \cdot |\mathcal{H}|^2 \cdot |s|^5)$ , i.e., our data complexity is cubic instead of quintic. We leave to future work a study of enumeration results for restricted classes of extraction grammars via Theorems 3 and 8.

## 7. Conclusions and Future Work

We have presented our formalism of annotated grammars and our results on the efficient enumeration of all annotations of an input string. Our results achieve output-linear delay, and cubic-time preprocessing if the grammar is unambiguous, quadratic-time if it is rigid, and linear-time for profiled-deterministic PDAnns.

The main question left open by our work is that of the precise complexity of this task, depending on the grammar formalism. For instance, can we improve the  $\mathcal{O}(n^3)$  algorithm to match the complexity of Valiant’s parser? For which grammar classes can we extend the linear-time preprocessing approach? We believe, however, that a complete classification is out of reach, given that classifying the fine-grained complexity of parsing is still open to a large extent even in the case of unannotated CFGs.

**Acknowledgements.** Amarilli was partially supported by the ANR project EQUUS ANR-19-CE48-0019 and by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – 431183758. Muñoz and Riveros were funded by ANID - Millennium Science Initiative Program - Code ICN17\_002.

## References

- [1] A. Abboud, A. Backurs, and V. V. Williams. If the current clique algorithms are optimal, so is Valiant’s parser. *SIAM J. Comput.*, 47(6), 2018.
- [2] A. V. Aho, J. E. Hopcroft, and J. D. Ullman. *The design and analysis of computer algorithms*. Addison-Wesley, 1974.
- [3] R. Alur and P. Madhusudan. Visibly pushdown languages. In *STOC*, 2004.
- [4] A. Amarilli, P. Bourhis, L. Jachiet, and S. Mengel. A circuit-based approach to efficient enumeration. In *ICALP*, 2017.
- [5] A. Amarilli, P. Bourhis, S. Mengel, and M. Niewerth. Constant-delay enumeration for nondeterministic document spanners. In *ICDT*, 2019.

- [6] A. Amarilli, P. Bourhis, S. Mengel, and M. Niewerth. Enumeration on trees with tractable combined complexity and efficient updates. In *PODS*, 2019.
- [7] A. Amarilli, P. Bourhis, S. Mengel, and M. Niewerth. Constant-delay enumeration for nondeterministic document spanners. *TODS*, 2020.
- [8] A. Amarilli, L. Jachiet, M. Muñoz, and C. Riveros. Efficient enumeration for annotated grammars. In *PODS*, 2022.
- [9] G. Bagan. MSO queries on tree decomposable structures are computable with linear delay. In *CSL*, 2006.
- [10] G. Bagan, A. Durand, and E. Grandjean. On acyclic conjunctive queries and constant delay enumeration. In *International Workshop on Computer Science Logic*. Springer, 2007.
- [11] C. Berkholz, J. Keppeler, and N. Schweikardt. Answering conjunctive queries under updates. In *PODS*, 2017.
- [12] J. Berstel. *Transductions and context-free languages*. Springer-Verlag, 2013.
- [13] N. Chomsky and M. P. Schützenberger. The algebraic theory of context-free languages. In *Studies in Logic and the Foundations of Mathematics*, volume 26. 1959.
- [14] B. Courcelle. Linear delay enumeration and monadic second-order logic. *Discret. Appl. Math.*, 157(12), 2009.
- [15] J. R. Driscoll, N. Sarnak, D. D. Sleator, and R. E. Tarjan. Making data structures persistent. In *STOC*, 1986.
- [16] A. Durand and E. Grandjean. First-order queries on structures of bounded degree are computable with constant delay. *ACM Transactions on Computational Logic (TOCL)*, 8(4), 2007.
- [17] R. Fagin, B. Kimelfeld, F. Reiss, and S. Vansummeren. Document spanners: A formal approach to information extraction. *JACM*, 62(2), 2015.
- [18] F. Florenzano, C. Riveros, M. Ugarte, S. Vansummeren, and D. Vrgoc. Constant delay algorithms for regular document spanners. In *PODS*, 2018.
- [19] F. Florenzano, C. Riveros, M. Ugarte, S. Vansummeren, and D. Vrgoc. Efficient enumeration algorithms for regular document spanners. *TODS*, 45(1), 2020.
- [20] S. Ginsburg and J. Ullian. Ambiguity in context free languages. *JACM*, 13(1), 1966.
- [21] O. Hartig and J. Pérez. Semantics and complexity of GraphQL. In *TheWebConf*, 2018.
- [22] M. Idris, M. Ugarte, and S. Vansummeren. The dynamic Yannakakis algorithm: Compact and efficient query processing under updates. In *SIGMOD*, 2017.

- [23] A. Kara, M. Nikolic, D. Olteanu, and H. Zhang. Trade-offs in static and dynamic evaluation of hierarchical queries. In *PODS*, 2020.
- [24] M. Kröll, R. Pichler, and S. Skritek. On the complexity of enumerating the answers to well-designed pattern trees. In *ICDT*, 2016.
- [25] M. Lange and H. Leiß. To CNF or not to CNF? An efficient yet presentable version of the CYK algorithm. *Informatica Didactica*, 8(2009), 2009.
- [26] H. A. Maurer. A direct proof of the inherent ambiguity of a simple context-free language. *JACM*, 16(2), 1969.
- [27] M. Muñoz and C. Riveros. Streaming enumeration on nested documents. In *ICDT*, 2022.
- [28] D. Nowotka and J. Srba. Height-deterministic pushdown automata. In *MFCS*, 2007.
- [29] L. Peterfreund. *The complexity of relational queries over extractions from text*. PhD thesis, Technion, 2019.
- [30] L. Peterfreund. Grammars for document spanners. In *ICDT*, 2021.
- [31] L. Peterfreund, B. ten Cate, R. Fagin, and B. Kimelfeld. Recursive programs for document spanners. In *ICDT*, 2019.
- [32] S. Schmitz. Can all unambiguous grammars be parsed in linear time? Theoretical Computer Science Stack Exchange, 2012. Version: 2012-03-02.
- [33] L. Segoufin. Enumerating with constant delay the answers to a query. In *ICDT*, 2013.

## A. Proofs of Section 3

### A.1. Proof of Theorem 1

We give here a self-contained argument for the proof of Theorem 1. A more detailed presentation can be found in [4] or [27].

We will first show how to build a structure that only represents sets that are non-empty and that cannot contain the empty string ( $\text{singleton}(\varepsilon)$ ). Then we will explain how we can extend this structure to support the empty set, and to support  $\text{singleton}(\varepsilon)$ .

**The DAG structure of enumerable sets.** As enumerable sets are a fully-persistent structure, we can represent our data structure as a global directed acyclic graph storing all enumerable sets constructed so far, to which we can only add new nodes pointing to existing nodes. In particular, for any node that is built, its set of descendants is immutable. Each node  $n$  will be associated to a set  $S(n)$ . In the global DAG, there are three kinds of nodes:

- *Singleton nodes*, a singleton node  $n$  carries a value  $x \in \mathcal{O}$  and as its name suggests,  $S(n) = \{x\}$ , i.e.  $n$  represents the set containing a single value which is the single-character string  $x$ .
- *Product nodes*, a product node has two outgoing edges. Intuitively, a product node  $n$  is such that  $S(n) = S(n_1) \otimes S(n_2)$ , where  $\otimes$  denotes the product operation on sets defined on the main text, and  $n_1$  and  $n_2$  are the two inputs.
- *Union nodes*, a union node has two or three outgoing edges. The first outgoing edge necessarily points to a node which is a singleton or product node. Intuitively, a union node  $n$  is such that  $S(n)$  is  $\cup_c S(c)$  where  $c$  ranges over the nodes to which  $n$  has an outgoing edge. In our construction we will make sure that union nodes with three children are never used directly to represent an enumerable sets but only used internally, in particular we never start the enumeration on such nodes.

**Implementing the various construction operations.** We now present how to implement the various operations for enumerable sets. For  $\text{singleton}$  and  $\text{prod}$  the operation is easy: for  $\text{singleton}(x)$ , we create a leaf node carrying the value  $\text{singleton}(x)$  and return it; for  $\text{prod}(n_1, n_2)$  we create a product node with  $n_1$  and  $n_2$  as outgoing edges and return it.

The union operation is the most complicated operation by far. Let  $n_1$  and  $n_2$  be the arguments. Remember that  $n_1$  and  $n_2$  cannot be union nodes that have three outgoing edges, because  $n_1$  and  $n_2$  are nodes that represent enumerable sets. We implement the union operator on  $n_1$  and  $n_2$  in the following way

- If one of the arguments is a singleton or product node, then the result is a union node with two children, having the singleton or product node as first child and the other argument as second child.

- If both arguments are union nodes  $n$  and  $n'$ , let  $n_1$  and  $n_2$  be the children of  $n$  with  $n_1$  a singleton or product node, and define likewise  $n'_1$  and  $n'_2$ . The result is a union node  $n''$  having as children  $n_1$  and a fresh union node  $n'''$  having as children  $n'_1$ ,  $n_2$ , and  $n'_2$ . Notice the returned node,  $n''$  has only two outgoing edges satisfying the condition that union nodes with three children are never used directly.

One can verify that the construction satisfies the invariants, i.e., that  $S(n)$  is as prescribed. This is clear for the product operator, and for the union operator the only difficult case is the union of two union nodes. Then, given that we had  $S(n) = S(n_1) \cup S(n_2)$ , and  $S(n') = S(n'_1) \cup S(n'_2)$ , the result is such that  $S(n'') = S(n_1) \cup S(n'_1) \cup S(n_2) \cup S(n'_2)$  which is correct. Furthermore, it is clear that we never return a union node with three children as the result of an operator, and that the unions that we create have the right number of children and that their first child is always a singleton or product node.

What is more, all operators are indeed in constant-time.

**Reduction to the enumeration of enumerable sets starting with a product node.** For the sake of simplicity, we will suppose in the rest of this section that the node to be enumerated is always a product node. If it is not we can easily create the product with a singleton node  $\$,$  enumerate the strings represented by this new node and remove from each string produced the last character (which will be the  $\$$ ).

**Auxiliary procedure: enumerating the exits of union nodes.** Given a node  $n,$  we say that an *exit* of  $n$  is a product or a singleton node reachable from  $n$  going only through union nodes. For our enumeration scheme, we will need a structure capable of enumerating the exits of a node with constant-delay between the any two exits, such that the structure can be initialized in constant time.

The set of exits for a node  $n$  can be expressed recursively: it is either  $n$  when  $n$  itself is a singleton or product node or it is the union of exits of the children of  $n$  otherwise. Our enumeration thus performs a depth first search but with a constant delay guarantee.

To do this, we will rely (as in a normal depth first search) on a stack  $S$  of nodes “to be enumerated”. At the beginning of the enumeration,  $S$  contains the node  $n,$  and the enumeration ends whenever  $S$  becomes empty. When looking for the next exit, we pop one node  $v$  from  $S.$  If  $v$  is a singleton or product node, the next exit is  $v.$  If  $v$  is a union node, we push back the children of  $v$  on the stack while making sure that the first child of  $v$  is on top. Then, we reiterate the procedure to find the next exit.

It is clear that this scheme enumerates all the exits and only once. To show that this procedure is indeed constant-delay, notice that, by design, the first child of a union node is never a union node itself. Therefore if we pop a union node, the next pop cannot be a union node again. Therefore once every two pops we produce at least one exit and therefore the procedure outputs the exits with constant delay.

**Enumeration state.** In the design of our enumeration algorithm, we must explain what is the *state* stored by the enumeration between any two outputs. This state is an *enumeration tree.* In an enumeration tree, the nodes are of one of the following two types: *concatenation* nodes (similar to the product nodes in our DAG) and *pointer* nodes, pointing to either a product or a singleton node in the DAG. For each concatenation node  $n$  we will store multiple pieces of information: two children, two auxiliary data structures

for exit enumeration, and as a pointer to the product node that  $n$  corresponds to. For a pointer node we will simply keep a pointer to the node in the DAG of enumerable sets, which is either a product node or a union node.

An enumeration-tree is said to be *output-ready* when no pointer nodes are pointing to product nodes. This means that the children of concatenation nodes are either concatenation nodes themselves or pointer nodes pointing to singleton nodes.

**Overall structure of the enumeration.** Each step of the enumeration is decomposed into three sub-phases: (1.) an *unfolding sub-phase* where we “unfold” the pointer nodes pointing to product nodes into concatenation nodes (resulting in an output-ready enumeration tree), then (2.) an *output sub-phase* where we compute the output, and finally (3.) a *pruning sub-phase* where we change the tree to point to the next solution by removing nodes where there are no more outputs. Each sub-phase will take a time proportional to the output currently being produced. After the pruning sub-phase we return the computed output. Thus, we ensure output-linear delay.

At the beginning of the enumeration for a enumerable set, remembering our assumption that the enumerable set corresponds to a product node  $n$ , the state of the enumeration is a pointer node pointing to  $n$ . The enumeration ends whenever  $n$  has no next solution (we explain below how this can be computed).

Note that we need to separate the pruning sub-phase of the current output from the next unfolding sub-phase, because the running time of the pruning sub-phase depends on the size of the current output, while the running time of the unfolding sub-phase of the next output depends on the size of the next output. Thus, without running the pruning sub-phase with the current output, we could violate the output-linear delay guarantee, because two successive outputs may have very different sizes.

**Unfolding sub-phase.** During the unfolding sub-phase, we explore the enumeration tree to look for a *pointer* node. Each time we find a pointer node  $p$  pointing to a product node  $n$  with children  $n_1$  and  $n_2$ , we replace  $p$  by a new concatenation node  $c$ . The two exit-enumeration structures  $S_1$  and  $S_2$  of the new concatenation node  $c$  enumerate the exits of the inputs  $n_1$  and  $n_2$  to the product node  $n$ . We use these structures to retrieve the first exits  $e_1$  and  $e_2$  of  $n_1$  and  $n_2$ , and the children  $c_1$  and  $c_2$  of  $c$  are pointer nodes pointing to these  $e_1$  and  $e_2$ . We recursively call the unfolding sub-phase into the two newly created nodes  $c_1$  and  $c_2$ . Clearly this unfolding sub-phase overall takes a time proportional to the size of the final tree. Further, at the end of this sub-phase, the enumeration tree is output-ready.

**Output sub-phase.** For the output sub-phase, we have a output buffer and we simply traverse our enumeration tree. Each time we encounter a pointer node, it points to a singleton node carrying the value  $x$ , and we append  $x$  to our buffer. This sub-phase clearly takes a time proportional to the size of the tree. Notice that the output produced in the buffer is also of size proportional to the tree: this is because the number of concatenation node corresponds to the number of pointer nodes (all of which point to singletons), minus one.

**Pruning sub-phase.** For the pruning sub-phase, we want to prepare for the next output. For this we will require to compute the Boolean information of whether a

node has a next solution. This can be computed from the output-ready tree recursively: for a pointer node, there is no next solution (because it points to a singleton), for a concatenation node there exists a next solution if either one of its children has a next solution or if one of the two exit-enumeration structures contains an exit to be enumerated.

To understand how we compute the next solution, note that for a product node  $n$  of  $n_1$  and  $n_2$  we will compute the elements of  $S(n)$  by computing first the concatenation of all strings of  $S(n_1)$  with the first string of  $S(n_2)$ , then with the second string of  $S(n_2)$ , etc., until we have covered all strings in  $S(n_2)$ .

Now, we can explain our recursive computation of the enumeration tree corresponding to the next solution: given a concatenation node  $c$  with children  $c_1$  and  $c_2$  and exit-enumeration structure  $S_1$  and  $S_2$  pointing to a product node  $n$  with children  $n_1$  and  $n_2$  there are several cases:

- When  $c_1$  has a next solution, we simply recurse into  $c_1$  to find the next solution.
- When  $c_1$  has no next solution but  $S_1$  has a next exit  $e_1$  to enumerate, we replace  $c_1$  with a pointer node to  $e_1$ .
- In the other cases, it means that we are done doing the enumeration of  $n_1$  for the current solution of  $n_2$ . We thus need to restart the enumeration of solutions to  $n_1$  which is done by resetting  $S_1$  to enumerate the exits of  $n_1$  and replacing  $c_1$  by the pointer node corresponding to the first exit enumerated by  $S_1$ . We also need to move the enumeration of  $n_2$  to its next solution. If  $c_2$  has a next solution we compute it recursively otherwise we replace  $c_2$  with the next exit in  $S_2$ .
- If  $c_2$  has no next solution and  $S_2$  has no next exit, then there is no next solution.

**Memory usage of the enumeration phase.** In Theorem 1, we claimed that our enumeration scheme uses a memory linear in the size of the input string. As we already explained, the number of nodes that we keep in an enumeration tree is at most linear in the size of the current output, but we also need to prove that the memory used by the auxiliary structures for exit enumeration stays linear in the size of the DAG of enumerable sets. For this, notice that any node  $n$  of the DAG cannot be stored twice. It cannot be stored within the same exit enumeration data structure as it would mean that the same solution will be enumerated twice, which is impossible by our assumption that unions are disjoint. Further, it cannot be stored in two different exit enumeration schemes, as otherwise there would be a product node where both inputs contain the letters of some element in  $S(n)$ , which is again forbidden.

**Adding the  $\emptyset$  operation.** To extend the data structure to support the  $\emptyset$  operator, we add a new node type to have a node representing the empty set. There will be only one node of this type. The  $\emptyset$  basic operation is simply to create the  $\emptyset$  node if it does not exist and return it, otherwise return the one that exists. We change the definition of the product operator to do the following on arguments  $n_1$  and  $n_2$ : if one of  $n_1$  and  $n_2$  is the  $\emptyset$  node, then return the  $\emptyset$  node; otherwise do as before. We change the definition of the union operator to do the following on arguments  $n_1$  and  $n_2$ : if one of  $n_1$  and  $n_2$

is the  $\emptyset$  node, then return the other argument; otherwise do as before. It is clear that these alternative definitions do not alter the semantics of nodes, and they guarantee that whenever we apply the definitions from before then the arguments are not  $\emptyset$  nodes.

For the enumeration phase, we add a new base case: if the node to enumerate is  $\emptyset$ , we immediately halt. Otherwise, it is clear that the  $\emptyset$  node can never be visited in the enumeration, because in our construction the  $\emptyset$  node never has an incoming edge.

**Adding the singleton( $\varepsilon$ ) operation.** Making a data structure that supports `singleton( $\varepsilon$ )` operator is done by storing, for each enumerable set, a pair of an enumerable set without support for empty string solutions (as defined above), and a Boolean indicating whether the empty string  $\varepsilon$  is captured. We now present the new operations on enumerable sets. We distinguish between the original functions and the new functions described here by adding an epsilon index to the new functions, e.g., `empty $_{\varepsilon}$` .

