

Generalizing input-driven languages: theoretical and practical benefits

Dino Mandrioli¹, Matteo Pradella^{1,2}

¹ DEIB – Politecnico di Milano, via Ponzio 34/5, Milano, Italy

² IEIIT – Consiglio Nazionale delle Ricerche, via Golgi 42, Milano, Italy
{dino.mandrioli, matteo.pradella}@polimi.it

Abstract. Regular languages (RL) are the simplest family in Chomsky’s hierarchy. Thanks to their simplicity they enjoy various nice algebraic and logic properties that have been successfully exploited in many application fields. Practically all of their related problems are decidable, so that they support automatic verification algorithms. Also, they can be recognized in real-time.

Context-free languages (CFL) are another major family well-suited to formalize programming, natural, and many other classes of languages; their increased generative power w.r.t. RL, however, causes the loss of several closure properties and of the decidability of important problems; furthermore they need complex parsing algorithms. Thus, various subclasses thereof have been defined with different goals, spanning from efficient, deterministic parsing to closure properties, logic characterization and automatic verification techniques.

Among CFL subclasses, so-called *structured* ones, i.e., those where the typical tree-structure is visible in the sentences, exhibit many of the algebraic and logic properties of RL, whereas deterministic CFL have been thoroughly exploited in compiler construction and other application fields.

After surveying and comparing the main properties of those various language families, we go back to *operator precedence languages* (OPL), an old family through which R. Floyd pioneered deterministic parsing, and we show that they offer unexpected properties in two fields so far investigated in totally independent ways: they enable parsing parallelization in a more effective way than traditional sequential parsers, and exhibit the same algebraic and logic properties so far obtained only for less expressive language families.

Keywords: regular languages, context-free languages, input-driven languages, visibly pushdown languages, operator-precedence languages, monadic second order logic, closure properties, decidability and automatic verification.

1 Introduction

Regular (RL) and context-free languages (CFL) are by far the most widely studied families of formal languages in the richest literature of the field. In Chomsky’s hierarchy, they are, respectively, in positions 2 and 3, 0 and 1 being recursively enumerable and context-sensitive languages.

Thanks to their simplicity, RL enjoy practically all positive properties that have been defined and studied for formal language families: they are closed under most algebraic operations, and most of their properties of interest (emptiness, finiteness, containment)

are decidable. Thus, they found fundamental applications in many fields of computer and system science: HW circuit design and minimization, specification and design languages (equipped with powerful supporting tools), automatic verification of SW properties, etc. One of their most relevant applications is now model-checking which exploits the decidability of the containment problem and important characterizations in terms of mathematical logics [12,23].

On the other hand, the typical linear structure of RL sentences makes them unsuitable or only partially suitable for application in fields where the data structure is more complex, e.g., is tree-like. For instance, in the field of compilation they are well-suited to drive lexical analysis but not to manage the typical nesting of programming and natural language features. The classical language family adopted for this type of modeling and analysis is the context-free one. The increased expressive power of CFL allows to formalize many syntactic aspects of programming, natural, and various other categories of languages. Suitable algorithms have been developed on their basis to parse their sentences, i.e., to build the structure of sentences as syntax-trees.

General CFL, however, lose various of the nice mathematical properties of RL: they are closed only under some of the algebraic operations, and several decision problems, typically the inclusion problem, are undecidable; thus, the automatic analysis and synthesis techniques enabled for RL are hardly generalized to CFL. Furthermore, parsing CFL may become considerably less efficient than recognizing RL: the present most efficient parsing algorithms of practical use for general CFL have an $O(n^3)$ time complexity.

The fundamental subclass of deterministic CFL (DCFL) has been introduced, and applied to the formalization of programming language syntax, to exploit the fact that in this case parsing is in $O(n)$. DCFL, however, do not enjoy enough algebraic and logic properties to extend to this class the successful applications developed for RL: e.g., although their equivalence is decidable, containment is not; they are closed under complement but not under union, intersection, concatenation and Kleene *.

From this point of view, *structured CFL* are somewhat in between RL and CFL. Intuitively, by structured CFL we mean languages where the structure of the syntax-tree associated with a given sentence is immediately apparent in the sentence. *Parenthesis languages (PL)* introduced in a pioneering paper by McNaughton [39] are the first historical example of such languages. McNaughton showed that they enjoy closure under Boolean operations (which, together with the decidability of the emptiness problem, implies decidability of the containment problem) and their generating grammars can be minimized in a similar way as *finite state automata (FSA)* are minimized (in fact an equivalent formalism for parenthesis languages are *tree automata* [45,13]).

Starting from PL various extensions of this family have been proposed in the literature, with the main goal of preserving most of the nice properties of RL and PL, yet increasing their generative power; among them *input-driven languages (IDL)* [40,48], later renamed *visibly pushdown languages (VPL)* [5] have been quite successful: they are closed under all traditional language operations (and therefore enjoy the consequent decidability properties). Also, they are characterized in terms of a *monadic second order (MSO)* logic by means of a natural extension of the classic characterization for RL originally and independently developed by Büchi, Elgot, and Trakhtenbrot [9,22,47].

For these reasons they are a natural candidate for extending model checking techniques from RL. To achieve such a goal in practice, however, MSO logic is not yet tractable due to the complexity of its decidability problems; thus, some research is going on to “pair” IDL with specification languages inspired by temporal logic as it has been done for RL [1].

On the other hand being IDL structured their parsing problem trivially scales down to the recognition problem but clearly they are not suitable to formalize languages whose structure is hidden from the surface sentences as, e.g., typical arithmetic expressions.

Rather recently, we resumed the study of an old class of languages which was interrupted a long time ago, namely *operator precedence languages (OPL)*. OPL and their generating grammars (OPG) have been introduced by Floyd [26] to build efficient deterministic parsers; indeed they generate a large and meaningful subclass of DCFL. In the past their algebraic properties, typically closure under Boolean operations [16], have been investigated with the main goal of designing inference algorithms for their languages [17]. After that, their theoretical investigation has been abandoned because of the advent of more powerful grammars, mainly LR ones [33,30], that generate all DCFL (although some deterministic parsers based on OPL’s simple syntax have been continuously implemented at least for suitable subsets of programming languages [29]).

The renewed interest in OPG and OPL has been ignited by two seemingly unrelated remarks: on the one hand we realized that they are a proper superclass of IDL and that all results that have been obtained for them (closures, decidability, logical characterization) extend naturally, but not trivially, to OPL; on the other hand new motivation for their investigation comes from their distinguishing property of *local parsability*: with this term we mean that their deterministic parsing can be started and led to completion from any position of the input string unlike what happens with general deterministic pushdown automata, which must necessarily operate strictly left-to-right from the beginning of the input. This property has a strong practical impact since it allows for exploiting modern parallel architectures to obtain a natural speed up in the processing of large tree-structured data. An automatic tool that generates parallel parsers for these grammars has already been produced and is freely available. The same local parsability property can also be exploited to incrementally analyze large structures without being compelled to reparse them from scratch after any modification thereof.

This renewed interest in OPL has also led to extend their study to ω -languages, i.e., those consisting of infinite strings: in this case too the investigation produced results that perfectly parallel the extension of other families, noticeably RL and IDL, from the finite string versions to the infinite ones.

In this paper we follow the above “story” since its beginning to these days and, for the first time, we join within the study of one single language family two different application domains, namely parallel parsing on the one side, and algebraic and logic characterization finalized to automatic verification on the other side. To the best of our knowledge, OPL is the largest family that enjoys all of such properties.

The paper is structured as follows: initially we resume some background on the two main families of formal languages. Section 2 introduces RL and various formalisms widely used in the literature to specify them; special attention is given to their logic

characterization which is probably less known to the wider computer science audience and is crucial for some fundamental results of this paper. Section 3 introduces CFL in a parallel way as for regular ones. Then, two sections are devoted to different subclasses of CFL that allowed to obtain important properties otherwise lacking in the larger original family: precisely, Section 4 deals with *structured CFL*, i.e., those languages whose tree-shaped structure is in some way immediately apparent from their sentences with no need to parse them; it shows that for such subclasses of CFL important properties of RL, otherwise lacking in the larger original family, still hold. Section 5, instead, considers the problem of efficient deterministic parsing for CFL. Finally, Section 6 “puts everything together” by showing that OPL on the one hand are significantly more expressive than traditional structured languages, but enjoy the same important properties as regular and structured CFL and, on the other hand, enable exploiting parallelism in parsing much more naturally and efficiently than for general deterministic CFL.

Since all results presented in this paper have already appeared in the literature, we based our presentation more on intuitive explanation and simple examples than on detailed technical constructions and proofs, to which appropriate references are supplied for the interested reader. However, to guarantee self-containedness, some fundamental definitions have been necessarily given in full formality. For convenience we do not add a final ‘s’ to acronyms when used as plurals so that, e.g., CFL denotes indifferently a single language, the language family and all languages in the family.

2 Regular Languages

The family of regular languages (RL) is one of the most important families in computer science. In the traditional Chomsky’s hierarchy it is the least powerful language family. Its importance stems from both its simplicity and its rich set of properties.

RL is a very robust family, as it enjoys many properties and practically all of them are decidable. It is defined through several different devices, both operational and descriptive. Among them we mention *Finite State Automata* (FSA), which are used for various applications, not only in computer science, *Regular Grammars*, *Regular Expressions*, often used in computing for describing the lexical elements of programming languages and in many programming environments for managing program sources, and various logic classifications that support automatic verification of their properties.

2.1 Regular Grammars and Finite State Automata

Chomsky introduced his seminal hierarchy in 1956. The interested reader may refer to a manual of formal languages, such as [30,14] for a more complete description.

Grammars are generative rewriting systems, and consist of a set of rewriting rules, also called productions. Intuitively, a production $x \rightarrow y$ means that a string x can be replaced with a string y in the rewriting.

Definition 1. A grammar G is a tuple (V_N, Σ, P, S) , where

- V_N , called nonterminal alphabet, and Σ , called terminal alphabet, are finite sets of characters such that $V_N \cap \Sigma = \emptyset$; $V_N \cup \Sigma$ is named V ;

- P is a finite set of productions or rules having the form $V^*V_NV^* \rightarrow V^*$; in a production $x \rightarrow y$, x is called the left-hand side (abbreviated as *lhs*), and y the right-hand side (abbreviated as *rhs*);
- $S \in V_N$ is the starting character (or axiom).

The following *naming conventions* are adopted for letters and strings, unless otherwise specified: lowercase Latin letters at the beginning of the alphabet a, b, \dots denote terminal characters; uppercase Latin letters A, B, \dots denote nonterminal characters; lowercase Latin letters at the end of the alphabet x, y, z, \dots denote terminal strings; lowercase Greek letters α, \dots, ω denote strings over V .

To define the language of a grammar, we need the concept of *derivation*. Grammars are a generative device, because with them one always starts with the axiom (usually the S character), then *derives* new strings by applying an iterative rewriting as specified by grammar's productions, until a string belonging to Σ^* is obtained. This intuition is formalized by the following definition.

Definition 2. Consider a grammar $G = (V_N, \Sigma, P, S)$. A string β is derivable in one step from a string α , denoted by $\alpha \Rightarrow_G \beta$ iff $\alpha = \alpha_1\alpha_2\alpha_3$, $\beta = \alpha_1\alpha'_2\alpha_3$, and $\alpha_2 \rightarrow \alpha'_2 \in P$.

The reflective and transitive closure of the relation \Rightarrow_G is denoted by $\overset{*}{\Rightarrow}_G$. The language $L(G)$ generated by G is the set $\{x \in \Sigma^* \mid S \overset{*}{\Rightarrow}_G x\}$.

Regular languages are defined by restricted grammars, called *regular grammars* (also *left-(or right-)linear grammars*).

Definition 3. A grammar $G = (V_N, \Sigma, P, S)$, such that $P \subseteq \{X \rightarrow Ya \mid X \in V_N, Y \in V_N \cup \{\varepsilon\}, a \in \Sigma\}$, is called *left-linear*. Symmetrically, if its rules have the form $X \rightarrow aY$, G is called *right-linear*.

A grammar is called *regular* if it is either *left-linear* or *right-linear*. A language is *regular* if it is generated by some regular grammar.

Example 1. Consider a right-linear grammar $G_1 = (V_N, \Sigma, P, S)$, where $V_N = \{S, B\}$, $\Sigma = \{a, b\}$, $P = \{S \rightarrow aS, S \rightarrow aB, S \rightarrow a, B \rightarrow bB, B \rightarrow b\}$.

For simplicity and shortness, it is customary to group together productions with the same lhs, writing all the corresponding right parts separated by $|$. In this case, $P = \{S \rightarrow aS \mid aB \mid a, B \rightarrow bB \mid b\}$. It is also customary to drop G from the relation \Rightarrow_G , if the grammar is clear from the context.

Here is an example derivation: $S \Rightarrow aS \Rightarrow aaS \Rightarrow aaaB \Rightarrow aaabB \Rightarrow aaabb$. Hence we see that $aaabb \in L(G_1)$. By simple reasoning we obtain $L(G_1) = \{a^+b^*\}$.

Probably the most known and used notation to define RL is that of *Finite State Automata*. An automaton is a simple operational device, which assumes a *state* from a finite set, and its behavior is based on the concept of *state transition*.

Definition 4. A Finite State Automaton (FSA) \mathcal{A} is a tuple $(Q, \Sigma, \delta, q_0, F)$, where

- Q is a finite set of states;
- Σ is the input alphabet;
- $\delta \subseteq Q \times \Sigma \times Q$ is the transition relation;

- $I \subseteq Q$ is the set of initial states;
- $F \subseteq Q$ is the set of final states.

If I is a singleton and δ is a function, i.e., $\delta : Q \times \Sigma \rightarrow Q$, \mathcal{A} is called *deterministic* (or DFSA).

Whereas grammars *generate* strings through derivations, automata *recognize* or *accept* them through sequences of *transitions* or *runs*. In the case of FSA their recognition is formalized as follows.

Definition 5. Let the FSA \mathcal{A} be $(Q, \Sigma, \delta, I, F)$. The transition sequence of \mathcal{A} starting from a state $q \in Q$, reading a string $x \in \Sigma^*$, and ending in a state q' is written $(q, x, q') \in \delta^*$ where the relation $\delta^* \subseteq Q \times \Sigma^* \times Q$ is inductively defined as follows:

- $(q, \varepsilon, q) \in \delta^*$;
- $(q, y, q') \in \delta^* \wedge (q', a, q'') \in \delta \implies (q, ya, q'') \in \delta^*$, for $q, q', q'' \in Q, a \in \Sigma, y \in \Sigma^*$.

The language accepted by \mathcal{A} is defined as: $L(\mathcal{A}) = \{x \mid \exists q_0 \in I, \exists q_F \in F : (q_0, x, q_F) \in \delta^*\}$.

Example 2. Figure 1 shows a classical representation of an example automaton, where Q is $\{q_0, q_1\}$, q_0 is marked as initial by the leftmost arrow; $F = \{q_1\}$ and the ball of the final state has a double line. Transitions are depicted through labeled arrows, e.g. the arrows connecting q_0 to q_0 itself and to q_1 mean that $(q_0, a, q_0) \in \delta$ and $(q_0, a, q_1) \in \delta$.

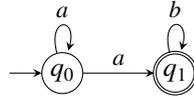


Fig. 1. The FSA \mathcal{A}_1 recognizing the language $L(G_1)$ of Example 1.

It is easy to see that $L(\mathcal{A}_1) = \{a^+b^*\}$. E.g., $(q_0, a, q_0) \in \delta$ implies $(q_0, a, q_0) \in \delta^*$ and $(q_0, a, q_1) \in \delta$ implies $(q_0, a, q_1) \in \delta^*$; from the same transitions we have $\{(q_0, aa, q_0), (q_0, aa, q_1), (q_0, aaa, q_0), (q_0, aaa, q_1)\} \subseteq \delta^*$; then, $(q_1, b, q_1) \in \delta$ and $(q_0, aaa, q_1) \in \delta^*$ imply $(q_0, aaab, q_1) \in \delta^*$; $(q_1, b, q_1) \in \delta$ and $(q_0, aaab, q_1) \in \delta^*$ imply $(q_0, aaabb, q_1) \in \delta^*$. Hence, $aaabb \in L(\mathcal{A}_1)$, since $q_1 \in F$.

By comparing the grammar of Example 1 with the automaton of Example 2 it is immediate to generalize their equivalence to the whole class RL: for any regular grammar an equivalent FSA can be built and conversely.

An important result on FSA is that they are *determinizable*, i.e., given a generic FSA, it is always possible to define an equivalent DFSA. The main idea is to note that the power set of a finite set is still finite, so if we have a FSA with a set of states Q , we can define another automaton with a set of states $\mathcal{P}(Q)$ so that the original transition relation can be made a function. More formally, given an FSA $(Q, \Sigma, \delta, I, F)$, a deterministic equivalent automaton is $\mathcal{A}_D = (\mathcal{P}(Q), \Sigma, \delta_D, \{I\}, F_D)$, where $\delta_D(q, a, \bigcup_{q' \in q, (q', a, q'') \in \delta} q'')$, $q \in \mathcal{P}(Q)$, $a \in \Sigma$, and $F_D = \{Q' \subseteq Q \mid Q' \cap F \neq \emptyset\}$.

We illustrate this construction with an example.

Example 3. The automaton $\mathcal{A}_1 = (Q, \Sigma, \delta, \{q_0\}, F)$ of Example 2 is not deterministic. A deterministic equivalent automaton is shown in Figure 2.

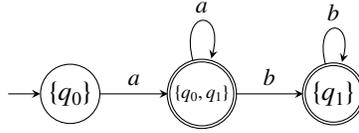


Fig. 2. The DFSA \mathcal{A}_2 recognizing the language $L(G_1)$ of Example 1.

RL enjoy many interesting closure properties, for instance they are closed w.r.t. all Boolean operations (union, intersection, and complement).

As far as complement is concerned, to build up the related automaton, one must start with a deterministic automaton (if not, we have just shown how to define an equivalent deterministic automaton). First, we modify the original automaton so that it reads any possible string – in this way, if the input string is not in the automaton’s language, it changes its state to an additional “error” state, if needed, so that the function δ is total. Then, we exchange the role between accepting and non-accepting states, so that a rejected string becomes part of the language, and vice versa. We illustrate this construction with our running example.

Example 4. We start with the automaton recognizing $L = a^+b^*$ of Figure 2, since it is deterministic; we also rename its states, for brevity.

Then, we introduce an error state, called q_e , and swap accepting with non-accepting roles: the resulting automaton is shown in Figure 3. It is easy to see that $L(\mathcal{A}_3) = \neg L(\mathcal{A}_2) = \neg L(G_1)$.

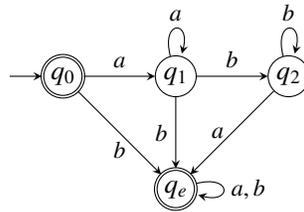


Fig. 3. The DFSA \mathcal{A}_3 recognizing the complement of language $L(G_1)$ of Example 1.

To prove that RL are closed w.r.t. intersection, we start from two FSA $\mathcal{A}' = (Q', \Sigma, \delta', I', F')$, $\mathcal{A}'' = (Q'', \Sigma, \delta'', I'', F'')$, and build a “product” automaton,

$$\mathcal{A}_p = (Q' \times Q'', \Sigma, \delta_p, I' \times I'', F' \times F''),$$

where $(\langle q_1, q'_1 \rangle, a, \langle q_2, q'_2 \rangle) \in \delta_p$, iff $(q_1, a, q_2) \in \delta'$, and $(q'_1, a, q'_2) \in \delta''$.

Closure w.r.t. union follows from De Morgan's laws.

Example 5. Let us consider the automaton of Figure 1 and apply the construction for intersection with the automaton of Figure 4. The resulting automaton is in Figure 5.

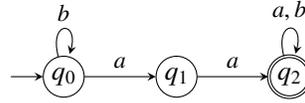


Fig. 4. The FSA \mathcal{A}_4 , recognizing language $b^*aa\{a,b\}^*$.

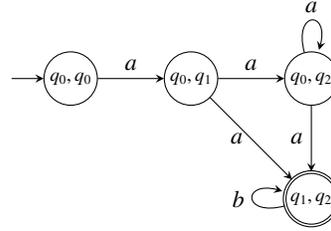


Fig. 5. FSA recognizing $(a^+b^*) \cap (b^*aa\{a,b\}^*)$, obtained from \mathcal{A}_1 and \mathcal{A}_4 .

Another fundamental property of RL is that their emptiness problem, i.e. whether or not $L = \emptyset$, is decidable. This in turn is an immediate consequence of the so-called *pumping lemma* which, informally, states that, due to the finiteness of the state space, the behavior of any FSA must eventually be periodic. More formally, this property is stated as follows:

Lemma 1. *Let \mathcal{A} be a FSA; for every string $x \in L(\mathcal{A})$, with $|x| > |Q|$, there exists a factorization $x = ywz$ such that $x = yw^n z \in L(\mathcal{A})$ for every $n \geq 0$.*

As a consequence of this lemma, if in a regular language there exists a sentence of any length, then there exists also a string of length $\leq |Q|$, hence the decidability of the emptiness problem and, in turn, of the containment problem (is $L(\mathcal{A}_1) \subseteq L(\mathcal{A}_2)$?) and of the equivalence problem, thanks to the closure w.r.t. Boolean operations. The lemma can also be exploited to show that no FSA can recognize $L = \{a^n b^n \mid n \geq 0\}$.

Other notable closures, not of major interest in this paper, are w.r.t. *concatenation*, *Kleene **, *string homomorphism*, and *inverse string homomorphism*; FSA are *minimizable*, i.e. given a FSA, there is an algorithm to build an equivalent automaton with the minimum possible number of states.

2.2 Logic characterization

From the very beginning of formal language and automata theory the investigation of the relations between defining a language through some kind of abstract machine and through a logic formalism has produced challenging theoretical problems and important applications in system design and verification. A well-known example of such an application is the classical Hoare's method to prove the correctness of a Pascal-like program w.r.t. a specification stated as a pair of pre- and post-conditions expressed through a first-order theory [32].

