



**HAL**  
open science

# Extracting Context-Free Grammars from Recurrent Neural Networks using Tree-Automata Learning and A\* Search

Benoît Barbot, Benedikt Bollig, Alain Finkel, Serge Haddad, Igor Khmelnitsky, Martin Leucker, Daniel Neider, Rajarshi Roy, Lina Ye

► **To cite this version:**

Benoît Barbot, Benedikt Bollig, Alain Finkel, Serge Haddad, Igor Khmelnitsky, et al.. Extracting Context-Free Grammars from Recurrent Neural Networks using Tree-Automata Learning and A\* Search. ICGI 2021 - 15th International Conference on Grammatical Inference, Aug 2021, New York City / Virtual, United States. pp.113-129. hal-03285433

**HAL Id: hal-03285433**

**<https://hal.science/hal-03285433v1>**

Submitted on 24 Aug 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Extracting Context-Free Grammars from Recurrent Neural Networks using Tree-Automata Learning and A\* Search\*

**Benoît Barbot**

*Université Paris-Est Créteil, France*

BENOIT.BARBOT@U-PEC.FR

**Benedikt Bollig**

**Alain Finkel**

**Serge Haddad**

**Igor Khmelnitsky**

*Université Paris-Saclay, CNRS, ENS Paris-Saclay, LMF, Gif-sur-Yvette, France*

BOLLIG@LSV.FR

ALAIN.FINKEL@LSV.FR

HADDAD@LSV.FR

KHMELNITSKY@LSV.FR

**Martin Leucker**

*Institute for Software Engineering and Programming Languages, Universität zu Lübeck, Germany*

LEUCKER@ISP.UNI-LUEBECK.DE

**Daniel Neider**

**Rajarshi Roy**

*Max Planck Institute for Software Systems, Kaiserslautern, Germany*

NEIDER@MPI-SWS.ORG

RAJARSHI@MPI-SWS.ORG

**Lina Ye**

*Université Paris-Saclay, CNRS, ENS Paris-Saclay, LMF, Inria, Gif-sur-Yvette, France  
CentraleSupélec, Université Paris-Saclay, France*

LINAYE@LRI.FR

## Abstract

This paper presents (i) an active learning algorithm for visibly pushdown grammars and (ii) shows its applicability for learning surrogate models of recurrent neural networks (RNNs) trained on context-free languages. Such surrogate models may be used for verification or explainability. Our learning algorithm makes use of the proximity of visibly pushdown languages and regular tree languages and builds on an existing learning algorithm for regular tree languages. Equivalence tests between a given RNN and a hypothesis grammar rely on a mixture of A\* search and random sampling. An evaluation of our approach on a set of RNNs from the literature shows good preliminary results.

## 1. Introduction

*Context-free languages* (CFLs), which are generated by context-free grammars (CFGs), abound in many application areas, for example when facing *formal* languages and applications such as programming languages and compilers, but especially also when processing natural language or controlled natural language. *Visibly pushdown languages* (VPLs), introduced by [Alur and Madhusudan \(2004, 2009\)](#), are a robust subclass of CFLs with interesting closure and decidability properties, as explained in further detail below—and are the class of languages studied in this paper. The idea is that the underlying pushdown automata are *input-driven* ([Mehlhorn, 1980](#)), i.e., every letter from the given alphabet is assigned a type among *push*, *pop*, and *internal* (we therefore deal with a *visibly pushdown alphabet*).

---

\* Extended version with appendix. This work was partly supported by the PHC PROCOPE 2020 project LeaRNNify (number 44707TK), funded by Campus France and DAAD, and by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) grant number 434592664.

Learning *representations* or *models* such as (neural) networks, grammars, or automata from given examples or by querying underlying systems is an important tool when working with such languages. It has been considered in the field of machine learning for sequence-processing tasks such as time-series prediction or sentiment analysis, but also in the field of grammatical inference (Vaandrager, 2017). While in the first setting, typically, a finite set of words is given as a training set from which a model such as a recurrent neural network (RNN) is derived, in the second setting, further queries to a so-called *minimally adequate teacher* (MAT) may be asked to shape the learning result. A prominent MAT learning algorithm is Angluin’s  $L^*$  for regular word languages (Angluin, 1987).

In this paper, as a first contribution, we present a novel learning algorithm for VPLs, given a minimally adequate teacher.

An important application area of such a learning algorithm, as pursued in this paper, is to derive so-called *surrogate models*, also known as approximation models, of recurrent neural networks (RNNs). RNNs play an important role in natural-language processing or time-series prediction, amongst others. While a neural network is often difficult to analyze and to understand, the surrogate model shares essential features of the underlying network but allows for simpler means for its analysis and explainability.

As a second contribution, we show that our algorithm can indeed be used for deriving a so-called visibly pushdown grammar (VPG) usable for explaining the language accepted by an underlying RNN. To this end, we perform queries to the network and infer an automaton model which is then translated into a grammar. The latter provides structural information of the underlying network which can hardly be obtained from the network directly.

**Our Approach.** As mentioned above, our learning algorithm is for the class of VPLs. Alur and Madhusudan (2004) established a close relationship between VPLs and regular tree languages. We exploit this relationship and use an existing learning algorithm for regular tree languages (Sakakibara, 1992; Drewes and Högberg, 2007) to derive a grammar-based representation of a VPL, resulting in a MAT learning algorithm.

This is similar to Sakakibara’s algorithm (Sakakibara, 1992), which infers CFGs in terms of tree automata learned using *structural* queries. In our case, we also adopt tree interpretations of the words that are queried, but with not exactly the same structure.

In fact, Kumar et al. (2006) and Isberner (2015) had already pointed out that it would be possible to use the algorithm of Sakakibara (1992) for learning regular tree languages to obtain a tree representation of a VPL, albeit mentioning two potential obstacles for this. First, the final *visibly pushdown automaton* is non-deterministic, requiring thus the exponential cost in obtaining a deterministic one. Furthermore, certain structural properties cannot be guaranteed that are expected from recursive programs. Our work focusing on practical learning of VPLs shows that these critical issues can be well handled by adapting the improved version of (Sakakibara, 1992) by Drewes and Högberg (2007) and by leveraging the computational power of RNNs.

One advantage of our algorithm as opposed to other algorithms for classes of CFLs is that it is easier to understand as it is based on the case for tree languages. Moreover, its correctness essentially follows from the correctness of the tree-learning algorithm so that, in principle, we can plug-in any other tree-automata learning algorithm having the same

interfaces. Another advantage is its extensibility to non-context-free languages, insofar as they have a representation as tree languages (Madhusudan and Parlato, 2011).

**Application to RNNs.** Our work is inspired by Yellin and Weiss (2021a), who infer CFGs from RNNs by extracting a sequence of deterministic finite automata (DFAs) using the algorithm proposed by Weiss et al. (2018), and exploiting the notion of pattern rule sets (PRSs), from which the CFG rules are derived. Experiments show that many interesting CFLs can be learned. There are nevertheless some difficulties to overcome. For example, a sequence of extracted DFAs often contains some noise, either from the RNN training or from the application of the  $L^*$  algorithm. Consequently, incorrect patterns are frequently inserted into the DFA sequence, which can deviate from the PRS. To handle this, a voting and threshold scheme has been proposed such that the languages of the given RNNs were mostly recovered in terms of CFGs, while several others were partially or incorrectly learned.

The class of VPLs is incomparable to the language class handled by Yellin and Weiss (2021a) (cf. Example 2 in Section 2). It must be fairly noted that our algorithm relies on a partitioning of the input alphabet into *push*, *pop*, and *internal* symbols, which is not required by Yellin and Weiss (2021a). However, it turns out that all the 15 benchmark languages considered by Yellin and Weiss (2021a) are VPLs.

In Yellin and Weiss (2021a), checking equivalence between the given RNN and a hypothesis grammar relies on an orthogonal learned abstraction of the RNN. In our case, the equivalence query relies on two complementary tests to look for words belonging to their symmetric difference, i.e., words in only one of the two corresponding languages.

Apart from two exceptions, the languages from Yellin and Weiss (2021a) are very well learned with our approach, even some of the languages that are only partially generalized by applying the other approach. This demonstrates that our algorithm may be a worthwhile alternative when dealing with structured data (annotated linguistic data, programs, XML documents, etc.), i.e., in presence of a visibly pushdown alphabet.

**Further Related Work.** There are a wide range of learning algorithms for regular languages. Let us mention some of them. Angluin (1982) used reversibility to identify a class of regular languages from positive data alone using DFAs. Then, Angluin (1987) showed that the class of all regular languages could be learned using the  $L^*$  algorithm in the MAT model, where the teacher can answer both membership queries and equivalence queries. Rivest and Schapire (1993) proposed binary search to determine a single suffix of a counterexample that causes refinement, while Kearns and Vazirani (1994) suggested constructing a discrimination tree instead of the observation table. Then, Balcázar et al. (1997) provided a unified view on these learning algorithms, resulting in the observation pack framework such that one can construct an automaton from them. Isberner et al. (2014) presented the TTT algorithm, which is extremely efficient, especially in presence of long counterexamples, thanks to a refined counterexample-analysis and redundancy-free organization of observations.

