

# CS 3719 (Theory of Computation and Algorithms) – Pushdown automata and context-free grammars

Antonina Kolokolova\*

January 24, 2019

## 1 Pushdown automata

We just proved that some languages are not regular. And, moreover, the examples we have seen are languages that can easily be computed by a simple algorithm in any modern programming language. Now a natural question is whether it is possible to add a little extra power to NFAs, so that the resulting model of computation would be able to handle languages such as  $\{0^n q^n\}$ . Indeed, it is possible to do so by giving NFAs a little “ability to count”, some unlimited memory. There are several ways of adding memory to NFAs, and we will look at two of them, Pushdown Automata and Turing machine (in short, Pushdown Automata have an access to an unlimited stack, and Turing machines to an unlimited tape.) In this lecture, we will look at Pushdown automata and analyze its power. Later we will show that Pushdown Automata still fall short of computing many languages that we view as easily computable by our usual algorithmic techniques.

Informally, a pushdown automaton is just an NFA with a stack. So the additional part of the description of such an automaton should include transitions that operate with the stack, as well as stack alphabet.

**Definition 7.** A pushdown automaton (PDA) is a 6-tuple  $(Q, \Sigma, \Gamma, \delta, q_0, F)$  where  $Q$  is the set of states,  $\Sigma$  the input alphabet,  $\Gamma$  the stack alphabet,  $q_0$  the start state,  $F$  is the set of finite states and the transition function  $\delta : Q \times (\Sigma \cup \{\epsilon\}) \times (\Gamma \cup \{\epsilon\}) \rightarrow \mathcal{P}(Q \times (\Gamma \cup \{\epsilon\}))$ .

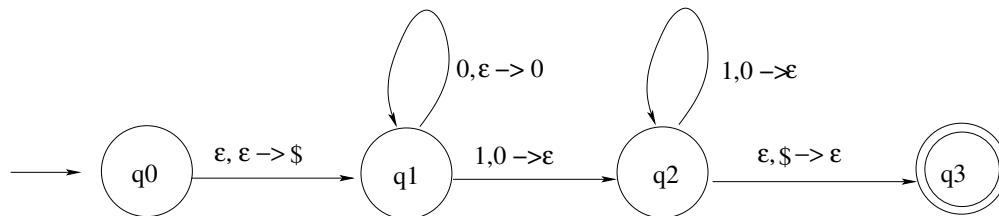
That is, each transition of a PDA pops a symbol (possibly  $\epsilon$ , corresponding to popping nothing) off the stack;  $\delta$  specifies which set of states to go from the current state on reading an input symbol and a stack symbol, and for every such state, which, if any, symbol to push onto the stack.

---

\*The material in this set of notes came from many sources, in particular “Introduction to Theory of Computation” by Sipser and course notes of U. of Toronto CS 364.

Note that this definition extends the definition of a non-deterministic finite automaton. It is possible to define deterministic pushdown automata by similarly extending the definition of a DFA. However, there are languages accepted by non-deterministic PDAs that no deterministic one can accept (such as a set of palindromes). The proof of this is beyond the scope of this course; however, it is good to remember this as a case where non-determinism actually results in a more powerful model of computation: it didn't for the finite automata. Later in the course we will talk about the P vs. NP problem, a major open problem in computer science which asks about the role of non-determinism for feasible computation.