Such a verification problem is undecidable if the involved formalisms have the computational power of Turing machines but may become decidable for less powerful formalisms as in the important case of RL. In particular, Büchi, in his seminal work [9], provided a complete characterization of regular languages in terms of a suitable logic, so that the decidability properties of this class of languages could be exploited to achieve automatic verification; later on, in fact a major breakthrough in this field has been obtained thanks to advent of model checking. The logic proposed by Büchi to characterize regular languages is a Monadic Second Order (MSO) one that is interpreted over the natural numbers representing the position of a character in the string: the string is accepted by a FSA iff it satisfies a sentence in the corresponding logic. The basic elements of the logic defined by Büchi are summarized here:

- First-order variables, denoted as lowercase letters at the end of the alphabet, x, y, \dots are interpreted over the natural numbers \mathbb{N} (these variables are written in boldface to avoid confusion with strings);
- Second-order variables, denoted as uppercase letters at the end of the alphabet, written in boldface, X, Y, \dots are interpreted over the power set of natural numbers $\mathcal{P}(\mathbb{N})$;
- For a given input alphabet Σ , the monadic predicate $a(\cdot)$ is defined for each $a \in \Sigma$: $a(x)$ evaluates to true in a string iff the character at position x is a ;
- The *successor* predicate is denoted by succ , i.e. $\text{succ}(x, y)$ means that $y = x + 1$.
- The well-formed formulas are defined by the following syntax:

$$\varphi := a(x) \mid x \in X \mid \text{succ}(x, y) \mid \neg\varphi \mid \varphi \vee \varphi \mid \exists x(\varphi) \mid \exists X(\varphi).$$

- They are interpreted in the natural way. Furthermore the usual predefined abbreviations are introduced to denote the remaining propositional connectives, universal quantifiers, arithmetic relations ($=, \neq, <, >$), sums and subtractions between first order variables and numeric constants. E.g. $x = y$ is an abbreviation for $\forall X(x \in X \iff y \in X)$; $x = z - 2$ stands for $\exists y(\text{succ}(z, y) \wedge \text{succ}(y, x))$
- A sentence is a closed formula of the MSO logic; also, the notation $w \models \varphi$ denotes that string w satisfies sentence φ . For a given sentence φ , the language $L(\varphi)$ is defined as³

$$L(\varphi) = \{w \in \Sigma^+ \mid w \models \varphi\}.$$

³ When specifying languages by means of logic formulas, the empty string must be excluded because formulas refer to string positions.

For instance formula $\forall \mathbf{x}, \mathbf{y}(a(\mathbf{x}) \wedge \text{succ}(\mathbf{x}, \mathbf{y}) \Rightarrow b(\mathbf{y}))$ defines the language of strings where every occurrence of character a is immediately followed by an occurrence of b .

Büchi's seminal result is synthesized by the following theorem.

Theorem 1. *A language L is regular iff there exists a sentence φ in the above MSO logic such that $L = L(\varphi)$.*

The proof of the theorem is constructive, i.e., it provides an algorithmic procedure that, for a given FSA builds an equivalent sentence in the logic, and conversely; next we offer an intuitive explanation of the construction, referring the reader to, e.g., [46] for a complete and detailed proof.

From FSA to MSO logic

The key idea of the construction consists in representing each state q of the automaton as a second order variable X_q , which is the set of all string's positions where the machine is in state q . Without lack of generality we assume the automaton to be deterministic, and that $Q = \{0, 1, \dots, m\}$, with 0 initial, for some m . Then we encode the definition of the FSA recognizing L as the conjunction of several clauses each one representing a part of the FSA definition:

- The transition $\delta(q_i, a) = q_j$ is formalized by $\forall \mathbf{x}, \mathbf{y}(\mathbf{x} \in X_i \wedge a(\mathbf{x}) \wedge \text{succ}(\mathbf{x}, \mathbf{y}) \Rightarrow \mathbf{y} \in X_j)$.
- The fact that the machine starts in state 0 is represented as $\exists \mathbf{z}(\nexists \mathbf{x}(\text{succ}(\mathbf{x}, \mathbf{z})) \wedge \mathbf{z} \in X_0)$.
- Being the automaton deterministic, for each pair of distinct second order variables X_i and X_j we need the subformula $\nexists \mathbf{y}(\mathbf{y} \in X_i \wedge \mathbf{y} \in X_j)$.
- Acceptance by the automaton, i.e. $\delta(q_i, a) \in F$, is formalized by: $\exists \mathbf{y}(\nexists \mathbf{x}(\text{succ}(\mathbf{y}, \mathbf{x})) \wedge \mathbf{y} \in X_i \wedge a(\mathbf{y}))$.
- Finally the whole language L is the set of strings that satisfy the global sentence $\exists X_0, X_1, \dots, X_m(\varphi)$, where φ is the conjunction of all the above clauses.

At this point it is not difficult to show that the set of strings satisfying the above global formula is exactly L .

From MSO logic to FSA

The construction in the opposite sense has been proposed in various versions in the literature. Here we summarize its main steps along the lines of [46]. First, the MSO sentence is translated into a standard form using only second-order variables, the \subseteq predicate, and variables W_a , for each $a \in \Sigma$, denoting the set of all the positions of the word containing the character a . Moreover, we use Succ, which has the same meaning of succ, but, syntactically, has second order variable arguments that are singletons. This simpler, equivalent logic, is defined by the following syntax:

$$\varphi := X \subseteq W_a \mid X \subseteq Y \mid \text{Succ}(X, Y) \mid \neg\varphi \mid \varphi \vee \varphi \mid \exists X(\varphi).$$

As before, we also use the standard abbreviations, e.g. $\wedge, \vee, =$. To translate first order variables to second order variables we need to state that a (second order) variable is a singleton. Hence we introduce the abbreviation: $\text{Sing}(X)$ for $\exists Y(Y \subseteq X \wedge Y \neq X \wedge \nexists Z(Z \subseteq X \wedge Z \neq Y \wedge Z \neq X))$. Then, $\text{Succ}(X, Y)$ is defined only for X and Y singletons.

The following step entails the inductive construction of the equivalent automaton. This is built by associating a single automaton to each elementary subformula and by composing them according to the structure of the global formula. This inductive approach requires to use open formulas. Hence, we are going to consider words on the alphabet $\Sigma \times \{0, 1\}^k$, so that X_1, X_2, \dots, X_k are the free variables used in the formula; 1 in the, say, j -th component means that the considered position belongs to X_j , 0 vice versa. For instance, if $w = (a, 0, 1)(a, 0, 0)(b, 1, 0)$, then $w \models X_2 \subseteq W_a, w \models X_1 \subseteq W_b$, with X_1 and X_2 singletons.

Formula transformation

1. First order variables are translated in the following way: $\exists x(\varphi(x))$ becomes $\exists X(\text{Sing}(X) \wedge \varphi'(X))$, where φ' is the translation of φ , and X is a fresh variable.
2. Subformulas having the form $a(x), \text{succ}(x, y)$ are translated into $X \subseteq W_a, \text{Succ}(X, Y)$, respectively.
3. The other parts remain the same.

Inductive construction of the automaton We assume for simplicity that $\Sigma = \{a, b\}$, and that $k = 2$, i.e. two variables are used in the formula. Moreover we use the shortcut symbol \circ to mean all possible values.

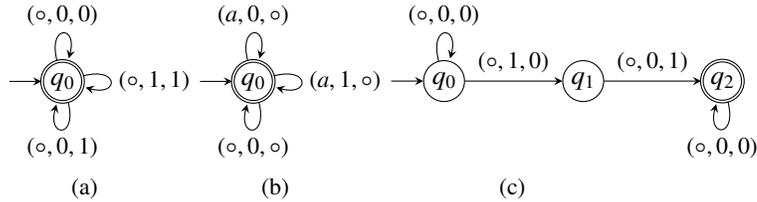


Fig. 6. Automata for the construction from MSO logic to FSA.

- The formula $X_1 \subseteq X_2$ is translated into an automaton that checks that there are 1's for the X_1 component only in positions where there are also 1's for the X_2 component (Figure 6 (a)).
- The formula $X_1 \subseteq W_a$ is analogous: the automaton checks that positions marked by 1 in the X_1 component must have symbol a (Figure 6 (b)).
- The formula $\text{Succ}(X_1, X_2)$ considers two singletons, and checks that the 1 for component X_1 is immediately followed by a 1 for component X_2 (Figure 6 (c)).
- Formulas inductively built with \neg and \vee are covered by the closure of regular languages w.r.t. complement and union, respectively.

- For \exists , we use the closure under alphabet projection: we start with an automaton with input alphabet, say, $\Sigma \times \{0, 1\}^2$, for the formula $\varphi(\mathbf{X}_1, \mathbf{X}_2)$; we need to define an automaton for the formula $\exists \mathbf{X}_1(\varphi(\mathbf{X}_1, \mathbf{X}_2))$. But in this case the alphabet is $\Sigma \times \{0, 1\}$, where the last component represents the only free remaining variable, i.e. \mathbf{X}_2 .

The automaton \mathcal{A}_\exists is built by starting from the one for $\varphi(\mathbf{X}_1, \mathbf{X}_2)$, and changing the transition labels from $(a, 0, 0)$ and $(a, 1, 0)$ to $(a, 0)$; $(a, 0, 1)$ and $(a, 1, 1)$ to $(a, 1)$, and those with b analogously. The main idea is that this last automaton nondeterministically “guesses” the quantified component (i.e. \mathbf{X}_1) when reading its input, and the resulting word $w \in (\Sigma \times \{0, 1\}^2)^*$ is such that $w \models \varphi(\mathbf{X}_1, \mathbf{X}_2)$. Thus, \mathcal{A}_\exists recognizes $\exists \mathbf{X}_1(\varphi(\mathbf{X}_1, \mathbf{X}_2))$.

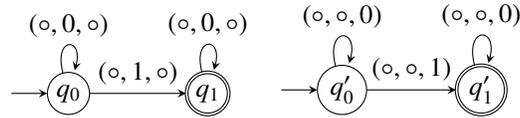
We refer the reader to the available literature for a full proof of equivalence between the logic formula and the constructed automaton. Here we illustrate the rationale of the above construction through the following example.

Example 6. Consider the language $L = \{a, b\}^* aa\{a, b\}^*$: it consists of the strings satisfying the formula:

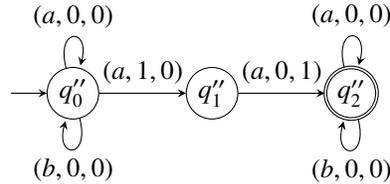
$$\varphi_L = \exists \mathbf{x} \exists \mathbf{y} (\text{succ}(\mathbf{x}, \mathbf{y}) \wedge a(\mathbf{x}) \wedge a(\mathbf{y})).$$

As seen before, first we translate this formula into a version using only second order variables: $\varphi'_L = \exists \mathbf{X}, \mathbf{Y} (\text{Sing}(\mathbf{X}) \wedge \text{Sing}(\mathbf{Y}) \wedge \text{Succ}(\mathbf{X}, \mathbf{Y}) \wedge \mathbf{X} \subseteq W_a \wedge \mathbf{Y} \subseteq W_a)$.

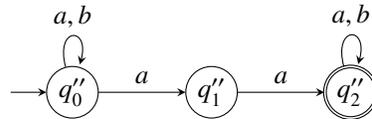
Here are the automata for $\text{Sing}(\mathbf{X})$ and $\text{Sing}(\mathbf{Y})$:



The above automata could also be obtained by expanding the definition of Sing and then projecting the quantified variables. By intersecting the automata for $\text{Sing}(\mathbf{X})$, $\text{Sing}(\mathbf{Y})$, $\text{Succ}(\mathbf{X}, \mathbf{Y})$ we obtain an automaton which is identical to the one we defined for translating formula $\text{Succ}(\mathbf{X}_1, \mathbf{X}_2)$, where here \mathbf{X} takes the role of \mathbf{X}_1 and \mathbf{Y} of \mathbf{X}_2 . Combining it with those for $\mathbf{X} \subseteq W_a$ and $\mathbf{Y} \subseteq W_a$ produces:



Finally, by projecting on the quantified variables \mathbf{X} and \mathbf{Y} we obtain the automaton for L .



The logical characterization of a class of languages, together with the decidability of the containment problem, is the main door towards automatic verification techniques: suppose in fact to have a logic formalism \mathfrak{L} recursively equivalent to an automaton family \mathfrak{A} ; then, one can use a language $L_{\mathfrak{L}}$ in \mathfrak{L} to specify the requirements of a given system and an abstract machine \mathcal{A} in \mathfrak{A} to implement the desired system: the correctness of the design defined by \mathcal{A} w.r.t. to the requirements stated by $L_{\mathfrak{L}}$ is therefore formalized by $L(\mathcal{A}) \subseteq L_{\mathfrak{L}}$, i.e., all behaviors realized by the machine are also satisfying the requirements. This is just the case with FSA and MSO logic for RL; in practice, however, the decision algorithms based on the above MSO logic have been proved of intractable complexity; thus, the great success of model-checking as the fundamental technique for automatic verification has been obtained by resorting to less powerful, typically temporal, logics which compensated the limited loss of generality by a complexity considered, and empirically verified in many practical experiments, as more tractable, despite a theoretical PSPACE-completeness [23].

3 Context-free Languages

Context-free languages (CFL), with their generating context-free grammars (CFG) and recognizing pushdown automata (PDA), are, together with RL, the most important chapter in the literature of formal languages. CFG have been introduced by Noam Chomsky in the 1950s as a formalism to capture the syntax of natural languages. Independently, essentially the same formalism has been developed to formalize the syntax of the first high level programming language, FORTRAN; in fact it is also referred to as Backus-Naur form (BNF) honoring the chief scientists of the team that developed FORTRAN and its first compiler. It is certainly no surprise that the same formalism has been exploited to describe the syntactic aspects of both natural and high level programming languages, since the latter ones have exactly the purpose to make algorithm specification not only machine executable but also similar to human description.

The distinguishing feature of both natural and programming languages is that complex sentences can be built by combining simpler ones in an a priori unlimited hierarchy: for instance a conditional sentence is the composition of a clause specifying a condition with one or two sub-sentences specifying what to do if the condition holds, and possibly what else to do if it does not hold. Such a typical nesting of sentences suggests a natural representation of their *structure* in the form of a tree shape. The possibility of giving a sentence a tree structure which hints at its semantics is a sharp departure from the rigid linear structure of regular languages. As an example, consider a simple arithmetic expression consisting of a sequence of operands with either a + or a * within any pair of them, as in $2 + 3 * 2 + 1 * 4$. Sentences of this type can be easily generated by, e.g., the following regular grammar:

$$\begin{aligned} S &\rightarrow 1 \mid 2 \mid \dots 0 \mid 1A \mid 2A \mid \dots 0A \\ A &\rightarrow +S \mid *S \end{aligned}$$

However, if we compute the value of the above expression by following the linear structure given to it by the grammar either by associating the sum and the multiply operations to the left or to the right we would obtain, respectively, 44 or 20 which is not the way we

learned to compute the value of the expression at school. On the contrary, we first compute the multiply operations and then the sum of the three partial results, thus obtaining 12; this, again, suggests to associate the semantics of the sentence—in this case the value of the expression—with a syntactic structure that is more appropriately represented by a tree, as suggested in Figure 7, than by a flat sequence of symbols.

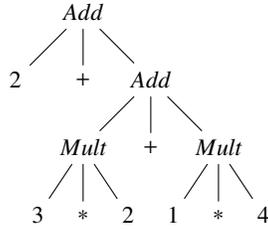


Fig. 7. A tree structure that shows the precedence of multiplication over addition in arithmetic expressions.

CFG are more powerful than regular ones in that not only they assign in a natural way a tree structure to their sentences but also the class of languages they generate strictly includes RL (see Section 3.2). CFG are defined as a special case of Chomsky’s general ones in the following way.

Definition 6. A grammar $G = (V_N, \Sigma, P, S)$ is context-free iff the lhs of all its productions is a single nonterminal.

Among all derivations associated with CFG it is often useful to distinguish the leftmost and rightmost ones:

Definition 7. A CFG derivation $S \xRightarrow{*}_G x$ is a leftmost, denoted as $S \xRightarrow{l}_G x$ (resp. rightmost, denoted as $S \xRightarrow{r}_G x$) one iff all of its immediate derivation steps are of the form $x A \beta \Rightarrow_G x \alpha \beta$ (resp. $\beta A x \Rightarrow_G \beta \alpha x$), with $x \in \Sigma^*$, $A \in V_N$, $\beta, \alpha \in V^*$.

In other words in a leftmost derivation, at every step the leftmost nonterminal character is the lhs of the applied rule. It is immediate to verify that for every derivation of a CFG there are an equivalent leftmost and an equivalent rightmost one, i.e., derivations that produce the same terminal string.⁴

Example 7. The following CF grammar GAE_1 generates the same numerical arithmetic expressions as above but assigns them the appropriate structure exemplified in Figure 7. In fact there is an obvious correspondence between grammar derivations and the trees whose root is the axiom and for every internal node, labeled by a nonterminal, its children are, in order, the rhs of a production whose lhs is the label of the father. In particular, a leftmost (resp. rightmost) derivation corresponds to a left (resp. right), top-down

⁴ This property does not hold for more general classes of grammars.

traversal of the tree; symmetrically a left bottom-up tree traversal produces the mirror image of a rightmost derivation. The tree that represents the derivation of a sentence x by a CFG is named *syntax-tree* or *syntactic tree* of x . Notice that in the following, for the sake of simplicity, in arithmetic expressions we will make use of a unique terminal symbol e to denote any numerical value.

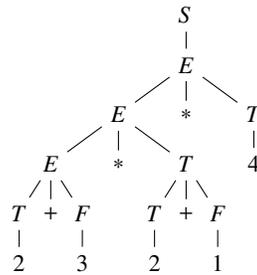
$$\begin{aligned}
 GAE_1 : S &\rightarrow E \mid T \mid F \\
 E &\rightarrow E + T \mid T * F \mid e \\
 T &\rightarrow T * F \mid e \\
 F &\rightarrow e
 \end{aligned}$$

It is an easy exercise augmenting the above grammar to let it generate more general arithmetic expressions including more arithmetic operators, parenthesized sub-expressions, symbolic operands besides numerical ones, etc.

Consider now the following slight modification GAE_2 of GAE_1 :

$$\begin{aligned}
 GAE_2 : S &\rightarrow E \mid T \mid F \\
 E &\rightarrow E * T \mid T + F \mid e \\
 T &\rightarrow T + F \mid e \\
 F &\rightarrow e
 \end{aligned}$$

GAE_1 and GAE_2 are equivalent in that they generate the same language; however, GAE_2 assigns to the string $2+3*2+1*4$ the tree represented in Figure 8: if we executed the operations of the string in the order suggested by the tree –first the lower ones so that their result is used as an operand for the higher ones– then we would obtain the value 60 which is not the “right one” 12.



The above examples, inspired by the intuition of arithmetic expressions, further emphasize the strong connection between the structure given to a sentence by the syntax-tree and the semantics of the sentence: as it is the case with natural languages too, an ambiguous sentence may exhibit different meanings.

The examples also show that, whereas a sentence of a regular language has a fixed –either right or left-linear structure, CFL sentences have a tree structure which in general, is not immediately “visible” by looking at the sentence, which is the frontier of its associated tree. As a consequence, in the case of CFL analyzing a sentence is often not only a matter of deciding whether it belongs to a given language, but in many cases it is also necessary to build the syntax-tree(s) associated with it: its structure often drives the evaluation of its semantics as intuitively shown in the examples of arithmetic expressions and systematically applied in all major applications of CFL, including designing their compilers. The activity of jointly accepting or rejecting a sentence and, in case of acceptance, building its associated syntax-tree(s) is usually called *parsing*.

3.1 Pushdown automata as parsers

It is well known that a typical data structure to build and visit trees is the *stack* or *pushdown store*. Thus, it is no surprise that the abstract machines adopted in formal language theory to analyze CFL are *pushdown automata*, i.e., devices supplied with a finite state control augmented with a LIFO auxiliary memory. Among the many versions of such automata proposed in the literature to parse CFL or subfamilies thereof we report here a classical one.

Definition 8. A *pushdown automaton (PDA)* \mathcal{A} is a tuple $(\Sigma, Q, \Gamma, \delta, q_0, Z_0, F)$ where

- Σ is the finite terminal alphabet;
- Q is the finite set of states;
- Γ is the finite stack alphabet;
- $\delta \subseteq_F Q \times \Gamma \times (\Sigma \cup \{\varepsilon\}) \times Q \times \Gamma^*$ is the transition relation, where the notation $A \subseteq_F B$ means that A is a finite subset of B ;
- $I \subseteq Q$ is the set of initial states;
- $Z_0 \in \Gamma$ is the initial stack symbol;
- $F \subseteq Q$ is the set of final or accepting states.

A PDA \mathcal{A} is *deterministic (DPDA)* iff

- I is a singleton;
- δ is a function $\delta : Q \times \Gamma \times (\Sigma \cup \{\varepsilon\}) \rightarrow Q \times \Gamma^*$;
- $\forall (q, A)$ if $\delta(q, A, a)$ is defined for some a , then $\delta(q, A, \varepsilon)$ is not defined; in other words, for any pair (q, A) the automaton may perform either an ε -move or a reading move but not both of them.

A *configuration* of a PDA is a triple (x, q, γ) where x is the string to be analyzed, q is the current state of the automaton, and γ is the content of the stack.

The transition relation between two configurations $c_1 \mapsto c_2$ holds iff $c_1 = (x_1, q_1, \gamma_1)$, $c_2 = (x_2, q_2, \gamma_2)$ and, either:

- $x_1 = ay, x_2 = y, \gamma_1 = A\beta, \gamma_2 = \alpha\beta$, and $(q_1, A, a, q_2, \alpha) \in \delta$, or
- $x_1 = x_2, \gamma_1 = A\beta, \gamma_2 = \alpha\beta$, and $(q_1, A, \varepsilon, q_2, \alpha) \in \delta$.

A string x is accepted by a PDA \mathcal{A} iff $c_0 = (x, q_0, Z_0) \xrightarrow{*} (\varepsilon, q_F, \gamma)$ for some $q_0 \in I$, $q_F \in F$, and $\gamma \in \Gamma^*$. The language accepted or recognized by \mathcal{A} is the set $L(\mathcal{A})$ of all strings accepted by \mathcal{A} .

A transition of the second type above, i.e. where no input character is consumed because the applied transition relation is of type $(q_1, A, \varepsilon, q_2, \alpha)$ is called an ε -move.

PDA as defined above are recursively equivalent to CFG, i.e. for any given CFG G a PDA \mathcal{A} can effectively be built such that $L(\mathcal{A}) = L(G)$ and conversely. DPDA define deterministic CFL (DCFL). Example 8 illustrates informally the rationale of this correspondence by presenting a PDA equivalent to grammar GAE_1 . It also provides a hint on how a PDA recognizing a CFL can be augmented as a pushdown transducer which also builds the syntax-tree of the input sentence (if accepted).