Some researchers adapted learning algorithms for regular languages to learn CFLs. For example, Clark and Eyraud (2007) presented an exact analogue of that proposed by Angluin (1982) for a limited class of CFLs by combining the correspondence of non-terminals to the syntactic congruence class with weak substitutability. Then, Clark (2010) expanded this approach by adopting an extended MAT to answer equivalence queries where the hypothesis may not be in the learnable class. Yoshinaka and Clark (2010) extended the syntactic

congruence to tuples of strings to learn efficiently some sorts of multiple CFGs. Even though the above algorithms for learning CFLs have shown some promising results, they are limited to some constrained class. The learnability of the whole class of CFLs is widely believed to be intractable (de la Higuera, 2005).

Since decades, some approaches have been developed to extract simpler and explainable surrogate models from a neural network to facilitate comprehension and verification (Thrun, 1994; Omlin and Giles, 1996). New algorithms for extracting (weighted or unweighted) DFAs from RNNs have been proposed recently, with promising applications in verification (Weiss et al., 2018; Mayr and Yovine, 2018; Ayache et al., 2018; Rabusseau et al., 2019; Weiss et al., 2019; Mayr et al., 2020). They may also turn out to be useful for generalizations to other, more complex classes of languages. Up to now, however, there has been little research on extracting CFGs from RNNs. With the exception of (Yellin and Weiss, 2021a), existing approaches rely on an RNN augmented with external stack memory, either continuous or discrete (Das et al., 1992; Sun et al., 1997). In such a hybrid system, besides the classical input symbols, the input includes also what is read from the top of the stack.

**Outline.** Section 2 recalls basic notions such as CFLs, CFGs, and VPLs. Trees and tree automata are presented in Section 3. In Section 4, we recall the tree-automata learning algorithm that we exploit, in Section 5, to learn grammars for VPLs. In Section 6, we apply our algorithm to inferring grammars from RNNs. We conclude in Section 7.

## 2. Context-Free and Visibly Pushdown Grammars

In this section, we recall standard concepts such as context-free languages and grammars. We also present their subclass of visibly pushdown languages.

### 2.1. Context-Free Languages and Grammars

Let  $\Sigma$  be an alphabet, i.e., a nonempty finite set. A word over  $\Sigma$  is a finite sequence  $w = a_1 \dots a_n$  of letters  $a_i \in \Sigma$ . The length  $|w|$  of  $w$  is  $n$ . The unique word of length 0 is the empty word, denoted by  $\varepsilon$ . By  $\Sigma^*$ , we denote the set of all finite words over  $\Sigma$ .

Any set  $L \subseteq \Sigma^*$  is called a language. For two languages  $L_1, L_2 \subseteq \Sigma^*$ , we let  $L_1 \oplus L_2$  denote their symmetric difference, i.e., the language  $(L_1 \setminus L_2) \cup (L_2 \setminus L_1)$ .

For a finite set  $U$ , we denote by  $\mathcal{P}(U)$  its powerset and by  $|U|$  its size.

**Definition 1** A context-free grammar (CFG) over  $\Sigma$  is a tuple  $G = (N, S, \rightarrow)$  where  $N$  is a finite set of nonterminal symbols with  $N \cap \Sigma = \emptyset$ ,  $S \in N$  is the start symbol, and  $\rightarrow \subseteq N \times (\Sigma \cup N)^*$  is the finite set of rules. A rule  $(A, w) \in \rightarrow$  is usually written as  $A \rightarrow w$ .

The language  $L(G) \subseteq \Sigma^*$  of  $G$  is defined using a global rewrite relation  $\Rightarrow \subseteq (\Sigma \cup N)^* \times (\Sigma \cup N)^*$  defined by  $uAv \Rightarrow uwv$  for all rules  $A \rightarrow w$  and  $u, v \in (\Sigma \cup N)^*$ . With this, we let  $L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$  where  $\Rightarrow^*$  denotes the reflexive transitive closure of the binary relation  $\Rightarrow$ .

We say that  $L \subseteq \Sigma^*$  is a *context-free language (CFL)* if there is a CFG  $G$  over  $\Sigma$  such that  $L(G) = L$ . It should be noted that the class of CFLs enjoys equivalent characterizations, e.g., via pushdown automata.

**Example 1** For  $n \geq 1$ , consider the grammar  $G_n$  given by  $S \rightarrow \varepsilon \mid p_i S q_i S$  (for all  $i \in \{1, \dots, n\}$ ) over the alphabet  $\Sigma_n = \{p_1, \dots, p_n, q_1, \dots, q_n\}$ . Hereby, as usual,  $\mid$  separates several possible right-hand sides of rules. Then,  $L(G_n)$  is the Dyck language of order  $n$  of well-bracketed words, where  $p_i$  is an opening and  $q_i$  its corresponding closing bracket.

## 2.2. Visibly Pushdown Languages and Their Grammars

The class of *visibly pushdown languages* has been introduced by [Alur and Madhusudan \(2004, 2009\)](#). It was originally defined in terms of visibly pushdown automata, but can be equivalently characterized by a subclass of CFGs. VPLs constitute a robust class that, unlike the class of CFLs, is closed under complement.

The idea is to assign to every letter from an alphabet a precise role. Speaking in terms of automata, every letter is either a push, a pop, or an internal symbol. This clearly is a restriction: A pushdown automaton for the CFL  $\{a^n b a^n \mid n \in \mathbb{N}\}$  has to perform a certain number of push operations while reading the first  $n$  occurrences of  $a$  before the  $b$ , and pop operations when reading the remaining letters  $a$ . On the other hand,  $\{a^n b^n \mid n \in \mathbb{N}\}$  can be recognized by a pushdown automaton where a stack symbol is pushed when reading an  $a$  and a stack symbol is popped when reading a  $b$ . Accordingly, a *visibly pushdown alphabet* is an alphabet  $\Sigma = \Sigma_{\text{push}} \uplus \Sigma_{\text{pop}} \uplus \Sigma_{\text{int}}$  that is partitioned into *push*, *pop*, and *internal letters*.

For the rest of the paper,  $\Sigma$  will always denote a given visibly pushdown alphabet.

**Definition 2** A visibly pushdown grammar (VPG) over  $\Sigma$  is a CFG  $(N, S, \rightarrow)$  such that every rule has one of the following forms (where  $A, B, C \in N$ ):  $A \rightarrow \varepsilon$ , or  $A \rightarrow cB$  with  $c \in \Sigma_{\text{int}}$ , or  $A \rightarrow aBbC$  with  $a \in \Sigma_{\text{push}}$  and  $b \in \Sigma_{\text{pop}}$ .

A language  $L \subseteq \Sigma^*$  is called a *visibly pushdown language (VPL)* over  $\Sigma$  if there is a VPG  $G$  over  $\Sigma$  such that  $L(G) = L$ .

**Example 2** For  $n \geq 1$ , consider again the grammar  $G_n$  from [Example 1](#). In fact,  $G_n$  is a VPG for  $\Sigma_{\text{push}} = \{p_1, \dots, p_n\}$ ,  $\Sigma_{\text{pop}} = \{q_1, \dots, q_n\}$ , and  $\Sigma_{\text{int}} = \emptyset$  so that  $L(G_n)$  is a VPL. Another example of a VPL is  $\{a^n x b^n \mid n \in \mathbb{N}\}$  where  $\Sigma_{\text{push}} = \{a\}$ ,  $\Sigma_{\text{pop}} = \{b\}$ , and  $\Sigma_{\text{int}} = \{x\}$ . This language is not captured by the PRS-formalism presented by [Yellin and Weiss \(2021a\)](#) (cf. ([Yellin and Weiss, 2021b, Section C.3](#))).

We observe that, due to the form of permitted rules, a VPL  $L$  can only contain words  $w \in \Sigma^*$  that are well-formed in a certain sense. The set  $\mathcal{W}_\Sigma$  of *well-formed* words over  $\Sigma$  is defined as the language  $L(G_\Sigma)$  of the “most permissive” VPG:  $G_\Sigma = (\{S\}, S, \rightarrow)$  with set of rules  $\{S \rightarrow \varepsilon\} \cup \{S \rightarrow cS \mid c \in \Sigma_{\text{int}}\} \cup \{S \rightarrow aSbS \mid a \in \Sigma_{\text{push}} \text{ and } b \in \Sigma_{\text{pop}}\}$ .

The general framework by [Alur and Madhusudan \(2004, 2009\)](#) can also cope with words that have unmatched push or pop positions. For simplicity, we restrict here to well-formed words. However, the algorithms can be extended straightforwardly to the general case.

With  $w = a_1 \dots a_n \in \mathcal{W}_\Sigma$ , we can associate a unique binary relation  $\sim \subseteq \{1, \dots, n\} \times \{1, \dots, n\}$  connecting a push with a unique pop position: For  $i, j \in \{1, \dots, n\}$ , we let  $i \sim j$  if  $i < j$ ,  $a_i \in \Sigma_{\text{push}}$ ,  $a_j \in \Sigma_{\text{pop}}$ , and  $a_{i+1} \dots a_{j-1}$  is well-formed. We call the pair  $(w, \sim)$  (with  $w \in \mathcal{W}_\Sigma$ ) a *nested word*. A nested word over  $\Sigma$  with  $\Sigma_{\text{push}} = \{a, b\}$ ,  $\Sigma_{\text{pop}} = \{\underline{a}, \underline{b}\}$ , and  $\Sigma_{\text{int}} = \{c\}$  is depicted in [Figure 1](#). We do not exploit nested words in this paper, but it is helpful to think of well-formed words as nested words when we encode them as trees.