For `empty $_{\varepsilon}$` , we return  $(\emptyset, \text{false})$ .

For `singleton $_{\varepsilon}$` ( $x$ ) we return the pair  $(\text{singleton}(x), \text{false})$  if  $x \neq \varepsilon$  and  $(\emptyset, \text{true})$  otherwise.

For `union $_{\varepsilon}$` (( $e_1, b_1$ ), ( $e_2, b_2$ )) we return  $(\text{union}(e_1, e_2), b_1 \vee b_2)$ .

For `prod $_{\varepsilon}$` (( $e_1, b_1$ ), ( $e_2, b_2$ )), the Boolean that we return is  $b_1 \wedge b_2$ , and the enumerable set that we return depends on the value of  $b_1$  and  $b_2$  we return `prod( $e_1, e_2$ )` when  $\neg b_1 \wedge \neg b_2$ , `union(prod( $e_1, e_2$ ),  $e_1$ )` when  $b_1 \wedge \neg b_2$ , `union(prod( $e_1, e_2$ ),  $e_2$ )` when  $b_2 \wedge \neg b_1$  and `union(prod( $e_1, e_2$ ), union( $e_1, e_2$ ))` when  $b_1 \wedge b_2$ .

To run the enumeration phase on a pair  $(e, b)$ , we start by outputting  $\varepsilon$  if  $b$  is set, and then we enumerate  $e$  as explained previously.

## A.2. Proof of Proposition 2

We show Proposition 2. The result is shown as in [25], though we give the proof in a self-contained fashion for convenience. We additionally show here that the transformation preserves rigidity (defined in Section 4) as this property is used in that section.

**Conditions 1 and 2: removing useless nonterminals.** We first perform a linear-time exploration from the terminals to mark the nonterminals  $X$  that can derive some string of terminals. The base case is if a nonterminal  $X$  has a rule  $X \rightarrow \alpha$  where  $\alpha$  only consists of terminals (in particular  $\alpha = \varepsilon$ ), then we mark it. The induction is that whenever a nonterminal  $X$  has a rule  $X \rightarrow \alpha$  where  $\alpha$  only consists of terminals and of marked nonterminals, then we mark  $X$ . At the end of this process, it is clear that any nonterminal that is not known to derive a string of terminals indeed does not derive any string, because any derivation of a string of terminals from a nonterminal  $X$  would witness that all nonterminals in this derivation, including  $X$ , should have been marked, which is impossible. Hence, we can remove the nonterminals that are not marked without changing the language or successful derivations of the grammar, and satisfy condition 1 in linear time.

Second, we perform a linear-time exploration from the start symbol  $S$  to mark the nonterminals  $X$  that can be reached in a derivation from  $S$ . The base case is that  $S$  is marked. The induction is that whenever a nonterminal  $Y$  occurs in the right-hand side of a rule having  $X$  as its left-hand side, and  $X$  is marked, then we mark  $Y$ . At the end

of the process, if a nonterminal  $X$  is not marked, then indeed there is no derivation from  $S$  that produces a string featuring  $X$ , as otherwise it would witness that  $X$  is marked, which is impossible. Hence, we can again remove the nonterminals that are not marked, the grammar and successful derivations are again unchanged, and we satisfy condition 2 in linear time.

As the transformations here only remove nonterminals and rules that cannot appear in a derivation, they clearly preserve unambiguity as well as rigidity.

**Condition 3: shape of rules.** We first ensure that every right-hand side of a rule is of size  $\leq 2$ . Given the annotated grammar  $\mathcal{G}$ , for every rule  $X \rightarrow \alpha$  where  $|\alpha| > 2$ , letting  $\alpha = \alpha_1 \cdots \alpha_n$ , we introduce  $n - 2$  fresh nonterminals  $X_{\alpha,1}, \dots, X_{\alpha,n-2}$ , and replace the rule by the following:  $X \rightarrow \alpha_1 X_{\alpha,1}$ ,  $X_{\alpha,1} \rightarrow \alpha_2 X_{\alpha,2}$ , ...,  $X_{\alpha,n-2} \rightarrow \alpha_{n-1} \alpha_n$ .

We make sure that the right-hand side of rules of size 2 consist only of nonterminals by introducing fresh intermediate nonterminals whenever necessary, which rewrite to the requisite terminal.

It is then clear that the result satisfies condition 3, and that there is a one-to-one correspondence between derivations in the original grammar and derivations in the rewritten grammar. To see this, note that there is an obvious one-to-one function which maps derivations from the original grammar into derivations in the new grammar, and that there is a slightly more involved function which receives a derivation in the new grammar, and builds a derivation in the original grammar by following the steps detailed above (and using the fact that each fresh nonterminal is associated to exactly one rule), which is also one-to-one. We conclude that the original grammar is unambiguous if and only if the new grammar is unambiguous.

The last point to check is that the arity-2 transformation preserves rigidity, i.e., if the original annotated grammar is rigid then so is the image of the transformation. Let  $X$  be some symbol of the original grammar  $\mathcal{G}$ , and  $w \in \Sigma^*$  be a string. Let us show that all derivations from the corresponding symbol  $X'$  of the rewritten grammar  $\mathcal{G}'$  have same shape. We do so by induction on the length of  $w$  and then on the topological order on nonterminals. The base case of  $w$  of length 0 is clear: the possible derivations are sequences of applications of rules of the form  $Y \rightarrow Z$  in a sequence of some fixed length, followed by a rule of the form  $Y \rightarrow \varepsilon$ , and what can happen in the rewritten grammar is the same.

For the inductive case, as  $\mathcal{G}$  is rigid, we know that there must be one fixed profile  $\pi \in \{0,1\}^k$  such that all derivations of  $w$  from  $X$  start by the application of a rule  $X \rightarrow \alpha$  where  $\alpha$  corresponds to profile  $\pi$ , i.e., it has length  $k$  and its  $i$ -th character is a nonterminal or terminal according to the value of the  $i$ -th bit of  $\pi$ . Otherwise the existence of two different right-hand-side profiles would contradict rigidity. Furthermore, by considering the possible sub-derivations from  $\alpha_1$  (including the empty derivation if  $\alpha_1$  is a terminal), we know that  $\alpha_1$  derives some fixed prefix of  $w$  and that all such derivations have the same sequence of profiles; otherwise we would witness a contradiction to rigidity. By applying the same argument successively to  $\alpha_2, \dots, \alpha_k$ , we deduce that there must be a partition of  $w = w_1 \cdots w_k$  such that, in all derivations of  $w$  from  $X$ , the derivation applies a rule with right-hand having profile  $\pi$  to produce some string  $\alpha_1 \cdots \alpha_k$ , and then each

$\alpha_i$  derives an annotation of  $w_i$  and for each  $i$  all possible derivations of some annotation of  $w_i$  by some  $i$ -th element in the right-hand side of such a rule has the same sequence of profiles.

As the string is nonempty we know that  $k > 0$ . Further, if  $k = 1$  then  $X$  and the productions involving  $X$  were not rewritten so we immediately conclude either with the case of a rule  $X \rightarrow \tau$  for a terminal  $\tau$  or by induction hypothesis on the nonterminals in the topological order for the case of a rule of the form  $X \rightarrow Y$ . Hence, we assume that  $k \geq 2$ .

We know by induction that, in the rewritten grammar, the derivation from  $X$  will start by rewriting  $X$  to  $Y_1 X_{\alpha,1}$ , the  $X_{\alpha,1}$  being itself rewritten to  $Y_2 X_{\alpha,2}$ , and so on, for some right-hand side  $\alpha$  of a rule  $X \rightarrow \alpha$  having profile  $\pi$ . Clearly each  $Y_i$  will have to derive an annotation of the  $w_i$  in the partitioning of  $w$ , as a derivation following a different partitioning would witness a derivation in the original grammar that contradicts rigidity. Now, the profile  $\pi$  indicates if each  $Y_i$  is a nonterminal of the initial grammar or a fresh nonterminal introduced to rewrite to a terminal. In the latter case, there is no possible deviation in profiles. In the former case, we conclude by induction hypothesis that each  $Y_i$  derives annotations of its  $w_i$  that all have the same profile, and we conclude that all derivations in the rewritten grammar indeed have the same profile, concluding the proof.

### A.3. Proof of Proposition 4

We know from [1] that for any  $c > 0$ , there exists a fixed grammar  $\mathcal{G}$  such that determining whether a string  $w$  is derived by  $\mathcal{G}$ , cannot be solved in time  $\mathcal{O}(|w|^{\omega-c})$ , unless the conjecture in graph algorithms mentioned in [1] is false.

We will see that this conditional lower bound translates to unambiguous annotated grammars. Indeed, we will show that for each grammar  $\mathcal{G}$  there exists an unambiguous annotated grammar  $\mathcal{G}'$  such that  $w$  is derived by  $\mathcal{G}$  if and only if  $[[\mathcal{G}']](w)$  is non-empty. Therefore after the preprocessing of  $w$  for  $\mathcal{G}'$ , we know in constant time whether  $w$  is derived by  $\mathcal{G}$  which proves that the preprocessing of  $\mathcal{G}'$  on  $w$  requires  $\mathcal{O}(|w|^{\omega-c})$  time, assuming the conjecture is true.

Now let us show how to translate a grammar  $\mathcal{G}$  into an unambiguous annotated grammar  $\mathcal{G}'$ . This can be challenging, because  $\mathcal{G}$  is not necessarily unambiguous: for this reason we need to define  $\mathcal{G}'$  intuitively by adding annotations that disambiguate the various possible derivations of  $\mathcal{G}$ , to guarantee that the result is unambiguous. As this is cumbersome to do on grammars, we use the correspondence between annotated grammars and pushdown annotators (Proposition 10), shown later in the article.

In this proof, we will use the notion of *pushdown automata* (PDA); see Definition 18 for the formal definition. Let us consider a PDA  $\mathcal{P}$  which is equivalent to  $\mathcal{G}$ . As is standard with PDAs, we can change the given definition to suppose without loss of generality that no transition in  $\mathcal{P}$  is an  $\varepsilon$ -transition. Specifically, we consider PDAs in a slightly different model where transitions are of the form  $(q_1, a, s_1, q_2, s_2) \in Q \times \Sigma \times \Gamma^+ \times Q \times \Gamma^+$ : such a transition means that in state  $q_1$ , when the top stack symbols are  $s_1$  and the next letter to read is  $a$ , the automaton can read the letter, move to state  $q_2$  and replace  $s_1$  by  $s_2$  on the stack. We create our unambiguous PDAnn  $\mathcal{P}'$  from  $\mathcal{P}$  by replacing each transition

$t = (q_1, a, s_1, q_2, s_2)$  to a set of transitions that first pop the symbols of  $s_1$  from the stack, then reads  $a$ , then pushes the symbols of  $s_2$  onto the stack. The first state of this transition is  $q_1$ , the last state is  $q_2$  but we make sure that each of the intermediate states are unique to  $t$ . Furthermore, the transition that reads the letter  $a$  outputs a symbol unique to the transition  $t$ . Therefore, by construction there is a bijection between runs of  $\mathcal{P}$  and runs of  $\mathcal{P}'$  and the PDAnn  $\mathcal{P}'$  is unambiguous because the run used for each output can be retrieved from that output.

We conclude by using Proposition 10 to obtain an equivalent annotated grammar  $\mathcal{G}'$ , which is also unambiguous. Thus, we know that on any unannotated string  $w$ , the set  $\llbracket \mathcal{G}' \rrbracket(w)$  is empty if  $\mathcal{G}$  does not derive  $w$ , and non-empty if it does. Thus, we know that, if we assume the conjecture is true, we cannot determine in  $\mathcal{O}(|w|^{\omega-c})$  whether  $\llbracket \mathcal{G}' \rrbracket(w)$  is empty or not. But if we have an algorithm to enumerate  $\llbracket \mathcal{G}' \rrbracket(w)$  with output-linear delay, as any output has size  $\mathcal{O}(|w|)$  in  $|w|$ , we can do this with a complexity linear in  $|w|$  which is that of the preprocessing of the enumeration algorithm. Thus, we conclude that the preprocessing conditionally requires  $\Omega(|w|^{\omega-c})$  time.

## B. Proofs of Section 4

### B.1. Proof of Theorem 5

In this appendix, we will use the notion of PDAnn introduced in Section 5, and we will use Proposition 10, which is also stated in Section 5 and proved in Appendix C.1.

To prove Theorem 5, we introduce a general-purpose normal form on PDAnn, where, intuitively, the only choices that can be made during a run are between the *types* of transition to apply. This is similar to Lemma 1 in [27].

**Definition 15.** *A PDAnn  $\mathcal{P}$  is deterministic-modulo-profile if it satisfies the following conditions:*

1. *for each state  $p$  there is at most one push transition that starts on  $p$ , formally  $|\{q, \gamma \in Q \times \Gamma \mid (p, q, \gamma) \in \Delta\}| \leq 1$*
2. *for each state  $p$  and stack symbol  $\gamma$  there is at most one pop transition that starts on  $p, \gamma$ , formally  $|\{q \in Q \mid (p, \gamma, q) \in \Delta\}| \leq 1$*
3. *for each state  $p$ , letter  $a$ , and output  $\sigma \in \Omega$ , there is at most one read-write transition that starts on  $p, a, \sigma$ , formally, we have  $|\{q \in Q \mid (p, (a, \sigma), q) \in \Delta\}| \leq 1$ .*
4. *for each state  $p$  and letter  $a$ , there is at most one read transition that starts on  $p, a$ , formally  $|\{q \in Q \mid (p, a, q) \in \Delta\}| \leq 1$ .*

**Lemma 2.** *Let  $\mathcal{P}$  be a PDAnn. We can build an equivalent PDAnn  $\mathcal{P}'$  which is deterministic-modulo-profile. The transformation takes exponential time, i.e., time  $\mathcal{O}(2^{|\mathcal{P}|^c})$  for some  $c > 0$ .*

*Further, on any string  $w$ , there is an accepting run of  $\mathcal{P}$  on  $w$  with profile  $\pi$  iff there is an accepting run of  $\mathcal{P}'$  on  $w$  with the same profile.*

*Proof.* The proof is similar to the determinization of visibly pushdown automata [27, Lemma 1].

Given  $\mathcal{P} = (Q, \Sigma, \Omega, \Gamma, \Delta, q_0, F)$ , we build  $\mathcal{P}' = (Q', \Sigma, \Omega, \Gamma', \Delta', S_I, F')$  as follows. We build  $Q' = 2^{Q \times Q}$ , intuitively denoting a set of pairs of states  $(p, q)$  of  $\mathcal{P}$  such that  $\mathcal{P}$  can be at state  $q$  at this point if it was at state  $p$  when the topmost stack symbol was pushed. We build  $\Gamma' = 2^{Q \times \Gamma \times Q}$ , intuitively specifying the sets of possible stack symbols and remembering the state just after the previous stack symbol was pushed and the state just after that symbol was pushed. We build  $S_I = \{(q_0, q_0)\}$ , meaning that initially we are at the initial state  $q_0$  and were here when the stack was initialized. We build  $F' = \{S \mid (q_0, q) \in S \text{ for some } q \in F\}$ , meaning that we accept when  $\mathcal{P}$  reaches a final state and we were at the initial state when the stack was initialized. Let  $\Delta'$  be defined as follows:

- The (unique) push transition from a state  $S \in Q'$  makes  $\mathcal{P}'$  push a stack symbol  $S'$  and move to a state  $T$ , intuitively defined as follows. For every pair  $(p, p')$  of  $S$  and push transition  $(p', q, \gamma) \in \Delta$  in the original PDAnn, we can move to state  $(q, q)$  and push on the stack the symbol  $(p, \gamma, q)$ . The stack symbol  $S'$  is the set of all possible stack symbols that can be pushed in this way, and  $T$  is the set of all possible states that can be reached in this way.

Formally, for every  $S \in Q'$  we include  $(S, S', T)$  in  $\Delta'$ , where:

$$T = \{(p, \gamma, q) \mid (p, p') \in S \text{ and } (p', q, \gamma) \in \Delta \text{ for some } p, p', q \in Q, \gamma \in \Gamma\},$$

$$S' = \{(q, q) \mid (p, p') \in S \text{ and } (p', q, \gamma) \in \Delta \text{ for some } p, p', q \in Q, \gamma \in \Gamma\}$$

- The (unique) pop transition from a state  $S \in Q'$  and topmost stack symbol  $T \in \Gamma'$  makes  $\mathcal{P}'$  move to a state  $T'$  intuitively defined as follows. For every pair  $(p', q')$  of  $S$ , we consider all triples  $(p, \gamma, p')$  of the topmost stack symbol  $T$ , and if the original PDAnn had a pop transition  $(q', \gamma, q) \in \Delta$ , then we can pop the topmost stack symbol and go to the state  $(p, q)$ . The new state  $T'$  is the set of all pairs  $(p, q)$  that can be reached in this way.

Formally, for every  $(S, T) \in Q' \times \Gamma'$  we include  $(S, T, S')$  in  $\Delta'$ , where:

$$S' = \{(p, q) \mid (p, \gamma, p') \in T \text{ and } (p', q') \in S \text{ and } (q', \gamma, q) \in \Delta \text{ for some } p, p', q, q' \in Q, \gamma \in \Gamma\},$$

- The (unique) read-write transition from a state  $S \in Q'$  on a letter  $a \in \Sigma$  and output  $\sigma \in \Omega$  makes  $\mathcal{P}'$  move to a state  $S'$  intuitively defined as follows: we consider all pairs  $(p, p')$  in  $S$  and all transitions from  $p'$  with  $a$  and  $\sigma$  in  $\mathcal{P}$  to some state  $q$ , and move to all possible pairs  $(p', q)$ .

Formally, for every  $(S, a, \sigma) \in Q' \times \Sigma \times \Omega$  we include  $(S, (a, \sigma), S')$  in  $\Delta'$ , where:

$$S' = \{(p, q) \mid (p, p') \in S \text{ and } (p', (a, \sigma), q) \in \Delta \text{ for some } p, p', q \in Q\}.$$

- The (unique) read transition from a state  $S \in Q'$  on a letter  $a \in \Sigma$  makes  $\mathcal{P}'$  move to a state  $S'$  intuitively defined as follows: we consider all pairs  $(p, p')$  in  $S$  and all transitions from  $p'$  with  $a$  in  $\mathcal{P}$  to some state  $q$ , and move to all possible pairs  $(p', q)$ .

Formally, for every  $(S, a) \in Q' \times \Sigma$  we include  $(S, a, S')$  in  $\Delta'$ , where:

$$S' = \{(p, q) \mid (p, p') \in S \text{ and } (p', a, q) \in \Delta \text{ for some } p, p', q \in Q\}.$$

It is clear by definition that  $\mathcal{P}'$  is deterministic-modulo-profile, and it is clear that the running time of the construction satisfies the claimed time bound.