Example 8. Consider the following PDA \mathcal{A} :

$$\begin{aligned} Q &= \{q_0, q_1, q_F\}, \Gamma = \{E, T, F, +, *, e, Z_0\}, \\ \delta &= \{(q_0, Z_0, \varepsilon, q_1, EZ_0), (q_0, Z_0, \varepsilon, q_1, TZ_0), (q_0, Z_0, \varepsilon, q_1, FZ_0)\} \cup \\ &\{(q_1, E, \varepsilon, q_1, E+T), (q_1, E, \varepsilon, q_1, T*F), (q_1, E, \varepsilon, q_1, e)\} \cup \\ &\{(q_1, T, \varepsilon, q_1, T*F), (q_1, T, \varepsilon, q_1, e)\} \cup \{(q_1, F, \varepsilon, q_1, e)\} \cup \\ &\{(q_1, +, +, q_1, \varepsilon), (q_1, *, *, q_1, \varepsilon), (q_1, e, e, q_1, \varepsilon)\} \cup \\ &\{(q_1, Z_0, \varepsilon, q_F, \varepsilon)\}. \end{aligned}$$

\mathcal{A} accepts, e.g., the input string $e + e * e$ through the following sequence of configuration transitions

$$\begin{aligned} (e + e * e, q_0, Z_0) &\vdash (e + e * e, q_1, EZ_0) \vdash (e + e * e, q_1, E + TZ_0) \vdash (e + e * e, q_1, e + \\ &TZ_0) \vdash (+ e * e, q_1, +TZ_0) \vdash (e * e, q_1, TZ_0) \vdash (e * e, q_1, T * FZ_0) \vdash (e * e, q_1, e * \\ &FZ_0) \vdash (* e, q_1, *FZ_0) \vdash (e, q_1, FZ_0) \vdash (e, q_1, eZ_0) \vdash (\varepsilon, q_1, Z_0) \vdash (\varepsilon, q_F, \varepsilon) \end{aligned}$$

In fact \mathcal{A} has been naturally derived from the original grammar, GAE_1 in this case, in such a way that its transitions perfectly parallel a grammar's leftmost derivation: at any step if the top of the stack stores a nonterminal it is nondeterministically replaced by the corresponding rhs of a production of which it is the lhs, in such a way that the leftmost character of the rhs is put on top of the stack; if it is a terminal it is compared with the current input symbol and, if they match the stack is popped and the reading head advances. Notice however, that, being \mathcal{A} nondeterministic, there are also several sequences of transitions beginning with the same c_0 that do not accept the input.

Strictly speaking the above automaton is not a real *parser* for the corresponding grammar since it simply accepts all and only the strings generated thereby. However, it is quite easy to make it a complete parser by adding a few translation components: precisely, a *pushdown transducer* is a pushdown automaton provided with an output alphabet O which augments the range of δ in such a way that at every transition the transducer appends to the output, initially empty, a string in O^* ; the output string (strings in case of nondeterminism) produced at the end of the computation applied to the input x is the translation $\tau(x)$ produced by the transducer.

For instance, it is not difficult to build a full parser for grammar GAE_1 by adding to the automaton of Example 8 an output alphabet consisting of labels that uniquely identify G 's productions: for every transition of the automaton that has been built in

one-to-one correspondence with G 's production set it suffices to augment δ by adding as output component of each triple the label of the corresponding production. Assume that productions $E \rightarrow E + T \mid T * F$ are labeled 1, 2, respectively; then δ is enriched accordingly with $\{(q_0, E, \varepsilon, q_0, T + E, 1), (q_0, E, \varepsilon, q_0, F * T, 2), \dots\} \subseteq \delta$. Conversely, transitions that simply verify that the element on top of the stack coincides with the current input character do not output anything. Clearly, the transducer built in this way from the original grammar for any string of the language delivers an output that corresponds to a top-down, leftmost visit of the syntax-tree(s) of the input string where any internal node is the label of the production used by the grammar to expand it. Thus, we can assume (nondeterministic) transducers as general parsers for CFG; if the original grammar is ambiguous the corresponding transducer produces several translations for the same input.

3.2 A comparative analysis of CFL and RL properties

CFG and PDA are more powerful formalisms than the corresponding regular ones. Such a greater power allows for their application in larger fields such as formalizing the syntax and semantics of programming languages and many other tree-structured data. The greater power of the CFL family than the RL one abides not only in the structure of their sentences but even in sets of sentences that constitute their languages. For instance, the language $L_1 = \{a^n b^n \mid n \geq 0\}$ is easily recognized by a PDA but no FSA can accept it due its finite memory formalized by the state space Q , whereas, intuitively, an unbounded memory is necessary to count the number of a 's to be able to compare their number with that of the following b 's; this claim can be easily proved formally by exploiting the fundamental pumping Lemma 1.

Not surprisingly, however, such an increased power of CFL w.r.t. RL comes at a price in terms of loss of various properties and of increased complexity of related analysis algorithms. Herewith we briefly examine which properties of RL still hold for CFL and which ones are lost.

The original pumping lemma for regular languages can be naturally extended to a more general version holding for CFL:

Lemma 2. *Let G be a CFG; there is a constant k , recursively computable as a function of G , such that, for every string $x \in L(G)$, with $|x| > k$, there exists a factorization $x = yw_1zw_2u$ such that $x = yw_1^nz_2w_2^nu \in L(G)$ for every $n \geq 0$.*

A first important and natural consequence of the CF version of the pumping lemma is that the emptiness problem for CFL is still decidable; this can be easily shown in a parallel way as for regular languages. Another application of the lemma analogous to that used for RL can show that various languages, such as, e.g., $L_2 = \{a^n b^n c^n \mid n \geq 0\}$ are not CF.

On the opposite side the fundamental effective equivalence between deterministic and nondeterministic FSA does not hold for PDA. A formal proof of this statement can be found in most texts on formal language theory, e.g. in [30]; here we only offer an intuitive explanation based on a simple counterexample. Consider the language $L = \{ww^R \mid w \in \{a, b\}^*\}$, where w^R denotes the mirror image of w . It is quite easy to build a nondeterministic PDA that accepts L : initially it pushes the read symbols onto the

stack up to a point when it nondeterministically changes state and, from that point on it compares the symbol on top of the stack with the read one: if it reaches the end of the input and empties the stack at the same time, through a series of successful comparisons, then that particular computation accepts the input string; of course there will also be many other computations that fail due to a wrong guess about the middle of the input string. It is also *intuitively* apparent that there is no way to decide deterministically where is the middle of the input sentence unless a device other than a PDA is adopted, e.g., a Turing machine.

This substantial difference in recognition power between general PDA and DPDAs has several consequences on the properties of the class of languages they accept, resp. CFL and DCFL. Here we list the most important ones:

- CFL are closed w.r.t. union but not w.r.t. complement and intersection. In fact closure w.r.t. union is easily shown by building, for two languages L_1, L_2 two corresponding grammars G_1 and G_2 with disjoint nonterminal alphabets; then simply adding a new axiom S and two productions $S \rightarrow S_1 \mid S_2$, we obtain a grammar that generates all and only the strings in $L_1 \cup L_2$.
On the other hand consider again the above language L_2 , which is not a CFL: it is the intersection of $L_3 = \{a^n b^n c^* \mid n \geq 0\}$ with $L_4 = \{a^* b^n c^n \mid n \geq 0\}$ which are clearly both DCFL. By De Morgan's laws, CFL are not closed w.r.t. complement as well.
- DCFL are closed under complement but neither under union nor under intersection. Closure under complement can be obtained by exploiting determinism and switching F and $Q \setminus F$ in a similar way as done for FSA, provided that a suitable normal form of the original DPDA is built before applying the switch of the accepting states. We refer the reader to the more specialized literature, e.g., [30], for the technicalities of the necessary normal form. The previous counterexample also shows the non-closure of DCFL w.r.t. intersection and therefore w.r.t. union, again thanks to De Morgan's laws.
- CFL are closed under concatenation, Kleene *, reversal or mirror image, homomorphism, whereas DCFL are not closed under anyone of these operations. We justify only the claim about homomorphism since the other operations are of minor interest for the purpose of this paper. Given an alphabet homomorphism $h : \Sigma_1 \rightarrow \Sigma_2^*$ with its natural extension $h : \Sigma_1^* \rightarrow \Sigma_2^*$, it is immediate, for a grammar G , to build a G' such that $L(G') = h(L(G))$: it suffices to replace every occurrence of an element $a \in \Sigma_1$ in G 's definition by $h(a)$, including the particular case $h(a) = \varepsilon$. On the other side consider the language $L = \{wcw^R \mid w \in \{a, b\}^*\}$: it is immediate to build a DPDA recognizing it thanks to the new c character which marks the center of the input string, but it has already been shown that $h(L)$, with h defined as $h(a) = a, h(b) = b, h(c) = \varepsilon$, cannot be recognized by any DPDA.
- The lack of the above closure properties for both CFL and DCFL, prevents their natural exploitation to obtain decidability of the inclusion property; in fact this problem has been proved to be undecidable in both cases [30].⁵

⁵ As side remark, the equivalence problem is undecidable for general CFL and has been proved to be decidable for DCFL after remaining open for a fairly long time [44].

3.3 Logic characterization

The lack of the above closure properties also hampers a natural extension of the logic characterization of regular languages to CFL: in particular the classic inductive construction outlined in Section 2.2 strongly hinges on the correspondence between logical connectives and Boolean operations on sub-languages. Furthermore, the linear structure of RL allows any move of a FSA to depend only on the current state which is associated with a position x of the input string and on the input symbol located at position $x + 1$; on the contrary the typical nested structure of CFL sentences imposes that the move of the PDA may depend on information stored in the stack, which in turn may depend on information read from the input much earlier than the current move.

Despite these difficulties some interesting results concerning a logic characterization of CFL have been obtained. In particular it is worth mentioning the characterization proposed in [35]. The key idea is to enrich the syntax of the second order logic with a *matching relation symbol* M which takes as arguments two string position symbols x and y : a matching relation interpreting M must satisfy the following axioms:

- $M(x, y) \Rightarrow x < y$: y always follows x ;
- $M(x, y) \Rightarrow \nexists z(z \neq x \wedge z \neq y \wedge (M(x, z) \vee M(z, y) \vee M(z, x) \vee M(y, z)))$: M is one-to-one;
- $\forall x, y, z, w((M(x, y) \wedge M(z, w) \wedge x < z < y) \Rightarrow x < w < y)$: M is *nested*, i.e., if we represent graphically $M(x, y)$ as an edge from x to y such edges cannot cross.

The matching relation is then used to represent the tree structure(s) associated with a CF language sentence: for instance consider the (ambiguous) grammar G_{amb}

$$\begin{aligned} G_{amb} : S &\rightarrow A_1 \mid A_2 \\ A_1 &\rightarrow aaA_1bb \mid aabb \\ A_2 &\rightarrow aA_3b \\ A_3 &\rightarrow aA_2b \mid ab \end{aligned}$$

G_{amb} induces a natural matching relation between the positions of characters in its strings. For instance Figure 9 shows the two relations associated with the string $aaaabbbb$.

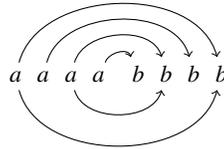


Fig. 9. Two matching relations for $aaaabbbb$, one is depicted above and the other below the string.

More generally we could state that for a grammar G and a sentence $x = a_1a_2\dots a_n \in L(G)$ with $\forall k, a_k \in \Sigma$, $(i, j) \in M$ iff $S \xrightarrow{*}_G a_1a_2\dots a_{i-1}Aa_{j+1}\dots a_n \xrightarrow{*}_G a_1a_2\dots a_n$. It is im-

mediate to verify, however, that such a definition in general does not guarantee the above axioms for the matching relation: think e.g., to the previous grammar GAE_1 which generates arithmetic expressions. For this reason [35] adopts the so-called *double Greibach normal form* (DGNF) which is an effective but non-trivial transformation of a generic CFG into an equivalent one where the first and last characters of any production are terminal.⁶ It is now immediate to verify that a grammar in DGNF does produce a matching relation on its strings that satisfies all of its axioms.

Thanks to the new relation and the corresponding symbol [35] defines a second order logic $CF\mathcal{L}$ that characterizes CFL: the sentences of $CF\mathcal{L}$ are first-order formulas prefixed by a single existential quantifier applied to the second order variable M representing the matching relation. Thus, intuitively, a $CF\mathcal{L}$ sentence such as $\exists M(\phi)$ where M occurs free in ϕ and ϕ has no free occurrences of first-order variables is satisfied iff there is a structure defined by a suitable matching relation such that the positions that satisfy the occurrences of M in ϕ also satisfy the whole ϕ . For instance the sentence

$$\exists M, z \left(\begin{array}{l} \left(\begin{array}{l} \exists x(\text{succ}(z, x) \wedge M(0, z) \wedge \\ \forall x, y(M(x, y) \Rightarrow a(x) \wedge b(y)) \wedge \\ \exists y \forall x \left(\begin{array}{l} (0 \leq x < y \Rightarrow a(x)) \wedge \\ (x \geq y \geq z \Rightarrow b(x)) \end{array} \right) \wedge \end{array} \right) \wedge \\ \left(\begin{array}{l} \forall x, y \left(M(x, y) \Rightarrow \begin{array}{l} (x > 0 \Rightarrow M(x-1, y+1)) \wedge \\ (x < y-2 \Rightarrow M(x+1, y-1)) \end{array} \right) \right) \wedge \\ \vee \\ \left(\begin{array}{l} \forall x, y \left(M(x, y) \Rightarrow \begin{array}{l} (x > 1 \Rightarrow M(x-2, y+2)) \wedge \\ \neg M(x-1, y+1)) \wedge \\ (x < y-4 \Rightarrow M(x+2, y-2)) \wedge \\ \neg M(x+1, y-1) \end{array} \right) \right) \end{array} \right) \end{array} \right)$$

is satisfied by all and only the strings of $L(G_{amb})$ with both the M relations depicted in Figure 9.

After having defined the above logic, [35] proved its equivalence with CFL in a fairly natural way but with a few non-trivial technicalities: with an oversimplification, from a given CFG in DGNF a corresponding logic formula is built inductively in such a way that $M(x, y)$ holds between the positions of leftmost and rightmost leaves of any subtree of a grammar's syntax-tree; conversely, from a given logic formula a tree-language, i.e., a set of trees, is built such that the frontiers of its trees are the sentences of a CFL. However, as the authors themselves admit, this result has a minor potential for practical applications due to the lack of closure under complementation. The need to resort to the DGNF puts severe constraints on the structures that can be associated with the strings, a priori excluding, e.g., linear structures typical of RL; nonetheless the introduction of the M relation opened the way for further important developments as we will show in the next sections.

3.4 Parsing

PDA are the natural abstract machines to recognize CFL as well as FSA are for RL; we have also seen that augmenting PDA with a suitable output device makes them push-

⁶ To be more precise, the normal form introduced in [35] is further specialized, but for our introductory purposes it is sufficient to consider any DGNF.

down transducers (PDT) which can be used as parsers, i.e. to build the syntax-tree(s) associated with a given string of the language. Unlike the case of RL, however, DPDA are not able to recognize all CFL. Thus, if we want to recognize or parse a generic CFL, in principle we must simulate all possible computations of a nondeterministic PDA (or PDT); this approach clearly raises a critical complexity issue. For instance, consider the analysis of any string in Σ^* by the PDA accepting $L = \{ww^R \mid w \in \{a,b\}^*\}$: in the worst case at each move the automaton “splits” its computation in two branches by guessing whether it reached the middle of the input sentence or not; in the first case the computation proceeds deterministically to verify that from that point on the input is the mirror image of what has been read so far, whereas the other branch of the computation proceeds still waiting for the middle of the string and splitting again and again at each step. Thus, the total number of different computations equals the length of the input string, say n , and each of them has in turn an $O(n)$ length; therefore, simulating the behavior of such a nondeterministic machine by means of a deterministic algorithm to check whether at least one of its possible computations accepts the input has an $O(n^2)$ total complexity.

The above example can be generalized in the following way: on the one hand we have

Statement 1 *Every CFL can be recognized in real-time, i.e. in a number of moves equal to the length of the input sentence, by a, possibly nondeterministic, PDA.*

One way to prove the statement is articulated in two steps:

1) First, an original CFG generating L is transformed into the *Greibach normal form (GNF)*⁷:

Definition 9. *A CFG is in Greibach normal form [30] iff the rhs of all its productions belongs to ΣV_N^* .*

The procedure given in [28] to transform any CF grammar into the normal form essentially is based on transforming any *left recursion*, i.e. a derivation such as $A \xRightarrow{*} A\alpha$ into an equivalent right one $B \xRightarrow{*} \alpha B$.

2) Starting from a grammar in GNF the procedure to build an equivalent PDA therefrom can be “optimized” by:

- restricting Γ to V_N only;
- when a symbol A is on top of the stack, a single move reads the next input symbol, say a , and replaces A in the stack with the string α of the rhs of a production $A \rightarrow a\alpha$, if any (otherwise the string is rejected).

Notice that such an automaton is real-time since there are no more ε -moves but, of course, it may still be nondeterministic. If the grammar in GNF is such that there are no two distinct productions of the type $A \rightarrow a\alpha$, $A \rightarrow a\beta$, then the automaton built in this way is a real-time DPDA that is able to build leftmost derivations of the grammar. In Section 5.1 we will go back to the issue of building deterministic parsers for CFL.

⁷ The previous DGNF is clearly a “symmetric variant” of the original GNF but is not due to the same author and serves different purposes.

Statement 1, on the other hand, leaves us with the natural question: “provided that purely nondeterministic machines are not physically available and at most can be approximated by parallel machines which however cannot exhibit an unbounded parallelism, how can we come up with some deterministic parsing algorithm for general CFL and which complexity can exhibit such an algorithm?”. In general it is well-known that simulating a nondeterministic device with complexity⁸ $f(n)$ by means of a deterministic algorithm may expose to the risk of even an exponential complexity $O(2^{f(n)})$. For this reason on the one hand many applications, e.g., compiler construction, have restricted their interest to the subfamily of DCFL; on the other hand intense research has been devoted to design efficient deterministic parsing algorithms for general CFL by departing from the approach of simulating PDA. The case of parsing DCFL will be examined in Section 5; parsing general CFL, instead is of minor interest in this paper, thus we simply mention the two most famous and efficient of such algorithms, namely the one due to Cocke, Kasami, and Younger, usually referred to as CKY and described, e.g., in [30], and the one by Earley reported in [14]; they both have an $O(n^3)$ time complexity.⁹

To summarize, in this section we have seen that CFL are considerably more general than RL but, on the one side they require parsing algorithms to assign a given sentence an appropriate (tree-shaped and usually not immediately visible) structure and, on the other side, they lose several closure and decidability properties typical of the simpler RL. Not surprisingly, therefore, much, largely unrelated, research has been devoted to face both such challenges; in both cases major successes have been obtained by introducing suitable subclasses of the general language family: on the one hand parsing can be accomplished for DCFL much more efficiently than for nondeterministic ones (furthermore in many cases, such as e.g. for programming languages, nondeterminism and even ambiguity are features that are better to avoid than to pursue); on the other hand various subclasses of CFL have been defined that retain some or most of the properties of RL yet increasing their generative power. In the next sections we resume, in the reverse order, the major results obtained on both sides, so far within different research areas and by means of different subfamilies of CFL. As anticipated in the introduction, however, we will see in the following sections that one of such families allows for major improvements on both sides.

4 Structured context-free languages

RL sentences have a fixed, right or left, linear structure; CFL sentences have a more general tree-structure, of which the linear one is a particular case, which normally is not immediately visible in the sentence and, in case of ambiguity, it may even happen that the same sentence has several structures. R. McNaughton, in his seminal paper [39], was probably the first one to have the intuition that, if we “embed the sentence structure in the sentence itself” in some sense making it visible from the frontier of the

⁸ As usual we assume as the complexity of a nondeterministic machine the length of the shortest computation that accepts the input or of the longest one if the string is rejected.

⁹ We also mention (from [30]) that in the literature there exists a variation of the CKY algorithm due to Valiant that is completely based on matrix multiplication and therefore has the same asymptotic complexity of this basic problem, at the present state of the art $O(n^{2.37})$.

syntax-tree (as it happens implicitly with RL since their structure is fixed a priori), then many important properties of RL still hold for such special class of CFL.

Informally, we name such CFL *structured* or *visible structure* CFL. The first formal definition of this class given by McNaughton is that of *parenthesis languages*, where each subtree of the syntax-tree has a frontier embraced by a pair of parentheses; perhaps the most widely known case of such a language, at least in the field of programming languages is the case of LISP. Subsequently several other equivalent or similar definitions of structured languages, have been proposed in the literature; not surprisingly, an important role in this field is played by *tree-languages* and their related recognizing machines, *tree-automata*. Next we browse through a selection of such language families and their properties, starting from the seminal one by McNaughton.

4.1 Parenthesis grammars and languages

Definition 10. For a given terminal alphabet Σ let $[,]$ be two symbols $\notin \Sigma$. A parenthesis grammar (PG) with alphabet $\Sigma \cup \{[,]\}$ is a CFG whose productions are of the type $A \rightarrow [\alpha]$, with $\alpha \in V^*$.

It is immediate to build a parenthesis grammar naturally associated with any CFG: for instance the following PG is derived from the GAE_1 of Example 7:

$$\begin{aligned} GAE_{[]} : S &\rightarrow [E] \mid [T] \mid [F] \\ E &\rightarrow [E + T] \mid [T * F] \mid [e] \\ T &\rightarrow [T * F] \mid [e] \\ F &\rightarrow [e] \end{aligned}$$

It is also immediate to realize that, whereas GAE_1 does not make immediately visible in the sentence $e + e * e$ that $e * e$ is the frontier of a subtree of the whole syntax-tree, $GAE_{[]}$ generates $[[[e] + [[e] * [e]]]]$ (but not $[[[[e] + [e]] * [e]]]]$), thus making the structure of the syntax-tree immediately visible in its parenthesized frontier.

Sometimes it is convenient, when building a PG from a normal one, to omit parentheses in the rhs of *renaming rules*, i.e., rules whose rhs reduces to a single nonterminal, since such rules clearly do not significantly affect the shape of the syntax-tree. In the above example such a convention would avoid the useless outermost pair of parentheses.

Given that the trees associated with sentences generated by PG are isomorphic to the sentences, the parsing problem for such languages disappears and scales down to a simpler recognition problem as it happens for RL. Thus, rather than using the full power of general PDA for such a job, *tree-automata* have been proposed in the literature as a recognition device equivalent to PG as well as FSA are equivalent to regular grammars. Intuitively, a tree-automaton (TA) traverses either top-down or bottom-up a labeled tree to decide whether to accept it or not, thus it is a tree-recognizer or tree-acceptor.

4.2 Tree Automata

Here we give an introductory view of the tree automata formalism: our formal definition is inspired by the traditional ones, e.g., in [45,13] but with minor adaptations for

convenience and homogeneity with the notation adopted in this paper; to avoid possible terminological confusion with other literature on tree-automata we give a different name to such automata.