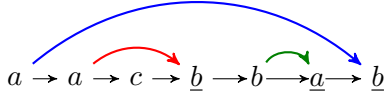


Figure 1: A nested word

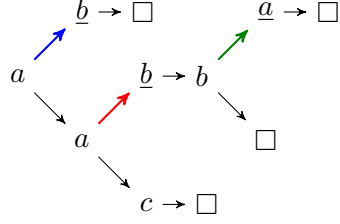


Figure 2: Its encoding as a tree

### 3. Trees and Regular Tree Languages

The reason why VPLs are so robust is that they are close to tree languages. In fact, nested words as introduced in the previous section can be represented as trees. Trees are defined over a *ranked alphabet*, i.e., an alphabet  $\Gamma = \Gamma_0 \uplus \Gamma_1 \uplus \dots \uplus \Gamma_{k_{max}}$  that is partitioned into letters of arity  $k \in \{0, \dots, k_{max}\}$  where  $k_{max} \in \mathbb{N}$  is the maximal arity. Unless otherwise stated, we let  $\Gamma$  be a fixed ranked alphabet.

A tree  $t$  over  $\Gamma$  is a term that is generated according to the grammar  $t ::= a(t_1, \dots, t_k)$ , where  $k$  ranges over  $\{0, \dots, k_{max}\}$  and  $a$  over  $\Gamma_k$ . Figure 2 depicts a syntax-tree-based representation of the tree  $a(a(c(\square()), \underline{b}(b(\square()), \underline{a}(\square()))), \underline{b}(\square()))$  over the ranked alphabet given by  $\Gamma_0 = \{\square\}$ ,  $\Gamma_1 = \{\underline{a}, \underline{b}, c\}$ , and  $\Gamma_2 = \{a, b\}$ .

The size  $|t|$  of  $t$  is the number of its nodes, i.e., the number of occurrences of symbols from  $\Gamma$ . Let  $\text{Trees}(\Gamma)$  denote the set of all trees over  $\Gamma$ .

The algorithm by [Drewes and Högberg \(2007\)](#), on which our approach is based, infers regular tree languages in terms of tree automata (later, when a tree automaton represents a VPL, we will be able to extract a corresponding VPG representation).

**Definition 3** A nondeterministic finite (bottom-up) tree automaton (NTA) over  $\Gamma$  is a tuple  $\mathcal{B} = (Q, \delta, F)$  where  $Q$  is the nonempty finite set of states,  $F \subseteq Q$  is the set of final states, and  $\delta : \bigcup_{k \in \{0, \dots, k_{max}\}} (\Gamma_k \times Q^k) \rightarrow \mathcal{P}(Q)$  is the transition function. We will write  $\delta(a(q_1, \dots, q_k))$  instead of  $\delta(a, q_1, \dots, q_k)$ .

We call  $\mathcal{B}$  deterministic (a DTA) if  $|\delta(a(q_1, \dots, q_k))| = 1$  for all arguments  $a, q_1, \dots, q_k$ . Then,  $\delta$  can also be seen as a total (i.e., complete) function  $\delta : \bigcup_{k \in \{0, \dots, k_{max}\}} (\Gamma_k \times Q^k) \rightarrow Q$ . We let  $\text{DTA}(\Gamma)$  denote the set of DTAs over  $\Gamma$ .

From  $\delta : \bigcup_{k \in \{0, \dots, k_{max}\}} (\Gamma_k \times Q^k) \rightarrow \mathcal{P}(Q)$ , we obtain a function  $\hat{\delta} : \text{Trees}(\Gamma) \rightarrow \mathcal{P}(Q)$  letting, for an arity  $k \in \{0, \dots, k_{max}\}$ ,  $a \in \Gamma_k$ , and  $t_1, \dots, t_k \in \text{Trees}(\Gamma)$ ,  $\hat{\delta}(a(t_1, \dots, t_k)) = \bigcup_{q_1 \in \hat{\delta}(t_1), \dots, q_k \in \hat{\delta}(t_k)} \delta(a(q_1, \dots, q_k))$ . We can now define the tree language recognized by  $\mathcal{B}$  as  $T(\mathcal{B}) = \{t \in \text{Trees}(\Gamma) \mid \hat{\delta}(t) \cap F \neq \emptyset\}$ .

We call a tree language  $T \subseteq \text{Trees}(\Gamma)$  *regular* if it is recognized by some NTA over  $\Gamma$ .

We now state some important and well-known facts about tree automata. For more details, we refer the reader to ([Comon et al., 2007](#)).

**Fact 1 (minimal DTA)** For every NTA  $\mathcal{B} = (Q, \delta, F)$ , there is a unique (up to isomorphism) minimal DTA  $\mathcal{B}' = (Q', \delta', F')$  such that  $T(\mathcal{B}') = T(\mathcal{B})$ . We can assume  $|Q'| \leq 2^{|Q|}$ .

The *index* of a regular tree language  $T$  is the number of states of the minimal DTA recognizing  $T$ .

While DTAs capture the class of regular tree languages, deterministic *top-down* finite tree automata (Comon et al., 2007), which we do not define here, are strictly less expressive.

**Fact 2 (membership and emptiness)** (i) Given an NTA  $\mathcal{B}$  and a tree  $t \in \text{Trees}(\Gamma)$ , one can decide in polynomial time whether  $t \in T(\mathcal{B})$ . For DTAs, there is a linear-time algorithm. (ii) For a given NTA  $\mathcal{B}$ , one can decide in polynomial time whether  $T(\mathcal{B}) \neq \emptyset$ .

## 4. Learning Deterministic Tree Automata

In her seminal work, Angluin (1987) provided the algorithm  $L^*$ , which can infer a deterministic finite automaton for a given regular word language that can only be accessed via two types of queries: *membership queries* (MQs) and *equivalence queries* (EQs). The algorithm has later been extended, first by Sakakibara (1992) to CFGs and then by Drewes and Högberg (2007) to tree automata. The latter algorithm, called  $TL^*$  in this paper, can infer a DTA over a fixed ranked alphabet  $\Gamma$  for a given (unknown) regular tree language  $T$ . Hereby,  $T$  can be accessed through membership queries and equivalence queries, which are implemented by “oracle” mappings  $\text{MQ}_{\text{tree}} : \text{Trees}(\Gamma) \rightarrow \{\text{yes}, \text{no}\}$  and  $\text{EQ}_{\text{tree}} : \text{DTA}(\Gamma) \rightarrow \{\text{yes}\} \cup \text{Trees}(\Gamma)$ :

- We say that  $\text{MQ}_{\text{tree}}$  is *sound* for  $T$  if, for all  $t \in \text{Trees}(\Gamma)$ ,  $\text{MQ}_{\text{tree}}(t) = \text{yes}$  iff  $t \in T$ .
- We say that  $\text{EQ}_{\text{tree}}$  is *counterexample-sound* for  $T$  if, for all  $\mathcal{B} \in \text{DTA}(\Gamma)$  and  $t \in \text{Trees}(\Gamma)$  such that  $\text{EQ}_{\text{tree}}(\mathcal{B}) = t$ , we have  $t \in T \oplus T(\mathcal{B})$  (i.e.,  $t$  is a *counterexample*).
- We call  $\text{EQ}_{\text{tree}}$  *equivalence-sound* for  $T$  if, for all  $\mathcal{B} \in \text{DTA}(\Gamma)$  such that  $\text{EQ}_{\text{tree}}(\mathcal{B}) = \text{yes}$ , we have  $T = T(\mathcal{B})$ .

These queries act as “oracles” and their answers are delivered instantaneously. Ideally, one assumes that  $\text{EQ}_{\text{tree}}$ , which checks the current *hypothesis* computed by the learning algorithm, is both counterexample- and equivalence-sound. In practice, this is not always the case. In fact, in our experiments, we will make weaker assumptions on  $\text{EQ}_{\text{tree}}$ .

The algorithm  $TL^*$  by Drewes and Högberg (2007) takes as input a ranked alphabet  $\Gamma$  and two functions  $\text{MQ}_{\text{tree}} : \text{Trees}(\Gamma) \rightarrow \{\text{yes}, \text{no}\}$  and  $\text{EQ}_{\text{tree}} : \text{DTA}(\Gamma) \rightarrow \{\text{yes}\} \cup \text{Trees}(\Gamma)$ . If  $TL^*(\Gamma, \text{MQ}_{\text{tree}}, \text{EQ}_{\text{tree}})$  terminates, it outputs a DTA over  $\Gamma$ .

**Fact 3 (Drewes and Högberg (2007))** Let  $T \subseteq \text{Trees}(\Gamma)$  be a regular tree language, say with index  $n$  (the minimal DTA for  $T$  has  $n$  states). Suppose  $\text{MQ}_{\text{tree}}$  is sound for  $T$  and that  $\text{EQ}_{\text{tree}}$  is both counterexample- and equivalence-sound for  $T$ . Then,  $TL^*(\Gamma, \text{MQ}_{\text{tree}}, \text{EQ}_{\text{tree}})$  terminates and outputs the unique minimal DTA  $\mathcal{B}$  with  $n$  states such that  $T(\mathcal{B}) = T$ . The overall running time is polynomial in  $|\Gamma|$ ,  $n^{k_{\max}}$ , and the maximal size of a counterexample returned by  $\text{EQ}_{\text{tree}}$ .