We now show that  $\mathcal{P}$  and  $\mathcal{P}'$  are equivalent.

Now, for the forward direction, let us first assume without loss of generality that whenever  $\mathcal{P}$  makes a push transition then the stack symbol that it pushes is annotated with the state reached just after the push. Then we will show that every instantaneous description that can be reached by  $\mathcal{P}$  can be reached by  $\mathcal{P}'$  by induction on the run. Specifically, we show by induction on the length of the run  $\rho$  the following claim: if  $\mathcal{P}$  has a run  $\rho$  on a string  $w$  that produces  $\mu$  from an initial state  $q_0 \in T$  to an instantaneous description  $(q, i), \alpha$ , with  $\alpha = \gamma_0 p_0, \dots, \gamma_m p_m$  being the sequence of the stack symbols and states annotating them, then  $\mathcal{P}'$  has a run  $\rho'$  on  $w$  from  $S_I$  to an instantaneous description  $(S, i), \alpha'$  with  $\alpha' = T_0 \dots T_m$  such that  $T_0$  contains  $(q_0, \gamma_0, p_0)$ ,  $T_1$  contains  $(p_0, \gamma_1, p_1)$ , ...,  $T_m$  contains  $(p_{m-1}, \gamma_m, p_m)$  and  $S$  contains  $(p_m, q)$ ; further  $\rho$  and  $\rho'$  have the same profile.

The base case of an empty run on a string is immediate: if  $\mathcal{P}$  has an empty run from an initial state  $q_0$ , then it reaches the instantaneous description with  $(q_0, 0)$  and the empty stack, and then  $\mathcal{P}'$  then has an empty run reaching the instantaneous description  $(S, 0)$  with the empty stack and  $S$  indeed contains  $(q_0, q_0)$ .

For the induction case, assume that  $\mathcal{P}$  has a non-empty run  $\rho_+$  on a string  $w$  that produces  $\mu$ . First, write  $\rho_+$  as a run  $\rho$  followed by one single transition of  $\mathcal{P}$ . We know  $\mathcal{P}$  has a run  $\rho$  on  $w$  which produces  $\mu$  from an initial state  $q_0$  to an instantaneous description  $(q, i), \alpha$ , with  $\alpha = \gamma_0 p_0, \dots, \gamma_m p_m$ . By the induction hypothesis, we know that  $\mathcal{P}'$  has a run  $\rho'$  on  $w$  from  $(q_0, q_0)$  to an instantaneous description  $(S, i), \alpha'$  with  $\alpha' = T_0 \dots T_m$  such that  $T_0$  contains  $(q_0, \gamma_0, p_0)$ ,  $T_1$  contains  $(p_0, \gamma_1, p_1)$ , ...,  $T_m$  contains  $(p_{m-1}, \gamma, p_m)$  and  $S$  contains  $(p_m, q)$ ; and  $\rho'$  and  $\rho$  have the same profile. We now distinguish on the type of the transition used to extend  $\rho$  to  $\rho_+$ .

If that transition is a read-write transition  $(q, (a, \sigma), q')$ , we consider the read-write transition of  $\mathcal{P}'$  labeled with  $(a, \sigma)$  from  $T$ , and call  $S'$  the state that  $\mathcal{P}'$  reaches. As  $(p_m, q) \in S$  and  $(q, (a, \sigma), q') \in \Delta$ , we know that  $(p_m, q') \in S'$ . Thus,  $\mathcal{P}'$  can read  $(a, \sigma)$  and reach a suitable state  $S'$  and position  $i+1$  and the stacks are unchanged so the claim is proven.

If that transition is a read transition  $(q, a, q')$ , we follow an analogous reasoning.

If that transition is a push transition  $(q, q', \gamma)$ , the position of  $\mathcal{P}$  is unchanged and the new stack is extended by  $\gamma$  annotated with state  $q'$ . Consider the push transition of  $\mathcal{P}'$  from  $q$ , and call  $S'$  the state reached and  $T = T_{m+1}$  the stack symbol that is pushed.

As  $(p_m, q) \in S$  and  $(q, q', \gamma) \in \Delta$ , we know that  $T$  contains  $(p_m, \gamma, q')$ , and  $S'$  contains  $(q', q')$ , which is what we needed to show.

If that transition is a pop transition,  $(q, \gamma_m, q')$ , the position of  $\mathcal{P}$  is unchanged and the topmost stack symbol is removed. Consider the topmost stack symbol  $T_m$  and the transition of  $\mathcal{P}'$  that pops it from  $S$ , and call  $S'$  the state that we reach. We know that  $S$  contains  $(p_m, q)$  and  $T_m$  contains  $(p_{m-1}, \gamma_m, p_m)$  and  $(q, \gamma_m, q') \in \Delta$ , so  $S'$  contains  $(p_{m-1}, q')$ , which is what we needed to show.

Note that, in all four cases, the profile of  $\rho_+$  and  $\rho'_+$  is the same, because this was true of  $\rho$  and  $\rho'$ , and the type of transition done to extend  $\rho'$  to  $\rho'_+$  is the same as the type of transition done to extend  $\rho$  to  $\rho_+$ .

The inductive claim is therefore shown, and thus if  $\mathcal{P}$  has a run  $\rho$  on some string  $w$  that produces  $\mu$  starting at some initial state  $q_0$  and ending at state  $q$ , then  $\mathcal{P}'$  has a run  $\rho'$  on  $w$  which produces  $\mu$  and ending at a state of the form  $(q_0, q)$  for  $q_0$  and having same profile. Thus, if  $\rho$  is accepting then  $q$  is final for  $\mathcal{P}$  and  $(q_0, q)$  is final for  $\mathcal{P}'$  so  $\rho'$  is accepting. This concludes the forward implication.

We now show the backward implication, and show it again by induction, again assuming that  $\mathcal{P}$  annotates the symbols of its stack with the state reached just after pushing them. We show by induction on the length of a run  $\rho'$  the following claim: if  $\mathcal{P}'$  has a run  $\rho'$  on a string  $w$  that produces  $\mu$  from its initial state to an instantaneous description  $(S, i), \alpha'$  with  $\alpha' = T_0, \dots, T_m$  being the sequence of the stack symbols, then for any choice of elements  $(q_0, \gamma_0, p_0) \in T_0, (p_0, \gamma_1, p_1) \in T_1, \dots, (p_{m-1}, \gamma_m, p_m) \in T_m$  and  $(p_m, q) \in S$  it holds that  $\mathcal{P}$  has a run  $\rho$  on  $w$  producing  $\mu$  from some initial state  $q_0$  to the instantaneous description  $(q, i), \alpha$  with  $\alpha = \gamma_0 q_0, \dots, \gamma_m q_m$  (writing next to each stack symbol the state that annotates it), and  $\rho'$  and  $\rho$  have the same profile.

The base case of an empty run on a string is again immediate: if  $\mathcal{P}'$  has an empty run from its initial state, then it reaches the instantaneous description with  $(S_I, 0)$  and empty stack, and then  $\mathcal{P}$  has an empty run from any initial state  $q_0$  to  $q_0$  so that indeed  $S_I$  contains  $(q_0, q_0)$ .

For the induction case, assume that  $\mathcal{P}'$  has a non-empty run  $\rho'_+$  on  $w$  which produces  $\mu$ . We write again  $\rho'_+$  as a run  $\rho'$  followed by one single transition of  $\mathcal{P}'$ . We know  $\mathcal{P}'$  has a run  $\rho'$  on  $w$  which produces  $\mu$  from the initial state  $S_I$  to an instantaneous description  $(S, i), \alpha'$ , with  $\alpha' = T_0 \dots T_m$ . By the induction hypothesis, we know that for any choice of elements  $(q_0, \gamma_0, p_0) \in T_0, (p_0, \gamma_1, p_1) \in T_1, \dots, (p_{m-1}, \gamma_m, p_m) \in T_m$  and  $(p_m, q) \in S$ , then  $\mathcal{P}$  has a run  $\rho$  on  $w$  which produces  $\mu$  from some initial state  $q_0$  to the instantaneous description  $(q, i), \alpha$  with  $\alpha = \gamma_0 q_0, \dots, \gamma_m q_m$ , and  $\rho$  and  $\rho'$  have the same profile. We now distinguish on the type of transition used to extend  $\rho'$  to  $\rho'_+$ .

If the last transition is a read-write transition  $(S, (a, \sigma), S')$  with  $S'$  defined as in the construction, consider any choice of  $(q_0, \gamma_0, p_0) \in T_0, (p_0, \gamma_1, p_1) \in T_1, \dots, (p_{m-1}, \gamma_m, p_m) \in T_m$  and  $(p_m, q') \in S'$ , and then there must be some state  $p''$  such that  $(p'', (a, \sigma), q) \in \Delta$  and  $(p_m, p'') \in S$ . Using the induction hypothesis but picking  $(p_m, p'') \in S$ , we obtain a run  $\rho$  of  $\mathcal{P}$  on  $w$  which produces  $\mu$ , with the correct stack and ending at position  $i$  on state  $p''$ , which we can extend by the read transition  $(p'', (a, \sigma), q)$  to reach state  $q$  at position  $i + 1$  without touching the stack, proving the result.

If the last transition is a read transition  $(S, a, S')$  with  $S'$  defined as in the construction,

we follow an analogous reasoning.

If the last transition is a push transition  $(S, S', T)$  with  $T$  defined as in the construction, consider any choice of  $(q_0, \gamma_0, p_0) \in T_0$ ,  $(p_0, \gamma_1, p_1) \in T_1$ , ...,  $(p_{m-1}, \gamma_m, p_m) \in T_m$ ,  $(p_m, \gamma_{m+1}, p_{m+1}) \in T_{m+1}$  and  $(p_{m+1}, q') \in S'$ . We know that we must have  $q' = p_{m+1}$ , and that there must be some state  $p''$  and push transition  $(p'', p_{m+1}, \gamma_{m+1})$  and pair  $(p_m, p'')$  in  $S$ . Using the induction hypothesis but picking  $(p_m, p'') \in S$ , we obtain a run  $\rho$  of  $\mathcal{P}$  on  $w$  which produces  $\mu$  with topmost stack symbol  $\gamma_m$ , ending at state  $p''$ , which we can extend with the push transition  $(p'', p_{m+1}, \gamma_{m+1})$  to obtain the desired stack and reach state  $p_{m+1} = q'$ , proving the result.

If the last transition is a pop transition  $(S, T, S')$  with  $S'$  defined as in the construction, consider any choice of  $(q_0, \gamma_0, p_0) \in T_0$ ,  $(p_0, \gamma_1, p_1) \in T_1$ , ...,  $(p_{m-2}, \gamma_{m-1}, p_{m-1}) \in T_{m-1}$ , and  $(p_{m-1}, q) \in S'$ . We know that there is a pair  $(p', q') \in S$  and a triple  $(p_{m-1}, \gamma_m, p')$  in  $T_m$  and a pop transition  $(q', \gamma_m, q)$  in  $\Delta$ . Applying the induction hypothesis, we get a run  $\rho$  of  $\mathcal{P}$  on  $w$  which produces  $\mu$  and with topmost stack symbol  $\gamma_m$  annotated with state  $p'$  and ending at state  $q'$ . The pop transition  $(q', \gamma_m, q)$  allows us to extend this run to reach state  $q$  and remove the topmost stack symbol, while the rest of the stack is correct, proving the result.

Again, we have ensured that  $\rho$  is extended to  $\rho_+$  with the same transition as the transition used to extend  $\rho'$  to  $\rho'_+$ , ensuring that  $\rho_+$  and  $\rho'_+$  have same profile. This concludes the proof of the backward induction, ensuring that if  $\mathcal{P}'$  has a run from  $S_I$  to some final state  $S$  reading a string  $w$  and producing  $\mu$ , and having  $(q_0, q_f)$  with  $q_f \in F$  in  $S$ , then  $\mathcal{P}$  has a run reading  $w$  which produces  $\mu$  going from  $q_0$  to the final state  $q_f$ . This concludes the backward implication and completes the proof.  $\square$

We can now show Theorem 5 via Proposition 10, using also the notion of *profiled PDAnn* defined in Section 5:

*Proof of Theorem 5.* Let  $\mathcal{G}$  be a rigid annotated grammar. Using Proposition 10, we transform it in polynomial time to a profiled PDAnn  $\mathcal{P}$ . Using Lemma 2, we build in exponential time an equivalent PDAnn  $\mathcal{P}'$  satisfying the conditions of the lemma.

We know that  $\mathcal{P}'$  is still profiled. Indeed, if we assume by contradiction that there is a string  $w$  on which  $\mathcal{P}'$  has two accepting runs with different profiles, then by the last condition of Lemma 2, the same is true of  $\mathcal{P}$ , contradicting the fact that  $\mathcal{P}$  is profiled.

Now, we claim that  $\mathcal{P}'$  is necessarily also unambiguous. To see why, consider two accepting runs  $\rho$  and  $\rho'$  of  $\mathcal{P}'$  on some string  $w$ . Since  $\mathcal{P}'$  is profiled,  $\rho$  and  $\rho'$  must have the same profile. But now, the conditions of Lemma 2 ensure that, knowing the input string  $w$  and profile, the runs  $\rho$  and  $\rho'$  are completely determined. Specifically, this is an immediate induction on the run. The base case is that there is only one initial state, so both  $\rho$  and  $\rho'$  must have the same initial state. Now, assuming by induction that the runs so far are identical and have the same stack, there are three cases. First, if the profile tells us that both runs make a push transition, the symbol pushed and state reached are determined by the last states of the runs so far, which are identical by inductive hypothesis. Second, if the profile tells us that both runs make a read-write transition (or read transition), the state reached is determined by the input and output

symbols (or just the input symbol), and by the last states of the run so far, which are identical by inductive hypothesis. Third, if the profile tells us that both runs make a pop transition, the state reaches is determined by the last state of the run so far, and the topmost symbol of the stack, which are identical by inductive hypothesis. This concludes the inductive proof.

Thus, for any two accepting runs  $\rho$  and  $\rho'$  on the string  $w$  which produce the same output, they must be identical. Thus,  $\mathcal{P}'$  is unambiguous. We use Theorem 10 to transform  $\mathcal{P}'$  back into an annotated grammar, which is still rigid and unambiguous, and equivalent to the original rigid annotated grammar  $\mathcal{G}$ . The overall complexity of the transformation is in  $\mathcal{O}((2^{(|\mathcal{G}|^c)^{c'})})^{c''}$  for some  $c, c', c'' > 0$ , so it is in  $\mathcal{O}(2^{|\mathcal{G}|^d})$  for some  $d > 0$  overall, and the time complexity is as stated.  $\square$

## B.2. Proof of Proposition 6

This proof is based on extending the definitions of unambiguity and rigidness of annotated grammars over unannotated context-free grammars. Indeed, an unambiguous annotated grammar with an empty output set is just an unambiguous CFG, and a rigid annotated grammar with an empty output set is a CFG for which every derivation of a given string  $w \in \Sigma^*$  has the same shape.

Consider the (unannotated) grammar  $\mathcal{G}'$  obtained from  $\mathcal{G}$  by removing all annotations on terminals, and making  $\Omega = \emptyset$ . It can be seen that  $L(\mathcal{G}') = L'$  since for each string  $w$ , if  $w \in L(\mathcal{G}')$ , then there is at least one  $\hat{w} \in L(\mathcal{G})$  with  $\text{str}(\hat{w}) = w$  and vice versa. Now, we claim that  $\mathcal{G}'$  is *rigid*, by extending the notion onto CFGs in the obvious way. To see this, consider a string  $w \in L(\mathcal{G}')$ ; all derivations of  $w$  by  $\mathcal{G}'$  correspond to derivations by  $\mathcal{G}$  of some  $\hat{w}$  such that  $\text{str}(\hat{w}) = w$ . Because  $\mathcal{G}$  is rigid, all these derivations have the same shape. Now, using Theorem 5, we can compute a rigid and unambiguous grammar  $\mathcal{G}''$  recognizing the same language over  $\Sigma^*$  as  $\mathcal{G}'$ , i.e.,  $L'$ . But as  $L'$  is a language without output, the unambiguity of  $\mathcal{G}''$  actually means that  $\mathcal{G}''$  is an unambiguous CFG. Hence,  $L'$  is recognized by an unambiguous grammar, concluding the proof.

## B.3. Proof of Proposition 7: Undecidability results on rigid grammars

We first show the undecidability of checking if an annotated grammar has an equivalent rigid annotated grammar:

**Proposition 16.** *Consider the problem, given an annotated grammar  $\mathcal{G}$ , of determining whether there exists some equivalent rigid annotated grammar equivalent to  $\mathcal{G}$ . This problem is undecidable.*

*Proof.* We reduce from the problem of deciding whether the language  $L_2$  of an input (unannotated) context-free grammar  $\mathcal{G}_2$  can be recognized by an unambiguous context-free grammar: this task is known to be undecidable [20]. Consider  $\mathcal{G}_2$  as an annotated grammar (with empty annotations). Let us show that  $L_2$  can be recognized by a rigid annotated grammar iff it can be recognized by an unambiguous context-free grammar, which concludes. For the forward direction, if  $L_2$  can be recognized by an unambiguous

context-free grammar, then that grammar is in particular rigid. For the backward direction, if  $L_2$  can be recognized by a rigid grammar, then Proposition 6 implies that  $L_2$  can also be recognized by an unambiguous context-free grammar. Thus, we have showed that the (trivial) reduction is correct.  $\square$

We next show that it is undecidable to check if an input annotated grammar is rigid:

**Proposition 17.** *Consider the problem, given an annotated grammar  $\mathcal{G}$ , of determining whether it is rigid. This problem is undecidable.*

*Proof.* We adapt the standard proof of undecidability [13, Ambiguity Theorem 2] for the problem of deciding, given an input unannotated grammar  $\mathcal{G}$ , if it is unambiguous. The reduction is from the Post Correspondence Problem (PCP), which is undecidable: we are given as input sequences  $\alpha_1, \dots, \alpha_n$  and  $\beta_1, \dots, \beta_n$  of strings over some alphabet  $\Sigma$ , and we ask whether there is a non-empty sequence of indices  $i_1, \dots, i_m$  of integers in  $[1, n]$  such that  $\alpha_{i_1} \dots \alpha_{i_m} = \beta_{i_1} \dots \beta_{i_m}$ . Given the input sequences  $\alpha_1, \dots, \alpha_n$  and  $\beta_1, \dots, \beta_n$  to the PCP, we consider the alphabet  $\Sigma' = \Sigma \cup \{1, \dots, n\}$ , and we consider the CFG having nonterminals  $S, S_1$ , and  $S_2$ , start symbol  $S$ , and rules  $S \rightarrow S_1, S \rightarrow S_2, S_1 \rightarrow \varepsilon, S_2 \rightarrow \varepsilon$ , and for each  $1 \leq i \leq n$  the productions  $S_1 \rightarrow \alpha_i S_1 i$  and  $S_2 \rightarrow \beta_i S_2 i$ .