Definition 11. A stencil of a terminal alphabet Σ is a string in $(\Sigma \cup \{N\})^*$. A stencil or tree alphabet is a finite set of stencils. The number of occurrences of the symbol N in a stencil st is called the stencil's arity and is denoted $n(st)$.

Definition 12. A top-down stencil automaton is a tuple $(Q, ST_\Sigma, \Delta, I)$, where

- Q is a finite set of states
- ST_Σ is a finite set of stencils of Σ ; $ST_{\Sigma,n}$ denotes the subset of ST_Σ of the stencils with arity n
- Δ is a collection of relations $\{\delta_n \subseteq Q \times ST_{(\Sigma,n)} \times Q^n\}$
- $I \subseteq Q$ is the set of initial states.

A top-down stencil automaton is deterministic iff

- For each n , δ_n is a function $\delta_n : Q \times ST_{(\Sigma,n)} \rightarrow Q^n$
- I is a singleton

Definition 13. A tree on a given pair (ST_Σ, Q) is a tree whose internal nodes are labeled by elements of Q , leaves are labeled by elements of Σ , and the homomorphism $h(q) = N$, with $q \in Q$ projects all strings of children of any node of the tree into an element of ST_Σ

A tree on a pair (ST_Σ, Q) is accepted by a top-down stencil automaton $\mathcal{A} = (Q, ST_\Sigma, \Delta, I)$ iff

- the root is labeled by an element $\in I$
- each pair (q_f, α) where q_f is the label of a father node, α is the string of its children in $(\Sigma \cup Q)^*$ is such that $(q_f, h(\alpha), q_1, \dots, q_n) \in \delta_n$, where n is the arity of $h(\alpha)$ and q_1, \dots, q_n are the states occurring in α , in the same order.

It is immediate to state a one-to-one, effective correspondence between parenthesis grammars and top-down stencil automata; in doing so it may be convenient to preliminarily eliminate from the grammar renaming rules [30] which do not have a real meaning in terms of sentence structure; such a normal form also requires that, instead of having just one axiom with possible renaming rules with it as lhs, S is a subset of V_N and there is a one-to-one correspondence between S and I .

Example 9. The following stencil automaton $\mathcal{A} = (Q, ST_\Sigma, \Delta, I)$ recognizes the syntax-trees of the grammar GAE_1 whose structure, up to the labeling of the internal nodes, is also described by the sentences of the parenthesis grammar $GAE_{[]}$.

- $Q = \{E, T, F\}$
- $ST_\Sigma = \{e, N + N, N * N\}$
- $\Delta = \{(E, e), (T, e), (F, e), (T, N * N, T, F), (E, N * N, T, F), (E, N + N, E, T)\}$
- $I = \{E, T, F\}$

If we compare the automata of Examples 8 and 9 we notice that both of them, resp., build and recognize the syntax-trees generated by grammar GAE_1 in top-down order; precisely, the computations of the automaton of Example 8 follow a leftmost derivation of the grammar, whereas the automaton of Example 9 simply states a partial order in the application of the δ relation going from father nodes to children. In both cases, however, it may be convenient to use automata that work in a symmetric way, i.e., resp. build or traverse the syntax-trees bottom up. Two important cases of bottom up parsers are described in Sections 5.2 and 6.2; here instead we introduce *bottom-up stencil automata* and compare them with the top-down ones.

Definition 14. A bottom-up stencil automaton is a tuple $(Q, ST_\Sigma, \Delta, F)$, where Q and ST_Σ are as in Definition 12,

- Δ is a collection of relations $\{\delta_n \subseteq Q^n \times ST_{(\Sigma, n)} \times Q\}$
- $F \subseteq Q$ is the subset of Q 's final states.

The automaton is deterministic iff for each n , δ_n is a function $\delta_n : Q^n \times ST_{(\Sigma, n)} \rightarrow Q$.

A tree on a pair (ST_Σ, Q) is accepted by a bottom-up stencil automaton $\mathcal{A} = (Q, ST_\Sigma, \Delta, F)$ iff

- each pair (α, q_f) where q_f is the label of a father node, α is the string of its children in $(\Sigma \cup Q)^*$ is such that $(q_1, \dots, q_n, h(\alpha), q_f) \in \delta_n$, where n is the arity of $h(\alpha)$ and q_1, \dots, q_n are the states occurring in α , in the same order;
- the root is labeled by an element $\in F$.

Structured CFL, of which parenthesis languages are a first major example, enjoy some properties that hold for RL but are lost by general CFL. In particular, we mention their closure under Boolean operations, which has several major benefits, including the decidability of the containment problem. The key milestones that allowed McNaughton to prove this closure are:

- The possibility to apply all operations within a *structure universe*, i.e., to a universe of syntax-trees rather than to the “flat” Σ^* ; in the case of parenthesis languages the natural universe associated with a given set of stencils is generated by the corresponding *stencil grammar*, i.e. the grammar with the only nonterminal N and one production for every stencil.
- The possibility of building a normal form of any parenthesis grammar *with no repeated rhs*, i.e., such that for every rhs there is only one production rewriting it. Such a construction has been defined by McNaughton by referring directly to the grammars but can be explained even more intuitively by noticing that, if we transform a bottom-up stencil automaton into a parenthesis grammar and conversely, the grammar has no repeated rhs iff the automaton is deterministic. In both cases the procedure to obtain such a normal form is a natural extension of the basic one applied to FSA and is based on defining a new set of states (resp. nonterminals) that is the power set of the original one; for instance, suppose that a grammar contains two rules $A \rightarrow [ab]$ and $B \rightarrow [ab]$; then they are “collapsed” into the unique one $AB \rightarrow [ab]$; then the procedure is iterated by replacing both A and B by the new nonterminal AB in every rhs where they occur and so on until no new nonterminal is

generated. It is then an easy exercise to prove the equivalence between the original grammar (resp. automaton) and the new one by means of a natural induction.

It is important to emphasize, however, that such a way to *determinize* a nondeterministic bottom-up stencil automaton does not work in the case of the top-down version of these automata: deterministic top-down tree (or stencil) automata in fact have been proven to be less powerful than their bottom-up counterpart. We refer the reader to a more specific literature, e.g. [13], for a proof of this claim; here we simply anticipate that something similar occurs for the parsing of DCFL.

On the basis of these two fundamental properties deriving effective closure w.r.t. Boolean operations within a given universe of stencils is a natural extension of the operations already outlined for RL (notice that RL are a special case of structured languages whose stencils are only linear, i.e., of the type aN, a (or Na, a) for $a \in \Sigma$).

- The complement language w.r.t. the universe of syntax-trees associated with a given set of stencils is obtained by
 - Adding a new conventional state q_{err} to a deterministic bottom up stencil automaton and completing each δ_n with it on the whole domain $Q^n \times ST_{(\Sigma, n)}$ in the same way as for RL.
 - Complementing the set of accepting states.
- The intersection between two languages sharing the same set of stencils is obtained by building a new set of states that is the cartesian product of the original ones and extending the δ function in the usual way.

As usual, an immediate corollary of these closure properties is the decidability of the containment problem for two languages belonging to the same universe.

We just mention another important result that is obtained as an extension of its version that holds for RL, namely, the possibility of minimizing the grammar generating a parenthesis language (resp. the automaton recognizing a tree language), w.r.t. the number of necessary nonterminals (resp. states); this type of results, however, is not in the scope of this paper; thus we refer the interested reader to the original paper by McNaughton or other subsequent literature.

On the basis of the important results obtained by McNaughton in his seminal paper, many other families of CFL have been defined in several decades of literature with the general goal of extending (at least some of the) closure and decidability properties, and logic characterizations that make RL such a “nice” class despite its limits in terms of generative power. Most of such families maintain the key property of being “structured” in some generalized sense w.r.t. parenthesis languages. In the following section we introduce so called *input-driven languages*, also known as *visibly pushdown languages* which received much attention in recent literature and exhibit a fairly complete set of properties imported from RL.

4.3 Input-driven or visibly pushdown languages

The concept of *input-driven CF language* has been introduced in [40] in the context of building efficient recognition algorithms for DCFL: according to [40] a DPDA is input-driven if the decision of the automaton whether to execute a *push move*, i.e. a move

where a new symbol is stored on top of the stack, or a *pop move*, i.e. a move where the symbol on top of the stack is removed therefrom, or a move where the symbol on top of the stack is only updated, depends exclusively on the current input symbol rather than on the current state and the symbol on top of the stack as in the general case. Later, several equivalent definitions of the same type of pushdown automata, whether deterministic or not, have been proposed in the literature; among them, here we choose an early one proposed in [4], which better fits with the notation adopted in this paper than the later one in [5].

Definition 15. *Let the input alphabet Σ be partitioned into three disjoint alphabets, $\Sigma = \Sigma_c \cup \Sigma_r \cup \Sigma_i$, named, respectively, the call alphabet, return alphabet, internal alphabet. A visibly pushdown automaton (VPA) over $(\Sigma_c, \Sigma_r, \Sigma_i)$ is a tuple $(Q, I, \Gamma, Z_0, \delta, F)$, where*

- Q is a finite set of states;
- $I \subseteq Q$ is the set of initial states;
- $F \subseteq Q$ is the set of final or accepting states;
- Γ is the finite stack alphabet;
- $Z_0 \in \Gamma$ is the special bottom stack symbol;
- δ is the transition relation, partitioned into three disjoint subrelations:
 - Call move: $\delta_c \subseteq Q \times \Sigma_c \times Q \times (\Gamma \setminus \{Z_0\})$,
 - Return move: $\delta_r \subseteq Q \times \Sigma_r \times \Gamma \times Q$,
 - Internal move: $\delta_i \subseteq Q \times \Sigma_i \times Q$.

A VPA is deterministic iff I is a singleton, and $\delta_c, \delta_r, \delta_i$ are functions:
 $\delta_c : Q \times \Sigma_c \rightarrow Q \times (\Gamma \setminus \{Z_0\})$, $\delta_r : Q \times \Sigma_r \times \Gamma \rightarrow Q$, $\delta_i : Q \times \Sigma_i \rightarrow Q$.

The automaton configuration, the transition relation between two configurations, the acceptance of an input string, and the language recognized by the automaton are then defined in the usual way: for instance if the automaton reads a symbol a in Σ_c while is in the state q and has C on top of the stack, it pushes onto the stack a new symbol D and moves to state q' provided that (q, a, q', D) belongs to δ_c ; notice that in this way the special symbol Z_0 can occur only at the bottom of the stack, during the computation. A language over an alphabet $\Sigma = \Sigma_c \cup \Sigma_r \cup \Sigma_i$ recognized by some VPA is called a *visibly pushdown language (VPL)*.

The following remarks help put IDL alias VPL in perspective with other families of CFL.

- Once PDA are defined in a standard form, with respect to the general one given in Definition 8 where their moves either push a new symbol onto the stack or remove it therefrom or leave the stack unaffected, the two definitions of IDL and VPL are equivalent. Both names for this class of languages are adequate: on the one side, the attribute input-driven emphasizes that the automaton move is determined exclusively on the basis of the current input symbol¹⁰; on the other side we can consider VPL as *structured languages* since the structure of their sentences is immediately *visible* thanks to the partitioning of Σ .

¹⁰ We will see, however, that the same term can be interpreted in a more general way leading to larger classes of languages.

- VPL generalize McNaughton’s parenthesis languages: open parentheses are a special case of Σ_c and closed ones of Σ_r ; further generality is obtained by the possibility of performing an unbounded number of internal moves, actually “purely finite state” moves between two matching call and return moves and by the fact that VPA can accept unmatched return symbols at the beginning of a sentence as well as unmatched calls at its end; a motivation for the introduction of such a generality is offered by the need of modeling systems where a computation containing a sequence of procedure calls is suddenly interrupted by a special event such as an exception or an interrupt.
- Being VPL essentially structured languages, their corresponding automata are just recognizers rather than real parsers.
- VPL are *real-time languages*; in fact VPA read one input symbol for each move. We have mentioned that this property can be obtained for nondeterministic PDA recognizing any CFL but not for deterministic ones.

VPL have obtained much attention since they represent a major step in the path aimed at extending many, if not all, of the important properties of RL to structured CFL. They are closed w.r.t. all major language operations, namely the Boolean ones, concatenation, Kleene * and others; this also implies the decidability of the containment problem, which, together with the characterization in terms of a MSO logic, again extending the result originally stated for RL, opens the door to applications in the field of automatic verification.

A key step to achieve such important results is the possibility of effectively determinizing nondeterministic VPA. The basic idea is similar to the classic one that works for RL, i.e. to replace the uncertainty on the current state of a nondeterministic automaton with the subset of Q containing all such possible states (see Example 3). Unlike the case of FSA however, when the automaton executes a return move it is necessary to “match” the current state with the one that was entered at the time of the corresponding call; to do so the key idea is to “pair” the set of states nondeterministically reached at the time of a return move with those that were entered at the time of the corresponding call; intuitively, the latter ones are memorized and propagated through the stack, whose alphabet is enriched with pairs of set states. As a result at the moment of the return it is possible to check whether some of the states memorized at the time of the call “match” with some of the states that can be currently held by the nondeterministic original automaton.

We do not go into the technical details of this construction, referring the reader to [4] for them; we just mention that, unlike the case of RL, the price to be paid in terms of size of the automaton to obtain a deterministic version from a nondeterministic one is $2^{O(s^2)}$, where s is the size of the original state set. In [5] the authors also proved that such a gap is unavoidable since there are VPL that are recognized by a nondeterministic VPA with a state set of cardinality s but are not recognized by any deterministic VPA with less than 2^{s^2} states. In Section 6.1 we will provide a similar proof of determinization for a more general class of automata.

Once a procedure to obtain a deterministic VPA from a nondeterministic one is available, closure w.r.t. Boolean operations follows quite naturally through the usual path already adopted for RL and parenthesis languages. Closure under other major lan-

guage operations such as concatenation and Kleene * is also obtained without major difficulties but we do not report on it since those operations are not of major interest for this paper. Rather, we wish here to go back to the issue of logical characterization.

The logic characterization of visibly pushdown languages We have seen in Section 3.3 that attempts to provide a logic characterization of general CFL produced only partial results due to the lack of closure properties and to the fact that CFL do not have an a priori fixed structure; in fact the characterization offered by [35] and reported here requires an existential quantification w.r.t. relation M that represents the structure of a string. Resorting to structured languages such as VPL instead allowed for a fairly natural generalization of the classical Büchi's result for RL.

The key steps to obtain this goal are [5]:

- Using the same relation M introduced in [35],¹¹ adding its symbol as a basic predicate to the MSO logic's syntax given in Section 2.2 for RL, and extending its interpretation in the natural way. This turns out to be simpler and more effective in the case of structured languages since, being such languages a priori unambiguous (the structure of the sentence is the sentence itself), there is only one such relation among the string positions and therefore there is no need to quantify it. Furthermore the relation is obviously one-to-one with a harmless exception due to the existence of unmatched closed parentheses at the beginning of the sentence and unmatched open ones at the end: in such cases conventional relations $M(-\infty, x)$, $M(x, +\infty)$ are stated.

- Repeating exactly the same path used for RL both in the construction of an automaton from a logic formula and in the converse one. This only requires the managing of the new M relation in both constructions; precisely:

In the construction from the MSO formula to VPA, the elementary automaton associated with the atomic formula $M(X, Y)$, where X and Y are the usual singleton second order variables for any pair of first order variables x and y , is represented by the diagram of Figure 10 where, like in the same construction for RL, \circ stands for any value of Σ for which the transition can be defined according to the alphabet partitioning, so that the automaton is deterministic, the second component of the triple corresponds to X , and the third to Y .¹²

In the construction from the VPA to the MSO formula, besides variables X_i for encoding states, we also need variables to encode the stack. We introduce variables C_A and R_A , for $A \in \Gamma$, to encode, respectively, the set of positions in which a *call* pushes A onto the stack, and in which a *return* pops A from the stack.

The following formula states that every pair (x, y) in M must belong, respectively, to exactly one C_A and exactly one R_A :

$$\forall x, y (M(x, y) \Rightarrow \bigvee_{A \in \Gamma} x \in C_A \wedge y \in R_A) \wedge \\ \forall x \bigwedge_{A \in \Gamma} \left(\begin{array}{c} x \in C_A \Rightarrow \neg \bigvee_{B \neq A} x \in C_B \\ \wedge \\ x \in R_A \Rightarrow \neg \bigvee_{B \neq A} x \in R_B \end{array} \right).$$

¹¹ Renamed *nesting relation* and denoted as \sim or ν in later literature.

¹² We use here the following notation for depicting VPA: a label a/B stands for a push move, a, B for a pop move, and a alone for an internal move.

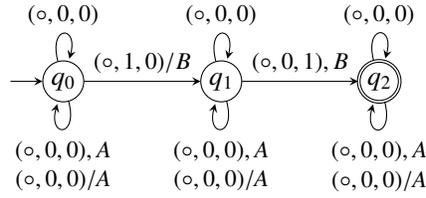


Fig. 10. VPA associated with $M(X, Y)$ atomic formula.

The following subformulas express the conditions for call and return transitions, respectively. They are added to the ones already described in Section 2.2 (without loss of generality, we assume that the original VPA is deterministic). Subformulas for internal transitions are almost identical to those for FSA.

$$\forall \mathbf{x}, \mathbf{y} \bigwedge_{0 \leq i \leq m} \bigwedge_{A \in \Gamma} \bigwedge_{a \in \Sigma} \left(\begin{array}{l} \mathbf{x} \in \mathbf{X}_i \wedge \text{succ}(\mathbf{x}, \mathbf{y}) \wedge \\ a(\mathbf{y}) \wedge \delta_c(i, a) = (j, A) \\ \Rightarrow \\ \mathbf{y} \in \mathbf{C}_A \wedge \mathbf{y} \in \mathbf{X}_j \end{array} \right)$$

$$\forall \mathbf{x}, \mathbf{y}, \mathbf{z} \bigwedge_{0 \leq i \leq m} \bigwedge_{A \in \Gamma} \bigwedge_{a \in \Sigma} \left(\begin{array}{l} \mathbf{y} \in \mathbf{X}_i \wedge \text{succ}(\mathbf{y}, \mathbf{z}) \wedge \\ M(\mathbf{x}, \mathbf{z}) \wedge \mathbf{z} \in \mathbf{R}_A \wedge a(\mathbf{z}) \\ \Rightarrow \\ \mathbf{z} \in \mathbf{X}_j \wedge j = \delta_r(i, a, A) \end{array} \right).$$

Example 10. Consider the alphabet $\Sigma = (\Sigma_c = \{a\}, \Sigma_r = \{b\}, \Sigma_i = \emptyset)$ and the VPA depicted in Figure 11.

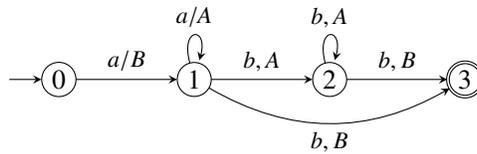


Fig. 11. A VPA recognizing $\{a^n b^n \mid n > 0\}$.

The MSO formula $\exists \mathbf{X}_0, \mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3, \mathbf{C}_A, \mathbf{C}_B, \mathbf{R}_A, \mathbf{R}_B (\varphi_{\mathcal{A}} \wedge \varphi_M)$ is built on the basis of such an automaton, where φ_M is the conjunction of the formulas defined above, and

$\varphi_{\mathcal{A}}$ is:

$$\begin{aligned}
& \exists z(\nexists x(\text{succ}(x, z)) \wedge z \in X_0) \wedge \\
& \exists y(\nexists x(\text{succ}(y, x)) \wedge y \in X_3) \wedge \\
& \forall x, y (x \in X_0 \wedge \text{succ}(x, y) \wedge a(y) \Rightarrow y \in C_B \wedge y \in X_1) \wedge \\
& \forall x, y (x \in X_1 \wedge \text{succ}(x, y) \wedge a(y) \Rightarrow y \in C_A \wedge y \in X_1) \wedge \\
& \forall x, y, z (y \in X_1 \wedge \text{succ}(y, z) \wedge M(x, z) \wedge z \in R_A \wedge b(z) \Rightarrow z \in X_2) \wedge \\
& \forall x, y, z (y \in X_2 \wedge \text{succ}(y, z) \wedge M(x, z) \wedge z \in R_A \wedge b(z) \Rightarrow z \in X_2) \wedge \\
& \forall x, y, z (y \in X_2 \wedge \text{succ}(y, z) \wedge M(x, z) \wedge z \in R_B \wedge b(z) \Rightarrow z \in X_3) \wedge \\
& \forall x, y, z (y \in X_1 \wedge \text{succ}(y, z) \wedge M(x, z) \wedge z \in R_B \wedge b(z) \Rightarrow z \in X_3).
\end{aligned}$$

Other studies aimed at exploiting less complex but also less powerful logics, such as first-order and temporal ones to support a more practical automatic verification of VPL [1] as it happened with great success with model-checking for RL; algorithmic model-checking, however, is not the main focus of this paper and we do not go deeper into this issue.

4.4 Other structured context-free languages

As we said, early work on parenthesis languages and tree-automata ignited many attempts to enlarge those classes of languages and to investigate their properties. Among them VPL have received much attention in the literature and in this paper, probably thanks to the completeness of the obtained results –closure properties and logic characterization. To give an idea of the vastness of this panorama and of the connected problems, and to help comparison among them, in this section we briefly mention a few more of such families with no attempt for exhaustiveness.

Balanced grammars *Balanced grammars* (BG) have been proposed in [7] as a first approach to model mark-up languages such as XML by exploiting suitable extensions of parenthesis grammars. Basically a (BG) has a partitioned alphabet exactly in the same way as VPL; on the other hand any production of a BG has the form $A \rightarrow aab$ where $a \in \Sigma_c$, $b \in \Sigma_r$, and α is a regular expression¹³ over $V_N \cup \Sigma_i$.

Since it is well-known that the use of regular expressions in the rhs of CF grammars can be replaced by a suitable expansion by using exclusively “normal” rules, we can immediately conclude that balanced languages, i.e. those generated by BG, are a proper subclass of VPL (e.g. they do not admit unmatched elements of Σ_c and Σ_r). Furthermore they are not closed under concatenation and Kleene * [7]; we are not aware of any logic characterization of these languages.

Height-deterministic languages Height-deterministic languages, introduced in [41], are a more recent and fairly general way of describing CFL in terms of their structure. In a nutshell the hidden structure of a sentence is made explicit by making ε -moves visible, in that the original Σ alphabet is enriched as $\Sigma \cup \{\varepsilon\}$; if the original input string in Σ^*

¹³ A regular expression over a given alphabet V is built on the basis of alphabet’s elements by applying union, concatenation, and Kleene * symbols; it is well-known that the class of languages definable by means of regular expressions is RL.

is transformed into one over $\Sigma \cup \{\varepsilon\}$ by inserting an explicit ε wherever a recognizing automaton executes an ε -move, we obtain a linear representation of the syntax-tree of the original sentence, so that the automaton can be used as a real parser. We illustrate such an approach by means of the following example.