Note that, in the next sections,  $k_{\max}$  will be fixed. However, unlike in the case of word automata, the size of a smallest counterexample tree returned for an equivalence query may be exponential in the size of the target automaton.



## 5. Learning Visibly Pushdown Grammars

In this section, we exploit tree-automata learning for the inference of VPLs in terms of VPGs. The derived algorithm will then be exploited to extract grammars from RNNs.

### 5.1. Encoding Nested Words as Trees

The main link between words and trees is provided by an encoding of well-formed words as trees over a suitable ranked alphabet (Alur and Madhusudan, 2004, 2009).

Let  $\Sigma = \Sigma_{\text{push}} \uplus \Sigma_{\text{pop}} \uplus \Sigma_{\text{int}}$  be a visibly pushdown alphabet. To encode words from  $\mathcal{W}_\Sigma$  as trees, we introduce a suitable ranked alphabet  $\Gamma = \Gamma_0 \uplus \Gamma_1 \uplus \Gamma_2$  letting  $\Gamma_0 = \{\square\}$ ,  $\Gamma_1 = \Sigma_{\text{pop}} \cup \Sigma_{\text{int}}$ , and  $\Gamma_2 = \Sigma_{\text{push}}$ . That is, the maximal arity is 2. For the rest of this section, we fix  $\Sigma$  and the associated ranked alphabet  $\Gamma$ .

To a well-formed word  $w \in \mathcal{W}_\Sigma$ , we inductively assign a (parse) tree  $\langle\langle w \rangle\rangle \in \text{Trees}(\Gamma)$  as follows: (i)  $\langle\langle \varepsilon \rangle\rangle = \square()$ . (ii) If  $w = aw_1bw_2$  such that  $a \in \Sigma_{\text{push}}$ ,  $b \in \Sigma_{\text{pop}}$ , and  $w_1$  and  $w_2$  are well-formed, then  $\langle\langle w \rangle\rangle = a(\langle\langle w_1 \rangle\rangle, b(\langle\langle w_2 \rangle\rangle))$ . (iii) If  $c \in \Sigma_{\text{int}}$  and  $w$  is well-formed, then  $\langle\langle cw \rangle\rangle = c(\langle\langle w \rangle\rangle)$ . The encoding of the word from Figure 1 is illustrated in Figure 2.

Given  $L \subseteq \mathcal{W}_\Sigma$ , we let  $\langle\langle L \rangle\rangle = \{\langle\langle w \rangle\rangle \mid w \in L\} \subseteq \text{Trees}(\Gamma)$ . Moreover, we let  $\mathcal{T}_\Gamma = \langle\langle \mathcal{W}_\Sigma \rangle\rangle$  be the set of trees that encode a well-formed word. Note that  $\langle\langle \cdot \rangle\rangle : \mathcal{W}_\Sigma \rightarrow \mathcal{T}_\Gamma$  is injective and, therefore, a bijection. Indeed, its inverse mapping, which we denote by  $\llbracket \cdot \rrbracket$ , is given by  $\llbracket \square() \rrbracket = \varepsilon$ ,  $\llbracket a(t_1, b(t_2)) \rrbracket = a\llbracket t_1 \rrbracket b\llbracket t_2 \rrbracket$  and  $\llbracket c(t) \rrbracket = c\llbracket t \rrbracket$ . For  $T \subseteq \mathcal{T}_\Gamma$ , let  $\llbracket T \rrbracket = \{\llbracket t \rrbracket \mid t \in T\}$ .

Let us state some known facts on the relation between VPGs and NTAs/DTAs due to Alur and Madhusudan (2004, 2009).

**Fact 4** *For every VPL  $L$  over  $\Sigma$ , there is an NTA (or DTA)  $\mathcal{B}$  over  $\Gamma$  such that  $T(\mathcal{B}) = \langle\langle L \rangle\rangle$ . In particular, there is a DTA  $\mathcal{B}_{\text{parse}}$  over  $\Gamma$  with a constant number of states such that  $T(\mathcal{B}_{\text{parse}}) = \mathcal{T}_\Gamma$ .*

As we will extract grammars from tree automata, the following is particularly important:

**Fact 5** *Let  $\mathcal{B}$  be an NTA over  $\Gamma$  such that  $T(\mathcal{B}) \subseteq \mathcal{T}_\Gamma$ . One can compute, in polynomial time, a VPG  $\text{nta2vpg}(\mathcal{B})$  over  $\Sigma$  such that  $L(\text{nta2vpg}(\mathcal{B})) = \llbracket T(\mathcal{B}) \rrbracket$ .*

We give the translation of an NTA into a VPG, as the latter will yield the representation of a VPL learned in terms of the NTA. Suppose  $\mathcal{B} = (Q, \delta, F)$  is an NTA over  $\Gamma$  such that  $T(\mathcal{B}) \subseteq \mathcal{T}_\Gamma$ . We define  $\text{nta2vpg}(\mathcal{B}) = (N, \mathcal{I}, \rightarrow)$  as follows. In fact, instead of just one start symbol, we assume a set of start symbols  $\mathcal{I} \subseteq N$ . This is no more expressive than having one single start symbol, as we can always introduce a fresh start symbol, leading to all the right hand sides of rules associated with symbols from  $\mathcal{I}$ . Intuitively, the grammar derives a run of the NTA top-down, where states are successively replaced with input letters. So we let  $N = Q$  and  $\mathcal{I} = F$ . Moreover, the set of rules contains (i)  $\hat{q} \rightarrow \varepsilon$  for all  $\hat{q} \in \delta(\square())$ ; (ii)  $\hat{q} \rightarrow cq$  for all  $c \in \Sigma_{\text{int}}$ ,  $q \in Q$ , and  $\hat{q} \in \delta(c(q))$ ; (iii)  $\hat{q} \rightarrow apbq$  for all  $a \in \Sigma_{\text{push}}$ ,  $b \in \Sigma_{\text{pop}}$ , and  $p, q, q', \hat{q} \in Q$  such that  $q' \in \delta(b(q))$  and  $\hat{q} \in \delta(a(p, q'))$ .

For completeness, let us mention some connections with visibly pushdown automata (VPAs), which are effectively equivalent to VPGs w.r.t. expressive power so that we could also learn VPAs instead of VPGs (cf. Appendix A or (Alur and Madhusudan, 2004, 2009)

<b>Algorithm 1</b> Implementing $\text{MQ}_{\text{tree}}$ in terms of $\text{MQ}_{\text{vpl}}$	<b>Algorithm 2</b> Implementing $\text{EQ}_{\text{tree}}$ in terms of $\text{EQ}_{\text{vpl}}$
<pre> 1  <math>\text{MQ}_{\text{tree}}(t)</math>: 2      <b>if</b> <math>t \in T(\mathcal{B}_{\text{parse}})</math> 3          <b>then return</b> <math>\text{MQ}_{\text{vpl}}(\llbracket t \rrbracket)</math> 4          <b>else return no</b> </pre>	<pre> 1  <math>\text{EQ}_{\text{tree}}(\mathcal{B})</math>: 2      <b>if</b> <math>T(\mathcal{B}) \subseteq T(\mathcal{B}_{\text{parse}})</math> 3          <b>then return</b> <math>\text{EQ}_{\text{vpl}}(\mathcal{B})</math> 4      <b>else</b> 5          pick <math>t \in T(\mathcal{B}) \setminus T(\mathcal{B}_{\text{parse}})</math> 6          <b>return</b> <math>t</math> </pre>

for the definition of VPAs). For an NTA  $\mathcal{B}$  over  $\Gamma$  such that  $T(\mathcal{B}) \subseteq \mathcal{T}_\Gamma$ , one can compute, in polynomial time, a VPA  $\mathcal{A}$  over  $\Sigma$  such that  $L(\mathcal{A}) = \llbracket T(\mathcal{B}) \rrbracket$ . Conversely, for a VPA  $\mathcal{A}$  over  $\Sigma$ , one can compute, in polynomial time, an NTA  $\mathcal{B}$  over  $\Gamma$  such that  $T(\mathcal{B}) = \langle\langle L(\mathcal{A}) \rangle\rangle$ . Hence, there is also a DTA for  $\langle\langle L(\mathcal{A}) \rangle\rangle$  of exponential size. In general, this exponential blow-up cannot be avoided even when we start from a deterministic VPA.

## 5.2. Learning VPLs in Terms of VPGs

Recall that  $\Sigma$  is a fixed visibly pushdown alphabet and  $\Gamma$  is the derived ranked alphabet.