We claim that this grammar is ambiguous iff there is a solution to the Post correspondence problem. Indeed, given any solution  $\alpha_{i_1} \dots \alpha_{i_m} = \beta_{i_1} \dots \beta_{i_m}$ , considering the string  $\alpha_{i_1} \dots \alpha_{i_m} i_m \dots i_1 = \beta_{i_1} \dots \beta_{i_m} i_m \dots i_1$ , we can parse it with one derivation featuring  $S_1$  and one derivation featuring  $S_2$ . Conversely, if we can parse a string  $w \in \Sigma^*$  with two different derivations, we know that there cannot be two different derivations featuring  $S_1$ . Indeed, reading the string from right to left uniquely identifies the possible derivations from  $S_1$ . The same argument applies to derivations featuring  $S_2$ . Hence, if the grammar is ambiguous, then there is exactly one derivation featuring  $S_1$  and exactly one derivation featuring  $S_2$ . These two derivations can be used to find a solution to the Post correspondence problem.

We now adapt this proof to show the undecidability of rigidity. We say that an input to the PCP is *trivial* if there is  $i$  such that  $\alpha_i = \beta_i$ . We can clearly decide in linear time, given the input to the PCP, if it is trivial. Hence, the PCP is also undecidable in the case where the PCP is non-trivial. Now, when doing the reduction above on a PCP instance that is not trivial, we observe that two derivations of the same string can never have the same sequences of shapes. Indeed, if we have two derivations of the same string, then as we explained one must feature  $S_1$  and the other must feature  $S_2$ , and they give a solution  $\alpha_{i_1} \dots \alpha_{i_m} = \beta_{i_1} \dots \beta_{i_m}$  to the PCP. Assume by contradiction that both derivations have the same sequences of shapes. Then, it means that we have  $|\alpha_{i_j}| = |\beta_{i_j}|$  for every  $1 \leq j \leq m$ . In particular we have  $|\alpha_{i_1}| = |\beta_{i_1}|$ , and so we know that  $\alpha_{i_1} = \beta_{i_1}$  and the PCP instance was trivial, a contradiction.

Hence, let us reduce from the PCP on non-trivial instances to the problem of deciding whether an input annotated grammar is not rigid. Given a non-trivial PCP instance, we construct  $\mathcal{G}$  as above, but seeing it as an annotated grammar with no outputs. Then  $\mathcal{G}$  is not rigid iff there is a string  $w$  such that the empty annotation of  $w$  has two derivations that do not have the same sequence of shapes. But this is equivalent to  $\mathcal{G}$

being unambiguous when seen as a CFG. Indeed, for the forward direction, if  $\mathcal{G}$  has two such derivations on a string  $w$  then clearly  $w$  witnesses that  $\mathcal{G}$  is ambiguous when seen as a CFG. Conversely, if  $\mathcal{G}$  is ambiguous when seen as a CFG, we have explained in the previous paragraph that the two derivations must have different sequences of shapes, so  $\mathcal{G}$  is not rigid. Hence, we conclude that there is a solution to the input non-trivial PCP instance iff  $\mathcal{G}$  is not rigid. This establishes that the problem is undecidable and concludes the proof.  $\square$

#### B.4. Proof of Proposition 9

Assume we have such an algorithm  $\mathcal{A}$ . Consider a procedure which receives an unambiguous unannotated CFG  $\mathcal{G}$  and an input string  $w$ , converts  $\mathcal{G}$  into an annotated grammar  $\mathcal{G}'$  with empty output set. Since  $\mathcal{G}'$  is unambiguous and rigid, we can run  $\mathcal{A}$  over  $\mathcal{G}'$  and  $w$ . If  $w \in L(\mathcal{G})$ , then  $\llbracket \mathcal{G}' \rrbracket(w) = \{\varepsilon\}$ , and if  $w \notin L(\mathcal{G})$ , then  $\llbracket \mathcal{G}' \rrbracket(w) = \emptyset$ . Thus, after the preprocessing phase of  $\mathcal{A}$  we need only to wait a constant amount of time to see if the string  $\varepsilon$  is given as output, or none is. We conclude that this procedure solves the problem with the same complexity as the preprocessing phase of  $\mathcal{A}$ .

### C. Proofs of Section 5

In this appendix, we will need to use the standard notion of a *pushdown automaton* (PDA), whose definition was omitted from the main text of the paper. We give it here:

**Definition 18.** A pushdown automaton (PDA) is a tuple  $\mathcal{A} = (Q, \Sigma, \Gamma, \Delta, q_0, F)$ , where  $Q$  is a finite set of states,  $\Sigma$  is the alphabet,  $\Gamma$  is a finite alphabet of stack symbols,  $q_0 \in Q$  is the initial state,  $F \subseteq Q$  are the final states. We assume that  $\Gamma$  is disjoint from  $\Sigma$ . Further,  $\Delta$  is a finite set of transitions of the following kind:

- Read transitions of the form  $(p, a, q) \in Q \times \Sigma \times Q$ , meaning that the automaton can go from state  $p$  to state  $q$  while reading the letter  $a$ ;
- Push transitions of the form  $(p, q, \gamma) \in Q \times Q \times \Gamma$ , meaning that the automaton can go from state  $p$  to state  $q$  while pushing the symbol  $\gamma$  on the stack;
- Pop transitions of the form  $(p, \gamma, q) \in Q \times \Gamma \times Q$ , meaning that, if the topmost symbol of the stack is  $\gamma$ , the automaton can go from  $p$  to  $q$  while removing this topmost symbol  $\gamma$ .

We omit the definition of the semantics of PDAs, which are standard, and allow us to define the *language*  $L(\mathcal{A})$  accepted by a PDA. It is also well-known that CFGs and PDAs have the same expressive power, i.e., given a CFG  $G$ , we can build in polynomial time a PDA  $\mathcal{A}$  which is *equivalent* in the sense that  $L(G) = L(\mathcal{A})$ , and vice-versa.

We will also need to use the standard notion of a *deterministic* PDA (with acceptance by final state). Formally:

**Definition 19.** Let  $\mathcal{A} = (Q, \Sigma, \Gamma, \Delta, q_0, F)$  be a PDA. For  $p \in Q$ , we define the next-transitions of  $p$  as the set  $\Delta(p)$  of all transitions in  $\Delta$  that start on  $p$ , i.e.,  $\Delta(p) = \{(p, x, y) \mid (p, x, y) \in \Delta\}$ . We say that a PDA  $\mathcal{A}$  is deterministic if for every state  $q \in Q$ , one of the following conditions hold:

- (a)  $\Delta(q) \subseteq Q \times \Sigma \times Q$  and  $|\{q' \mid (q, a, q') \in \Delta\}| \leq 1$  for each  $a \in \Sigma$ . Informally, all applicable transitions are read transitions, and there is at most one such applicable transition for each letter.
- (b)  $\Delta(q) \subseteq Q \times (Q \times \Gamma)$ , and  $|\Delta(q)| \leq 1$ . Informally, all applicable transitions are push transitions, and there is at most one such transition from  $q$ .
- (c)  $\Delta(q) \subseteq (Q \times \Gamma) \times Q$  and  $|\{q' \mid (q, \gamma, q') \in \Delta\}| \leq 1$  for each  $\gamma \in \Gamma$ . Informally, all applicable transitions from  $q$  are pop transitions, and there is at most one such applicable transition for each stack symbol.

It is clear that the definition ensures that, on every input string  $w$ , a deterministic PDA  $\mathcal{A}$  has at most one run accepting  $w$ , so that we can check in linear time in  $\mathcal{A}$  and  $w$  if  $w \in L(\mathcal{A})$ . Further, it is known that deterministic PDAs are strictly less expressive than general PDAs.

### C.1. Proof of Proposition 10

Let us first give the formal definitions needed for the statement of the result. We say that two annotated grammars  $\mathcal{G}$  and  $\mathcal{G}'$  are *equivalent* if they define the same functions, i.e.,  $\llbracket \mathcal{G} \rrbracket = \llbracket \mathcal{G}' \rrbracket$ . We define equivalence in the same way for two PDAnns, or for an annotated grammar and a PDAnn.

We first show one direction:

**Proposition 20.** *For any annotated grammar  $\mathcal{G}$ , we can build an equivalent PDAnn  $\mathcal{P}$  in polynomial time. Further, if  $\mathcal{G}$  is unambiguous then so is  $\mathcal{P}$ . Moreover, if  $\mathcal{G}$  is rigid, then  $\mathcal{P}$  is profiled.*

*Proof.* This is a standard transformation. Let  $\mathcal{G} = (V, \Sigma, \Omega, P, S)$ . We build a PDAnn  $\mathcal{P} = (Q, \Sigma, \Omega, \Gamma, \Delta, q_0, F)$  as follows: For every rule  $X \rightarrow \alpha$  of  $\mathcal{G}$  and position  $0 \leq i \leq |\alpha|$ , the PDAnn  $\mathcal{P}$  has a state  $(X, \alpha, i)$  in  $Q$ , plus a special state  $q_0$  which is the only initial and only final state. Also,  $\Gamma = Q$ . The set  $\Delta$  has the following transitions:

- A push transition  $(q_0, (S, \alpha, 0), q_0)$  and a pop transition  $((S, \alpha, |\alpha|), q_0, q_0)$ , for every rule  $S \rightarrow \alpha$ .
- For each state  $(X, \alpha, |\alpha|)$  for a production  $X \rightarrow \alpha$ , a pop transition reading a state from the stack and moving to that state.
- For each state  $(X, \alpha, i)$  where the  $(i + 1)$ -th element of  $\alpha$  (numbered from 1) is a nonterminal  $Y$ , for every rule  $Y \rightarrow \beta$ , a pop transition pushing  $(X, \alpha, i + 1)$ , and moving to  $(Y, \beta, 0)$ .

- For each state  $(X, \alpha, i)$  where the  $(i + 1)$ -th element of  $\alpha$  (numbered from 1) is a terminal  $\tau$ , a read transition moving to  $(X, \alpha, i + 1)$  reading the symbol of  $\tau$  and outputting the annotation of  $\tau$  (if any).

We will show a function that maps a given leftmost derivation  $S \Rightarrow_{\mathcal{G}} \gamma_1 \Rightarrow_{\mathcal{G}} \dots \Rightarrow_{\mathcal{G}} \gamma_m = \hat{w}$  into a run in  $\mathcal{P}$ . To do this, we convert this sequence of productions into a sequence of strings which has the same size as the run (minus one). These strings serve as an intermediate representation of both the derivation and the run. The process is essentially to simulate the run in  $\mathcal{P}$ .

- First, we reduce the derivation into a sequence of productions  $X_1 \Rightarrow \gamma_1, X_2 \Rightarrow \gamma_2, \dots, X_m \Rightarrow \gamma_m$  which uniquely defines the derivation.
- The alphabet in which we represent strings that produce other strings include two special markers  $\downarrow$  and  $\uparrow$ .
- We start on the string  $\downarrow S \uparrow$ .
- If the current string is  $\hat{u} \downarrow X\beta$ , and it is the  $i$ -th one that has reached a string of this form, then it must hold that  $X = X_i$ . We follow it by  $\hat{u} \downarrow \gamma_i \uparrow \beta$ .
- If the current string is  $\hat{u} \downarrow \tau\beta$ , for some terminal  $\tau$ , we follow it by  $\hat{u}\tau \downarrow \beta$ .
- If the current string is  $\hat{u} \downarrow \uparrow \beta$ , then we follow it by  $\hat{u} \downarrow \beta$ .
- If the current string is  $\hat{u} \downarrow$ , there is no follow up.

Interestingly, this function is completely reversible, since to obtain a sequence of productions from a sequence of strings in this model, all we need to do is to remove the markers  $\downarrow$  and  $\uparrow$  and eliminate the duplicate strings that appear. We will borrow the name **plain** to talk about the function which receives a string and returns one which deletes all markers. It is obvious that the resulting derivation is the original one.

Furthermore, and more interestingly, we can extend the function **shape** to receive one of these strings and return a string in the alphabet  $\{0, 1, \downarrow, \uparrow\}$ . For two derivations that have the same shape, the resulting sequences have the same shape as well.

This sequence of strings represents a run in  $\mathcal{P}$  almost verbatim, and we only need to adapt it into a sequence of pushes, pops and reads: We make a run  $\rho$  which starts on  $q_0$ , pushes  $(X, \gamma_1, 0)$  to the stack, and moves to the state  $(X, \gamma_1, 0)$ . This pairs exactly to the strings  $\downarrow S$  and  $\downarrow \gamma_1 \uparrow$ , which are the first two in the sequence. Then, we read the sequence of strings in order. If the current string is  $\hat{u} \downarrow X\beta$ , and this is the  $i$ -th time a string of this form is seen, then the current state must be  $(Y, \alpha_1 X_i \alpha_2, k)$ , where  $|\alpha_1| = k$ ; we push  $(Y, \alpha_1 X_i \alpha_2, k + 1)$  onto the stack, and move to the state  $(X_i, \gamma_i, 0)$ . If the current string is  $\hat{u} \downarrow a\beta$  for some  $a \in \Sigma$ , and the current state is  $(X, \gamma, k)$ , we read  $\tau$ , and move to the state  $(X, \gamma, k)$ . If the current string is  $\hat{u} \downarrow (a, \sigma)\beta$  for some  $a \in \Sigma$ , and the current state is  $(X, \gamma, k)$ , we read  $a$ , output  $\sigma$ , and move to the state  $(X, \gamma, k)$ . If the current string is  $\hat{u} \downarrow \uparrow \beta$ , we pop the topmost state from the stack and we move into that state. It is straightforward to see that this run represents exactly the leftmost

derivation  $S \Rightarrow_{\mathcal{G}}^* \hat{w}$ , and that for each annotated string  $\hat{w} \in L(\mathcal{G})$  if and only if there is a run of  $\mathcal{P}$  over  $w$  that produces  $\mu = \text{ann}(\hat{w})$  as output.

This function is also reversible. Consider a run of  $\mathcal{P}$  over a string  $w$  which produces  $\mu$  as output. This run must start on  $q_0$ , and then push  $q_0$  and move onto a state  $(S, \alpha, 0)$  for some rule  $S \rightarrow \alpha$ . Thus, our first two strings in the sequence are  $\downarrow S \uparrow$  and  $\downarrow \alpha \uparrow$ . If the current state is  $(X, \alpha, k)$  and the next transition is to push  $(X, \alpha, k+1)$  onto the stack to move into the state  $(Y, \gamma, 0)$ , then the current string is of the form  $\hat{u} \downarrow X\beta$ , so we follow it by the string  $\hat{u} \downarrow \gamma \uparrow \beta$ . If the next transition is a pop, then the current string is  $\hat{u} \downarrow \beta$ , so we follow it by  $\hat{u} \downarrow \beta$ . If the current transition is a read, then the current string is  $\hat{u} \downarrow a\beta$  for  $a \in \Sigma$ , so we follow it by  $\hat{u} \downarrow \beta$ . If the current transition is a read-write, then the current string is  $\hat{u} \downarrow (a, \sigma)\beta$  for  $(a, \sigma) \in \Sigma \times \Omega$ , so we follow it by  $\hat{u}(a, \sigma) \downarrow \beta$ . It can easily be seen that using the original function over this resulting sequence would give the original sequence back. We point that these two reversible functions mean that there is a one to one correspondence between derivations of  $S \Rightarrow_{\mathcal{G}}^* \hat{w}$  and accepting runs of  $\mathcal{P}$  over  $w$  with output  $\mu = \text{ann}(\hat{w})$ .

Similarly to the observation we made before, we notice that if we start on a sequence in the intermediate model, the profile of the resulting run  $\rho$  is fully given by the shape of the sequence (at each step, the size of the stack will be equal to the number of markers  $\uparrow$  present in the string).

Now assume that  $\mathcal{G}$  is unambiguous. Seeing that  $\mathcal{P}$  is unambiguous as well comes straightforwardly from the fact that the functions presented above are bijective.

Assume now that  $\mathcal{G}$  is rigid. Let  $w$  be an unannotated string and consider two runs  $\rho_1$  and  $\rho_2$  of  $\mathcal{P}$  over  $w$  which output  $\mu_1$  and  $\mu_2$  respectively. Convert these two runs into sequences  $\mathcal{S}_1$  and  $\mathcal{S}_2$  in the intermediate model. Note that if we convert these two sequences into derivations  $S \Rightarrow_{\mathcal{G}}^* \mu_i(w)$ , they will have the same shape. We can apply the functions above to obtain the two runs  $\rho_1$  and  $\rho_2$  back, and note that they have the same profile. We conclude that if  $\mathcal{G}$  is rigid, then  $\mathcal{P}$  is profiled.  $\square$

We then show another direction:

**Proposition 21.** *For any PDAnn  $\mathcal{P}$ , we can build an equivalent annotated grammar  $\mathcal{G}$  in polynomial time. Further, if  $\mathcal{P}$  is unambiguous then so is  $\mathcal{G}$ . Moreover, if  $\mathcal{P}$  is profiled, then  $\mathcal{G}$  is rigid.*

*Proof.* This is again a standard transformation. We first transform the input PDAnn  $\mathcal{P} = (Q, \Sigma, \Omega, \Gamma, \Delta, q_0, F)$  to accept by empty stack, i.e., to accept iff the stack is empty. To do this, we build an equivalent PDAnn  $\mathcal{P}' = (Q', \Sigma, \Omega, \Gamma', \Delta', q'_0, F')$  where  $Q' = Q \cup \{q'_0, q_e, q_f\}$ ,  $\Gamma' = \Gamma \cup \{\gamma_0\}$ ,  $F = \{q_f\}$ , and we add the following transitions to  $\Delta$  to obtain  $\Delta'$ : A push transition  $(q'_0, q_0, \gamma_0)$ , a pop transition  $(q, \gamma_0, q_f)$  for every  $q \in F$  (for runs that accept at a point where the stack is already empty), plus a pop transition  $(q, \gamma, q_e)$  for any other  $\gamma \in \Gamma$ , and a pop transition  $(q_e, q_0, q_f)$ .

This clearly ensures that there is a bijection between the accepting runs of  $\mathcal{P}$  and those of  $\mathcal{P}'$ : given an accepting run  $\rho$  of  $\mathcal{P}$ , the bijection maps it to an accepting run of  $\mathcal{P}'$  by extending it with a push transition at the beginning, and pop transitions at the

end. Further, all accepting runs in  $\mathcal{P}'$  now finish with an empty stack, more specifically a run is accepting iff it finishes with an empty stack.