Example 11. Consider the language $L = L_1 \cup L_2$, with $L_1 = \{a^n b^n c^* \mid n \geq 0\}$, $L_2 = \{a^* b^n c^n \mid n \geq 0\}$ which is well-known to be inherently ambiguous since a string such as $a^n b^n c^n$ can be parsed both according to L_1 's structure and according to L_2 's one. A possible nondeterministic PDA \mathcal{A} recognizing L could act as follows:

- \mathcal{A} pushes all a 's until it reaches the first b ; at this point it makes a nondeterministic choice:
 - in one case it makes a "pause", i.e., an ε -move and enters a new state, say q_1 ;
 - in the other case it directly enters a different state, say q_2 with no "pause";
- from now on its behavior is deterministic; precisely:
 - in q_1 it checks that the number of b 's equals the number of a 's and then accepts any number of c 's;
 - in q_2 it pushes the b 's to verify that their number equals that of the c 's.

Thus, the two different behaviors of \mathcal{A} when analyzing a string of the type $a^n b^n c^n$ would result into two different strings in the extended alphabet $\Sigma \cup \{\varepsilon\}$ ¹⁴, namely $a^n \varepsilon b^n c^n$ and $a^n b^n c^n$; it is now simple, if needed, to state a one-to-one correspondence between strings augmented with explicit ε and the different syntax-trees associated with the original input: in this example, $a^n \varepsilon b^n c^n$ corresponds to the structure of Figure 12 (a) and $a^n b^n c^n$ to that of Figure 12 (b). It is also easy to build other nondeterministic PDA recognizing L that "produce" different strings associated with different structures.

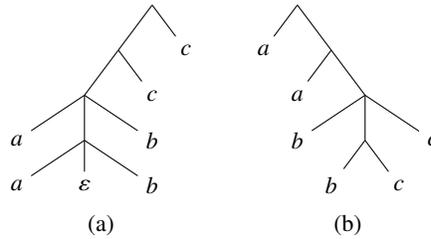


Fig. 12. Different structures for $a^n \varepsilon b^n c^n$ (a) and $a^n b^n c^n$ (b), for $n = 2$.

Once the input string structures are made visible by introducing the explicit ε character, the characteristics of PDA, of their subclasses, and of the language subfamilies they recognize, are investigated by referring to the *heights of their stack*. Precisely:

¹⁴ This could be done explicitly by means of a nondeterministic transducer that outputs a special marker in correspondence of an ε -move, but we stick to the original [41] formalization where automata are used exclusively as acceptors without resorting to formalized transducers.

- Any PDA is put into a *normalized form*, where
 - the δ relation is *complete*, i.e., it is defined in such a way that for every input string, whether accepted or not, the automaton scans the whole string:

$$\forall x \exists c (c_0 = (x, q_0, Z_0) \vdash^* c = (\varepsilon, q, \gamma));$$
 - every element of δ is exclusively in one of the forms: (q, A, o, q', AB) , or (q, A, o, q', A) , or $(q, A, o, q', \varepsilon)$, where $o \in (\Sigma \cup \{\varepsilon\})$;
 - for every $q \in Q$ all elements of δ having q as the first component either have $o \in \Sigma$ or $o = \varepsilon$, but not both of them.
- For any normalized \mathcal{A} and word $w \in (\Sigma \cup \{\varepsilon\})^*$, $\mathcal{N}(\mathcal{A}, w)$ denotes the set of all stack heights reached by \mathcal{A} after scanning w .
- \mathcal{A} is said *height-deterministic* (HPDA) iff $\forall w \in (\Sigma \cup \{\varepsilon\})^*$, $|\mathcal{N}(\mathcal{A}, w)| = 1$.
- The family of height-deterministic PDA is named HPDA; similarly, HDPDA denotes the family of deterministic height-deterministic PDA, and HRDPDA that of deterministic, real-time (i.e., those that do not perform any ε -move) height-deterministic PDA. The same acronyms with a final L replacing the A denote the corresponding families of languages.

It is immediate to realize (Lemma 1 of [41]) that every PDA admits an equivalent normalized one. Example 11 provides an intuitive explanation that every PDA admits an equivalent HPDA (Theorem 1 of [41]); thus HPDL = CFL; also, any (normalized) DPDA is, by definition an HPDA and therefore a deterministic HPDA; thus HDPDL = DCFL. Finally, since every deterministic VPA is already in normalized form and is a real-time machine, $VPL \subset HRDPDL$: the strict inclusion follows from the fact that $L = \{a^n b a^n\}$ cannot be recognized by a VPA since the same character a should belong both to Σ_c and to Σ_r .

Coupling the extension of the alphabet from Σ to $\Sigma \cup \{\varepsilon\}$ with the set $\mathcal{N}(\mathcal{A}, w)$ allows us to consider HPDL as a generalized kind of structured languages. As an intuitive explanation, let us go back to Example 11 and consider the two behaviors of \mathcal{A} when parsing the string $aabbcc$ once it has been “split” into $aabbcc$ and $aa\varepsilon bbcc$; the stack heights $\mathcal{N}(\mathcal{A}, w)$ for all their prefixes are, respectively: $(1, 2, 3, 4, 3, 2)$ and $(1, 2, 2, 1, 0, 0, 0)$ (if we do not count the bottom of the stack Z_0). In general, it is not difficult to associate every sequence of stack lengths during the parsing of an input string (in $(\Sigma \cup \{\varepsilon\})^*$!) with the syntax-tree visited by the HPDA.

As a consequence, the following fundamental definition of *synchronization* between HPDA can be interpreted as a form of structural equivalence.

Definition 16. *Two HPDA \mathcal{A} and \mathcal{B} are synchronized, denoted as $\mathcal{A} \sim \mathcal{B}$, iff $\forall w \in (\Sigma \cup \{\varepsilon\})^*$, $\mathcal{N}(\mathcal{A}, w) = \mathcal{N}(\mathcal{B}, w)$.*

It is immediate to realize that synchronization is an equivalence relation and therefore to associate an equivalence class $[\mathcal{A}]_{\sim}$ with every HPDA; we also denote as \mathcal{A} -HDPL the class of languages recognized by automata in $[\mathcal{A}]_{\sim}$. Then, in [41] the authors show that:

- For every *deterministic* HPDA \mathcal{A} the class \mathcal{A} -HDPL is a Boolean algebra.¹⁵

¹⁵ If the automaton is not deterministic only closures under union and intersection hold.

- Real-time HPDA can be determinized (with a complexity of the same order as for VPA), so that the class of real-time HPDL coincides with HRDPDL.

On the other hand neither HRDPDL nor HDPDL are closed under concatenation and Kleene * [15] so that the gain obtained in terms of generative power w.r.t. VPL has a price in terms of closure properties. We are not aware of logic characterizations for this class of structured languages.

Let us also mention that other classes of structured languages based on some notion of synchronization have been studied in the literature; in particular [41] compare their class with those of [25] and [11]. Finally, we acknowledge that our survey does not consider some logic characterization of tree or even graph languages which refer either to very specific families (such as, e.g. star-free tree languages [42]) and/or to an alphabet of basic elements, e.g., arcs connecting tree or graph nodes [10], which departs from the general framework of string (structured) languages.

5 Parsing context-free languages deterministically

We have seen in Section 3.4 that general CFL, having a hidden and sometimes ambiguous structure, need not only recognizing mechanisms but more complex parsing algorithms that, besides deciding whether a string belongs to the language or not, also produce the syntax-tree(s) formalizing its structure and possibly driving its semantics; parsing is the core of any compiler or interpreter. We noticed that not all CFL can be recognized –and therefore parsed– by DPDA: whereas any CFL can be recognized by a nondeterministic PDA that operates in real-time, the best deterministic algorithms to parse general CFL, such as CKY and Early’s ones have a $O(n^3)$ complexity, which is considered not acceptable in many fields such as programming language compilation.

For this and other reasons many application fields restrict their attention to DCFL; DPDA can be easily built, with no loss of generality, in such a way that they can operate in linear time, whether as pure recognizers or as parsers and language translators: it is sufficient, for any DPDA, to effectively transform it into an equivalent *loop-free* one, i.e. an automaton that does not perform more than a constant number, say k , of ε -moves before reading a character from the input or popping some element from the stack (see, e.g., [30]). In such a way the whole input x is analyzed in at most $h \cdot |x|$ moves, where h is a function of k and the maximum length of the string that the automaton can push on the stack in a single move. Notice, however, that in general it is not possible to obtain a DPDA that recognizes its language in real-time. Consider, for instance, the language $L = \{a^m b^n c^n d^m \mid m, n \geq 1\} \cup \{a^m b^+ e d^m \mid m \geq 1\}$: intuitively, a DPDA recognizing L must first push the as onto the stack; then, it must also store on the stack the subsequent bs since it must be ready to compare their number with the following cs , if any; after reading the bs , however, if it reads the e it must necessarily pop all bs by means of n ε -moves before starting the comparison between the as and the ds .

Given that, in general, it is undecidable to state whether a CFL is deterministic or not [31], the problem of automatically building a deterministic automaton, *if any*, from a given CFG is not trivial and deserved much research. In this section we will briefly recall two major approaches to the problem of deterministic parsing. We do not go deep into their technicalities, however, because the families of languages they can analyze are

not of much interest from the point of view of algebraic and logical characterizations. It is Section 6, instead, where we introduce a class of languages that allows for highly efficient deterministic parsing *and* exhibits practically all desirable algebraic and logical properties.

5.1 Top-down deterministic parsing

In Section 3.4 we introduced the GNF for CFG and observed that, if in such grammars there are no two distinct productions of the type $A \rightarrow a\alpha$, $A \rightarrow a\beta$, then the automaton built therefrom is deterministic. The above particular case has been generalized leading to the definition of *LL grammars*, so called because they allow to build deterministic parsers that scan the input Left-to-right and produce the Leftmost derivation thereof. Intuitively, an *LL(k)* grammar is such that, for any leftmost derivation

$$S\#^k \Rightarrow^* xA\alpha\#^k$$

it is possible to decide deterministically which production to apply to rewrite nonterminal A by “looking ahead” at most k terminal characters of the input string that follows x .¹⁶ Normally, the practical application of LL parsing is restricted to the case $k = 1$ to avoid memorizing and searching too large tables. In general this choice allows to cover a large portion of programming language syntax even if it is well-known that not all DCFL can be generated by LL grammars. For instance no LL grammar can generate the deterministic language $\{a^n b^n \mid n \geq 1\} \cup \{a^n c^n \mid n \geq 1\}$ since the decision on whether to join the a ’s with the b ’s or with the c ’s clearly requires an unbounded look-ahead.

Example 12. The following grammar is a simple transformation of GAE_1 in LL form; notice that the original left recursion $E \Rightarrow E + T$ has been replaced by a right one $E \Rightarrow^* T + E$ and similarly for nonterminal T .

$$\begin{aligned} GAE_{LL} : S &\rightarrow E\# \\ E &\rightarrow TE' \\ E' &\rightarrow +E \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *T \mid \varepsilon \\ F &\rightarrow e. \end{aligned}$$

5.2 Bottom-up deterministic parsing

So far the PDA that we used to recognize or parse CFL are working in a top-down manner by trying to build leftmost derivation(s) that produce the input string starting from grammar’s axiom. We also mentioned, however, that trees can be traversed also in a bottom-up way; a typical way of doing so is visiting them in *leftmost post-order*, i.e. scanning their frontier left-to right and, as soon as a string of children is identified, writing them followed by their father, then proceeding recursively until the root

¹⁶ The “tail” of $k\#$ characters is a simple trick to allow for the look ahead when the reading head of the automaton is close to the end of the input.

is reached. For instance such a visit of the syntax-tree by which GAE_1 generates the string $e + e * e$ is $eE + eT * eFTES$ which is the reverse of the rightmost derivation $S \Rightarrow E \Rightarrow E + T \Rightarrow E + T * F \Rightarrow E + T * e \Rightarrow E + e * e \Rightarrow e + e * e$.

Although Definition 8 introduces PDA in such a way that they are naturally oriented towards building leftmost derivations, it is not difficult to let them work in a bottom-up fashion: the automaton starts reading the input left-to-right and pushes the read character on top of the stack (formally, since at every transition the character on top of the stack is *replaced* by a string, it rewrites the existing symbol and adds the read character over it); as soon as –by means of its finite memory control device– it realizes that a whole rhs is on top of the stack it replaces it by the corresponding lhs (again, to be fully consistent with Definition 8, this operation should not be formalized as a single transition but could be designed as a sort of “macro-move” consisting of several micro-steps that remove one character per step and whose last one also pushes the lhs character).¹⁷ This type of operating by a PDA is called *shift-reduce parsing* since it consists in *shifting* characters from the input to the stack and *reducing* them from a rhs to the corresponding lhs.

Not surprisingly, the “normal” behavior of such a bottom-up parser is, once again, nondeterministic: in our example once the rhs e is identified, the automaton must apply a reduction either to F or to T or to E . Even more critical is the choice that the automaton must take after having read the substring $e+e$ and having (if it did the correct reductions) $E + T$ on top of the stack: in this case the string could be the rhs of the rule $E \rightarrow E + T$, and in that case the automaton should reduce it to E but the T on the top could also be the beginning of another rhs, i.e., $T * F$, and in such a case the automaton should go further by shifting more characters before doing any reduction; this is just the opposite situation of what happens with VPL, where the current input symbol allows the machine to decide how to progress.

In fact, the previous type of nondeterminism, i.e., the choice whether to reduce e to E or T or F , could be resolved as we did with structured languages, by eliminating repeated rhs through a classical power set construction. In this case, instead, the automaton must build by itself the unknown structure and this may imply a lack of knowledge about having reached or not a complete rhs. “Determinizing” bottom-up parsers has been an intense research activity during the 1960’s, as well as for their top-down counterpart. In Section 6 we thoroughly examine one of the earliest practical deterministic bottom-up parsers and the class of languages they can recognize, namely Floyd’s *operator precedence languages*. The study of these languages, however, has been abandoned after a while due to advent of a more powerful class of grammars –the LR ones, defined and investigated by Knuth [33], whose parsers proceed Left-to-right as well as the LL ones but produce (the reverse of) Rightmost derivations. LR grammars in fact, unlike LL and operator precedence ones, generate all DCFL.

¹⁷ A formal definition of PDA operating in the above way is given, e.g., in [41] and reported in Section 4.4

6 Operator precedence languages

Operator precedence grammars (OPG) have been introduced by R. Floyd in his pioneering paper [26] with the goal of building efficient, deterministic, bottom-up parsers for programming languages. In doing so he was inspired by the hidden structure of arithmetic expressions which suggests to “give precedence” to, i.e. to execute first multiplicative operations w.r.t. the additive ones, as we illustrated through Example 7. Essentially, the goal of deterministic bottom-up parsing is to unambiguously decide when a complete rhs has been identified so that we can proceed with replacing it with the unique corresponding lhs with no risk to apply some roll-back if another possible reduction was the right one. Floyd achieved such a goal by suitably extending the notion of precedence between arithmetic operators to all grammar terminals, in such a way that a complete rhs is enclosed within a pair (*yields-precedence*, *takes-precedence*). OPG obtained a considerable success thanks to their simplicity and to the efficiency of the parsers obtained therefrom; incidentally, some independent studies also uncovered interesting algebraic properties ([16]) which have been exploited in the field of grammar inference ([17]). As we anticipated in the introduction, however, the study of these grammars has been dismissed essentially because of the advent of other classes, such as the LR ones, which can generate all DCFL; OPG instead do not have such power as we will see soon, although they are able to cover most syntactic features of normal programming languages, and parsers based on OPG are still in practical use (see, e.g., [29]).

Only recently we renewed our interest in this class of grammars and languages for two different reasons that are the object of this survey. On the one side in fact, OPL, despite being apparently not structured, since they require and have been motivated by parsing, have shown rather surprising relations with various families of structured languages; as a consequence it has been possible to extend to them all the language properties investigated in the previous sections. On the other side, OPG enable parallelizing their parsing in a natural and efficient way, unlike what happens with other parsers which are compelled to operate in a strict left-to-right fashion, thus obtaining considerable speed-up thanks to the wide availability of modern parallel HW architectures.

Therefore, after having resumed the basic definitions and properties of OPG and their languages, we show, in Section 6.1, that they considerably increase the generative power of structured languages but, unlike the whole class of DCFL, they still enjoy all algebraic and logic properties that we investigated for such smaller classes. In Section 6.2 we show how parallel parsing is much better supported by this class of grammars than by the presently used ones.

Definition 17. *A grammar rule is in operator form if its rhs has no adjacent nonterminals; an operator grammar (OG) contains only such rules.*

Notice that the grammars considered in Example 7 are OG. Furthermore any CF grammar can be recursively transformed into an equivalent OG one [30].

Next, we introduce the notion of precedence relations between elements of Σ : we say that *a* is *equal in precedence* to *b* iff the two characters occur consecutively, or at

most with one nonterminal in between, in some rhs of the grammar; in fact, when evaluating the relations between terminal characters for OPG, nonterminals are inessential, or “transparent”. a yields precedence to b iff a can occur at the immediate left of a subtree whose leftmost terminal character is b (again whether there is a nonterminal character at the left of b or not is inessential). Symmetrically, a takes precedence over b iff a occurs as the rightmost terminal character of a subtree and b is its following terminal character. These concepts are formally defined as follows.

Definition 18. For an OG G and a nonterminal A , the left and right terminal sets are

$$\begin{aligned}\mathcal{L}_G(A) &= \{a \in \Sigma \mid A \xRightarrow{*} Baa\} \\ \mathcal{R}_G(A) &= \{a \in \Sigma \mid A \xRightarrow{*} \alpha aB\} \text{ where } B \in V_N \cup \{\varepsilon\}.\end{aligned}$$

The grammar name G will be omitted unless necessary to prevent confusion.

For an OG G , let α, β range over $(V_N \cup \Sigma)^*$ and $a, b \in \Sigma$. The three binary operator precedence (OP) relations are defined as follows:

- equal in precedence: $a \doteq b \iff \exists A \rightarrow \alpha A b b \beta, B \in V_N \cup \{\varepsilon\}$,
- takes precedence: $a \succ b \iff \exists A \rightarrow \alpha D b \beta, D \in V_N \text{ and } a \in \mathcal{R}_G(D)$,
- yields precedence: $a \prec b \iff \exists A \rightarrow \alpha a D \beta, D \in V_N \text{ and } b \in \mathcal{L}_G(D)$.

For an OG G , the operator precedence matrix (OPM) $M = OPM(G)$ is a $|\Sigma| \times |\Sigma|$ array that, for each ordered pair (a, b) , stores the set M_{ab} of OP relations holding between a and b .

For the grammar GAE_1 of Example 7 the left and right terminal sets of nonterminals E, T and F are, respectively:

$\mathcal{L}(E) = \{+, *, e\}$, $\mathcal{L}(T) = \{*, e\}$, $\mathcal{L}(F) = \{e\}$, $\mathcal{R}(E) = \{+, *, e\}$, $\mathcal{R}(T) = \{*, e\}$, and $\mathcal{R}(F) = \{e\}$.

| | | | |
|---|---|----|---|
| | + | * | e |
| + | > | << | < |
| * | > | > | < |
| e | > | > | |

Fig. 13. The OPM of the GAE_1 of Example 7.

Figure 13 displays the OPM associated with the grammar of GAE_1 of Example 7 where, for an ordered pair (a, b) , a is one of the symbols shown in the first column of the matrix and b one of those occurring in its first row. Notice that, unlike the usual arithmetic relations denoted by similar symbols, the above precedence relations do not enjoy any of the transitive, symmetric, reflexive properties.

Definition 19. An OG G is an operator precedence or Floyd grammar (OPG) iff $M = OPM(G)$ is a conflict-free matrix, i.e., $\forall a, b, |M_{ab}| \leq 1$. An operator precedence language (OPL) is a language generated by an OPG.

An OPG is in Fischer normal form (FNF) iff it is invertible, i.e., no two rules have the same rhs, has no empty rules, i.e., rules whose rhs is ε , except possibly $S \rightarrow \varepsilon$, and no renaming rules, i.e., rules whose rhs is a single nonterminal, except possibly those whose lhs is S .

For every OPG an equivalent one in FNF can effectively be built [24,30]. A FNF grammar (manually) derived from GAE_1 of Example 7 is GAE_{FNF} :

$$\begin{aligned} S &\rightarrow E \mid T \mid F \\ E &\rightarrow E + T \mid E + F \mid T + T \mid F + F \mid F + T \mid T + F \\ T &\rightarrow T * F \mid F * F \\ F &\rightarrow e \end{aligned}$$

We can now see how the precedence relations of an OPG can drive the deterministic parsing of its sentences: consider again the sentence $e + e * e$; add to its boundaries the conventional symbol $\#$ which implicitly yields precedence to any terminal character and to which every terminal character takes precedence, and evaluate the precedence relations between pairs of consecutive symbols; they are displayed below:

$$\# \langle e \rangle + \langle e \rangle * \langle e \rangle \#.$$

The three occurrences of e enclosed within the pair (\langle, \rangle) , are the rhs of production $F \rightarrow e$; thanks to the fact that the grammar is in FNF there is no doubt on its corresponding lhs; therefore they can be reduced to F . Notice that such a reduction could be applied *in any order, possibly even in parallel*; this feature will be exploited later in Section 6.2 but for the time being let us consider a traditional bottom-up parser proceeding rigorously left-to-right. Thus, the first reduction produces the string $F + e * e$; if we now recompute the precedence relations on the new string, due to the “transparency” of nonterminals, we obtain $\# \langle F + \langle e \rangle * \langle e \rangle \#$.

At this point the next rhs to be reduced is the second occurrence of e ; after a third similar reduction the original string and the corresponding precedence relations, are reduced to $\# \langle F + \langle F * F \rangle \#$ where the (only) next rhs to be reduced is $F * F$.

Notice that there is no doubt on whether the first nonterminal F should be part of a rhs *beginning* with F or of another one *ending* with F , such as $\# \langle F + F \langle * F \rangle \#$. In fact, if the rhs to be reduced were just $*F$, its corresponding lhs would be a nonterminal adjacent to the F at its left, thus contradicting the hypothesis of the grammar being an OG. At this point the bottom-up shift-reduce algorithm continues deterministically until the axiom is reached and a syntax-tree of the original sentence –represented by the mirror image of the rightmost derivation– is built.