We now present an algorithm, called  $\text{VPL}^*$  in the following, that learns a VPL  $L \subseteq \mathcal{W}_\Sigma$  in terms of a DTA for the tree language  $\langle\langle L \rangle\rangle \subseteq \mathcal{T}_\Gamma$  that can then be translated into a VPG according to Fact 5. In particular, the equivalence query will take a DTA as argument, rather than a VPA. Essentially, we rely on the algorithm  $\text{TL}^*$ . However, equivalence and membership queries are now answered w.r.t. the VPL  $L$ . More precisely, we deal with a mapping  $\text{MQ}_{\text{vpl}} : \mathcal{W}_\Sigma \rightarrow \{\text{yes}, \text{no}\}$  and a partial mapping  $\text{EQ}_{\text{vpl}} : \text{DTA}(\Gamma) \rightarrow \{\text{yes}\} \cup \mathcal{T}_\Gamma$  whose domain is the set of DTAs  $\mathcal{B} \in \text{DTA}(\Gamma)$  such that  $T(\mathcal{B}) \subseteq \mathcal{T}_\Gamma$ :

- We call  $\text{MQ}_{\text{vpl}}$  *sound* for  $L$  if, for all  $w \in \mathcal{W}_\Sigma$ , we have  $\text{MQ}_{\text{vpl}}(w) = \text{yes}$  iff  $w \in L$ .
- We say that  $\text{EQ}_{\text{vpl}}$  is *counterexample-sound* for  $L$  if, for all  $\mathcal{B} \in \text{DTA}(\Gamma)$  such that  $T(\mathcal{B}) \subseteq \mathcal{T}_\Gamma$  and all  $t \in \mathcal{T}_\Gamma$ ,  $\text{EQ}_{\text{vpl}}(\mathcal{B}) = t$  implies  $\llbracket t \rrbracket \in L \oplus \llbracket T(\mathcal{B}) \rrbracket$ .
- We say that  $\text{EQ}_{\text{vpl}}$  is *equivalence-sound* for  $L$  if, for all  $\mathcal{B}$  over  $\Gamma$  such that  $T(\mathcal{B}) \subseteq \mathcal{T}_\Gamma$ ,  $\text{EQ}_{\text{vpl}}(\mathcal{B}) = \text{yes}$  implies  $L = \llbracket T(\mathcal{B}) \rrbracket$ .

Our algorithm  $\text{VPL}^*$  for learning VPLs uses  $\text{TL}^*$  as a black-box. Therefore, we define a mapping  $\text{MQ}_{\text{tree}} : \text{Trees}(\Gamma) \rightarrow \{\text{yes}, \text{no}\}$  and a mapping  $\text{EQ}_{\text{tree}} : \text{DTA}(\Gamma) \rightarrow \{\text{yes}\} \cup \text{Trees}(\Gamma)$  that implement the membership and equivalence queries for tree languages, respectively (cf. Algorithms 1 and 2). The algorithm  $\text{VPL}^*$  (Algorithm 3) then simply calls  $\text{TL}^*$  with parameters  $(\Gamma, \text{MQ}_{\text{tree}}, \text{EQ}_{\text{tree}})$  and translates the resulting DTA into a VPG.

**Algorithm 1.** Membership query  $\text{MQ}_{\text{tree}}(t)$  with  $t \in T(\mathcal{B}_{\text{parse}}) = \mathcal{T}_\Gamma$  is answered in terms of  $\text{MQ}_{\text{vpl}}(\llbracket t \rrbracket)$  (line 3). If, on the other hand,  $t \notin T(\mathcal{B}_{\text{parse}})$ , the query returns **no** (line 4).

**Algorithm 2.** Recall that we are looking for a tree automaton for the language  $T = \langle\langle L \rangle\rangle$ , which is included in  $T(\mathcal{B}_{\text{parse}})$ . We will, therefore, first check whether this inclusion also applies to the current hypothesis DTA  $\mathcal{B}$ , i.e., whether  $T(\mathcal{B}) \subseteq T(\mathcal{B}_{\text{parse}})$ . If not,

---

**Algorithm 3** VPL\*

```
1  $\mathcal{B} \leftarrow \text{TL}^*(\Gamma, \text{MQ}_{\text{tree}}, \text{EQ}_{\text{tree}})$  /* MQtree and EQtree from Algorithms 1 and 2 */
2 return nta2vpg( $\mathcal{B}$ )
```

---

then we can find a tree  $t \in T(\mathcal{B}) \setminus T(\mathcal{B}_{\text{parse}})$ , which serves as a counterexample to the equivalence query (line 5). So suppose that  $T(\mathcal{B}) \subseteq T(\mathcal{B}_{\text{parse}})$ . Let us assume that  $\text{EQ}_{\text{vpl}}$  is both counterexample- and equivalence-sound. If it returns a tree  $t = \text{EQ}_{\text{vpl}}(\mathcal{B})$ , then  $\llbracket t \rrbracket \in L \oplus \llbracket T(\mathcal{B}) \rrbracket$  so that  $t$  can indeed be used to refine the hypothesis  $\mathcal{B}$ . If, on the other hand,  $\text{EQ}_{\text{vpl}}(\mathcal{B}) = \text{yes}$ , then  $L = \llbracket T(\mathcal{B}) \rrbracket$ , i.e.,  $T(\mathcal{B}) = \langle\langle L \rangle\rangle$ , so that we can return  $\mathcal{B}$  as a suitable tree-language representation. Algorithm 3 then returns the VPG  $G = \text{nta2vpg}(\mathcal{B})$ . According to Fact 5, we have  $L(G) = \llbracket T(\mathcal{B}) \rrbracket = L$ .

**Theorem 4** *Let  $L$  be a VPL and  $\hat{\mathcal{B}}$  be the minimal DTA such that  $T(\hat{\mathcal{B}}) = \langle\langle L \rangle\rangle$ . Assume  $\text{MQ}_{\text{vpl}}$  is sound for  $L$  and that  $\text{EQ}_{\text{vpl}}$  is both counterexample- and equivalence-sound for  $L$ . Then, VPL\* (Algorithm 3) terminates and eventually returns a VPG  $G$  of size polynomial in the size of  $\hat{\mathcal{B}}$  such that  $L(G) = L$ . The overall running time is polynomial in  $|\Sigma|$ , the index of  $\hat{\mathcal{B}}$ , and the maximal size of a counterexample returned in lines 3 and 6 of Algorithm 2.*

See Appendix B for the proof of Theorem 4. Note that, like in  $\text{TL}^*$ , a smallest counterexample tree returned for an equivalence query may be of exponential size. Also note that the size of the returned VPG  $G$  is at most exponential in the size of a minimal (non-deterministic) VPA recognizing  $L$ .

## 6. Experiments

We applied Algorithm 3 to recurrent neural networks (RNNs) in order to extract VPGs. We implemented it in Python 3.6, using the Numpy library.<sup>1</sup> Since we are comparing our extractions to those done in (Yellin and Weiss, 2021a, [https://github.com/tech-srl/RNN\\_to\\_PRS\\_CFG](https://github.com/tech-srl/RNN_to_PRS_CFG)), we use the 15 RNNs they trained, and a modified version of the interface they wrote to communicate with these RNNs. All benchmarks were performed on a computer equipped by Intel i5-8250U CPU with 4 cores, 16GB of memory, and Ubuntu Linux 18.03.

The experiments done for this paper are *preliminary*, yet they show promise and motivate a more in-depth investigation. We describe some ideas for some more in-depth experimentation in the end of the section.

**Recurrent Neural Networks.** RNNs can be seen as language acceptors. For the purpose of this paper, it is enough to think of an RNN  $R$  as an infinite automaton with infinite state space  $Q$  (e.g.,  $Q = \mathbb{R}^{\text{dim}}$  for some dimension  $\text{dim} \geq 1$ ), initial state  $q_0 \in Q$ , transition function  $\delta : Q \times \Sigma \rightarrow Q$ , and a mapping  $\text{score} : Q \rightarrow \mathbb{R}$  (e.g., indicating a probability of acceptance or of an “end of sequence” token). As usual,  $\delta$  is extended to  $\hat{\delta} : Q \times \Sigma^* \rightarrow Q$  over words by applying  $\delta$  letter by letter. Then,  $R$  computes a (score) function  $R : \Sigma^* \rightarrow \mathbb{R}$  by  $R(w) = \text{score}(\hat{\delta}(q_0, w))$ . Moreover, given a threshold  $\tau \in \mathbb{R}$ , one can associate with  $R$  a language letting  $L(R) = \{w \in \Sigma^* \mid R(w) \geq \tau\}$  or using different threshold criteria.

---

1. The code is available here: [https://github.com/LeaRNNify/VPA\\_learning](https://github.com/LeaRNNify/VPA_learning)

Table 1: Definition of some CFLs ( $X$  and  $Y$  are finite sets of words)

$\mathcal{L}(X, Y):$	<i>RE-Dyck</i> ( $X, Y$ ):	<i>Dyck</i> <sub>1</sub> $\triangleleft X$ :
$S \rightarrow \varepsilon \mid xSy$	$S \rightarrow xAy$	$S \rightarrow p_1Aq_1$
(for all $x \in X$ and $y \in Y$ )	$A \rightarrow xAy \mid AA \mid \varepsilon$	$A \rightarrow p_1Aq_1 \mid AA \mid \varepsilon \mid x$
	(for all $x \in X$ and $y \in Y$ )	(for all $x \in X$ )
<i>Dyck</i> <sub><math>n</math></sub> :	<i>Alternating</i> :	<i>Dyck</i> <sub>2</sub> $\triangleleft X$ :
$S \rightarrow p_iAq_i$	$S \rightarrow A \mid B$	$S \rightarrow p_iAq_i$
$A \rightarrow p_iAq_i \mid AA \mid \varepsilon$	$A \rightarrow (B) \mid \varepsilon$	$A \rightarrow p_iAq_i \mid AA \mid \varepsilon \mid x$
(for all $i \in \{1, \dots, n\}$ )	$B \rightarrow [A] \mid \varepsilon$	(for all $i \in \{1, 2\}$ and $x \in X$ )

In fact, [Yellin and Weiss \(2021a\)](#) use *language-model RNNs* where, in addition, every letter gets a dedicated score in a given state. They define the semantics  $L(R)$  as the *locally  $\tau$ -truncated support* where “acceptance” is subject to the condition that the score of every letter (including “end of sequence”) has a score greater than  $\tau$  (cf. [Hewitt et al., 2020](#)).