Now, we can perform the transformation. The nonterminals of the grammar are triples of the form  $(q, \gamma, q')$  for states  $q$  and  $q'$  and a stack symbol  $\gamma$ . Intuitively,  $(p, \gamma, q')$  will derive the strings that can be read by the PDAnn starting from state  $p$ , reaching some other state  $q$  with the same stack, not seeing the stack at all in the process, and then popping  $\gamma$  to reach  $q'$ .

The production rules are the following:

- A rule  $S \rightarrow (q_0, \gamma_0, q_f)$ .
- A rule  $(p, \gamma, q') \rightarrow (q, \gamma', r)(r, \gamma, q')$  for every nonterminal  $(p, \gamma, q')$ , push transition  $(p, q, \gamma') \in \Delta$  and state  $r$ .
- A rule  $(p, \gamma, q) \rightarrow \varepsilon$  for each pop transition  $(p, \gamma, q) \in \Delta$ .
- A rule  $(p, \gamma, q') \rightarrow (a, \sigma)(q, \gamma, q')$  for each read-write transition  $(p, (a, \sigma), q)$ , and a rule  $(p, \gamma, q') \rightarrow a(q, \gamma, q')$  for each read transition  $(p, a, q)$ , for each nonterminal  $(p, \gamma, q')$ .

As we did in Proposition 20, we will show a function which receives an accepting run  $\rho$  over  $w$  in  $\mathcal{P}$  with output  $\mu$  and outputs a leftmost derivation  $S = \alpha_1 \Rightarrow_{\mathcal{G}} \alpha_2 \Rightarrow_{\mathcal{G}} \dots \Rightarrow_{\mathcal{G}} \alpha_m = \mu(w)$ . The way we do this is quite straightforward: There is a one-to-one correspondence between snapshots in the run to each  $\alpha_i$ . Indeed, it can be seen that  $\alpha_i = \hat{u}(q_1, \gamma_1, q_2)(q_2, \gamma_2, q_3) \dots (q_k, \gamma_k, q_f)$  for some string  $\hat{u} \in (\Sigma \cup (\Sigma \times \Omega))^*$ , some states  $q_1, \dots, q_k$  and stack symbols  $\gamma_1, \dots, \gamma_k$ . Moreover, the  $i$ -th stack in the run is equal to  $\gamma_1 \gamma_2 \dots \gamma_k$ , whereas each state  $q_j$  is the first state that is reached after popping the respective  $\gamma_{j-1}$ . We see that this function is fully reversible, as each production corresponds unequivocally to a transition in particular. This implies that  $\mathcal{P}$  is unambiguous if, and only if  $\mathcal{G}$  is unambiguous.

For the next part of the proof, we bring attention to the fact that there are exactly four possible shapes on the right sides of the rules in  $\mathcal{G}$ . Each of these directly map to some type of transition, be it the initial push transition  $(q'_0, q_0, \gamma_0)$ , a different push transition, a pop transition, or a read (or read-write) transition. To be precise, these shapes are the strings 1, 11,  $\varepsilon$  and 01 respectively. From here it can be easily seen that, while comparing a run  $\rho$  to its respective derivation  $S \Rightarrow_{\mathcal{G}}^* \hat{w}$ , each production in the run immediately tells which type of transition was taken, and each transition in the run immediately tells which rule (and therefore, rule shape) was used. Therefore, each derivation shape maps to exactly one stack profile and vice versa, from which we conclude that  $\mathcal{P}$  is profiled if, and only if,  $\mathcal{G}$  is rigid.  $\square$

## C.2. Proof of Lemma 1

We will show a linear-time reduction to enumeration for an I/O-deterministic VPT [27]. We will state the necessary preliminaries to use the result presented there. We will also

adapt the models slightly to fit our results better while also keeping the results trivially equivalent.

A structured alphabet is a triple  $(\Sigma^<, \Sigma^>, \Sigma^|)$  consisting of three disjoint sets  $\Sigma^<$ ,  $\Sigma^>$ , and  $\Sigma^|$  that contain *open*, *close*, and *neutral* symbols respectively. The set of *well-nested strings* over  $\Sigma$ , denoted as  $\Sigma^{<*>}$ , is defined as the closure of the following rules:  $\Sigma^| \cup \{\varepsilon\} \subseteq \Sigma^{<*>}$ , if  $w_1, w_2 \in \Sigma^{<*>} \setminus \{\varepsilon\}$  then  $w_1 \cdot w_2 \in \Sigma^{<*>}$ , and if  $w \in \Sigma^{<*>}$  and  $a \in \Sigma^<$  and  $b \in \Sigma^>$  then  $a \cdot w \cdot b \in \Sigma^{<*>}$ .

A *Visibly Pushdown Transducer* (VPT) is a tuple  $\mathcal{T} = (Q, \hat{\Sigma}, \Gamma, \Omega, \Delta, I, F)$  where  $Q$  is a state set,  $\hat{\Sigma}$  is a structured alphabet,  $\Gamma$  is set of a stack symbols,  $\Omega$  is the output alphabet,  $I$  is the set of initial states,  $F$  is the set of final states, and  $\Delta \subseteq (Q \times (\Sigma^< \cup \Sigma^< \times \Omega) \times Q \times \Gamma) \cup (Q \times (\Sigma^> \cup \Sigma^> \times \Omega) \times \Gamma \times Q) \cup (Q \times (\Sigma^| \cup \Sigma^| \times \Omega) \times Q)$  is the transition relation. A run  $\rho$  of  $\mathcal{T}$  over a well-nested string  $w = a_1 a_2 \dots a_n \in (\Sigma \cup \Sigma \times \Omega)^{<*>}$  is a sequence of the form  $\rho = (q_1, \sigma_1) \xrightarrow{s_1} \dots \xrightarrow{s_n} (q_{n+1}, \sigma_{n+1})$  where  $q_i \in Q$ ,  $\sigma_i \in \Gamma^*$ ,  $q_1 \in I$ ,  $\sigma_1 = \varepsilon$ , each  $s_i$  is either equal to  $a_i$ , or to  $a_i/\sigma_i$  for some  $\sigma_i \in \Omega$ , and for every  $i \in [1, n]$  the following holds:

1. If  $s_i = a \in \Sigma^<$ , then  $(q_i, a, q_{i+1}, \gamma) \in \Delta$ , and if  $s_i = a/\sigma$ , for  $a \in \Sigma^<$ , then  $(q_i, (a, \sigma), q_{i+1}, \gamma) \in \Delta$ , for some  $\gamma \in \Gamma$  with  $\sigma_{i+1} = \gamma \sigma_i$ ,
2. If  $s_i = a \in \Sigma^>$ , then  $(q_i, a, \gamma, q_{i+1}) \in \Delta$ , and if  $s_i = a/\sigma$ , for  $a \in \Sigma^>$ , then  $(q_i, (a, \sigma), \gamma, q_{i+1}) \in \Delta$ , for some  $\gamma \in \Gamma$  with  $\sigma_i = \gamma \sigma_{i+1}$ , and
3. If  $s_i = a \in \Sigma^|$ , then  $(p_i, a, q_{i+1}) \in \Delta$ , and if  $s_i = a/\sigma$  for  $a \in \Sigma^|$ , then  $(p_i, (a, \sigma), q_{i+1}) \in \Delta$  with  $\sigma_i = \sigma_{i+1}$ . We say that the run is accepting if  $q_{n+1} \in F$ .

Given a VPT  $\mathcal{T}$  and a  $w \in \Sigma^{<*>}$ , we define the set  $\llbracket \mathcal{T} \rrbracket(w)$  of all outputs of  $\mathcal{T}$  over  $w$  as:  $\llbracket \mathcal{T} \rrbracket(w) = \{\text{ann}(\rho) \mid \rho \text{ is an accepting run of } \mathcal{T} \text{ over } w\}$ . where  $\text{ann}$  for runs in VPT is defined analogously to PDAnn. That is, if  $\rho = (q_1, \sigma_1) \xrightarrow{s_1} \dots \xrightarrow{s_n} (q_{n+1}, \sigma_{n+1})$ , then  $\text{ann}(\rho) = \omega_1 \dots \omega_n$  where  $\omega_i = (\sigma, i)$  if  $s_i = a/\sigma$  and  $\omega_i = \varepsilon$  otherwise.

We say that  $\mathcal{T}$  is *unambiguous* if for every  $w$  and  $\mu$  there is at most one accepting run  $\rho$  of  $\mathcal{T}$  which produces  $\mu$ . In [27] these VPT are called input-output unambiguous.

The theorem we use can be stated as follows:

**Theorem 22.** ([27], Theorem 3) *There is an algorithm that receives an unambiguous VPT  $\mathcal{T} = (Q, \hat{\Sigma}, \Gamma, \Omega, \Delta, q_0, F)$  and an input string  $w$ , and enumerates the set  $\llbracket \mathcal{T} \rrbracket(w)$  with output-linear delay after a preprocessing phase that takes  $\mathcal{O}(|Q|^2 \cdot |\Delta| \cdot |w|)$  time.*

The rest of the proof will consist on showing a linear-time reduction from the problem of enumerating the set  $\llbracket \mathcal{P} \rrbracket_\pi(w)$  for an unambiguous PDAnn  $\mathcal{P}$  and input string  $w$  to the problem of enumerating the set  $\llbracket \mathcal{T} \rrbracket(w')$  for an unambiguous VPT  $\mathcal{T}$ , and input string  $w'$ .

Let  $\mathcal{P} = (Q, \Sigma, \Omega, \Gamma, \Delta, q_0, F)$  and let  $w \in \Sigma^*$  be an input string. Consider the structured alphabet  $\hat{\Sigma} = (\{<\}, \{>\}, \Sigma)$  for some  $<, > \notin \Sigma$ . Assume  $\pi = \pi_1, \dots, \pi_m$ . We construct a well-nested string  $w' = b_1 \dots b_{m-1}$  where  $b_i = <$  if  $\pi_i > \pi_{i+1}$ ,  $b_i = >$  if  $\pi_i < \pi_{i+1}$ , and  $b_i = w_j$  otherwise, where  $i$  is the  $j$ -th index in which  $\pi_i = \pi_{i+1}$ . We also build a table  $\text{Ind}$  such that  $\text{Ind}(i) = j$  for each of the indices in the third case. We build a VPT

$\mathcal{T} = (Q, \hat{\Sigma}, \Gamma, \Omega, \Delta', I, F)$  where  $I = \{q_0\}$  and we get  $\Delta'$  by replacing every push transition  $(p, q, \gamma) \in \Delta$  by  $(p, \prec, q, \gamma)$  and every pop transition  $(p, \gamma, q) \in \Delta$  by  $(p, \succ, q, \gamma)$ . Note that read and read-write transitions are untouched.

Let  $w$  be an input string, and let  $\mu$  be an output. Consider the output  $\mu'$  which is obtained by shifting the indices in  $\mu$  to those that correspond in  $w'$ . We argue that for each run  $\rho$  of  $\mathcal{P}$  over  $w$  with profile  $\pi$  which produces  $\mu$ , there is exactly one run  $\rho'$  of  $\mathcal{T}$  over  $w'$  which produces  $\mu'$ , and vice versa. We see this by a straightforward induction argument on the size of the run. This immediately implies that for each output  $\mu \in \llbracket \mathcal{T} \rrbracket(w)$  there exists exactly one output  $\mu' \in \llbracket \mathcal{P} \rrbracket_\pi(w)$ , which has its indices shifted as we mentioned. The algorithm then consists on simulating the procedure from Theorem 22 over  $\mathcal{T}$  and  $w'$ , and before producing an output  $\mu$ , we replace the indices to the correct ones following the table `Ind`. The time bounds are unchanged since the table `Ind` has linear size in  $m$ , and replacing the index on some output  $\mu$  can be done linearly on  $|\mu|$ . We conclude that there is algorithm that enumerates the set  $\llbracket \mathcal{P} \rrbracket_\pi(w)$  with output-linear delay after a preprocessing that takes  $\mathcal{O}(|Q|^2 \cdot |\Delta| \cdot |\pi|)$  time.

### C.3. Proof of Proposition 11

For this result, we use the notion of PDA (Definition 18) and deterministic PDA (Definition 19) that were omitted from the main text.

As we have done in previous proofs, the strategy consists on starting with a profiled-deterministic PDAnn  $\mathcal{P}$ , and building a PDAnn  $\mathcal{P}'$  by eliminating the output symbols from each transition. This PDAnn behaves almost identically to a pushdown automaton  $\mathcal{A}$  in the sense that if  $w \in L(\mathcal{A})$ , then  $\llbracket \mathcal{P} \rrbracket(w) = \{\varepsilon\}$ , and that if  $w \notin L(\mathcal{A})$  then  $\llbracket \mathcal{P} \rrbracket(w) = \emptyset$ . Whenever this holds, we say that the PDAnn and the pushdown automaton are *equivalent*. It is simple to see that for this  $\mathcal{A}$  it holds that  $L(\mathcal{A}) = L'$ . To conclude the proof, we must show that  $\mathcal{A}$  can be made deterministic. Without loss of generality, we remove from  $\mathcal{A}$  all *inaccessible states*, i.e., all states for which there is no run that goes to the state.

First, we will prove that the PDAnn  $\mathcal{P}'$  is profiled-deterministic. Let  $w$  be a string in  $\Sigma^*$  and let  $\rho'_1$  and  $\rho'_2$  be two partial runs of  $\mathcal{P}'$  over  $w$  with the same profile, and with last configurations  $(q, i)$  and  $(q', i)$ . There clearly exist partial runs  $\rho_1$  and  $\rho_2$  of  $\mathcal{P}$  over  $w$  with the same profile, which can be obtained by replacing each transition by one of the transitions in  $\mathcal{P}$  it was replaced by. Since  $\mathcal{P}$  is profile-deterministic, then one of the following must hold in  $\mathcal{P}$ : (1)  $\Delta(q) \cup \Delta(q') \subseteq Q \times (\Sigma \cup \Sigma \times \Omega) \times Q$ , i.e., all transitions from  $q$  and  $q'$  are read or read-write transitions; (2)  $\Delta(q) \cup \Delta(q') \subseteq Q \times (Q \times \Gamma)$ , i.e., all transitions from  $q$  and  $q'$  are push transitions; or (3)  $\Delta(q) \cup \Delta(q') \subseteq (Q \times \Gamma) \times Q$ , i.e., all transitions from  $q$  and  $q'$  are pop transitions. Note that if (2) or (3) hold, then in the new PDAnn  $\mathcal{P}$  the condition holds again in  $\mathcal{P}'$  trivially since none of the transitions in  $\Delta(q)$  and  $\Delta(q')$  was changed. Moreover, if (1) holds, then it can be seen that all of the transitions that belonged in  $Q \times (\Sigma \times \Omega) \times Q$  now belong in  $Q \times \Sigma \times Q$ , which also leaves the condition unchanged in  $\mathcal{P}'$ . We conclude that  $\mathcal{P}'$  is profiled-deterministic.

The next step is to use Lemma 2 from  $\mathcal{P}'$  to obtain an equivalent PDAnn  $\mathcal{P}''$  which is deterministic-modulo-profile. We will argue that if we start with  $\mathcal{P}'$ , which was profiled-

deterministic, then the resulting  $\mathcal{P}''$  is equivalent to a pushdown automaton  $\mathcal{A}'$  which is also deterministic. Let  $w$  be an input string in  $\Sigma^*$  and let  $\rho''$  be a partial run of  $\mathcal{A}$  over  $w$  with last configuration  $(S, i)$  and with topmost symbol on the stack  $T$ . Let us recall what  $\mathcal{P}''$  being deterministic-modulo-profile entails that the following conditions hold:

1. There is at most one push transition that starts on  $S$ ; formally, we have:

$$|\{S', T \in Q'' \times \Gamma'' \mid (S, S', T) \in \Delta\}| \leq 1.$$

2. There is at most one pop transition that starts on  $S, T$ ; formally, for each  $\gamma$ , we have:

$$|\{S' \in Q \mid (S, \gamma, S') \in \Delta\}| \leq 1.$$

3. For each letter  $a$ , and output  $\sigma \in \Omega$ , there is at most one read-write transition that starts on  $S, a, \sigma$ ; formally, we have

$$|\{S' \in Q'' \mid (S, (a, \sigma), S') \in \Delta''\}| \leq 1.$$

4. For each letter  $a$ , there is at most one read transition that starts on  $S, a$ ; formally, we have:

$$|\{q \in Q'' \mid (S, a, S') \in \Delta''\}| \leq 1.$$

We will show that at most one of these conditions holds. Recall that in the transformation, the states of  $\mathcal{P}''$  are sets which contain pairs of states  $(p, q) \in Q' \times Q'$ , and the stack symbols are triples  $(p, \gamma, q) \in Q' \times \Gamma' \times Q'$ . Now, recall the claim that was proven in the lemma, on the backwards direction:

If  $\mathcal{P}''$  has a run  $\rho''$  on a string  $w$ , producing output  $\mu$ , from its initial state to an instantaneous description  $(S, i), \alpha'$  with  $\alpha' = T_0, \dots, T_m$  being the sequence of the stack symbols, then for any choice of elements  $(q_0, \gamma_0, p_0) \in T_0, (p_0, \gamma_1, p_1) \in T_1, \dots, (p_{m-1}, \gamma_m, p_m) \in T_m$  and  $(p_m, q) \in S$  it holds that  $\mathcal{P}'$  has a run  $\rho'$  on  $w$  producing output  $\mu$  from some initial state  $q_0$  to the instantaneous description  $(q, i), \alpha$  with  $\alpha = \gamma_0 q_0, \dots, \gamma_m q_m$  (writing next to each stack symbol the state that annotates it), and  $\rho''$  and  $\rho'$  have the same profile.

Since  $\mathcal{P}'$  is profiled-deterministic, then each run  $\rho'^+$  which continues  $\rho'$  by one step must have the same shape. This implies that exactly one of the following conditions must hold:

- The last transition in  $\rho'^+$  is a read or read-write transition. Therefore, all transitions from  $q$  are either read or read-write transitions.
- The last transition in  $\rho'^+$  is a push transition. Therefore, all transitions from  $q$  are pop transitions.
- The last transition in  $\rho'^+$  is a pop transition. Therefore, all transitions from  $q$  are pop transitions.