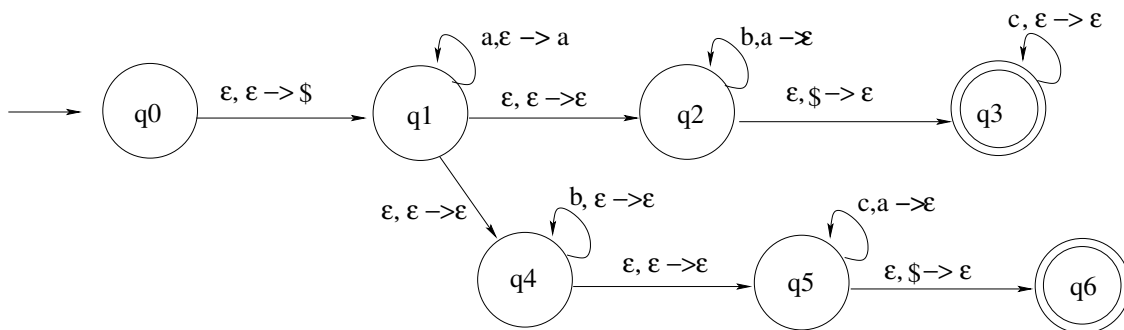
**Example 1.** Recall the language  $\{0^n 1^n\}$  which we have shown to be non-regular in previous lecture. The following PDA accepts this language.



Here, the tape alphabet is  $\Gamma = \{0, \$\}$ , where  $\$$  is a special symbol used as an empty stack marker. Since we never need to put 1 on the stack in this PDA, it is OK not to have 1 in  $\Gamma$ , although it is common to take  $\Sigma \subseteq \Gamma$ .

Let us do one more example of a Pushdown Automaton. In this case, we will consider a language where non-determinism is really unavoidable:  $\{a^i b^j c^k \mid i = j \vee i = k\}$ .

**Example 2.** The following PDA recognizes the language  $\{a^i b^j c^k \mid i = j \vee i = k\}$ .



Here, the tape alphabet is  $\Gamma = \{a, \$\}$ . Note the non-deterministic choice this PDA makes: in the state  $q_1$  it forks between the automaton matching letters  $b$  to letters  $a$  on the stack, or matching the  $c$ 's.

## 2 Context-Free Grammars

We have shown earlier that regular languages can be described by regular expressions. It is natural to ask is there a similar description for the languages computed by pushdown automata. Indeed there is, and it is a natural formalism that has been used for a long time on its own: context-free grammars. Noam Chomsky defined them as one of the types in his hierarchy; the context-free grammar syntax has been used in the programming language community to describe the syntax of programming languages since Algol (known there as Backus-Naur Form).

When you think of a grammar the first thing that might come to mind is studying a natural language such as French in school, and all these rules about nouns and verbs. Indeed, it is possible to do much of natural language processing (although not everything) using context-free grammar formalism.

**Example 3.** Consider the following rules:

$\langle sentence \rangle$	$- \rangle \langle nounphrase \rangle \langle verbphrase \rangle$
$\langle verbphrase \rangle$	$- \rangle \langle verb \rangle \mid \langle verb \rangle \langle nounphrase \rangle$
$\langle nounphrase \rangle$	$- \rangle \text{Jane} \mid \text{the assignment}$
$\langle verb \rangle$	$- \rangle \text{solved} \mid \text{did} \mid \text{decided}$

For example, a sentence “Jane did the assignment” can be generated by this grammar, and also “Jane solved”. On the other hand, this grammar can generate “The assignment solved Jane”, which might not be a desired result.

Now that we have seen an example, let’s define formally what is a context-free grammar.

**Definition 8.** A context-free grammar (CFG) is a 4-tuple  $(V, \Sigma, R, S)$ , where

- 1)  $V$  is a finite set of variables, with  $S \in V$  the start variable.
- 2)  $\Sigma$  is a finite set of **terminals** (disjoint from the set of variables).
- 3)  $R$  is a finite set of rules, with each rule consisting of a variable followed by  $- \rangle$  followed by a string of variables and terminals.

A grammar is a set of substitution rules, where a variable at the head of a rule can be substituted by the string in the body of a rule whenever that variable occurs. Let  $A \rightarrow w$  be a rule of the grammar, where  $w$  is a string of variables and terminals. Then  $A$  can be

replaced in another rule by  $w$ :  $uAv$  in a body of another rule can be replaced by  $uwv$  (we say  $uAv$  yields  $uwv$ , denoted  $uAv \Rightarrow uwv$ ). If there is a sequence  $u = u_1, u_2, \dots, u_k = v$  such that for all  $i$ ,  $1 \leq i < k$ ,  $u_i \Rightarrow u_{i+1}$  then we say that  $u$  derives  $v$  (denoted  $v \stackrel{*}{\Rightarrow} u$ .)