Several well-known architectures are available to effectively represent RNNs, such as (simple) Elman RNNs, LSTM ([Hochreiter and Schmidhuber, 1997](#)), and GRUs ([Cho et al., 2014](#)). Generally, depending on the architecture, the expressive power of RNNs goes beyond the regular languages. So it is worthwhile to study extraction methods for classes of CFLs.

**Methodology and Results.** The 15 CFLs considered by [Yellin and Weiss \(2021a\)](#) are given in Table 2, together with the CFGs from Table 1. For conciseness, they are defined in terms of general CFGs. However, it turns out that all of them are VPLs. In most cases, there is arguably a canonical partition of the alphabet into a visibly pushdown alphabet. For all these VPLs, we considered the RNNs provided by [Yellin and Weiss \(2021a\)](#), which were trained on sample sets generated by a probabilistic version of a corresponding CFG.

In our experiments, we used a Kearns-Vazirani variation of  $TL^*$  (cf. [Drewes et al. \(2011\)](#)). A query  $MQ_{\text{vpl}}(w)$  for a well-formed word  $w$  was answered according to the given RNN  $R$ , i.e.,  $MQ_{\text{vpl}}(w) = \text{yes}$  iff  $w \in L(R)$ . To answer a query  $EQ_{\text{vpl}}(\mathcal{B})$ , we used two independent subroutines that look for counterexample words (of length under 30):

- (i) We chose 1500 random words in the current hypothesis language  $\llbracket T(\mathcal{B}) \rrbracket$ .
- (ii) We noticed that the 15 languages represented by the RNNs from ([Yellin and Weiss, 2021a](#)) are very sparse (similarly to a lot of other examples of RNN languages from the literature). Therefore, taking only random samples from the RNN would most likely produce an empty language. In order to avoid this, we implemented a type of  $A^*$  exploration (cf. [Russell and Norvig, 2020](#)) in the rooted directed tree of all words  $\Sigma^*$ , where each vertex is a word  $w \in \Sigma^*$  and its children are  $wa$  for  $a \in \Sigma$ . This word-exploration technique relies on an evaluation function  $f : \Sigma^* \rightarrow \mathbb{R}$  where the higher the score of a word the higher its priority to be explored first. The function we chose for this depended on two things: 1) The average score given to the word by the RNN and its neighborhood (where the assumption is that the higher the score the closer we are to a word in the RNN language), and 2) The length of the word, where

we preferred shorter words. The function we chose is

$$f(w) = \frac{1}{|w|^2} \sum_{w' \in \Sigma^* \text{ s.t. } |w'| \leq d} R(ww')$$

where  $R(\cdot)$  is the score given to the word by the RNN, and the size of the neighborhood was chosen to be  $d = 4$ . Using this type of exploration, we generated a set  $P$  (performed once in the beginning of the run) of positive examples from the RNN language (only *well-formed words*; timeout of 60 seconds).

Note that  $\text{EQ}_{\text{vpl}}$  is counterexample-sound for  $L(R)$  but not necessarily equivalence-sound. It was sufficiently precise on our small set of examples, which leads us to believe that there is a reasonable chance that a further investigation might reveal that it is well performing in a more general environment. Though the given trained RNNs have imperfections, the intended languages are learned in most cases. Table 2 indicates the time needed to learn a VPG, averaging across five runs, and the number of rules extracted. In most runs, the extracted VPGs are equivalent to the respective CFGs the RNNs were trained on. Exceptions are  $L_{14}$  and  $L_{15}$  for which we obtain grammars approximating the respective languages. This happens due to structural errors in the given RNNs w.r.t. the target languages.

To give a (successful) example, Table 3 depicts the grammar that was output for  $L_{10}$ .

**Fixing Mistakes Variation.** The previous result can easily be ruined by a wrong sample of words. For example, we could pick a word that is in the RNN language but not in the original language. To mitigate this problem, one can do the following: Denote by  $P$  the set of positive examples generated from the RNN, let  $H$  be the current hypothesis grammar, and let  $\text{pos}(H) = \frac{|P \cap L(H)|}{|P|}$ . Assume that  $H$  comes with a counterexample  $w_c$  and a new hypothesis  $H'$ . If  $\text{pos}(H') < \text{pos}(H)$ , then we keep refining both of them, but making sure that  $w_c$  cannot be used as a counterexample for  $H$ . In the end, we return the hypothesis which is “closest” to the RNN language, i.e., the one with largest  $\text{pos}(H)$ . For example, by increasing the sampling length from the RNN ( $30 \rightarrow 40$ ) and the size of the sample set ( $1500 \rightarrow 2000$ ), we ruined (most of the RNN have some errors in them, it is just a matter of time to find them) the extraction of language  $L_8$ , but using the procedure above, we manage to fixed this issue.

**Agnostic Learning.** Some criticism might be given to the fact that we assume the visibly pushdown alphabet to be known. To solve this issue, we generated a set  $P$  of positive words from the RNN like before. Using  $P$ , we examined all the possible visibly pushdown alphabets (there may be several), picked the best suited alphabet (with the least number of internal symbols), and continued learning. We succeeded in 8 of the 13 languages that were successful in the non-agnostic case.

**Further experiments and optimizations.** As mentioned in the beginning of the section, there is a strong need for further experimentation and evaluation of this algorithm. Here is a preliminary list of ideas towards these goals:

- Larger pool of RNNs representing different visibly pushdown languages.
- In this paper, the number of queries taken while checking equality was chosen to be very small. This was done in order to reproduce the exact CFLs, since  $L^*$  type of algorithms are inherently sensitive to errors. Even one error can produce a completely different VPG. If

Table 2: Results for learning RNNs

	Language	Visibly Pushdown Alphabet			#Rules	Time
		Push	Pop	Int		
$L_1$	$\mathcal{L}(\{a\}, \{b\})$	$\{a\}$	$\{b\}$		3	1s
$L_2$	$\mathcal{L}(\{a, b\}, \{c, d\})$	$\{a, b\}$	$\{c, d\}$		9	23s
$L_3$	$\mathcal{L}(\{ab, cd\}, \{ef, gh\})$	$\{a, b, c, d\}$	$\{e, f, g, h\}$		13	74s
$L_4$	$\mathcal{L}(\{ab\}, \{cd\})$	$\{a, b\}$	$\{c, d\}$		4	1s
$L_5$	$\mathcal{L}(\{abc\}, \{def\})$	$\{a, b, c\}$	$\{d, e, f\}$		5	1s
$L_6$	$\mathcal{L}(\{ab, c\}, \{de, f\})$	$\{a, c\}$	$\{d, f\}$	$\{b, e\}$	10	49s
$L_7$	$Dyck_2$	$\{p_1, p_2\}$	$\{q_1, q_2\}$		19	69s
$L_8$	$Dyck_3$	$\{p_1, p_2, p_3\}$	$\{q_1, q_2, q_3\}$		28	74s
$L_9$	$Dyck_4$	$\{p_1, \dots, p_4\}$	$\{q_1, \dots, q_4\}$		37	79s
$L_{10}$	$RE\text{-}Dyck(\{(abcd), \{wxyz\})$	$\{(, a, b, c, d\}$	$\{w, x, y, z, \}$		10	7s
$L_{11}$	$RE\text{-}Dyck(\{ab, c\}, \{de, f\})$	$\{a, c\}$	$\{d, f\}$	$\{b, e\}$	27	59s
$L_{12}$	<i>Alternating</i>	$\{(, []$	$\{), \}\}$		5	2s
$L_{13}$	$Dyck_1 \triangleleft \{a, b, c\}$	$\{p_1\}$	$\{q_1\}$	$\{a, b, c\}$	19	66s
$L_{14}$	$Dyck_2 \triangleleft \{a, b, c\}$	$\{p_1, p_2\}$	$\{q_1, q_2\}$	$\{a, b, c\}$	–	65s
$L_{15}$	$Dyck_1 \triangleleft \{abc, d\}$	$\{p_1\}$	$\{q_1\}$	$\{a, b, c, d\}$	–	51s

Table 3: Learned VPG for  $L_{10}$  with start symbol  $A_1$ 

$A_1 \rightarrow (A_2)A_0$	$A_2 \rightarrow aA_3zA_0$	$A_4 \rightarrow cA_5xA_0$	$A_5 \rightarrow dA_1wA_0$	$A_6 \rightarrow (A_2)A_1$
$A_0 \rightarrow \varepsilon$	$A_3 \rightarrow bA_4yA_0$	$A_5 \rightarrow dA_0wA_0$	$A_5 \rightarrow dA_6wA_0$	$A_6 \rightarrow (A_2)A_6$

the goal had been to produce a VPG whose language is statistically close to the original one, then one could have taken a more probabilistic approach using, for example, the Chernoff-Hoeffding bound (cf. (Khmelnitsky et al., 2021)).