This first introduction to OPG allows us to draw some early important properties thereof:

- In some sense OPL are *input-driven* even if they do not match exactly the definition of these languages: in fact, the decision of whether to apply a push operation (at the

beginning of a rhs) or a shift one (while scanning a rhs) or a pop one (at the end of a rhs) depends only on terminal characters but not on a *single* one, as a look-ahead of one more terminal character is needed.¹⁸

- The above characteristic is also a major reason why OPL, though allowing for efficient deterministic parsing of various practical programming languages [29,18,26], do not cover the whole family DCFL; consider in fact the language $L = \{0a^n b^n \mid n \geq 0\} \cup \{1a^n b^{2n} \mid n \geq 0\}$: a DPDA can easily “remember” the first character in its state; then push all the *a*’s onto the stack and, when it reaches the *bs* decide whether to scan one or two *bs* for every *a* depending on its memory of the first read character. On the contrary, it is clear that any grammar generating *L* would necessarily exhibit some precedence conflict.¹⁹
- We like to consider OPL as *structured languages* in a generalized sense since, once the OPM is given, the structure of their sentences is immediately defined and univocally determinable as it happens. e.g., with VPL once the alphabet is partitioned into call, return, and internal alphabet. However, we would be reluctant to label OPL as *visible* since there is a major difference between parenthesis-like terminals which make the language sentence isomorphic to its syntax-tree, and precedence relations which help *building* the tree but are computed only during the parsing. Indeed, not all of them are immediately visible in the original sentence: e.g., in some cases such as in the above sentence $\# \langle F + \langle F * F \rangle \#$ precedence relations are not even matched so that they can be assimilated to real parentheses only when they mark a complete rhs. In summary, we would consider that OPL are structured (by the OPM) as well as PL (by explicit parentheses), VPL (by alphabet partitioning), and other families of languages; however, we would intuitively label them at most as “semi-visible” since making their structure visible requires some partial parsing, though not necessarily a complete recognition.

6.1 Algebraic and logic properties of operator precedence languages

OPL enjoy all algebraic and logic properties that have been illustrated in the previous sections for much less powerful families of structured languages.

As a first step we introduce the notion of a *chain* as a formal description of the intuitive concept of “semi-visible structure”. To illustrate the following definitions and properties we will continue to make use of examples inspired by arithmetic expressions but we will enrich such expressions with, possibly nested, explicit parentheses as the visible part of their structure. For instance the following grammar GAE_P is a natural enrichment of previous GAE_1 to generate arithmetic expressions that involve parenthesized subexpressions (we use the slightly modified symbols ‘(’ and ‘)’’ to avoid

¹⁸ As it happens in other deterministic parsers such as LL or LR ones (see Section 5).

¹⁹ The above *L* is instead LL (see Section 5.1); on the contrary, the language $\{a^n b^n \mid n \geq 1\} \cup \{a^n c^n \mid n \geq 1\}$ is OPL but not LL; thus, OPL and LL languages are uncomparable.

overloading with other uses of parentheses).

$$\begin{aligned}
 GAE_P : S &\rightarrow E \mid T \mid F \\
 E &\rightarrow E + T \mid T * F \mid e \mid \langle E \rangle \\
 T &\rightarrow T * F \mid e \mid \langle E \rangle \\
 F &\rightarrow e \mid \langle E \rangle
 \end{aligned}$$

Definition 20 (Operator precedence alphabet). An operator precedence (OP) alphabet is a pair (Σ, M) where Σ is an alphabet and M is a conflict-free operator precedence matrix, i.e. a $|\Sigma \cup \{\#\}|^2$ array that associates at most one of the operator precedence relations: $\doteq, <, >$ with each ordered pair (a, b) . As stated above the delimiter $\#$ yields precedence to other terminals and other terminals take precedence over it (with the special case $\# \doteq \#$ for the final reduction of renaming rules.) Since such relations are stated once and forever, we do not explicitly display them in OPM figures.

If $M_{ab} = \{\circ\}$, with $\circ \in \{<, \doteq, >\}$, we write $a \circ b$. For $u, v \in \Sigma^*$ we write $u \circ v$ if $u = xa$ and $v = by$ with $a \circ b$.

Definition 21 (Chains). Let (Σ, M) be a precedence alphabet.

- A simple chain is a word $a_0 a_1 a_2 \dots a_n a_{n+1}$, written as ${}^{a_0}[a_1 a_2 \dots a_n]^{a_{n+1}}$, such that: $a_0, a_{n+1} \in \Sigma \cup \{\#\}$, $a_i \in \Sigma$ for every $i : 1 \leq i \leq n$, $M_{a_0 a_{n+1}} \neq \emptyset$, and $a_0 < a_1 \doteq a_2 \dots a_{n-1} \doteq a_n > a_{n+1}$.
- A composed chain is a word $a_0 x_0 a_1 x_1 a_2 \dots a_n x_n a_{n+1}$, with $x_i \in \Sigma^*$, where ${}^{a_0}[a_1 a_2 \dots a_n]^{a_{n+1}}$ is a simple chain, and either $x_i = \varepsilon$ or ${}^{a_i}[x_i]^{a_{i+1}}$ is a chain (simple or composed), for every $i : 0 \leq i \leq n$. Such a composed chain will be written as ${}^{a_0}[x_0 a_1 x_1 a_2 \dots a_n x_n]^{a_{n+1}}$.
- The body of a chain ${}^a[x]^b$, simple or composed, is the word x .

Example 13. Figure 14 (a) depicts the $OPM(GAE_P)$, whereas Figure 14 (b) represents the “semi-visible” structure induced by the operator precedence alphabet of grammar GAE_P for the expression $\#e + e * \langle e + e \rangle \#$: $\#[x_0 + x_1]^\#$, $+[y_0 * y_1]^\#$, $*[\langle w_0 \rangle]^\#$, $\langle [z_0 + z_1] \rangle^\flat$ are composed chains and $\#[e]^+$, $+[e]^*$, $\langle [e]^+ \rangle$, $+[e]^\flat$ are simple chains.

Definition 22 (Compatible word). A word w over (Σ, M) is compatible with M iff the two following conditions hold:

- For each pair of letters c, d , consecutive in w , $M_{cd} \neq \emptyset$;
- for each substring x of $\#w\#$ such that $x = a_0 x_0 a_1 x_1 a_2 \dots a_n x_n a_{n+1}$, if $a_0 < a_1 \doteq a_2 \dots a_{n-1} \doteq a_n > a_{n+1}$ and, for every $0 \leq i \leq n$, either $x_i = \varepsilon$ or ${}^{a_i}[x_i]^{a_{i+1}}$ is a chain (simple or composed), then $M_{a_0 a_{n+1}} \neq \emptyset$.

For instance, the word $e + e * \langle e + e \rangle$ is compatible with the operator precedence alphabet of grammar GAE_P , whereas $e + e * \langle e + e \rangle \langle e + e \rangle$ is not.

Thus, given an OP alphabet, the set of possible chains over Σ^* represents the universe of possible structured strings compatible with the given OPM.

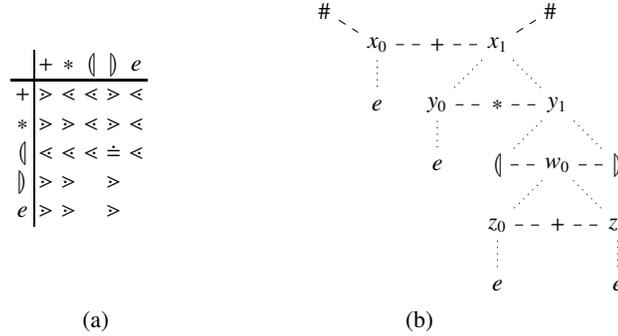


Fig. 14. OPM of grammar GAE_P (a) and structure of the chains in the expression $\#e + e * (e + e)\#$ (b).

Operator precedence automata Despite abstract machines being the classical way to formalize recognition and parsing algorithms for any family of formal languages, and despite OPL having been invented just with the purpose of supporting deterministic parsing, their theoretical investigation has been abandoned before a family of automata completely equivalent to their generative grammars appeared in the literature. Only recently, when the numerous still unexplored benefits obtainable from this family appeared clear to us, we filled up this hole with the herewith resumed definition ([37]). The formal model presented in this paper is a “traditional” left-to-right automaton, although, as we already anticipated and will thoroughly exploit in the next Section 6.2 a distinguishing feature of OPL is that their parsing can be started in arbitrary positions with no harm nor loss of efficiency. This choice is explained by the need to extend and to compare results already reported for other language families. The original slightly more complicated version of this model was introduced in [38].

Definition 23 (Operator precedence automaton). An operator precedence automaton (OPA) is a tuple $\mathcal{A} = (\Sigma, M, Q, I, F, \delta)$ where:

- (Σ, M) is an operator precedence alphabet,
- Q is a set of states (disjoint from Σ),
- $I \subseteq Q$ is the set of initial states,
- $F \subseteq Q$ is the set of final states,
- $\delta \subseteq Q \times (\Sigma \cup Q) \times Q$ is the transition relation, which is the union of three disjoint relations:

$$\delta_{shift} \subseteq Q \times \Sigma \times Q, \quad \delta_{push} \subseteq Q \times \Sigma \times Q, \quad \delta_{pop} \subseteq Q \times Q \times Q.$$

An OPA is deterministic iff

- I is a singleton
- All three components of δ are functions:

$$\delta_{shift} : Q \times \Sigma \rightarrow Q, \quad \delta_{push} : Q \times \Sigma \rightarrow Q, \quad \delta_{pop} : Q \times Q \rightarrow Q.$$

We represent an OPA by a graph with Q as the set of vertices and $\Sigma \cup Q$ as the set of edge labelings. The edges of the graph are denoted by different shapes of arrows to distinguish the three types of transitions: there is an edge from state q to state p labeled by $a \in \Sigma$ denoted by a dashed (respectively, normal) arrow iff $(q, a, p) \in \delta_{\text{shift}}$ (respectively, $\in \delta_{\text{push}}$) and there is an edge from state q to state p labeled by $r \in Q$ and denoted by a double arrow iff $(q, r, p) \in \delta_{\text{pop}}$.

To define the semantics of the automaton, we need some new notations.

We use letters p, q, p_i, q_i, \dots to denote states in Q . Let Γ be $\Sigma \times Q$ and let $\Gamma' = \Gamma \cup \{Z_0\}$ be the *stack alphabet*; we denote symbols in Γ' as $[a, q]$ or Z_0 . We set $\text{symbol}([a, q]) = a$, $\text{symbol}(Z_0) = \#$, and $\text{state}([a, q]) = q$. Given a string $\Pi = \pi_n \dots \pi_2 \pi_1 Z_0$, with $\pi_i \in \Gamma'$, $n \geq 0$, we set $\text{symbol}(\Pi) = \text{symbol}(\pi_n)$, including the particular case $\text{symbol}(Z_0) = \#$.

As usual, a *configuration* of an OPA is a triple $c = \langle w, q, \Pi \rangle$, where $w \in \Sigma^* \#$, $q \in Q$, and $\Pi \in \Gamma^* Z_0$.

A *computation* or *run* of the automaton is a finite sequence of *moves* or *transitions* $c_1 \mapsto c_2$; there are three kinds of moves, depending on the precedence relation between the symbol on top of the stack and the next symbol to read:

push move: if $\text{symbol}(\Pi) < a$ then $\langle ax, p, \Pi \rangle \mapsto \langle x, q, [a, p]\Pi \rangle$, with $(p, a, q) \in \delta_{\text{push}}$;

shift move: if $a \doteq b$ then $\langle bx, q, [a, p]\Pi \rangle \mapsto \langle x, r, [b, p]\Pi \rangle$, with $(q, b, r) \in \delta_{\text{shift}}$;

pop move: if $a > b$ then $\langle bx, q, [a, p]\Pi \rangle \mapsto \langle bx, r, \Pi \rangle$, with $(q, p, r) \in \delta_{\text{pop}}$.

Observe that shift and pop moves are never performed when the stack contains only Z_0 .

Push and shift moves update the current state of the automaton according to the transition relations δ_{push} and δ_{shift} , respectively: push moves put a new element on top of the stack consisting of the input symbol together with the current state of the automaton, whereas shift moves update the top element of the stack by *changing its input symbol only*. Pop moves remove the element on top of the stack, and update the state of the automaton according to δ_{pop} on the basis of the pair of states consisting of the current state of the automaton and the state of the removed stack symbol; notice that in this moves the input symbol is used only to establish the $>$ relation and it remains available for the following move.

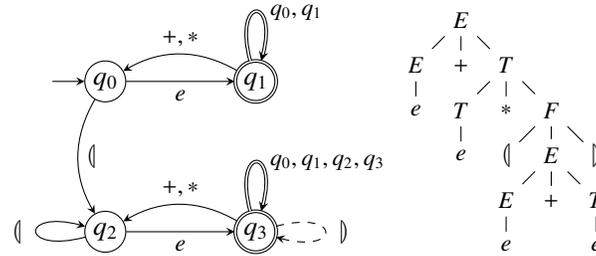
The language accepted by the automaton is defined as:

$$L(\mathcal{A}) = \left\{ x \mid \langle x\#, q_I, Z_0 \rangle \xrightarrow{*} \langle \#, q_F, Z_0 \rangle, q_I \in I, q_F \in F \right\}.$$

Example 14. The OPA depicted in Figure 15 accepts the language of arithmetic expressions generated by grammar GAE_P . The same figure also shows the syntax-tree of the sentence $e + e * (e + e)$ and an accepting computation on this input.

Notice the similarity of the above definition of OPA with that of VPA (Definition 15) and with the normalized form for PDA given in Section 4.4.

Showing the equivalence between OPG and OPAs, though somewhat intuitive, requires to overtake a few non-trivial technical difficulties, mainly in the path from OPG to OPAs. Here we offer just an informal description of the rationale of the two constructions and an illustrating example; the full proof of the equivalence between OPG and OPA can be found in [37].



| input | state | stack |
|----------------------|-------|--|
| $e + e * (e + e) \#$ | q_0 | Z_0 |
| $+e * (e + e) \#$ | q_1 | $[e, q_0]Z_0$ |
| $+e * (e + e) \#$ | q_1 | Z_0 |
| $e * (e + e) \#$ | q_0 | $[+, q_1]Z_0$ |
| $*(e + e) \#$ | q_1 | $[e, q_0][+, q_1]Z_0$ |
| $*(e + e) \#$ | q_1 | $[+, q_1]Z_0$ |
| $(e + e) \#$ | q_0 | $[\ast, q_1][+, q_1]Z_0$ |
| $e + e) \#$ | q_2 | $[(, q_0][\ast, q_1][+, q_1]Z_0$ |
| $+e) \#$ | q_3 | $[e, q_2][(, q_0)[\ast, q_1][+, q_1]Z_0$ |
| $+e) \#$ | q_3 | $[(, q_0)[\ast, q_1][+, q_1]Z_0$ |
| $e) \#$ | q_2 | $[+, q_3][(, q_0)[\ast, q_1][+, q_1]Z_0$ |
| $) \#$ | q_3 | $[e, q_2][+, q_3][(, q_0)[\ast, q_1][+, q_1]Z_0$ |
| $) \#$ | q_3 | $[+, q_3][(, q_0)[\ast, q_1][+, q_1]Z_0$ |
| $) \#$ | q_3 | $[(, q_0)[\ast, q_1][+, q_1]Z_0$ |
| $\#$ | q_3 | $[(, q_0)[\ast, q_1][+, q_1]Z_0$ |
| $\#$ | q_3 | $[\ast, q_1][+, q_1]Z_0$ |
| $\#$ | q_3 | $[+, q_1]Z_0$ |
| $\#$ | q_3 | Z_0 |

Fig. 15. Automaton and example of computation for the language of Example 14. Recall that shift, push and pop transitions are denoted by dashed, normal and double arrows, respectively.

For convenience and with no loss of generality, let G be an OPG with no empty rules, except possibly $S \rightarrow \varepsilon$, and no renaming rules, except possibly those whose lhs is S , an OPA \mathcal{A} equivalent to G is built in such a way that a successful computation thereof corresponds to building bottom-up the mirror image of a rightmost derivation of G : \mathcal{A} performs a push transition when it reads the first terminal of a new rhs; it performs a shift transition when it reads a terminal symbol inside a rhs, i.e. a leaf with some left sibling leaf; it performs a pop transition when it completes the recognition of a rhs, then it guesses (nondeterministically, if there are several rules with the same rhs) the nonterminal at the lhs.

Each state of \mathcal{A} contains two pieces of information: the first component represents the prefix of the rhs under construction, whereas the second component is used to recover the rhs *previously under construction* (see Figure 16) whenever all rhs's nested below have been completed. Without going into the details of the construction and the

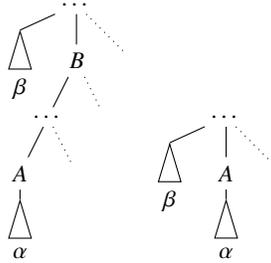


Fig. 16. When parsing α , the prefix previously under construction is β .

formal equivalence proof between G and \mathcal{A} , we further illustrate the rationale of the construction through the following Example.

Example 15. Consider again grammar GAE_p . Figure 17 shows the first part of an accepting computation of the automaton derived therefrom. Consider, for instance, step 3 of the computation: at this point the automaton has already reduced (nondeterministically) the first e to E and has pushed the following $+$ onto the stack, paired with the state from which it was coming; thus, its new state is $\langle E+, \varepsilon \rangle$; at step 6, instead, the state is $\langle T*, E+ \rangle$ because at this point the automaton has built the $T*$ part of the current rhs and remembers that the prefix of the suspended rhs is $E+$. Notice that the computation partially shown in Figure 17 is equal to that of Figure 15 up to a renaming of the states; in fact the shape of syntax-trees and consequently the sequence of push, shift and pop moves in OPL depends only on the OPM, not on the visited states.

The converse construction from OPAs to OPG is somewhat simpler; in essence, from a given OPA $\mathcal{A} = (\Sigma, M, Q, I, F, \delta)$ a grammar G is built whose nonterminals are 4-tuples $(a, q, p, b) \in \Sigma \times Q \times Q \times \Sigma$, written as $\langle^a p, q^b \rangle$. G 's rules are built on the basis of \mathcal{A} 's chains as follows (particular cases are omitted for simplicity):

| step | input | state | stack |
|------|----------------------|--|---|
| 0 | $e + e * (e + e) \#$ | $\langle \varepsilon, \varepsilon \rangle$ | Z_0 |
| 1 | $+e * (e + e) \#$ | $\langle e, \varepsilon \rangle$ | $[e, \langle \varepsilon, \varepsilon \rangle] Z_0$ |
| 2 | $+e * (e + e) \#$ | $\langle E, \varepsilon \rangle$ | Z_0 |
| 3 | $e * (e + e) \#$ | $\langle E+, \varepsilon \rangle$ | $[+, \langle E, \varepsilon \rangle] Z_0$ |
| 4 | $*(e + e) \#$ | $\langle e, \varepsilon \rangle$ | $[e, \langle E+, \varepsilon \rangle][+, \langle E, \varepsilon \rangle] Z_0$ |
| 5 | $*(e + e) \#$ | $\langle T, E+ \rangle$ | $[+, \langle E, \varepsilon \rangle] Z_0$ |
| 6 | $(e + e) \#$ | $\langle T*, E+ \rangle$ | $[*, \langle T, E+ \rangle][+, \langle E, \varepsilon \rangle] Z_0$ |

Fig. 17. Partial accepting computation of the automaton built from grammar GAE_P .

- for every simple chain ${}^{a_0}[a_1 a_2 \dots a_n]^{a_{n+1}}$, if there is a sequence of \mathcal{A} 's transitions that, while reading the body of the chain starting from q_0 leaves \mathcal{A} in q_{n+1} , include the rule

$$\langle {}^{a_0} q_0, q_{n+1} {}^{a_{n+1}} \rangle \longrightarrow a_1 a_2 \dots a_n$$

- for every composed chain ${}^{a_0}[x_0 a_1 x_1 a_2 \dots a_n x_n]^{a_{n+1}}$, add the rule

$$\langle {}^{a_0} q_0, q_{n+1} {}^{a_{n+1}} \rangle \longrightarrow \Lambda_0 a_1 \Lambda_1 a_2 \dots a_n \Lambda_n$$

if there is a sequence of \mathcal{A} 's transitions that, while reading the body of the chain starting from q_0 leaves \mathcal{A} in q_{n+1} , and, for every $i = 0, 1, \dots, n$, $\Lambda_i = \varepsilon$ if $x_i = \varepsilon$, otherwise $\Lambda_i = \langle {}^{a_i} q_i, q_i' {}^{a_{i+1}} \rangle$ if $x_i \neq \varepsilon$ and there is a path leading from q_i to q_i' when traversing x_i .

Since the size of G 's nonterminal alphabet is bounded, the above procedure eventually terminates when no new rules are added to P .²⁰

We have seen that a fundamental issue to state the properties of most abstract machines is their determinizability: in the cases examined so far we have realized that the positive basic result holding for RL extends to the various versions of structured CFL, though at the expenses of more intricate constructions and size complexity of the deterministic versions obtained from the nondeterministic ones, but not to the general CF family. Having OPL been born just with the motivation of supporting deterministic parsing, and being they structured as well, it is not surprising to find that for any non-deterministic OPA with s states an equivalent deterministic one can be built with $2^{O(s^2)}$ states, as it happens for the analogous construction for VPL: in [37] besides giving a detailed construction for the above result, it is also noticed that the construction of an OPA from an OPG is such that, if the OPG is in FNF, then the obtained automaton is already deterministic since the grammar has no repeated rhs. As a consequence, producing a deterministic OPA from an OPG by first putting the OPG into FNF produces an automaton of an exponentially smaller size than the other way around.

²⁰ The above claim can be easily proved if the OPM has no circularities in the \doteq relation, since this implies an upper bound to the length of P 's rhs; in the (seldom) case where this hypothesis is not verified other ones can be exploited (see [37] for a more detailed analysis.)

Operator precedence vs other structured languages A distinguishing feature of OPL is that, on the one side they have been defined to support deterministic parsing, i.e., the construction of the syntax-tree of any sentence which is not immediately visible in the sentence itself but, on the other side, they can still be considered as structured in the sense that their syntax-trees are univocally determined once an OPM is fixed, as it happens when we enclose grammar's rhs within parentheses or we split the terminal alphabet into $\Sigma_c \cup \Sigma_r \cup \Sigma_i$. It is therefore natural to compare their generative power with that of other structured languages.

On this respect, the main result of this section is that *OPL strictly include VPL*, which in turn strictly include parenthesis languages and the languages generated by balanced grammars (discussed in Section 4.4).