Assume bullet point 1 holds. Note that there are no read-write transitions in  $\mathcal{P}'$  so there are only read transitions. From here, we prove that only (4) is true simply by inspecting the transformation in the lemma; if (1) held, then there would be a push transition from  $q$  in  $\mathcal{P}'$ , if (2) held then there would be a pop transition from  $q$  and  $\gamma$  in  $\mathcal{P}'$ , and (3) never holds. Now, assume bullet point 2 holds. From here, we prove that only (1) can be true; if (2) held, then there would be a pop transition from  $q$  and  $\gamma$  in  $\mathcal{P}'$ , if (4) held, then there would be a read transition from  $q$  in  $\mathcal{P}'$ , and again, (3) is never true. Lastly, assume bullet point 3 holds. From here, we prove that only (2) can be true; if (1) held, then there would be a push transition from  $q$  in  $\mathcal{P}'$ , if (4) held, then there would be a read transition from  $q$  in  $\mathcal{P}'$ , and yet again, (3) is never true. We conclude that from the 4 points, at most one of these can be true at the same time.

Now we prove that the equivalent PDA  $\mathcal{A}$  is deterministic. Let  $q$  be a state of  $\mathcal{A}$ . As all states of  $\mathcal{A}$  are accessible, pick  $\rho''$  to be a run that reaches state  $q$ . We have argued that at most one of the points in the list above is true of  $\mathcal{P}'$ , and it cannot be point (3). Now, we see that (a) is equivalent to (4), that (b) is equivalent to (1) and (c) is equivalent to (2). Since only one of the conditions among (1), (2) or (4) can be true, the same holds for (a), (b) and (c), from which we conclude that  $\mathcal{A}$  is deterministic. This completes the proof.

#### C.4. Proof of Proposition 12

Consider a profiled-deterministic PDAnn  $\mathcal{P}$ . To prove that it is profiled, consider an input string  $w \in \Sigma^*$ . We will prove by a simple induction argument that any two runs of  $\mathcal{P}$  over  $w$  have the same profile. The base case is trivial since the run is of length 0, and the profile up to now is composed simply of the stack size 0. Assume now that for each pair of runs  $\rho$  and  $\rho'$  of  $\mathcal{P}$  over  $w$  of size  $k$ , that they have the same profile. We will show that for every pair of runs  $\rho_1$  and  $\rho_2$  over  $w$  of size  $k + 1$ , they have the same profile as well. Note that the runs  $\rho_1^-$  and  $\rho_2^-$  that are obtained by removing the last step have the same profile, by the hypothesis. From the definition of profiled-deterministic it can be directly seen that if (1) the last transition in  $\rho_1$  is a read or read-write transition, then for the runs  $\rho_1^-$  and  $\rho_2^-$ , the only choices are read or read-write transitions, from which we deduce that the last transition in  $\rho_2$  is a read or read-write transition as well, if (2) the last transition in  $\rho_1$  is a push transition, then for the run  $\rho_1^-$  and  $\rho_2^-$  the only choice are push transitions, and therefore the last transition in  $\rho_2$  has to be a push transition as well, and if (3) the last transition in  $\rho_1$  is a pop transition, then for  $\rho_1^-$  and  $\rho_2^-$  the only choices are pop transitions, so the last transition in  $\rho_2$  must be a pop transition as well. We obtain that  $\rho_1$  and  $\rho_2$  have the same profile, and from the induction argument, we conclude that  $\mathcal{P}$  is profiled.

Now, consider a profiled-deterministic PDAnn  $\mathcal{P}$  and an input string  $w$ . We will prove that the unique profile of accepting runs of  $\mathcal{P}$  over  $w$  can be computed in linear time in  $|w|$ . The way we do this is by using the pushdown automaton  $\mathcal{A}$  that was constructed in Proposition 11. By inspecting the proof, it can be seen that the unique profile of  $\mathcal{P}$  over  $w$  is maintained throughout the construction. Indeed, the first construction simply removes the output symbols, which does not affect the profile, and the second construction has

an invariant that keeps the profile intact as well. Therefore, by running the automaton  $\mathcal{A}$  over  $w$ , and storing the stack sizes at each step, we obtain a profile  $\pi$  which is exactly the same profile of the accepting runs of  $\mathcal{P}$  over  $w$ . To finish the proof, we only need to argue that this profile has linear size on  $|w|$  (from a data complexity perspective). This follows from the fact that any run of a deterministic pushdown automaton  $\mathcal{A}$  over a string  $w$  has  $\mathcal{O}(f(\mathcal{A}) \times |w|)$  length, for some computable function  $f$ . This can be seen from a counting argument: (1) There is a maximum stack size  $k$  that can be reached in an accepting run of  $\mathcal{P}$  over  $w$  from an empty stack through  $\varepsilon$ -transitions, which is given by the number of states in  $\mathcal{A}$ . Otherwise, there are two configurations which are reachable from one another in a way such the stack, as it was at the first configuration, is not seen. This implies that there is a loop, and since  $\mathcal{A}$  is deterministic,  $\mathcal{A}$  does not accept  $w$ . (2) From a given stack, the maximum numbers of steps that can be taken without reading from  $w$ , and without seeing the topmost symbol on the stack is given by the number of possible stacks of size  $k$ . (3) Between a read (or read-write) transition and the next one, the maximum height difference is  $k$ , and if we move out of a read (or read-write) transition with a certain stack, from (2) we can see that we can only do a fixed number of steps before consuming some symbol from this stack, and therefore, the number of steps is bounded by a factor depending on  $\mathcal{A}$  multiplied by the size of the stack up until this point, which is linear on the number of symbols in  $w$  read so far. We conclude that  $w'$  has size linear on  $w$ , from a data complexity point of view.

## D. Proofs of Section 6

### D.1. Formal definitions on extraction grammars

We first give formal definitions that were omitted from the main text of the paper for lack of space. We first formally define *ref-words*, i.e., strings with the special *variable operations*

**Definition 23** (Ref-words). *For the set of variables  $\mathcal{X}$  we define the variable operations of  $\mathcal{X}$  by  $C_{\mathcal{X}} = \bigcup_{x \in \mathcal{X}} \{\vdash_x, \dashv_x\}$ . A ref-word is a string over  $\Sigma \cup C_{\mathcal{X}}$ , and we let  $\text{plain} : (\Sigma \cup C_{\mathcal{X}})^* \rightarrow \Sigma^*$  be the morphism over ref-words that removes variable operations. A ref-word is valid if each variable in  $\mathcal{X}$  is opened exactly once and then closed exactly once. A valid ref-word  $r$  then defines a mapping  $\eta_r$  over  $\text{plain}(r)$  in the following way: for each variable  $x \in \mathcal{X}$ , there is a unique factorization  $r = r_x^p \cdot \vdash_x \cdot r_x \cdot \dashv_x \cdot r_x^s$ , and we set  $\eta^r(x) := [i, j]$  for  $i := |\text{plain}(r_x^p)| + 1$  and  $j := i + |\text{plain}(r_x)|$ .*

We then define *extraction grammars* from [30]: they are simply CFGs over ref-words:

**Definition 24** (Extraction grammar). *An extraction grammar is a tuple  $\mathcal{H} = (V, \Sigma, \mathcal{X}, P, S)$  where  $V$  is a finite set of nonterminals,  $\Sigma$  is the alphabet,  $\mathcal{X}$  is a set of variables,  $P$  is a finite set of production rules of the form  $A \rightarrow \alpha$  with  $A \in V$  and  $\alpha \in (V \cup \Sigma \cup C_{\mathcal{X}})^*$ , and  $S \in V$  is the start symbol. We assume that  $V$ ,  $\Sigma$ , and  $\mathcal{X}$  are pairwise disjoint. Note the difference with annotated grammars: the annotations must correspond to variable*

operations (in order to capture spans), and they are represented as separate terminals instead of annotating existing terminals, which we believe makes the design of enumeration algorithms less convenient.

The semantics of extraction grammars is similar to that for annotated grammars. It is defined through derivations. Specifically, the rules  $P$  define the (left) derivation relation  $\Rightarrow_{\mathcal{H}} \subseteq (V \cup \Sigma \cup \mathcal{C}_{\mathcal{X}})^* \times (V \cup \Sigma \cup \mathcal{C}_{\mathcal{X}})^*$  such that  $uA\beta \Rightarrow_{\mathcal{H}} u\alpha\beta$  iff  $u \in (\Sigma \cup \mathcal{C}_{\mathcal{X}})^*$ ,  $A \in V$ ,  $\alpha, \beta \in (V \cup \Sigma \cup \mathcal{C}_{\mathcal{X}})^*$ , and  $A \rightarrow \alpha \in P$ . We denote by  $\Rightarrow_{\mathcal{H}}^*$  the reflexive and transitive closure of  $\Rightarrow_{\mathcal{H}}$ . Then the language defined by  $\mathcal{H}$  is the set of ref-words  $L(\mathcal{H}) = \{w \in (\Sigma \cup \mathcal{C}_{\mathcal{X}})^* \mid S \Rightarrow_{\mathcal{H}}^* w\}$ . This language defines a spanner  $\llbracket \mathcal{H} \rrbracket$  as follows: for every document  $d \in \Sigma^*$ ,

$$\llbracket \mathcal{H} \rrbracket(d) = \{\eta^r \mid r \in L(\mathcal{H}), r \text{ is valid, and } \text{plain}(r) = d\}.$$

Note that ref-words that are not valid are ignored. An extraction grammar  $\mathcal{H}$  is called functional if every  $r \in L(\mathcal{H})$  is valid, and it is called unambiguous if for every  $r \in L(\mathcal{H})$  there exists exactly one derivation of  $r$  from  $S$ .

We can now formally define the *equivalence* between an extraction grammar and an annotated grammar. To do so, we first explain how we can translate mappings to annotations:

**Definition 25** (Output associated to a mapping). *Given a set  $\mathcal{X}$  of variables, the corresponding set of annotations  $\Omega_{\mathcal{X}}$  will be the powerset of  $\mathcal{C}_{\mathcal{X}}$ . Now, given a mapping  $\eta$  on a document  $d$  and variables  $\mathcal{X}$  to an annotation, we let  $\mathcal{I} = \bigcup_{x \in \mathcal{X}} \{i, j \mid \eta(x) = [i, j]\}$  be the set of indices which appear in some span of  $\eta$ . Further, for each  $k \in \mathcal{I}$ , let  $S_k = \{\vdash_x \mid \exists j. \eta(x) = [k, j]\} \cup \{\dashv_x \mid \exists i. \eta(x) = [i, k]\}$ . We now define the output  $\text{out}(\eta) = (i_1, S_{i_1}) \dots (i_m, S_{i_m})$  where  $\mathcal{I} = \{i_1, \dots, i_m\}$  and  $i_1 < \dots < i_m$ , namely, we group the captures for each position as a set and use this set as the annotation. Note that the largest index that appears in the annotation can be  $|d| + 1$  because of the range of spans.*

We can now define *equivalence* between extraction grammars and annotation grammars:

**Definition 26** (Equivalent annotated grammar). *We say that an extraction grammar  $\mathcal{H}$  over variables  $\mathcal{X}$  has an equivalent annotated grammar  $\mathcal{G}$  if  $\mathcal{G}$  is over the set of annotations  $\Omega_{\mathcal{X}}$  and over the alphabet  $\Sigma \cup \{\#\}$  for  $\#$  a fresh symbol, and if for every document  $d \in \Sigma^*$  and every mapping  $\eta$  of  $d$  over  $\mathcal{X}$ , we have  $\eta \in \llbracket \mathcal{H} \rrbracket(d)$  iff  $\text{out}(\eta) \in \llbracket \mathcal{G} \rrbracket(d \cdot \#)$ . The  $\#$ -symbol at the end is used because of the difference in the indexing of spans (from 1 to  $|d| + 1$ ) and annotations (from 1 to  $d$ ).*

## D.2. Proof of Proposition 13

Recall that, in the statement of this result, the formal notion of an *equivalent annotated grammar* is the one defined above (Definition 26). Recall also the formal definition of *ref-words* (Definition 23).

Let  $r$  be a ref-word in  $\Sigma \cup \mathcal{C}_{\mathcal{X}}$  and let  $\hat{w}$  be an annotated string in  $\Sigma \cup (\Sigma \times \Omega_{\mathcal{X}})$ . We say that  $r$  and  $\hat{w}$  are *equivalent* if  $\text{plain}(r \cdot \#) = \text{str}(\hat{w})$  and  $\text{out}(\eta^r) = \text{ann}(\hat{w})$ . For

example, the ref-word  $r_1 = \vdash_x \mathbf{a} \mathbf{a} \dashv_x \vdash_y \mathbf{b} \mathbf{b} \dashv_y \mathbf{b}$  is equivalent to  $\hat{w}_1 = (\mathbf{a}, \{\vdash_x\}) \mathbf{a} (\mathbf{b}, \{\dashv_x, \vdash_y\}) \mathbf{b} \mathbf{b} (\mathbf{b}, \{\dashv_y\}) \#$ .

The overall strategy of this proof is going to be to construct an annotated grammar  $\mathcal{G}$  in a way such that for every ref-word  $r \in L(\mathcal{H})$  there exists an equivalent annotated string  $\hat{w} \in L(\mathcal{G})$ , and vice versa. It is clear that this implies that  $\mathcal{G}$  and  $\mathcal{H}$  are equivalent.

The way we build  $\mathcal{G}$  will look like we are “pushing” the variable operations to the next terminal to the right. We will do this process one variable operation at a time.

First, we need to define an intermediate model between those of extraction grammars and annotated grammars. We define *extraction grammars with annotations* as a straightforward extension of extraction grammars which allow annotations on terminals that are not variable operations. For the set of variables  $\mathcal{X}$ , recall from Definition 23 that we define the variable operations of  $\mathcal{X}$  by  $\mathbf{C}_{\mathcal{X}} = \{\vdash_x, \dashv_x \mid x \in \mathcal{X}\}$ , and recall from Definition 25 that we define  $\Omega_{\mathcal{X}} = 2^{\mathbf{C}_{\mathcal{X}}}$ . An *extraction grammar with annotations* is a tuple  $\mathcal{F} = (V, \Sigma, \mathcal{X}, P, S)$  where  $V$  is a finite set of nonterminal symbols,  $\Sigma$  is an alphabet and  $\mathcal{X}$  is a set of variables, such that  $V$ ,  $\Sigma$ ,  $\Sigma \times \Omega_{\mathcal{X}}$ , and  $\mathbf{C}_{\mathcal{X}}$  are pairwise disjoint,  $P$  is a finite set of rules of the form  $A \rightarrow \alpha$  with  $A \in V$  and  $\alpha \in (V \cup \Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}})^*$ , and  $S \in V$  is the start symbol. As in the other models, the semantic of extraction grammars is defined through derivations. Specifically, the set  $P$  defines the (left) derivation relation  $\Rightarrow_{\mathcal{F}} \subseteq (V \cup \Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}})^* \times (V \cup \Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}})^*$  such that  $\hat{u}A\beta \Rightarrow_{\mathcal{F}} \hat{u}\alpha\beta$  iff  $\hat{u} \in (\Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}})^*$ ,  $A \in V$ ,  $\alpha, \beta \in (V \cup \Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}})^*$ , and  $A \rightarrow \alpha \in P$ . We denote by  $\Rightarrow_{\mathcal{F}}^*$  the reflexive and transitive closure of  $\Rightarrow_{\mathcal{F}}$ . Then the language defined by  $\mathcal{F}$  is  $L(\mathcal{F}) = \{\hat{w} \in (\Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}})^* \mid S \Rightarrow_{\mathcal{F}}^* \hat{w}\}$ . In addition, we assume that no string  $\hat{w}$  in  $L(\mathcal{F})$  has a variable operation as its last symbol.

An extraction grammar with annotations generates strings over  $\Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}}$ , which we now refer to as *annotated ref-words*, and each annotated ref-word defines an output. We will define the semantics of extraction grammars with annotations recursively by using the semantics of annotated grammars as a starting point, that is, by extending the function  $\text{ann}$  to receive strings over  $\Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}}$ . In particular, for an annotated ref-word  $\hat{r} \in (\Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}})^*$  the result of  $\text{ann}(\hat{r})$  stays the same. For a string  $\hat{r} = \hat{u}\kappa a \hat{v} \in (\Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}})^*$  where  $\kappa \in \mathbf{C}_{\mathcal{X}}$  and  $a \in \Sigma$ , we define  $\text{ann}(\hat{r}) = \text{ann}(\hat{u}(a, \{\kappa\})\hat{v})$ , and for a string  $\hat{r} = \hat{u}\kappa(a, \sigma)\hat{v} \in (\Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}})^*$  where  $\sigma \in \Omega_{\mathcal{X}}$ , we define  $\text{ann}(\hat{r}) = \text{ann}(\hat{u}(a, \sigma \cup \{\kappa\})\hat{v})$ .

Further, we extend the function  $\text{str}$  to receive strings over  $\Sigma \cup (\Sigma \times \Omega_{\mathcal{X}}) \cup \mathbf{C}_{\mathcal{X}}$  as  $\text{str}(\hat{r}) = \text{str}(\text{plain}(\hat{r}))$ . Therefore, for an extraction grammar with annotations  $\mathcal{F}$  and a string  $w \in \Sigma^*$  we define the function  $\llbracket \mathcal{F} \rrbracket$  as:  $\llbracket \mathcal{F} \rrbracket(w) := \{\text{ann}(\hat{r}) \mid \hat{r} \in L(\mathcal{F}) \wedge \text{str}(\hat{r}) = w\}$ . We maintain the notions of equivalency between annotated grammars, extraction grammars, and extraction grammars with annotations. Likewise, we define equality between annotated strings, ref-words and annotated ref-words in the obvious way. Further, we note that any extraction grammar  $\mathcal{H} = (V, \Sigma, \mathcal{X}, P, S)$  is equivalent to the extraction grammar with annotations  $\mathcal{F} = (V', \Sigma \cup \{\#\}, \mathcal{X}, P', S')$  where  $V' = V \cup \{S'\}$  for some  $S' \notin V$ , and  $P' = P \cup \{S' \rightarrow S\#\}$ . It is obvious that  $\mathcal{F}$  is unambiguous if and only if  $\mathcal{H}$  is unambiguous.