- Testing this technique on RNNs that represent a language which is not necessarily a CFL (e.g., the Amazon sentiment analysis, which was previously unlearnable with this type of method due to the sparseness of the language).

- Currently, in order to learn a VPG representing the given language, we learn a DTA and then translate it into a VPG. Instead, one may try to use an approach to learning the VPG/VPA (well-formed and not well-formed) directly.

## 7. Conclusion

We presented an algorithm to learn VPLs in the MAT framework. As an application, we focused on the extraction of grammars from RNNs. Our experiments suggest that the algorithm is a suitable alternative to current approaches when we deal with structured data.

Learning VPLs has potential applications in formal verification (Alur and Madhusudan, 2004, 2009), so it would be worthwhile to conduct an evaluation in that domain, too.

**Acknowledgment.** We would like to thank the reviewers for their helpful remarks, which led to an improved presentation of the paper.

## References

- Rajeev Alur and P. Madhusudan. Visibly pushdown languages. In László Babai, editor, *Proceedings of the 36th Annual ACM Symposium on Theory of Computing, Chicago, IL, USA, June 13-16, 2004*, pages 202–211. ACM, 2004. doi: 10.1145/1007352.1007390. URL <https://doi.org/10.1145/1007352.1007390>.
- Rajeev Alur and P. Madhusudan. Adding nesting structure to words. *J. ACM*, 56(3):16:1–16:43, 2009. doi: 10.1145/1516512.1516518. URL <https://doi.org/10.1145/1516512.1516518>.
- Dana Angluin. Inference of reversible languages. *J. ACM*, 29(3):741–765, 1982.
- Dana Angluin. Learning regular sets from queries and counterexamples. *Inf. Comput.*, 75(2):87–106, 1987.
- Stéphane Ayache, Rémi Eyraud, and Noé Goudian. Explaining black boxes on sequential data using weighted automata. In Olgierd Unold, Witold Dyrka, and Wojciech Wieczorek, editors, *Proceedings of the 14th International Conference on Grammatical Inference, ICGI 2018, Wrocław, Poland, September 5-7, 2018*, volume 93 of *Proceedings of Machine Learning Research*, pages 81–103. PMLR, 2018.
- José L. Balcázar, Josep Díaz, and Ricard Gavaldà. Algorithms for learning finite automata from queries: A unified view. In *Advances in Algorithms, Languages, and Complexity - In Honor of Ronald V. Book*, pages 53–72. Kluwer, 1997.
- Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proc. EMNLP*, pages 1724–1734. ACL, 2014.
- Alexander Clark. Distributional learning of some context-free languages with a minimally adequate teacher. In José M. Sempere and Pedro García, editors, *Proc. of ICGI 2010*, volume 6339 of *LNCS*, pages 24–37. Springer, 2010. doi: 10.1007/978-3-642-15488-1\_4. URL [https://doi.org/10.1007/978-3-642-15488-1\\_4](https://doi.org/10.1007/978-3-642-15488-1_4).
- Alexander Clark and Rémi Eyraud. Polynomial identification in the limit of substitutable context-free languages. *J. Mach. Learn. Res.*, 8:1725–1745, 2007.
- H. Comon, M. Dauchet, F. Jacquemard, R. Gilleron, C. Löding, D. Lugiez, S. Tison, and M. Tommasi. Tree automata techniques and applications. Available on: <http://www.grappa.univ-lille3.fr/tata>, 2007. release October, 12th 2007.

- Sreerupa Das, C. Lee Giles, and Guo-Zheng Sun. Using prior knowledge in a {NNPDA} to learn context-free languages. In *Advances in Neural Information Processing Systems 5, [NIPS Conference, Denver, Colorado, USA, November 30 - December 3, 1992]*, pages 65–72. Morgan Kaufmann, 1992.
- Colin de la Higuera. A bibliographical study of grammatical inference. *Pattern Recognit.*, 38(9):1332–1348, 2005.
- Frank Drewes and Johanna Högberg. Query learning of regular tree languages: How to avoid dead states. *Theory Comput. Syst.*, 40(2):163–185, 2007. doi: 10.1007/s00224-005-1233-3. URL <https://doi.org/10.1007/s00224-005-1233-3>.
- Frank Drewes, Johanna Högberg, and Andreas Maletti. MAT learners for tree series: an abstract data type and two realizations. *Acta Informatica*, 48(3):165–189, 2011. doi: 10.1007/s00236-011-0135-x.
- John Hewitt, Michael Hahn, Surya Ganguli, Percy Liang, and Christopher D. Manning. Rnns can generate bounded hierarchical languages with optimal memory. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 1978–2010. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.emnlp-main.156.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, 1997.
- Malte Isberner. *Foundations of active automata learning: an algorithmic perspective*. PhD thesis, Technical University Dortmund, Germany, 2015. URL <http://hdl.handle.net/2003/34282>.
- Malte Isberner, Falk Howar, and Bernhard Steffen. The TTT algorithm: A redundancy-free approach to active automata learning. In Borzoo Bonakdarpour and Scott A. Smolka, editors, *Proc. of RV 2014*, volume 8734 of *LNCS*, pages 307–322. Springer, 2014. doi: 10.1007/978-3-319-11164-3\_26. URL [https://doi.org/10.1007/978-3-319-11164-3\\_26](https://doi.org/10.1007/978-3-319-11164-3_26).
- Michael J. Kearns and Umesh V. Vazirani. *An Introduction to Computational Learning Theory*. MIT Press, 1994.
- Igor Khmelnitsky, Daniel Neider, Rajarshi Roy, Xuan Xie, Benoît Barbot, Benedikt Bollig, Alain Finkel, Serge Haddad, Martin Leucker, and Lina Ye. Property-directed verification and robustness certification of recurrent neural networks. In *Proceedings of the 19th International Symposium on Automated Technology for Verification and Analysis (ATVA 2021)*, 2021. To appear.
- Viraj Kumar, P. Madhusudan, and Mahesh Viswanathan. Minimization, learning, and conformance testing of boolean programs. In *Proc. of CONCUR 2006*, volume 4137 of *LNCS*, pages 203–217. Springer, 2006. doi: 10.1007/11817949\_14. URL [https://doi.org/10.1007/11817949\\_14](https://doi.org/10.1007/11817949_14).



- P. Madhusudan and Gennaro Parlato. The tree width of auxiliary storage. In Thomas Ball and Mooly Sagiv, editors, *Proc. of POPL 2011*, pages 283–294. ACM, 2011. doi: 10.1145/1926385.1926419. URL <https://doi.org/10.1145/1926385.1926419>.
- Franz Mayr and Sergio Yovine. Regular inference on artificial neural networks. In *Proc. of CD-MAKE 2018*, volume 11015 of *LNCS*, pages 350–369. Springer, 2018.
- Franz Mayr, Ramiro Visca, and Sergio Yovine. On-the-fly black-box probably approximately correct checking of recurrent neural networks. In *Proc. of CD-MAKE 2020*, volume 12279 of *LNCS*, pages 343–363. Springer, 2020. doi: 10.1007/978-3-030-57321-8\_19. URL [https://doi.org/10.1007/978-3-030-57321-8\\_19](https://doi.org/10.1007/978-3-030-57321-8_19).
- Kurt Mehlhorn. Pebbling mountain ranges and its application of defl-recognition. In J. W. de Bakker and Jan van Leeuwen, editors, *Automata, Languages and Programming, 7th Colloquium, Noordwijkerhout, The Netherlands, July 14-18, 1980, Proceedings*, volume 85 of *LNCS*, pages 422–435. Springer, 1980. doi: 10.1007/3-540-10003-2\_89. URL [https://doi.org/10.1007/3-540-10003-2\\_89](https://doi.org/10.1007/3-540-10003-2_89).
- Christian W. Omlin and C. Lee Giles. Extraction of rules from discrete-time recurrent neural networks. *Neural Networks*, 9(1):41–52, 1996. doi: 10.1016/0893-6080(95)00086-0.
- Guillaume Rabusseau, Tianyu Li, and Doina Precup. Connecting weighted automata and recurrent neural networks through spectral learning. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS 2019, 16-18 April 2019, Naha, Okinawa, Japan*, volume 89 of *Proceedings of Machine Learning Research*, pages 1630–1639. PMLR, 2019.
- Ronald L. Rivest and Robert E. Schapire. Inference of finite automata using homing sequences. *Inf. Comput.*, 103(2):299–347, 1993.
- Stuart J. Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach (4th Edition)*. Pearson, 2020.
- Yasubumi Sakakibara. Efficient learning of context-free grammars from positive structural examples. *Inf. Comput.*, 97(1):23–60, 1992.
- Guo-Zheng Sun, C. Lee Giles, and Hsing-Hen Chen. The neural network pushdown automaton: Architecture, dynamics and training. In *Adaptive Processing of Sequences and Data Structures, International Summer School on Neural Networks, "E.R. Caianiello", Vietri sul Mare, Salerno, Italy, September 6-13, 1997, Tutorial Lectures*, volume 1387 of *LNCS*, pages 296–345. Springer, 1997.
- Sebastian Thrun. Extracting rules from artificial neural networks with distributed representations. In *Advances in Neural Information Processing Systems 7, [NIPS Conference, Denver, Colorado, USA, 1994]*, pages 505–512. MIT Press, 1994.
- Frits W. Vaandrager. Model learning. *Commun. ACM*, 60(2):86–95, 2017. doi: 10.1145/2967606.