This result was originally proved in [15]. To give an intuitive explanation of this claim consider the following remarks:

- Sequences $\in \Sigma_i^*$ can be assimilated to regular “sublanguages” with a linear structure; if we conventionally assign to them a left-linear structure, this can be obtained through an OPM where every character, but those in Σ_c , takes precedence over all elements in Σ_i ; by stating instead that all elements of Σ_c yield precedence to the elements in Σ_i we obtain that after every call the OPA pushes and pops all subsequent elements of Σ_i , as an FSA would do without using the stack.
- All elements of Σ_r produce a pop from the stack of the corresponding element of Σ_c , if any; thus we obtain the same effect by letting them take precedence over all other terminal characters.
- A VPA performs a push onto its stack when (and only when) it reads an element of Σ_c , whereas an OPA pushes the elements to which some other element yields precedence; thus, it is natural to state that whenever on top of the stack there is a call symbol, possibly after having visited a subtree whose result is stored as the state component in the top of the stack together with the terminal symbol, such a symbol yields precedence to the following call (roughly, open parentheses yield precedence to other open parentheses and closed parentheses take precedence over other closed parentheses).
- Once the whole subsequence embraced by two matching call and return is scanned and possibly reduced, the two terminals are faced, with the possible occurrence of an irrelevant nonterminal in between, and therefore the call must be equal in precedence to the return.
- Finally, the usual convention that # yields precedence to everything and everything takes precedence over # enables the managing of possible unmatched returns at the beginning of the sentence and unmatched calls at its end.

In summary, for every VPA \mathcal{A} with a given partitioned alphabet Σ , an OPM such as the one displayed in Figure 18 and an OPA \mathcal{A}' defined thereon can be built such that $L(\mathcal{A}') = L(\mathcal{A})$.

In [15] it is also shown the converse property, i.e., that whenever an OPM is such that the terminal alphabet can be partitioned into three disjoint sets $\Sigma_c, \Sigma_r, \Sigma_i$ such that the OPM has the shape of Figure 18, any OPL defined on such an OPM is also a VPL. Strict inclusion of VPL within OPL follows from the fact that VPL are real-time whereas OPL

| | | | | |
|------------|------------------|------------------|------------------|---|
| | Σ_c | Σ_r | Σ_i | Legend Σ_c denotes “calls” Σ_r denotes “returns” Σ_i denotes internal characters |
| Σ_c | \leq | \equiv | \leq | |
| Σ_r | \triangleright | \triangleright | \triangleright | |
| Σ_i | \triangleright | \triangleright | \triangleright | |

Fig. 18. A partitioned matrix, where $\Sigma_c, \Sigma_r, \Sigma_i$ are set of terminal characters. A precedence relation in position $\Sigma_\alpha, \Sigma_\beta$ means that relation holds between all symbols of Σ_α and all those of Σ_β .

include also non-real-time languages (see Section 5); there are also real-time OPL²¹ such as

$$L = \{b^n c^n \mid n \geq 1\} \cup \{f^n d^n \mid n \geq 1\} \cup \{e^n (fb)^n \mid n \geq 1\}$$

that are not VPL. In fact, strings of type $b^n c^n$ impose that b is a call and c a return; for similar reasons, f must be a call and d a return. Strings of type $e^n (fb)^n$ impose that at least one of b and f must be a return, a contradiction for a VP alphabet. In conclusion we have the following result:

Theorem 2. *VPL are the subfamily of OPL whose OPM is a partitioned matrix, i.e., a matrix whose structure is depicted in Figure 18.*

As a corollary OPL also strictly include balanced languages and parenthesis languages. OPL are instead uncomparable with HRDPDL: we have already seen that the language $L_1 = \{a^n b a^n\}$ is an HRDPDL but it is neither a VPL nor an OPL since it necessarily requires a conflict $a < a$ and $a > a$; conversely, the previous language $L_2 = \{a^m b^n c^n d^m \mid m, n \geq 1\} \cup \{a^m b^+ e d^m \mid m \geq 1\}$ can be recognized by an OPA but by no HRDPDA (see Section 5).

The increased power of OPL over other structured languages goes far beyond the mathematical containment properties and opens several application scenarios that are hardly accessible by “traditional” structured languages. The field of programming languages was the original motivation and source of inspiration for the introduction of OPL; arithmetic expressions, used throughout this paper as running examples, are just a small but meaningful subset of such languages and we have emphasized from the beginning that their partially hidden structure cannot be “forced” to the linearity of RL, nor can always be made explicit by the insertion of parentheses.

VPL too have been presented as an extension of parenthesis languages with the motivation that not always calls, e.g. procedure calls, can be matched by corresponding returns: a sudden closure, e.g. due to an exception or an interrupt or an unexpected end may leave an open chain of suspended calls. Such a situation, however, may need a generalization that cannot be formalized by the VPL formalism, since in VPL unmatched

²¹ When we say that an OPL L is real-time we mean, as usual, that there is an abstract machine, in particular a DPDA, recognizing it that performs exactly $|x|$ moves for every input string x ; this is not to say that an OPA accepting L operates in real-time, since OPA’s pop moves are defined as ε moves.

calls can occur only at the end of a string.²² Imagine, for instance, that the occurrence of an interrupt while serving a chain of calls imposes to abort the chain to serve immediately the interrupt; after serving the interrupt, however, the normal system behavior may be resumed with new calls and corresponding returns even if some previous calls have been lost due to the need to serve the interrupt with high priority. Various, more or less sophisticated, policies can be designed to manage such systems and can be adequately formalized as OPL. The next example describes a first simple case of this type; other more sophisticated examples of the same type and further ones inspired by different application fields can be found in [37].

Example 16 (Managing interrupts). Consider a software system that is designed to serve requests issued by different users but subject to interrupts. Precisely, assume that the system manages “normal operations” according to a traditional LIFO policy, and may receive and serve some interrupts denoted as *int*.

We model its behavior by introducing an alphabet with two pairs of calls and returns: *call* and *ret* denote the call to, and return from, a normal procedure; *int*, and *serve* denote the occurrence of an interrupt and its serving, respectively. The occurrence of an interrupt provokes discarding possible pending *calls* not already matched by corresponding *rets*; furthermore when an interrupt is pending, i.e., not yet served, calls to normal procedures are not accepted and consequently corresponding returns cannot occur; interrupts however, can accumulate and are served themselves along a LIFO policy. Once all pending interrupts have been served the system can accept new *calls* and manage them normally.

Figure 19 (a) shows an OPM that assigns to sequences on the above alphabet a structure compatible with the described priorities. Then, a suitable OPA can specify further constraints on such sequences; for instance the automaton of Figure 19 (b) restricts the set of sequences compatible with the matrix by imposing that all *int* are eventually served and the computation ends with no pending *calls*; furthermore unmatched *serve* and *ret* are forbidden. E.g., the string *call.call.ret.int.serve.call.ret* is accepted through the sequence of states $q_0 \xrightarrow{\text{call}} q_1 \xrightarrow{\text{call}} q_1 \xrightarrow{\text{ret}} q_1 \xrightarrow{q_1} q_1 \xrightarrow{q_0} q_0 \xrightarrow{\text{int}} q_2 \xrightarrow{\text{serve}} q_2 \xrightarrow{q_0} q_0 \xrightarrow{\text{call}} q_1 \xrightarrow{\text{ret}} q_1 \xrightarrow{q_0} q_0$; on the contrary, a sequence beginning with *call.serve* would not be accepted.

Closure and decidability properties Structured languages with *compatible structures* often enjoy many closure properties typical of RL; noticeably, families of structured languages are often Boolean algebras. The intuitive notion of compatible structure is formally defined for each family of structured languages; for instance two VPL have compatible structure iff their tri-partition of Σ is the same; two height-deterministic PDA languages (HPDL) have compatible structure if they are synchronized. In the case of OPL, the notion of structural compatibility is naturally formalized through the OPM.

²² Recently, such a weakness of VPL has been acknowledged in [3] where the authors introduced *colored VPL* to cope with the above problem; the extended family, however, still does not reach the power of OPL ([3]).

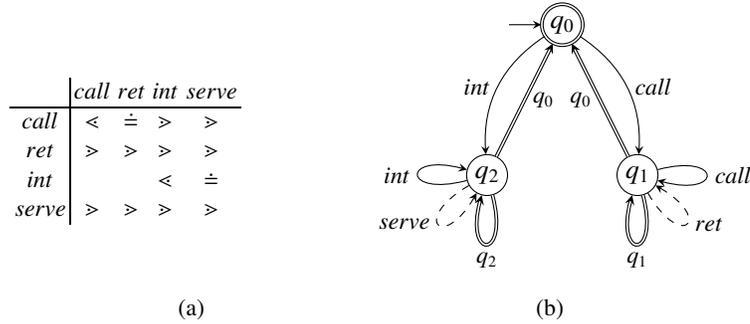


Fig. 19. OPM (a) and automaton (b) for the language of Example 16.

Definition 24. Given two OPM M_1 and M_2 , we define set inclusion and union:

$$M_1 \subseteq M_2 \text{ if } \forall a, b : (M_1)_{ab} \subseteq (M_2)_{ab}$$

$$M = M_1 \cup M_2 \text{ if } \forall a, b : M_{ab} = (M_1)_{ab} \cup (M_2)_{ab}.$$

Two matrices are compatible if their union is conflict-free. A matrix is total (or complete) if it contains no empty cell.

The following theorem has been proved originally in [16] by exploiting some standard forms of OPG that have been applied to grammar inference problems [17].

Theorem 3. For any conflict-free OPM M the OPL whose OPM is contained in M are a Boolean algebra. In particular, if M is complete, the top language of its associated algebra is Σ^* with the structure defined by M .

Notice however, that the same result could be proved in a simpler and fairly standard way by exploiting OPA and their traditional composition rules (which pass through determinization to achieve closure under complement). As usual in such cases, thanks to the decidability of the emptiness problem for general CFL, a major consequence of Boolean closures is the following corollary.

Corollary 1. The inclusion problem between OPL with compatible OPM is decidable.

Closure under concatenation and Kleene * has been proved more recently in [15]; whereas such closures are normally easily proved or disproved for many families of languages, the technicalities to achieve this result are not trivial for OPL; however we do not go into their description since closure or non-closure w.r.t. these operations is not of major interest in this paper.

Logic characterization Achieving a logic characterization of OPL has probably been the most difficult job in the recent revisit of these languages and posed new challenges w.r.t. the analogous path previously followed for RL and then for CFL [35] and VPL [5]. In fact we have seen that moving from linear languages such as RL to tree-shaped ones such as CFL led to the introduction of the relation M between the positions of leftmost and rightmost leaves of any subtree (generated by a grammar in DGNF); the obtained characterization in terms of first-order formulas existentially quantified w.r.t. the M relation (which is a representation of the sentence structure) however, was suffering from the lack of closure under complementation of CFL [35]; the same relation instead proved more effective for VPL thanks to the fact that they are structured and enjoy all necessary closures.

To the best of our knowledge, however, *all previous characterizations* of formal languages in terms of logics that refer to string positions (for instance, there is significant literature on the characterization of various subclasses of RL in terms of first-order or temporal logics, see, e.g., [23]) *have been given for real-time languages*. This feature is the key that allows, in the exploitation of MSO logic, to state a natural correspondence between automaton's state q_i and second-order variable X_i in such a way that the value of X_i is the set of positions where the state visited by the automaton is q_i .

OPL instead include also DCFL that are not real-time and, as a consequence, there are positions where the recognizing OPA traverses different configurations with different states. As a further consequence, the M relation adopted for CFL and VPL is not anymore a one-to-one relation since the same position may be the position of the left/rightmost leaf of several subtrees of the whole syntax-tree; this makes formulas such as the key ones given in Section 4.3 meaningless.

The following key ideas helped overtaking the above difficulties:

- A new relation μ replaces the original M adopted in [35] and [5]; μ is based on the look-ahead-look-back mechanism which drives the (generalized) input-driven parsing of OPL based on precedence relations: thus, whereas in $M(x, y)$ x, y denote the positions of the extreme leaves of a subtree, in $\mu(x, y)$ they denote the position of the *context* of the same subtree, i.e., respectively, of the character that yields precedence to the subtree's leftmost leaf, and of the one over which the subtree's rightmost leaf takes precedence. The new μ relation is not one-to-one as well, but, unlike the original M , its parameters x, y are not "consumed" by a pop transition of the automaton and remain available to be used in further automaton transitions of any type. In other words, μ holds between the positions 0 and $n + 1$ of every chain (see Definition 21). For instance, Figure 20 displays the μ relation, graphically denoted by arrows, holding for the sentence $e + e * (e + e)$ generated by grammar GAE_P : we have $\mu(0, 2)$, $\mu(2, 4)$, $\mu(5, 7)$, $\mu(7, 9)$, $\mu(5, 9)$, $\mu(4, 10)$, $\mu(2, 10)$, and $\mu(0, 10)$. Such pairs correspond to contexts where a reduce operation is executed during the parsing of the string (they are listed according to their execution order).

In general $\mu(x, y)$ implies $y > x + 1$, and a position x may be in relation μ with more than one position and vice versa. Moreover, if w is compatible with M , then $\mu(0, |w| + 1)$.

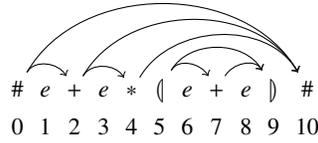


Fig. 20. The string $e + e * (e + e)$, with positions and relation μ .

Example 17. The following sentence of the MSO logic enriched with the μ relation defines, within the universe of strings compatible with the OPM of Figure 14(a), the language where parentheses are used only when they are needed (i.e. to give precedence to $+$ over $*$).

$$\forall \mathbf{x} \forall \mathbf{y} \left(\begin{array}{c} \mu(\mathbf{x}, \mathbf{y}) \wedge ((\mathbf{x} + 1) \wedge) (\mathbf{y} - 1) \\ \Rightarrow \\ (*(\mathbf{x}) \vee *(\mathbf{y})) \wedge \\ \mathbf{x} + 1 < \mathbf{z} < \mathbf{y} - 1 \wedge +(\mathbf{z}) \wedge \\ \exists \mathbf{z} \left(\begin{array}{c} \neg \exists \mathbf{u} \exists \mathbf{v} \left(\begin{array}{c} \mathbf{x} + 1 < \mathbf{u} < \mathbf{z} \wedge ((\mathbf{u}) \wedge) \\ \mathbf{z} < \mathbf{v} < \mathbf{y} - 1 \wedge ((\mathbf{v}) \wedge) \\ \mu(\mathbf{u} - 1, \mathbf{v} + 1) \end{array} \right) \end{array} \right) \end{array} \right)$$

- Since in every position there may be several states held by the automaton while visiting that position, instead of associating just one second-order variable to each state of the automaton we define three different sets of second-order variables, namely, $\mathbf{A}_0, \mathbf{A}_1, \dots, \mathbf{A}_N, \mathbf{B}_0, \mathbf{B}_1, \dots, \mathbf{B}_N$ and $\mathbf{C}_0, \mathbf{C}_1, \dots, \mathbf{C}_N$. Set \mathbf{A}_i contains those positions of word w where state q_i may be assumed after a shift or push transition, i.e. after a transition that “consumes” an input symbol. Sets \mathbf{B}_i and \mathbf{C}_i encode a pop transition concluding the reading of the body of a chain ${}^a[w_0 a_1 w_1 \dots a_l w_l]^{a_{l+1}}$ in a state q_i : set \mathbf{B}_i contains the position of symbol a that precedes the corresponding push, whereas \mathbf{C}_i contains the position of a_i , which is the symbol on top of the stack when the automaton performs the pop move relative to the whole chain.

Figure 21 presents such sets for the example automaton of Figure 15, with the same input as in Figure 20. Notice that each position, except the last one, belongs to exactly one \mathbf{A}_i , whereas it may belong to several \mathbf{B}_i and at most one \mathbf{C}_i .

We can now outline how an OPA can be derived from an MSO logic formula making use of the new symbol μ and conversely.

From MSO formula to OPA The construction from MSO logic to OPA essentially follows the lines given originally by Büchi, and reported in Section 2.2: once the original alphabet has been enriched and the formula has been put in the canonical form in the same way as described in Section 2.2, we only need to define a suitable automaton fragment to be associated with the new atomic formula $\mu(\mathbf{X}_i, \mathbf{X}_j)$; then, the construction of the global automaton corresponding to the global formula proceeds in the usual inductive way.

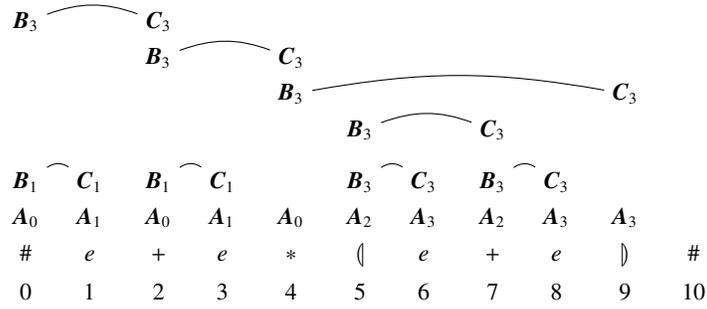


Fig. 21. The string of Figure 20 with B_i , A_i , and C_i evidenced for the automaton of Figure 15. Pop moves of the automaton are represented by linked pairs B_i, C_i .

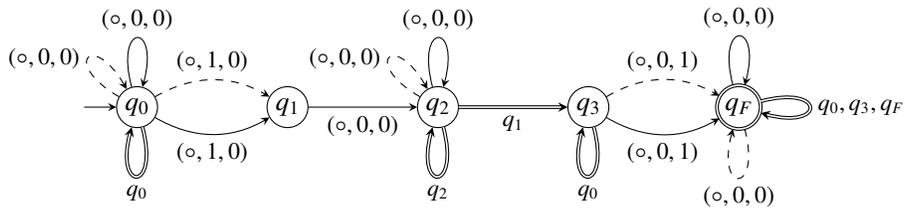


Fig. 22. OPA for atomic formula $\mu(X, Y)$.

Figure 22 represents the OPA for atomic formula $\psi = \mu(X, Y)$. As before, labels are triples belonging to $\Sigma \times \{0, 1\}^2$, where the first component encodes a character $a \in \Sigma$, the second the positions belonging to X (with 1) or not (with 0), while the third component is for Y . The symbol \circ is used as a shortcut for any value in Σ compatible with the OPM, so that the resulting automaton is deterministic.

The semantics of μ requires for $\mu(X, Y)$ that there must be a chain $^a[w_2]^b$ in the input word, where a is the symbol at the only position in X , and b is the symbol at the only position in Y . By definition of chain, this means that a must be read, hence in the position represented by X the automaton performs either a push or a shift move (see Figure 22, from state q_0 to q_1), as pop moves do not consume input. After that, the automaton must read w_2 . In order to process the chain $^a[w_2]^b$, reading w_2 must start with a push move (from state q_1 to state q_2), and it must end with one or more pop moves, before reading b (i.e. the only position in Y – going from state q_3 to q_F).

This means that the automaton, after a generic sequence of moves corresponding to visiting an irrelevant (for $\mu(X, Y)$) portion of the syntax-tree, when reading the symbol at position X performs either a push or a shift move, depending on whether X is the position of a leftmost leaf of the tree or not. Then it visits the subsequent subtree ending with a pop labeled q_1 ; at this point, if it reads the symbol at position Y , it accepts anything else that follows the examined fragment.

It is interesting to compare the diagram of Figure 22 with those of Figure 6 (c) and of Figure 10: the first one, referring to RL, uses two consecutive moves; the second one, referring to VPL, may perform an unbounded number of internal moves and of matching call-return pairs between the call-return pair in positions x, y ; the OPA does the same as the VPA but needs a pair of extra moves to take into account the look-ahead-look-back implied by precedence relations.

From the OPA \mathcal{A} to the MSO formula In this case the overall structure of the logic formula is the same as in the previous cases for RL and VPL, i.e., an existential quantification over second-order variables, which represent states through sets of positions within the string, of a global formula that formalizes a) the constraints imposed by automaton's transitions, b) the fact that in position 0 the automaton must be in an initial state, and c) that at the end of the string it must be in an accepting state. The complete formalization of the δ transition relation as a collection of formulas relating the various variables A_i, B_i, C_i , however, is much more involved than in the two previous cases. Here we only provide a few meaningful examples of such formulas, just to give the essential ideas of how they have been built; their complete set can be found in [37] together with the equivalence proof. Without loss of generality we *assume that the OPA is deterministic*.

Preliminarily, we introduce some notation to make the following formulas more understandable:

- When considering a chain $^a[w]^b$ we assume $w = w_0 a_1 w_1 \dots a_\ell w_\ell$, with $^a[a_1 a_2 \dots a_\ell]^b$ being a simple chain (any w_g may be empty). We denote by s_g the position of symbol a_g , for $g = 1, 2, \dots, \ell$ and set $a_0 = a$, $s_0 = 0$, $a_{\ell+1} = b$, and $s_{\ell+1} = |w| + 1$.
- $x < y$ states that the symbol in position x yields precedence to the one in position y and similarly for the other precedence relations

- The fundamental abbreviation

$$\text{Tree}(\mathbf{x}, \mathbf{z}, \mathbf{v}, \mathbf{y}) := \mu(\mathbf{x}, \mathbf{y}) \wedge \begin{pmatrix} (\mathbf{x} + 1 = \mathbf{z} \vee \mu(\mathbf{x}, \mathbf{z})) \wedge \\ \neg \exists \mathbf{t} (\mathbf{z} < \mathbf{t} < \mathbf{y} \wedge \mu(\mathbf{x}, \mathbf{t})) \wedge \\ (\mathbf{v} + 1 = \mathbf{y} \vee \mu(\mathbf{v}, \mathbf{y})) \wedge \\ \neg \exists \mathbf{t} (\mathbf{x} < \mathbf{t} < \mathbf{v} \wedge \mu(\mathbf{t}, \mathbf{y})) \end{pmatrix}$$

is satisfied, for every chain ${}^a[w]^b$ embraced within positions \mathbf{x} and \mathbf{y} by a (unique, maximal) \mathbf{z} such that $\mu(\mathbf{x}, \mathbf{z})$, if $w_0 \neq \varepsilon$, $\mathbf{z} = \mathbf{x} + 1$ if instead $w_0 = \varepsilon$; symmetrically for \mathbf{y} and \mathbf{v} . In particular, if w is the body of a simple chain, then $\mu(0, \ell + 1)$ and $\text{Tree}(0, 1, \ell, \ell + 1)$ are satisfied; if it is the body of a composed chain, then $\mu(0, |w| + 1)$ and $\text{Tree}(0, s_1, s_\ell, s_{\ell+1})$ are satisfied. If $w_0 = \varepsilon$ then $s_1 = 1$, and if $w_\ell = \varepsilon$ then $s_\ell = |w|$. In the example of Figure 20 relations $\text{Tree}(2, 3, 3, 4)$, $\text{Tree}(2, 4, 4, 10)$, $\text{Tree}(4, 5, 9, 10)$, $\text{Tree}(5, 7, 7, 9)$ are satisfied, among others.