We proceed as follows. First we convert  $\mathcal{H}$  into a functional extraction grammar (see Definition 24). As detailed in Peterfreund’s work [30, Propositions 10 and 12], this

takes running time  $\mathcal{O}(3^{2k}|\mathcal{H}|^2)$ , and if the initial grammar is unambiguous then so is the resulting grammar. Note that this implies that every ref-word  $r$  which is derivable from  $S$  contains each variable operation at most once. Hence we can build an equivalent extraction grammar with annotations  $\mathcal{F} = (V, \Sigma, \mathcal{X}, P, S)$  using the technique above. We then convert  $\mathcal{F}$  into a version of CNF which is slightly more restrictive than arity-two normal form: We allow rules of the form  $X \rightarrow YZ$ ,  $X \rightarrow \varepsilon$  and  $X \rightarrow \tau$ , for nonterminals  $X, Y$  and  $Z$  and a terminal  $\tau$ , but rules of the form  $X \rightarrow Y$  are not permitted. Converting to this formalism can be done in linear time in  $|\mathcal{F}|$  while preserving unambiguity, e.g., by transforming rules of the form  $X \rightarrow Y$  to  $X \rightarrow EY$  for some fresh nonterminal  $E$  with a rule  $E \rightarrow \varepsilon$ , and otherwise applying our result on arity-2 normal form (Proposition 2). We pick an order over the variable operations in  $\mathcal{C}_{\mathcal{X}}$  and for each  $\kappa \in \mathcal{C}_{\mathcal{X}}$  we do the following:

Define a function  $\text{proc}_{\kappa}$  that receives an annotated ref-word  $\hat{r} \in (\Sigma \cup (\Sigma \times 2^{\mathcal{C}_{\mathcal{X}}}) \cup \mathcal{C}_{\mathcal{X}})^*$  and:

1. if  $\hat{r} \in (\Sigma \cup (\Sigma \times 2^{\mathcal{C}_{\mathcal{X}}}) \cup (\mathcal{C}_{\mathcal{X}} \setminus \{\kappa\}))^*$ , then  $\text{proc}_{\kappa}(\hat{r}) = \hat{r}$ ,
2. if  $\hat{r} = \hat{u}\kappa\beta a\hat{v}$  for some  $\hat{u}, \hat{v} \in (\Sigma \cup (\Sigma \times 2^{\mathcal{C}_{\mathcal{X}}}) \cup (\mathcal{C}_{\mathcal{X}} \setminus \{\kappa\}))^*$ ,  $\beta \in (\mathcal{C}_{\mathcal{X}} \setminus \{\kappa\})^*$  and  $a \in \Sigma$ , then  $\text{proc}_{\kappa}(\hat{r}) = \hat{u}\kappa\beta(a, \{\kappa\})\hat{v}$ ,
3. if  $\hat{r} = \hat{u}\kappa\beta(a, T)\hat{v}$ , with  $T \subseteq \mathcal{C}_{\mathcal{X}} \setminus \{\kappa\}$ , then  $\text{proc}_{\kappa}(\hat{r}) = \hat{u}\kappa\beta(a, T \cup \{\kappa\})\hat{v}$ , and
4.  $\text{proc}_{\kappa}$  is undefined in any other case.

It is straightforward to see that the annotated ref-words  $\hat{r}$  and  $\hat{t} = \text{proc}_{\kappa}(\hat{r})$  are equivalent whenever  $\text{proc}_{\kappa}(\hat{r})$  is defined.

We build an extraction grammar with annotations  $\mathcal{F}' = (V', \Sigma, \mathcal{X}, P', S')$  where  $V' = V_{\text{out}} \cup V_{\text{in}} \cup V_{\text{left}} \cup V_{\text{mid}} \cup V_{\text{right}} \cup \{S'\}$ , and  $V_{\text{scr}} = \{A^{\text{scr}} \mid A \in P\}$  for  $\text{scr} \in \{\text{out}, \text{in}, \text{left}, \text{mid}, \text{right}\}$ , and  $P'$  is defined by the following rules:

- For each rule  $S \rightarrow AB$  in  $P$ , we add the rules  $S' \rightarrow A^{\text{out}}B^{\text{out}}$ ,  $S' \rightarrow A^{\text{in}}B^{\text{out}}$  and  $S' \rightarrow A^{\text{left}}B^{\text{right}}$  to  $P'$ .
- For each rule  $A \rightarrow BC$ , we add the rules  $A^{\text{out}} \rightarrow B^{\text{out}}C^{\text{out}}$ ,  $A^{\text{in}} \rightarrow B^{\text{in}}C^{\text{out}}$ ,  $A^{\text{in}} \rightarrow B^{\text{out}}C^{\text{in}}$ ,  $A^{\text{in}} \rightarrow B^{\text{left}}C^{\text{right}}$ ,  $A^{\text{left}} \rightarrow B^{\text{out}}C^{\text{left}}$ ,  $A^{\text{left}} \rightarrow B^{\text{left}}C^{\text{mid}}$ ,  $A^{\text{mid}} \rightarrow B^{\text{mid}}C^{\text{mid}}$ ,  $A^{\text{right}} \rightarrow B^{\text{right}}C^{\text{out}}$  and  $A^{\text{right}} \rightarrow B^{\text{mid}}C^{\text{right}}$  to  $P'$ .
- For each rule  $A \rightarrow a$ ,  $a \in \Sigma$ , we add  $A^{\text{out}} \rightarrow a$  and  $A^{\text{right}} \rightarrow (a, \{\kappa\})$ .
- For each rule  $A \rightarrow (a, T)$ ,  $a \in \Sigma$  and  $T \subseteq \mathcal{C}_{\mathcal{X}} \setminus \{\kappa\}$ , we add  $A^{\text{out}} \rightarrow (a, T)$  and  $A^{\text{right}} \rightarrow (a, T \cup \{\kappa\})$ .
- For each rule  $A \rightarrow \kappa'$ ,  $\kappa' \in \mathcal{C}_{\mathcal{X}} \setminus \{\kappa\}$ , we add the rules  $A^{\text{out}} \rightarrow \kappa'$  and  $A^{\text{mid}} \rightarrow \kappa'$ .
- For each rule  $A \rightarrow \kappa$ , we add the rule  $A^{\text{left}} \rightarrow \varepsilon$ .
- For each rule  $A \rightarrow \varepsilon$ , we add the rules  $A^{\text{out}} \rightarrow \varepsilon$  and  $A^{\text{mid}} \rightarrow \varepsilon$ .

Let  $\hat{\Sigma}_\kappa = \Sigma \cup (\Sigma \times 2^{\mathcal{C}_\mathcal{X}}) \cup \mathcal{C}_\mathcal{X} \setminus \{\kappa\}$ . For each nonterminal  $A \in V$  these hold:

$$\begin{aligned}
L(A^{\text{out}}) &= \{\hat{w} \mid A \Rightarrow_{\mathcal{F}}^* \hat{w}, \text{ where } \hat{w} \in \hat{\Sigma}_\kappa^*\} \\
L(A^{\text{in}}) &= \{\text{proc}_\kappa(\hat{w}) \mid A \Rightarrow_{\mathcal{F}}^* \hat{w}, \text{ where } \hat{w} = \hat{u}\kappa\beta a\hat{v} \text{ or } \hat{w} = \hat{u}\kappa\beta(a, T)\hat{v}, \\
&\quad \hat{u}, \hat{v} \in \hat{\Sigma}_\kappa^*, \beta \in (\mathcal{C}_\mathcal{X} \setminus \{\kappa\})^*, a \in \Sigma, T \subseteq \mathcal{C}_\mathcal{X} \setminus \{\kappa\}\} \\
L(A^{\text{left}}) &= \{\hat{w}\beta \mid A \Rightarrow_{\mathcal{F}}^* \hat{w}\kappa\beta, \text{ where } \hat{w} \in \hat{\Sigma}_\kappa^*, \beta \in (\mathcal{C}_\mathcal{X} \setminus \{\kappa\})^*\} \\
L(A^{\text{right}}) &= \{\beta(a, \{\kappa\})\hat{w} \mid A \Rightarrow_{\mathcal{F}}^* \beta a\hat{w}, \text{ where } \hat{w} \in \hat{\Sigma}_\kappa^*, \beta \in (\mathcal{C}_\mathcal{X} \setminus \{\kappa\})^*\} \cup \\
&\quad \{\beta(a, T \cup \{\kappa\})\hat{w} \mid A \Rightarrow_{\mathcal{F}}^* \beta(a, T)\hat{w}, \text{ where } \hat{w} \in \hat{\Sigma}_\kappa^*, \beta \in (\mathcal{C}_\mathcal{X} \setminus \{\kappa\})^*, T \subseteq \mathcal{C}_\mathcal{X} \setminus \{\kappa\}\} \\
L(A^{\text{mid}}) &= \{\beta \mid A \Rightarrow_{\mathcal{F}}^* \beta, \text{ where } \beta \in (\mathcal{C}_\mathcal{X} \setminus \{\kappa\})^*\}
\end{aligned}$$

These equalities are given without proof since they are not used, and are just for illustrating the idea behind the proof.

For the rest of our proof we will represent derivations  $X \Rightarrow^* \delta$  as the sequence of productions  $X_1 \Rightarrow \gamma_1, X_2 \Rightarrow \gamma_2, \dots, X_m \Rightarrow \gamma_m$ , where  $X_1 = X$  and  $\gamma_m = \delta$ , which uniquely determines the derivation by doing it in the leftmost way. We use this representation to state exactly how derivations in  $\mathcal{F}$  are translated to derivations in  $\mathcal{F}'$  and vice versa.

Another notion we need to address is how, in a given derivation  $X \Rightarrow^* \delta$ , instances of nonterminals are located with respect to each other. By an *instance of a nonterminal* (or just *instance*), we mean an  $X_i$  along with some specific derivation  $X_i \Rightarrow \gamma_i$  in the sequence. For some instances  $X_i$  and  $X_j$ , we say that  $X_i$  is a *descendant* of  $X_j$  if  $X_i = X_j$ , or if  $Y \Rightarrow X_i Z$ , or  $Y \Rightarrow Z X_i$  for some  $Y$  which is a descendant of  $X_j$ . We say that  $X_i$  is *to the left of*  $X_j$  (or  $X_j$  is *to the right of*  $X_i$ ) if there is a derivation  $X \Rightarrow YZ$  in the sequence such that  $X_i$  is a descendant of  $Y$  and  $X_j$  is a descendant of  $Z$ .

We note a few things in our construction: (1) Each  $X \in V_{\text{out}}$  only produces terminals (which do not include  $\kappa$ ) and nonterminals in  $V_{\text{out}}$ ; furthermore, every rule  $X \rightarrow YZ$  in  $P$  is copied into  $P'$  as  $X^{\text{out}} \rightarrow Y^{\text{out}}Z^{\text{out}}$ . (2) Nonterminals in  $V_{\text{in}}$  do not produce any terminals or  $\varepsilon$  directly, so they need to derive into some  $X \in V_{\text{in}}$  and some  $Y \in V_{\text{right}}$  to derive some string. (3) As with  $V_{\text{out}}$ , each  $X \in V_{\text{mid}}$  only produces terminals in  $\mathcal{C}_\mathcal{X} \setminus \{\kappa\}$  and nonterminals in  $V_{\text{mid}}$ . (4) Each  $X \in V_{\text{left}}$  (resp.  $V_{\text{right}}$ ) produces exactly one nonterminal  $X' \in V_{\text{left}}$  (resp.  $V_{\text{right}}$ ), or  $\varepsilon$  (resp.  $(a, T)$  for some  $a \in \Sigma$  and  $T \subseteq \mathcal{C}_\mathcal{X}$  such that  $\kappa \in T$ ); this, as a consequence, means that on each derivation from  $\mathcal{F}'$  where the first production is not  $S' \Rightarrow X^{\text{out}}Y^{\text{out}}$  there is exactly one derivation  $X \Rightarrow \varepsilon$  such that  $X \in V_{\text{left}}$  (resp., exactly one derivation  $X \Rightarrow (a, T)$ , such that  $X \in V_{\text{right}}$ ).

From point (1) we see that each annotated ref-word  $\hat{r} \in L(\mathcal{F})$  such that  $\hat{r} \in \hat{\Sigma}_\kappa^*$  (this is, which does not mention  $\kappa$  at all) can be derived by  $S'$  starting by  $S' \Rightarrow A^{\text{out}}B^{\text{out}}$ ,  $S' \Rightarrow a$ ,  $S' \Rightarrow \varepsilon$  or  $S' \Rightarrow \kappa'$ .

On the other hand, each annotated ref-word  $\hat{r} \in L(\mathcal{F}')$  which does not have  $\kappa$  on any annotation set was necessarily derived through rules of the form  $X^{\text{out}} \rightarrow Y^{\text{out}}Z^{\text{out}}$  which correspond to the rule  $X \rightarrow YZ$  in  $P$ , so we deduce that  $\hat{r} \in L(\mathcal{F})$ .

We shall now prove that for any string  $\hat{r} = \hat{u}\kappa\beta a\hat{v}$ , or  $\hat{r} = \hat{u}\kappa\beta(a, T)\hat{v}$ , where  $\hat{u}, \hat{v} \in \hat{\Sigma}_\kappa^*$ ,  $\beta \in \mathcal{C}_\mathcal{X} \setminus \{\kappa\}$ ,  $a \in \Sigma$  and  $T \subseteq \mathcal{C}_\mathcal{X} \setminus \{\kappa\}$  such that  $\hat{r} \in L(\mathcal{F})$ , it holds that  $\text{proc}(\hat{r}) \in L(\mathcal{F}')$ .

W.l.o.g., let  $\hat{r} = \hat{u}\kappa\beta a\hat{v}$  and consider some leftmost derivation of  $\hat{r}$  from  $\mathcal{F}$ :

$$\begin{aligned}
S &\Rightarrow_{\mathcal{F}}^* \hat{u}_1 A \delta \\
&\Rightarrow_{\mathcal{F}} \hat{u}_1 B C \delta \\
&\Rightarrow_{\mathcal{F}}^* \hat{u}_1 \hat{u}_2 \kappa \beta_1 C \delta \\
&\Rightarrow_{\mathcal{F}}^* \hat{u}_1 \hat{u}_2 \kappa \beta_1 \beta_2 a \hat{v}_1 \delta \\
&\Rightarrow_{\mathcal{F}}^* \hat{u}_1 \hat{u}_2 \kappa \beta_1 \beta_2 a \hat{v}_1 \hat{v}_2 = \hat{r},
\end{aligned}$$

where we have that  $\beta = \beta_1 \beta_2$ ,  $\hat{u} = \hat{u}_1 \hat{u}_2$  and  $\hat{v} = \hat{v}_1 \hat{v}_2$ . Note that this is an arbitrary derivation, and we are merely identifying these nonterminals  $A$ ,  $B$  and  $C$ . We also identify the nonterminals  $D$ , which produces  $\kappa$ , and  $E$ , which produces  $a$ . For the rest of the current part of the proof, we only refer to the *instances* of these nonterminals. Using this, we build a derivation from  $\mathcal{F}'$  step by step:

1. We have  $S' \Rightarrow_{\mathcal{F}'}^* \hat{u}_1 A^{\text{in}} \delta'$ , where  $\delta'$  is obtained by replacing each nonterminal  $X$  in  $\delta$  by  $X^{\text{out}}$ . We get this by starting with the derivation  $S \Rightarrow_{\mathcal{F}}^* \hat{u}_1 A \delta$ , and replacing  $X \Rightarrow_{\mathcal{F}} YZ$  by  $X^{\text{in}} \Rightarrow_{\mathcal{F}'} Y^{\text{in}} Z^{\text{out}}$  if  $A$  is a descendant of  $Y$ , by  $X^{\text{in}} \Rightarrow_{\mathcal{F}'} Y^{\text{out}} Z^{\text{in}}$  if  $Z$  is, or by  $X^{\text{out}} \Rightarrow_{\mathcal{F}'} Y^{\text{out}} Z^{\text{out}}$  if none is. We also replace each  $X \Rightarrow_{\mathcal{F}} \tau$ , for  $\tau \in \Sigma \cup (\Sigma \times 2^{\mathcal{C}\mathcal{X}}) \cup \mathcal{C}\mathcal{X} \setminus \{\kappa\} \cup \{\varepsilon\}$ , by  $X^{\text{out}} \Rightarrow_{\mathcal{F}'} \tau$ . If  $A = S$ , we replace  $A$  it by  $S'$ .
2. We have the rule  $A^{\text{in}} \rightarrow B^{\text{left}} C^{\text{right}}$  which was added to  $P'$ .
3. We have  $B^{\text{left}} \Rightarrow_{\mathcal{F}'}^* \hat{u}_2 \beta_1$ . We get this by starting from  $B \Rightarrow_{\mathcal{F}}^* \hat{u}_2 \kappa \beta_1$ , and we replace  $X \Rightarrow_{\mathcal{F}} YZ$  by  $X^{\text{left}} \Rightarrow_{\mathcal{F}'} Y^{\text{left}} Z^{\text{mid}}$  if  $E$  is a descendant of  $D$ , by  $X^{\text{left}} \Rightarrow_{\mathcal{F}'} Y^{\text{out}} Z^{\text{left}}$  if  $Z$  is, by  $X^{\text{out}} \Rightarrow_{\mathcal{F}'} Y^{\text{out}} Z^{\text{out}}$  if  $X$  is to the left of  $D$ , and by  $X^{\text{mid}} \Rightarrow_{\mathcal{F}'} Y^{\text{mid}} Z^{\text{mid}}$  if it is to the right. We also replace  $X \Rightarrow_{\mathcal{F}} \tau$  by  $X^{\text{out}} \Rightarrow_{\mathcal{F}'} \tau$  if  $X$  is to the left of  $D$ , and by  $X^{\text{mid}} \Rightarrow_{\mathcal{F}'} \tau$  if it is to the right. Lastly, we replace  $D \Rightarrow_{\mathcal{F}} \kappa$  by  $D^{\text{left}} \Rightarrow_{\mathcal{F}'} \varepsilon$ .
4. We have  $C^{\text{right}} \Rightarrow_{\mathcal{F}'}^* \beta_2(a, \{\kappa\}) \hat{v}_1$ . We get this by starting from  $C \Rightarrow_{\mathcal{F}}^* \beta_2 a \hat{v}_1$ , and we replace  $X \Rightarrow_{\mathcal{F}} YZ$  by  $X^{\text{right}} \Rightarrow_{\mathcal{F}'} Y^{\text{mid}} Z^{\text{right}}$  if  $E$  is descendant of  $Z$ , or by  $X^{\text{right}} \Rightarrow_{\mathcal{F}'} Y^{\text{right}} Z^{\text{out}}$  if  $Y$  is, by  $X^{\text{mid}} \Rightarrow_{\mathcal{F}'} Y^{\text{mid}} Z^{\text{mid}}$  if  $X$  is to the left of  $E$ , and by  $X^{\text{right}} \Rightarrow_{\mathcal{F}'} Y^{\text{right}} Z^{\text{right}}$  if it is to the right. We also replace  $X \Rightarrow_{\mathcal{F}} \tau$  by  $X^{\text{mid}} \Rightarrow_{\mathcal{F}'} \tau$  if  $X$  is to the left of  $E$ , and by  $X^{\text{out}} \Rightarrow_{\mathcal{F}'} \tau$  if it is to the right. Lastly, we replace  $E \Rightarrow_{\mathcal{F}} a$  by  $E^{\text{right}} \Rightarrow_{\mathcal{F}'} (a, \{\kappa\})$ .
5. We have  $\delta' \Rightarrow_{\mathcal{F}'}^* \hat{v}_2$ , which we obtain from  $\delta \Rightarrow_{\mathcal{F}}^* \hat{v}_2$  by replacing each  $X \Rightarrow_{\mathcal{F}} YZ$  by  $X^{\text{out}} \Rightarrow_{\mathcal{F}'} Y^{\text{out}} Z^{\text{out}}$ , and each  $X \Rightarrow_{\mathcal{F}} \tau$  by  $X^{\text{out}} \Rightarrow_{\mathcal{F}'} \tau$ .