- Gail Weiss, Yoav Goldberg, and Eran Yahav. Extracting automata from recurrent neural networks using queries and counterexamples. In *Proc. of ICML 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 5244–5253. PMLR, 2018.
- Gail Weiss, Yoav Goldberg, and Eran Yahav. Learning deterministic weighted automata with queries and counterexamples. In *Proc. of NeurIPS 2019*, pages 8558–8569, 2019.
- Daniel M. Yellin and Gail Weiss. Synthesizing context-free grammars from recurrent neural networks. In Jan Friso Groote and Kim Guldstrand Larsen, editors, *Proc. of TACAS 2021, Part I*, volume 12651 of *LNCS*, pages 351–369. Springer, 2021a. doi: 10.1007/978-3-030-72016-2\_19. URL [https://doi.org/10.1007/978-3-030-72016-2\\_19](https://doi.org/10.1007/978-3-030-72016-2_19).
- Daniel M. Yellin and Gail Weiss. Synthesizing context-free grammars from recurrent neural networks (extended version). *CoRR*, abs/2101.08200, 2021b. URL <https://arxiv.org/abs/2101.08200>.
- Ryo Yoshinaka and Alexander Clark. Polynomial time learning of some multiple context-free languages with a minimally adequate teacher. In *Proc. of FG 2010*, volume 7395 of *LNCS*, pages 192–207. Springer, 2010. doi: 10.1007/978-3-642-32024-8\_13. URL [https://doi.org/10.1007/978-3-642-32024-8\\_13](https://doi.org/10.1007/978-3-642-32024-8_13).

## Appendix A. Visibly Pushdown Automata

Though we are principally interested in inferring grammars, we give here the definition of visibly pushdown automata, which also constitute a characterization of the class of VPLs.

**Definition 5** A visibly pushdown automaton (VPA) over  $\Sigma$  is a tuple  $\mathcal{A} = (Q, \mathcal{S}, \delta, \iota, F)$  containing a finite set of control states  $Q$ , a nonempty finite set of stack symbols  $\mathcal{S}$ , an initial state  $\iota$ , and a set of final states  $F \subseteq Q$ . Moreover,  $\delta = (\delta_{\text{push}}, \delta_{\text{pop}}, \delta_{\text{int}})$  is a collection of transition functions  $\delta_{\text{push}} : Q \times \Sigma_{\text{push}} \rightarrow \mathcal{P}(Q \times \mathcal{S})$ ,  $\delta_{\text{pop}} : Q \times \Sigma_{\text{pop}} \times \mathcal{S} \rightarrow \mathcal{P}(Q)$ , and  $\delta_{\text{int}} : Q \times \Sigma_{\text{int}} \rightarrow \mathcal{P}(Q)$ . We call  $\mathcal{A}$  deterministic if all transition functions map all arguments to singleton sets.

A VPA recognizes a language  $L(\mathcal{A}) \subseteq \Sigma^*$ . Intuitively, it is the language of an infinite automaton whose states (we actually say configurations) are pairs  $(q, \sigma)$  where  $q \in Q$  is the current control state and  $\sigma \in \mathcal{S}^*$  is the current stack contents. With this, in the infinite automaton, we have a transition  $(q, \sigma) \xrightarrow{a} (q', \sigma')$  if there is  $A \in \mathcal{S}$  such that one of the following holds:

- $a \in \Sigma_{\text{push}}$  and  $(q', A) \in \delta_{\text{push}}(q, a)$  and  $\sigma' = \sigma \cdot A$
- $a \in \Sigma_{\text{pop}}$  and  $q' \in \delta_{\text{pop}}(q, a, A)$  and  $\sigma = \sigma' \cdot A$
- $a \in \Sigma_{\text{int}}$  and  $q' \in \delta_{\text{int}}(q, a)$  and  $\sigma' = \sigma$

We call  $(q, \sigma)$  a final configuration if  $q \in F$  and  $\sigma = \varepsilon$ . Moreover,  $(\iota, \varepsilon)$  is the only initial configuration. Finally, we define  $L(\mathcal{A})$  to be the language recognized by this infinite automaton in the expected way.

**Fact 6 (Alur and Madhusudan (2004, 2009))** Let  $L \subseteq \Sigma^*$ . Then,  $L$  is a VPL over  $\Sigma$  iff there is a VPA  $\mathcal{A}$  over  $\Sigma$  such that  $L(\mathcal{A}) = L$ .

## Appendix B. Proof of Theorem 4

**Proof** By Fact 5, we have to show that calling  $\text{TL}^*(\Gamma, \text{MQ}_{\text{tree}}, \text{EQ}_{\text{tree}})$  returns, in polynomial time, a DTA  $\mathcal{B}$  such that  $T(\mathcal{B}) = T(\hat{\mathcal{B}})$ . We will show that  $\text{MQ}_{\text{tree}}$  is sound for  $T(\hat{\mathcal{B}})$  and  $\text{EQ}_{\text{tree}}$  is counterexample- and equivalence-sound for  $T(\hat{\mathcal{B}})$ . By Fact 3, this implies that  $\text{TL}^*(\Gamma, \text{MQ}_{\text{tree}}, \text{EQ}_{\text{tree}})$  returns a DTA  $\mathcal{B}$  such that  $T(\mathcal{B}) = T(\hat{\mathcal{B}})$ . The running time is due to the fact that all additional operations in Algorithms 1 and 2 can be performed in polynomial time wrt. the input parameters (cf. Fact 2).

To show that  $\text{MQ}_{\text{tree}}$  is sound for  $T(\hat{\mathcal{B}})$ , let  $t \in \text{Trees}(\Gamma)$ . Assume  $\text{MQ}_{\text{tree}}(t) = \text{yes}$ . By Algorithm 1, this implies  $t \in T(\mathcal{B}_{\text{parse}})$  and  $\text{MQ}_{\text{vpl}}(\llbracket t \rrbracket) = \text{yes}$ . As  $\text{MQ}_{\text{vpl}}$  is sound for  $L$ , we have  $\llbracket t \rrbracket \in L$ . Since  $T(\hat{\mathcal{B}}) = \langle\langle L \rangle\rangle$ , we get  $t \in T(\hat{\mathcal{B}})$ . Conversely, assume  $\text{MQ}_{\text{tree}}(t) = \text{no}$ . If  $t \notin T(\mathcal{B}_{\text{parse}})$ , then  $t \notin T(\hat{\mathcal{B}})$ . So suppose  $t \in T(\mathcal{B}_{\text{parse}})$  and  $\text{MQ}_{\text{vpl}}(\llbracket t \rrbracket) = \text{no}$ . As  $\text{MQ}_{\text{vpl}}$  is sound for  $L$ , we have  $\llbracket t \rrbracket \notin L$ , which implies  $t \notin T(\hat{\mathcal{B}})$ .

Let us show that  $\text{EQ}_{\text{tree}}$  is counterexample-sound for  $T(\hat{\mathcal{B}})$ . Suppose  $\mathcal{B} \in \text{DTA}(\Gamma)$  and  $t \in \text{Trees}(\Gamma)$  such that  $\text{EQ}_{\text{tree}}(\mathcal{B}) = t$ . There are two cases. First, suppose  $t \in T(\mathcal{B}) \setminus T(\mathcal{B}_{\text{parse}})$ . As  $T(\hat{\mathcal{B}}) \subseteq T(\mathcal{B}_{\text{parse}})$ , we have  $t \in T(\mathcal{B}) \setminus T(\hat{\mathcal{B}})$  and, hence,  $t \in T(\mathcal{B}) \oplus T(\hat{\mathcal{B}})$ . Second, assume

$T(\mathcal{B}) \subseteq T(\mathcal{B}_{\text{parse}})$ . As  $\text{EQ}_{\text{vpl}}$  is counterexample-sound for  $L$ , this implies  $\llbracket t \rrbracket \in L \oplus \llbracket T(\mathcal{B}) \rrbracket$ . Due to  $T(\hat{\mathcal{B}}) = \langle\langle L \rangle\rangle$ , we get  $t \in T(\hat{\mathcal{B}}) \oplus T(\mathcal{B})$ .

Finally, we show that  $\text{EQ}_{\text{tree}}$  is equivalence-sound for  $T(\hat{\mathcal{B}})$ . Suppose  $\mathcal{B} \in \text{DTA}(\Gamma)$  such that  $\text{EQ}_{\text{tree}}(\mathcal{B}) = \text{yes}$ . Then,  $T(\mathcal{B}) \subseteq T(\mathcal{B}_{\text{parse}})$  and  $\text{EQ}_{\text{vpl}}(\mathcal{B}) = \text{yes}$ . As  $\text{EQ}_{\text{vpl}}$  is equivalence-sound for  $L$ , we get  $L = \llbracket T(\mathcal{B}) \rrbracket$ , which implies  $T(\hat{\mathcal{B}}) = T(\mathcal{B})$ . ■