- The shortcut $Q_i(\mathbf{x}, \mathbf{y})$ is used to represent that \mathcal{A} is in state q_i when at position \mathbf{x} and the next position to read, possibly after scanning a chain, is \mathbf{y} . Since the automaton is not real time, we must distinguish between push and shift moves (case $\text{Succ}_i(\mathbf{x}, \mathbf{y})$), and pop moves (case $\text{Next}_i(\mathbf{x}, \mathbf{y})$).

$$\begin{aligned} \text{Succ}_k(\mathbf{x}, \mathbf{y}) &:= \mathbf{x} + 1 = \mathbf{y} \wedge \mathbf{x} \in \mathbf{A}_k \\ \text{Next}_k(\mathbf{x}, \mathbf{y}) &:= \mu(\mathbf{x}, \mathbf{y}) \wedge \mathbf{x} \in \mathbf{B}_k \wedge \\ &\quad \exists \mathbf{z}, \mathbf{v} (\text{Tree}(\mathbf{x}, \mathbf{z}, \mathbf{v}, \mathbf{y}) \wedge \mathbf{v} \in \mathbf{C}_k) \\ Q_i(\mathbf{x}, \mathbf{y}) &:= \text{Succ}_i(\mathbf{x}, \mathbf{y}) \vee \text{Next}_i(\mathbf{x}, \mathbf{y}). \end{aligned}$$

E.g., with reference to Figures 20 and 21, $\text{Succ}_2(5, 6)$, $\text{Next}_3(5, 9)$, and $\text{Next}_3(5, 7)$ hold.

We can now show a meaningful sample of the various formulas that code the automaton's transition relation.

- The following formula states that if \mathcal{A} is in position \mathbf{x} and state q_i and performs a push transition by reading the character in position \mathbf{y} , it goes to state q_k according to the transition relation δ .

$$\varphi_{\text{push_fw}} := \forall \mathbf{x}, \mathbf{y} \bigwedge_{i=0}^N \bigwedge_{k=0}^N \bigwedge_{c \in \Sigma} \left(\begin{array}{l} \mathbf{x} \leq \mathbf{y} \wedge c(\mathbf{y}) \wedge Q_i(\mathbf{x}, \mathbf{y}) \wedge \\ \delta_{\text{push}}(q_i, c) = q_k \\ \Rightarrow \mathbf{y} \in \mathbf{A}_k \end{array} \right)$$

Notice that the original formula given in Section 2.2 for RL can be seen as a particular case of the above one.

- Conversely, if \mathcal{A} is in state q_k after a push starting from position \mathbf{x} and reading character c , in that position it must have been in a state q_i such that $\delta(q_i, c) = q_k$:

$$\varphi_{\text{push_bw}} := \forall \mathbf{x}, \mathbf{y} \bigwedge_{k=0}^N \bigwedge_{c \in \Sigma} \left(\begin{array}{l} \mathbf{x} \leq \mathbf{y} \wedge c(\mathbf{y}) \wedge \mathbf{y} \in \mathbf{A}_k \wedge \\ (\mathbf{x} + 1 = \mathbf{y} \vee \mu(\mathbf{x}, \mathbf{y})) \\ \Rightarrow \\ \bigvee_{i=0}^N (Q_i(\mathbf{x}, \mathbf{y}) \wedge \delta_{\text{push}}(q_i, c) = q_k) \end{array} \right)$$

- The formulas coding the shift transitions are similar to the previous ones and therefore omitted.

- To define $\varphi_{\delta_{\text{pop}}}$ we introduce the shortcut $\text{Tree}_{i,j}(\mathbf{x}, \mathbf{z}, \mathbf{v}, \mathbf{y})$, which represents the fact that \mathcal{A} is ready to perform a pop transition from state q_i having on top of the stack state q_j ; such pop transition corresponds to the reduction of the portion of string between positions \mathbf{x} and \mathbf{y} (excluded).

$$\text{Tree}_{i,j}(\mathbf{x}, \mathbf{z}, \mathbf{v}, \mathbf{y}) := \text{Tree}(\mathbf{x}, \mathbf{z}, \mathbf{v}, \mathbf{y}) \wedge Q_i(\mathbf{v}, \mathbf{y}) \wedge Q_j(\mathbf{x}, \mathbf{z}).$$

Formula $\varphi_{\delta_{\text{pop}}}$ is thus defined as the conjunction of three formulas. As before, the forward (abbreviated with the subscript f_w) formula gives the sufficient condition for two positions to be in the sets \mathbf{B}_k and \mathbf{C}_k , when performing a pop move, and the backward formulas state symmetric necessary conditions.

$$\begin{aligned} \varphi_{\text{pop-}f_w} &:= \forall \mathbf{x}, \mathbf{z}, \mathbf{v}, \mathbf{y} \bigwedge_{i=0}^N \bigwedge_{j=0}^N \bigwedge_{k=0}^N \left(\begin{array}{l} \text{Tree}_{i,j}(\mathbf{x}, \mathbf{z}, \mathbf{v}, \mathbf{y}) \wedge \\ \delta_{\text{pop}}(q_i, q_j) = q_k \\ \Rightarrow \\ \mathbf{x} \in \mathbf{B}_k \wedge \mathbf{v} \in \mathbf{C}_k \end{array} \right) \\ \varphi_{\text{pop-}bwB} &:= \forall \mathbf{x} \bigwedge_{k=0}^N \left(\begin{array}{l} \mathbf{x} \in \mathbf{B}_k \Rightarrow \\ \exists \mathbf{y}, \mathbf{z}, \mathbf{v} \bigvee_{i=0}^N \bigvee_{j=0}^N \text{Tree}_{i,j}(\mathbf{x}, \mathbf{z}, \mathbf{v}, \mathbf{y}) \wedge \\ \delta_{\text{pop}}(q_i, q_j) = q_k \end{array} \right) \\ \varphi_{\text{pop-}bwC} &:= \forall \mathbf{v} \bigwedge_{k=0}^N \left(\begin{array}{l} \mathbf{v} \in \mathbf{C}_k \Rightarrow \\ \exists \mathbf{x}, \mathbf{y}, \mathbf{z} \bigvee_{i=0}^N \bigvee_{j=0}^N \text{Tree}_{i,j}(\mathbf{x}, \mathbf{z}, \mathbf{v}, \mathbf{y}) \wedge \\ \delta_{\text{pop}}(q_i, q_j) = q_k \end{array} \right) \end{aligned}$$

6.2 Local parsability for parallel parsers

Let us now go back to the original motivation that inspired R. Floyd when he invented the OPG family, namely supporting efficient, deterministic parsing. In the introductory part of this section we noticed that the mechanism of precedence relations isolates grammar's rhs from their context so that they can be reduced to the corresponding lhs independently from each other. This fact guarantees that, in whichever order such reductions are applied, at the end a complete grammar derivation will be built; such a derivation corresponds to a visit of the syntax-tree, not necessarily leftmost or rightmost, and its *construction has no risk of applying any back-track* as it happens instead in nondeterministic parsing. We call this property *local parsability property*, which intuitively can be defined as the possibility of applying deterministically a bottom-up, shift-reduce parsing algorithm by inspecting only an a priori bounded portion of any string containing a rhs. Various formal definitions of this concept have been given in the literature, the first one probably being the one proposed by Floyd himself in [27]; a fairly general definition of local parsability and a proof that OPG enjoy it can be found in [6].

Local parsability, however, has the drawback that it loses chances of deterministic parsing when the information on how to proceed with the parsing is arbitrarily far from the current position of the parser. We therefore have a trade-off between the size of the family of recognizable languages, which in the case of LR grammars is the whole DCFL class (see Section 5.2), and the constraint of proceeding rigorously left-to-right for the parser. So far this trade-off has been normally solved in favor of the generality in the absence of serious counterparts in favor of the other option. We argue however, that

the massive advent of parallel processing, even in the case of small architectures such as those of tablets and smartphones, could dramatically change the present state of affairs. On the one side parallelizing parsers such as LL or LR ones requires reintroducing a kind of nondeterministic *guess* on the state of the parser in a given position, which in most cases voids the benefits of exploiting parallel processors (see [6] for an analysis of previous literature on various attempts to develop parallel parsers); on the contrary, OPL are from the beginning oriented toward parallel analysis whereas their previous use in compilation shows that they can be applied to a wide variety of practical languages, and further more as suggested by other examples given here and in [37].

Next we show how we exploited the local parsability property of OPG to realize a complete and general parallel parser for these grammars. A first consequence of the basic property is the following statement.

Statement 2 *For every substring $a\delta b$ of $\gamma a\delta b\eta \in V^*$ derivable from S , there exists a unique string α , called the irreducible string, deriving δ such that $S \xRightarrow{*} \gamma\alpha b\eta \xRightarrow{*} \gamma a\delta b\eta$, and the precedence relations between the consecutive terminals of $a\alpha b$ do not contain the pattern $\langle (\neq)^* \rangle$. Therefore there exists a factorization $a\alpha b = \zeta\theta$ into two possibly empty factors such that the left factor does not contain \langle and the right factor does not contain \rangle .*

On the basis of the above statement, the original parsing algorithm is generalized in such a way that it may receive as input a portion of a string, not necessarily enclosed within the delimiters #, and produces as output two stacks, one that stores the substring ζ and one that stores θ as defined in Statement 2, and a partial derivation of $a\alpha b = \zeta\theta \xRightarrow{*} a\delta b$. For instance, with reference to the grammar GAE_{FNF} , which is a FNF of GAE_1 , if we supply to such a generalized parser the partial string $+e * e * e + e$ we obtain $\zeta = +T+$, $\theta = e$ and $+T+ \xRightarrow{*} +e * e * e +$ since $+ \rangle +$ and $+ \langle e$. We call S^L and S^R the two stacks produced by this partial parsing.

At this point it is fairly easy to let several such generalized parsers work in parallel:

- Suppose to use k parallel processors, also called *workers*; then split the input into k chunks; given that an OP parser needs a look-ahead-look-back of one character, the chunks must overlap by one character for each consecutive pair. For instance, the global input string $\#e + e + e * e * e + *e + e\#$, with $k = 3$ could be split as shown below:

$$\# \overbrace{e + e}^1 + \overbrace{e * e * e +}^2 e \overbrace{*e + e\#}^3$$

where the unmarked symbols $+$ and e are shared by the adjacent segments. The splitting can be applied arbitrarily, although in practice it seems natural to use segments of approximately equal length and/or to apply some heuristic criterion (for instance, if possible one should avoid particular cases where only \langle or \rangle relations occur in a single chunk so that the parser could not produce any reduction).

- Each chunk is submitted to one of the workers which produces a partial result in the form of the pair (S^L, S^R) (notice that some of those partial stacks may be empty).
- The partial results are concatenated into a new string $\in V^*$ and the process is iterated until a short enough single chunk is processed and the original input string is

accepted or rejected. In practice it may be convenient to build the new segments to be supplied to the workers by facing an \mathcal{S}^R with the following \mathcal{S}^L so that the likelihood of applying many new reductions in the next pass is increased. For instance the $\zeta = +T+$ part produced by the parsing of the second chunk could be paired with the $\mathcal{S}^R = \#E+$ part obtained from the parsing of the first chunk, producing the string $\#E + T+$ to be supplied to a worker for the new iteration. Some experience shows that quite often optimal results in terms of speed-up are obtained with 2, at most 3 passes of parallel parsing.

[6] describes in detail *PAPAGENO*, a *PARallel PARser GENeratOr* built on the basis of the above algorithmic schema. It has been applied to several real-life data definition, or programming, languages including JSON, XML, and Lua and different HW architectures. The paper also reports on the experimental results in terms of the obtained speed-up compared with standard sequential parser generators as Bison. *PAPAGENO* is freely available at <https://github.com/PAPAGENO-devels/papageno> under GNU license.

7 Concluding remarks

The main goal of this paper has been to show that an old-fashioned and almost abandoned family of formal languages indeed offers considerable new benefits in apparently unrelated application fields of high interest in modern applications, i.e., automatic property verification and parallelization. In the first field OPL significantly extend the generative power of the successful class of VPL still maintaining all of their properties: to the best of our knowledge, OPL are the largest class of languages closed under all major language operations and provided with a complete classification in terms of MSO logic.

Various other results about this class of languages have been obtained or are under development, which have not been included in this paper for length limits. We mention here just the most relevant or promising ones with appropriate references for further reading.

- The theory of OPL for languages of finite length strings has been extended in [37] to so called ω -languages, i.e. languages of infinite length strings: the obtained results perfectly parallel those originally obtained by Büchi and others for RL and subsequently extended to other families, noticeably VPL [5]; in particular, ω -OPL lose determinizability in case of Büchi acceptance criterion as it happens for RL and VPL.
- Some investigation is going on to devise more tractable automatic verification algorithms than those allowed by the full characterization of these languages in terms of MSO logic. On this respect, the state of the art is admittedly still far from the success obtained with model checking exploiting various forms of temporal logics for FSA and several extensions thereof such as, e.g., timed automata [2]. Some interesting preliminary results have been obtained for VPL by [1] and for a subclass of OPL in [36].
- The local parsability property can be exploited not only to build parallel parsers but also to make them *incremental*, in such a way that when a large piece of text or

software code is locally modified its analysis should not be redone from scratch but only the affected part of the syntax-tree is “plugged” in the original one with considerable saving; furthermore incremental and/or parallel parsing can be naturally paired with incremental and/or parallel semantic processing, e.g. realized through the classic schema of *attribute evaluation* [34,14]. Some early results on incremental software verification by exploiting the locality property are reported in [8]. We also mention ongoing work on parallel XML-based query processing.

- A seminal paper by Schützenberger [43] introduced the concept of *weighted languages* as RL where each word is given a weight in a given algebra which may represent some “attribute” of the word such as importance or probability. Later, these weighted languages too have been characterized in terms of MSO logic [20] and such a characterization has also been extended to VPL [21] and ω -VPL [19]. Our two research groups are both confident that those results can also be extended to *weighted OPL* and are starting a joint investigation on this promising approach.

Acknowledgments. We acknowledge the contribution to our research given by Alessandro Barenghi, Stefano Crespi Reghizzi, Violetta Lonati, Angelo Morzenti, and Federica Panella.

References

1. R. Alur, M. Arenas, P. Barceló, K. Etessami, N. Immerman, and L. Libkin. First-order and temporal logics for nested words. *Logical Methods in Computer Science*, 4(4), 2008.
2. R. Alur and D. L. Dill. A Theory of Timed Automata. *Theor. Comput. Sci.*, 126(2):183–235, 1994.
3. R. Alur and D. Fisman. Colored nested words. In *Language and Automata Theory and Applications - 10th International Conference, LATA 2016, Prague, Czech Republic, March 14-18, 2016, Proceedings*, pages 143–155, 2016.
4. R. Alur and P. Madhusudan. Visibly Pushdown Languages. In *STOC: ACM Symposium on Theory of Computing (STOC)*, 2004.
5. R. Alur and P. Madhusudan. Adding nesting structure to words. *J. ACM*, 56(3), 2009.
6. A. Barenghi, S. Crespi Reghizzi, D. Mandrioli, F. Panella, and M. Pradella. Parallel parsing made practical. *Sci. Comput. Program.*, 112(3):195–226, 2015. DOI: 10.1016/j.scico.2015.09.002.
7. J. Berstel and L. Boasson. Balanced Grammars and Their Languages. In W. B. et al., editor, *Formal and Natural Computing*, volume 2300 of *LNCS*, pages 3–25. Springer, 2002.
8. D. Bianculli, A. Filieri, C. Ghezzi, and D. Mandrioli. Syntactic-semantic incrementality for agile verification. *Sci. Comput. Program.*, 97:47–54, 2015.
9. J. R. Büchi. Weak Second-Order Arithmetic and Finite Automata. *Mathematical Logic Quarterly*, 6(1-6):66–92, 1960.
10. D. Caucal. On infinite transition graphs having a decidable monadic theory. *Theor. Comput. Sci.*, 290(1):79–115, 2003.
11. D. Caucal. Synchronization of Pushdown Automata. In O. H. Ibarra and Z. Dang, editors, *Developments in Language Theory*, volume 4036 of *LNCS*, pages 120–132. Springer, 2006.
12. E. M. Clarke, E. A. Emerson, and A. P. Sistla. Automatic verification of finite-state concurrent systems using temporal logic specifications. *ACM Trans. Program. Lang. Syst.*, 8:244–263, April 1986.
13. H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, 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.

14. S. Crespi Reghizzi, L. Breveglieri, and A. Morzenti. *Formal Languages and Compilation, Second Edition*. Texts in Computer Science. Springer, 2013.
15. S. Crespi Reghizzi and D. Mandrioli. Operator Precedence and the Visibly Pushdown Property. *J. Comput. Syst. Sci.*, 78(6):1837–1867, 2012.
16. S. Crespi Reghizzi, D. Mandrioli, and D. F. Martin. Algebraic Properties of Operator Precedence Languages. *Information and Control*, 37(2):115–133, May 1978.
17. S. Crespi Reghizzi, M. A. Melkanoff, and L. Lichten. The Use of Grammatical Inference for Designing Programming Languages. *Commun. ACM*, 16(2):83–90, 1973.
18. K. De Bosschere. An Operator Precedence Parser for Standard Prolog Text. *Softw., Pract. Exper.*, 26(7):763–779, 1996.
19. M. Droste and S. Dück. Weighted automata and logics for infinite nested words. *CoRR*, abs/1506.07031, 2015.
20. M. Droste and P. Gastin. Weighted automata and weighted logics. *Theor. Comput. Sci.*, 380(1-2):69–86, 2007.
21. M. Droste and B. Pibaljomme. Weighted nested word automata and logics over strong bimonoids. *Int. J. Found. Comput. Sci.*, 25(5):641, 2014.
22. C. C. Elgot. Decision problems of finite automata design and related arithmetics. *Trans. Am. Math. Soc.*, 98(1):21–52, 1961.
23. E. A. Emerson. Temporal and modal logic. In *Handbook of Theoretical Computer Science, Volume B: Formal Models and Semantics (B)*, pages 995–1072. 1990.
24. M. J. Fischer. Some properties of precedence languages. In *STOC '69: Proc. first annual ACM Symp. on Theory of Computing*, pages 181–190, New York, NY, USA, 1969. ACM.
25. D. Fisman and A. Pnueli. Beyond Regular Model Checking. *FSTTCS: Foundations of Software Technology and Theoretical Computer Science*, 21, 2001.
26. R. W. Floyd. Syntactic Analysis and Operator Precedence. *J. ACM*, 10(3):316–333, 1963.
27. R. W. Floyd. Bounded context syntactic analysis. *Commun. ACM*, 7(2):62–67, 1964.
28. S. A. Greibach. A new normal-form theorem for context-free phrase structure grammars. *J. ACM*, 12(1):42–52, 1965.
29. D. Grune and C. J. Jacobs. *Parsing techniques: a practical guide*. Springer, New York, 2008.
30. M. A. Harrison. *Introduction to Formal Language Theory*. Addison Wesley, 1978.
31. J. Hartmanis and J. E. Hopcroft. What makes some language theory problems undecidable. *J. Comput. Syst. Sci.*, 4(4):368–376, 1970.
32. C. A. R. Hoare. An axiomatic basis for computer programming. *Commun. ACM*, 12(10):576–580, 1969.
33. D. E. Knuth. On the translation of languages from left to right. *Information and Control*, 8(6):607–639, 1965.
34. D. E. Knuth. Semantics of context-free languages. *Mathematical Systems Theory*, 2(2):127–145, 1968.
35. C. Lautemann, T. Schwentick, and D. Thérien. Logics for context-free languages. In L. Pacholski and J. Tiuryn, editors, *Computer Science Logic, 8th International Workshop, CSL '94, Kazimierz, Poland, September 25-30, 1994, Selected Papers*, volume 933 of *Lecture Notes in Computer Science*, pages 205–216. Springer, 1994.
36. V. Lonati, D. Mandrioli, F. Panella, and M. Pradella. First-order logic definability of free languages. In L. D. Beklemishev and D. V. Musatov, editors, *Computer Science - Theory and Applications - 10th International Computer Science Symposium in Russia, CSR 2015, Listvyanka, Russia, July 13-17, 2015, Proceedings*, volume 9139 of *Lecture Notes in Computer Science*, pages 310–324. Springer, 2015.
37. V. Lonati, D. Mandrioli, F. Panella, and M. Pradella. Operator precedence languages: Their automata-theoretic and logic characterization. *SIAM J. Comput.*, 44(4):1026–1088, 2015.

38. V. Lonati, D. Mandrioli, and M. Pradella. Precedence Automata and Languages. In *6th Int. Computer Science Symposium in Russia (CSR)*, volume 6651 of *LNCS*, pages 291–304. 2011.
39. R. McNaughton. Parenthesis Grammars. *J. ACM*, 14(3):490–500, 1967.
40. K. Mehlhorn. Pebbling mountain ranges and its application of DCFL-recognition. In *Automata, languages and programming (ICALP-80)*, volume 85 of *LNCS*, pages 422–435, 1980.
41. D. Nowotka and J. Srba. Height-Deterministic Pushdown Automata. In L. Kucera and A. Kucera, editors, *MFCS 2007, Český Krumlov, Czech Republic, August 26-31, 2007, Proceedings*, volume 4708 of *LNCS*, pages 125–134. Springer, 2007.
42. A. Pothhoff and W. Thomas. Regular tree languages without unary symbols are star-free. In Z. Ésik, editor, *Fundamentals of Computation Theory, 9th International Symposium, FCT '93, Szeged, Hungary, August 23-27, 1993, Proceedings*, volume 710 of *Lecture Notes in Computer Science*, pages 396–405. Springer, 1993.
43. M. P. Schützenberger. On the definition of a family of automata. *Information and Control*, 4(2-3):245–270, 1961.
44. G. Sénizergues. The equivalence problem for deterministic pushdown automata is decidable. In *Proceedings of the 24th International Colloquium on Automata, Languages and Programming, ICALP '97*, pages 671–681, London, UK, UK, 1997. Springer-Verlag.
45. J. Thatcher. Characterizing derivation trees of context-free grammars through a generalization of finite automata theory. *Journ. of Comp. and Syst.Sc.*, 1:317–322, 1967.
46. W. Thomas. Handbook of theoretical computer science (vol. B). chapter Automata on infinite objects, pages 133–191. MIT Press, Cambridge, MA, USA, 1990.
47. B. A. Trakhtenbrot. Finite automata and logic of monadic predicates (in Russian). *Doklady Akademii Nauk SSR*, 140:326–329, 1961.
48. B. von Braunmühl and R. Verbeek. Input-driven languages are recognized in log n space. In *Proceedings of the Symposium on Fundamentals of Computation Theory, Lect. Notes Comput. Sci. 158*, pages 40–51. Springer, 1983.