In the end, we get the following leftmost derivation from  $\mathcal{F}'$ :

$$\begin{aligned}
S &\Rightarrow_{\mathcal{F}'}^* \hat{u}_1 A^{\text{in}} \delta' \text{ (or } \hat{u}_1 S' \delta') \\
&\Rightarrow_{\mathcal{F}'} \hat{u}_1 B^{\text{left}} C^{\text{right}} \delta' \\
&\Rightarrow_{\mathcal{F}'}^* \hat{u}_1 \hat{u}_2 \beta_1 C^{\text{right}} \delta' \\
&\Rightarrow_{\mathcal{F}'}^* \hat{u}_1 \hat{u}_2 \beta_1 \beta_2(a, \{\kappa\}) \hat{v}_1 \delta' \\
&\Rightarrow_{\mathcal{F}'}^* \hat{u}_1 \hat{u}_2 \beta_1 \beta_2(a, \{\kappa\}) \hat{v}_1 \hat{v}_2 = \text{proc}_{\kappa}(\hat{r}),
\end{aligned}$$

which proves that  $\text{proc}(\hat{r}) \in L(\mathcal{F}')$ .

We will prove that for every  $\hat{s} = \hat{u}(a, T)\hat{v} \in L(\mathcal{F}')$  where  $\kappa \in T$  there is  $\hat{r}$  such that  $\text{proc}_\kappa(\hat{r}) = \hat{s}$  in a similar way. We argue that any leftmost derivation that produces  $\hat{s}$  has the following form:

$$\begin{aligned}
S' &\Rightarrow_{\mathcal{F}'}^* \hat{u}_1 A^{\text{in}} \delta_1 \quad (\text{or } \hat{u}_1 S' \delta_1) \\
&\Rightarrow_{\mathcal{F}'} \hat{u}_1 B^{\text{left}} C^{\text{right}} \delta_1 \\
&\Rightarrow_{\mathcal{F}'}^* \hat{u}_1 \hat{u}_2 D^{\text{left}} \delta_2 C^{\text{right}} \delta_1 \\
&\Rightarrow_{\mathcal{F}'} \hat{u}_1 \hat{u}_2 \delta_2 C^{\text{right}} \delta_1 \\
&\Rightarrow_{\mathcal{F}'}^* \hat{u}_1 \hat{u}_2 \beta_1 C^{\text{right}} \delta_1 \\
&\Rightarrow_{\mathcal{F}'}^* \hat{u}_1 \hat{u}_2 \beta_1 \beta_2(a, T) \hat{v}_1 \delta_1 \\
&\Rightarrow_{\mathcal{F}'}^* \hat{u}_1 \hat{u}_2 \beta_1 \beta_2(a, T) \hat{v}_1 \hat{v}_2 = \hat{s},
\end{aligned}$$

where  $\delta \in V_{\text{out}}^*$ ,  $\delta' \in V_{\text{mid}}^*$ ,  $\beta_1, \beta_2 \in (\mathcal{C}_{\mathcal{X}} \setminus \{\kappa\})^*$ ,  $\hat{u} = \hat{u}_1 \hat{u}_2$  and  $\hat{v} = \hat{v}_1 \hat{v}_2$ . The reasoning goes as follows:

- We know that  $S' \Rightarrow_{\mathcal{F}'}^* \hat{u}\beta(a, T)\hat{v} \in L(\mathcal{F}')$ . If  $X \Rightarrow_{\mathcal{F}'} (a, T)$  and  $\kappa \in T$ , then  $X \in V_{\text{right}}$ .
- From the way  $\mathcal{F}'$  was built, there must be a production  $X \Rightarrow_{\mathcal{F}'} YZ$  in  $S' \Rightarrow_{\mathcal{F}'}^* \hat{s}$  such that  $X \in V_{\text{in}}$  (or  $X = S'$ ),  $Y \in V_{\text{left}}$  and  $Z \in V_{\text{right}}$ , as it is the only way to derive a nonterminal in  $V_{\text{right}}$ . Let  $A^{\text{in}}$  (or  $S'$ ),  $B^{\text{left}}$  and  $C^{\text{right}}$  be these  $X$ ,  $Y$  and  $Z$  respectively.
- Seeing the rules in  $P'$  we note that every string of terminals that is derivable from  $B^{\text{left}}$  is of the form  $\hat{w}\beta$ , where  $\hat{w} \in \hat{\Sigma}_\kappa$  and  $\beta \in (\mathcal{C}_{\mathcal{X}} \setminus \{\kappa\})^*$ . Furthermore, this string satisfies that there is a production  $X \Rightarrow_{\mathcal{F}'} \varepsilon$  for some  $X \in V_{\text{left}}$  such that this  $\varepsilon$  is exactly at the left of where  $\beta$  begins. Let  $\hat{u}_2$  be this  $\hat{w}$ , let  $D^{\text{left}}$  be this  $X$ , and let  $\beta_1$  be this  $\beta$ .
- Likewise, we note that  $C^{\text{right}}$  always derives a string of terminals of the form  $\beta(a', T')\hat{w}$  for some  $\beta \in (\mathcal{C}_{\mathcal{X}} \setminus \{\kappa\})^*$  and  $\hat{w} \in \hat{\Sigma}_\kappa$ . Let  $\beta_2$  be this  $\beta$  and let  $\hat{v}_1$  be this  $\hat{w}$ .
- Lastly, let  $S' \Rightarrow_{\mathcal{F}'}^* \hat{u}_1 A^{\text{in}} \delta_1$  (or  $\hat{u}_1 S' \delta_1$ ) be the one that derives  $\hat{s}$ . From the rules in  $P'$ , we note that  $\delta_1$  is composed solely of nonterminals in  $V_{\text{out}}$ .

An important point that can be seen from this reasoning is that for each instance  $X \neq S'$  that appears in the derivation  $S' \Rightarrow_{\mathcal{F}'}^* \hat{s}$ , we can deduce the set  $V_{\text{scr}}$  for which  $X \in V_{\text{scr}}$ , among the options  $\text{scr} \in \{\text{out}, \text{in}, \text{left}, \text{mid}, \text{right}\}$ , by seeing its position in the derivation. To be precise, this is given from how  $X$  relates to the instances  $D^{\text{left}} \Rightarrow \varepsilon$ , and to  $E^{\text{right}} \Rightarrow (a, T)$ , for the nonterminal  $E^{\text{right}} \in V_{\text{right}}$  that satisfies this. (1) If  $X$  is to the left of  $D^{\text{left}}$ , then  $X \in V_{\text{out}}$ , (2) if  $D^{\text{left}}$  is a descendant of  $X$ , but  $E^{\text{right}}$  is not, then  $X \in V_{\text{left}}$ , (3) if  $X$  is to the right of  $D^{\text{left}}$ , and is to the left of  $E^{\text{right}}$ , then  $X \in V_{\text{mid}}$ , (4) if  $E^{\text{right}}$  is a descendant of  $X$ , but  $D^{\text{left}}$  is not, then  $X \in X_{\text{right}}$ , (5) if  $X$  is to the right of

$E^{\text{right}}$ , then  $X \in V_{\text{out}}$ , and (6) if both  $D^{\text{left}}$  and  $E^{\text{right}}$  are descendants of  $X$ , then  $X \in V_{\text{in}}$ . We bring attention to the fact that in this paragraph we referred only to the instances of  $D^{\text{left}}$  and  $E^{\text{right}}$  on the derivations mentioned above.

Another, more important point, is this reasoning gives us the derivation presented above. This derivation is translated into the following derivation in  $\mathcal{F}$ :

$$\begin{aligned}
S &\Rightarrow_{\mathcal{F}}^* \hat{u}_1 A \delta'_1 \quad (\text{or } \hat{u}_1 S \delta'_1) \\
&\Rightarrow_{\mathcal{F}} \hat{u}_1 B C \delta'_1 \\
&\Rightarrow_{\mathcal{F}}^* \hat{u}_1 \hat{u}_2 D \delta'_2 C \delta'_1 \\
&\Rightarrow_{\mathcal{F}} \hat{u}_1 \hat{u}_2 \kappa \delta'_2 C \delta'_1 \\
&\Rightarrow_{\mathcal{F}}^* \hat{u}_1 \hat{u}_2 \kappa \beta_1 C \delta'_1 \\
&\Rightarrow_{\mathcal{F}} \hat{u}_1 \hat{u}_2 \kappa \beta_1 \beta_2 (a, T \setminus \{\kappa\}) \hat{v}_1 \delta'_1, \quad \text{or} \\
&\quad \hat{u}_1 \hat{u}_2 \kappa \beta_1 \beta_2 a \hat{v}_1 \delta'_1, \\
&\Rightarrow_{\mathcal{F}}^* \hat{u}_1 \hat{u}_2 \kappa \beta_1 \beta_2 (a, T \setminus \{\kappa\}) \hat{v}_1 \hat{v}_2 = \hat{r}, \quad \text{or} \\
&\quad \hat{u}_1 \hat{u}_2 \kappa \beta_1 \beta_2 a \hat{v}_1 \hat{v}_2 = \hat{r},
\end{aligned}$$

Where  $\delta'_1$  and  $\delta'_2$  are obtained by replacing each  $X^{\text{mid}}$  by  $X$  in  $\delta_1$  and  $\delta_2$ , respectively. It is direct to see that this is a valid derivation since for every production  $X^x \Rightarrow_{\mathcal{F}'} Y^y Z^z$  there exists a valid production  $X \Rightarrow_{\mathcal{F}} YZ$ , for any  $x, y, z \in \{\text{in, out, left, mid, right}\}$ . Furthermore, for the production  $D^{\text{left}} \Rightarrow_{\mathcal{F}'} \varepsilon$  there exists  $D \Rightarrow_{\mathcal{F}} \kappa$ , and for  $C^{\text{right}} \Rightarrow_{\mathcal{F}'} (a, T)$  there exists  $C \Rightarrow_{\mathcal{F}} (a, T \setminus \{\kappa\})$  if  $T \neq \{\kappa\}$ , and  $C \Rightarrow_{\mathcal{F}} a$  if  $T = \{\kappa\}$ . Further, note that  $\text{proc}_{\kappa}(\hat{r}) = \hat{s}$ . We conclude that  $\hat{r} \in L(\mathcal{F})$  for some  $\hat{r}$  such that  $\text{proc}_{\kappa}(\hat{r}) = \hat{s}$ .

From the arguments above, we obtain that for each annotated ref-word  $\hat{r} \in L(\mathcal{F})$  there exists an equivalent annotated ref-word  $\hat{t} \in L(\mathcal{F})$ , given by  $\hat{t} = \text{proc}_{\kappa}(\hat{r})$ . Furthermore, we showed that for each annotated ref-word  $\hat{t} \in L(\mathcal{F}')$  there exists an equivalent  $\hat{r} \in L(\mathcal{F})$ . This implies that  $\mathcal{F}$  and  $\mathcal{F}'$  are equivalent.

Now, assume that  $\mathcal{F}$  is unambiguous. We will prove that  $\mathcal{F}'$  is unambiguous as well. Consider an annotated ref-word  $\hat{t} \in L(\mathcal{F}')$  and consider two sequences  $\mathcal{S}_1$  and  $\mathcal{S}_2$  which define the derivation  $S' \Rightarrow_{\mathcal{F}'}^* \hat{t}$ . We showed above how to translate these sequences into sequences  $\mathcal{S}'_1$  and  $\mathcal{S}'_2$  which define the derivation  $S \Rightarrow_{\mathcal{F}}^* \hat{r}$ , for some  $\hat{r}$  such that  $\hat{t} = \text{proc}_{\kappa}(\hat{r})$ . Since  $\mathcal{F}$  is unambiguous, these sequences are equal. Assume now that  $\mathcal{S}_1$  and  $\mathcal{S}_2$  are not equal, but since their translations into  $\mathcal{F}$  are the same, then it must be that for some production  $X \Rightarrow_{\mathcal{F}} YZ$  or  $X \Rightarrow_{\mathcal{F}} \tau$ , there must be two productions  $X^{x_1} \Rightarrow_{\mathcal{F}'} Y^{y_1} Z^{z_1}$  and  $X^{x_2} \Rightarrow_{\mathcal{F}'} Y^{y_2} Z^{z_2}$ , or  $X^{x_1} \Rightarrow_{\mathcal{F}'} \tau$  and  $X^{x_2} \Rightarrow_{\mathcal{F}'} \tau$  at the same position, for some  $(x_1, y_1, z_1) \neq (x_2, y_2, z_2)$ . We note that this is not possible since we argued that for a given derivation  $S' \Rightarrow_{\mathcal{F}'}^* \hat{r}$ , the set in which each nonterminal instance belongs, among  $V_{\text{out}}, V_{\text{in}}, V_{\text{left}}, V_{\text{mid}}, V_{\text{right}}$ , is fixed by its relation to certain instances of  $D^{\text{left}}$  and  $E^{\text{right}}$ . We conclude that  $\mathcal{F}'$  is unambiguous.

Assume  $\mathcal{F}'$  is unambiguous. We will prove that  $\mathcal{F}$  is unambiguous as well. Likewise, consider an annotated ref-word  $\hat{r} \in L(\mathcal{F})$ , and consider two sequences  $\mathcal{S}_1$  and  $\mathcal{S}_2$  which define the derivation  $S \Rightarrow_{\mathcal{F}}^* \hat{t}$ . We showed above how to convert these sequences into  $\mathcal{S}'_1$

and  $\mathcal{S}'_2$  which define the derivation  $S' \Rightarrow_{\mathcal{F}'}^* \text{proc}_\kappa(\hat{r})$ . Since  $\mathcal{F}'$  is unambiguous, then it must hold that  $\mathcal{S}'_1 = \mathcal{S}'_2$ . Note that the translation we showed consisted in replacing productions of the form  $X \Rightarrow_{\mathcal{F}} YZ$  by  $X^x \Rightarrow_{\mathcal{F}'} Y^y Z^z$ ,  $X \Rightarrow_{\mathcal{F}} \kappa$  by  $X^{\text{left}} \Rightarrow_{\mathcal{F}'} \varepsilon$ ,  $X \Rightarrow_{\mathcal{F}} a$  by  $X^{\text{right}} \Rightarrow_{\mathcal{F}'} (a, \{\kappa\})$  (or  $X \Rightarrow_{\mathcal{F}} (a, T)$  by  $X^{\text{right}} \Rightarrow_{\mathcal{F}'} (a, T \cup \{\kappa\})$ ) for a single fixed production, and  $X \Rightarrow_{\mathcal{F}} \tau$  by  $X^x \Rightarrow_{\mathcal{F}'} \tau$  in any other case, for some  $x, y, z \in \{\text{out}, \text{in}, \text{left}, \text{mid}, \text{right}\}$ . Therefore, the sequence  $S_1$  from which  $S'_1$  was obtained is uniquely defined, from which we deduce that  $S_1 = S_2$ , and we conclude that  $\mathcal{F}$  is unambiguous.

At the end of the procedure, we obtain an extraction grammar with annotations  $\mathcal{F}^\dagger = (V^\dagger, \Sigma, \mathcal{X}, P^\dagger, S^\dagger)$  such that there are no rules of the form  $X \rightarrow \kappa$  in  $P^\dagger$ , for any  $\kappa \in C_{\mathcal{X}}$ . From this, we obtain the annotated grammar  $\mathcal{G} = (V^\dagger, \Sigma, \Omega_{\mathcal{X}}, P^\dagger, S^\dagger)$  which is equivalent to  $\mathcal{H}$ . Furthermore,  $\mathcal{G}$  is unambiguous if and only if  $\mathcal{H}$  is unambiguous.

With respect to the running time of building  $\mathcal{G}$ , note that in each iteration of the algorithm, by starting on an extraction grammar with annotations  $\mathcal{F}$  with a set of rules  $P$ , the resulting  $\mathcal{F}'$  has a set of rules  $P'$  with a size of  $9|P|$ . Since this step is repeated twice for each variable  $x \in \mathcal{X}$  (once for each variable operation), the total running time is  $\mathcal{O}(9^{2|\mathcal{X}|}(3^{2|\mathcal{X}|}|\mathcal{H}|^2)) = \mathcal{O}(9^{3|\mathcal{X}|}|\mathcal{H}|^2)$ .

### D.3. Expressiveness examples

We complete Proposition 13 to give more intuition about the conciseness of extraction grammars vs annotated grammars, and the difference in expressiveness.

We first give a simple example to show that, with our notion of equivalence (Definition 26), we may indeed need an exponential number of symbols in the annotation set, implying that extraction grammars are in some cases exponentially more concise:

**Example 27.** Consider the following functional extraction grammar  $\mathcal{H}$  with  $n$  variables  $x_1, \dots, x_n$  and alphabet  $\{a\}$ :

$$\begin{aligned} \mathcal{H}: \quad A_1 &\rightarrow \vdash_{x_1} \neg_{x_1} A_2 \mid \vdash_{x_1} A_2 \neg_{x_1} \\ A_2 &\rightarrow \vdash_{x_2} \neg_{x_2} A_3 \mid \vdash_{x_2} A_3 \neg_{x_2} \\ &\vdots \\ A_n &\rightarrow \vdash_{x_n} \neg_{x_n} a \mid \vdash_{x_n} a \neg_{x_n} \end{aligned}$$

For the document  $a$ , this extraction grammar will output all possible combinations depending on whether  $\neg_{x_i}$  is at the beginning or end of  $a$  for each  $i \leq n$ . Thus, an equivalent annotated grammar will need to consider all possible subsets of  $\{\neg_x \mid x \in \mathcal{X}\}$  as possible annotations of the character  $a$ , which will require an exponential number of rules.

We then illustrate why annotated grammars are in fact *strictly* more expressive: in addition to capturing all extraction grammars (Proposition 13), annotated grammars can express functions that do not correspond to an annotation grammar.

**Example 28.** Consider a singleton annotation set  $\Omega = \{\sigma\}$ , a singleton alphabet  $\Sigma = \{a\}$ , and the annotated grammar with start symbol  $S$  and production  $S \rightarrow a(a, \sigma)S|a|\varepsilon$ . For each string of  $\Sigma^*$ , it produces one output where every other character is annotated. This

*cannot be expressed by an extraction grammar, as such a grammar fixes a finite set  $\mathcal{X}$  of variables independently from the input document, and each variable is mapped to only one span.*