**Definition 9.** If  $G$  is a context-free grammar, then the language of  $G$  is the set of all strings of terminals that can be generated from the start variable:  $\mathcal{L}(G) = \{w \in \Sigma^* \mid S \stackrel{*}{\Rightarrow} w\}$ . A language is called a context-free language (CFL) if there exists a CFG generating it.

In the above example, variables are  $\langle sentence \rangle$ ,  $\langle nounphrase \rangle$ ,  $\langle verbphrase \rangle$ , and  $\langle verb \rangle$ , with the start variable  $S = \langle sentence \rangle$ , the  $\Sigma = \{ \text{Jane, the assignment, solved, did, decided} \}$  and  $R$  consists of the rules listed above. To be more precise, there are more rules than 4: we used a shortcut “—” for several different rules starting with the same variable. E.g, we should have had two rules  $\langle nounphrase \rangle \rightarrow \text{Jane}$ , and  $\langle nounphrase \rangle \rightarrow \text{the assignment}$ .

To see how a CFG can generate a nonregular language, consider the language  $\{0^n 1^n\}$  for which we constructed a PDA in the last lecture.

**Example 4.** The following CFG generates the language  $\{0^n 1^n\}$ :

$$A \rightarrow 0A1 \mid \epsilon$$

A derivation of a string 000111 in this grammar is:

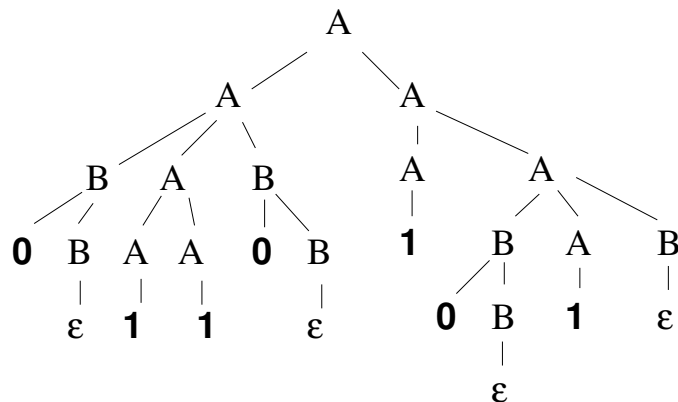
$$A \rightarrow 0A1 \rightarrow 00A11 \rightarrow 000A111 \rightarrow 000111$$

The last step is because adding or removing an empty substring  $\epsilon$  does not change a string.

We often visualize a derivation as a tree (a parse tree), with variables as internal nodes and terminals as leaves.

**Example 5.** Consider the following grammar. Here is a parse tree for a string 0110101. Note that there are several possible parse trees for this string. In this case, we say that a grammar is *ambiguous*.

$$\begin{aligned} A &\rightarrow AA \\ A &\rightarrow BAB \\ B &\rightarrow OB \\ B &\rightarrow \epsilon \\ A &\rightarrow 1 \end{aligned}$$



## 2.1 Arithmetic expressions

A canonical example of use of context-grammars is in parsing. In particular, here we will see how to parse an arithmetic expression using context-free grammars.

**Example 6.** Consider the following grammar  $G_1$ :

$$EXPR \rightarrow EXPR + EXPR | EXPR * EXPR | (EXPR) | x | y | z | 0 | 1 | 2 | 3$$

This grammar generates arithmetic expressions such as  $x + 2 * y$ . However, there is a problem: it generates it ambiguously, ignoring precedence rules. So evaluating the expression according to the tree might give different answers, depending whether the first rule applied was multiplication or addition.

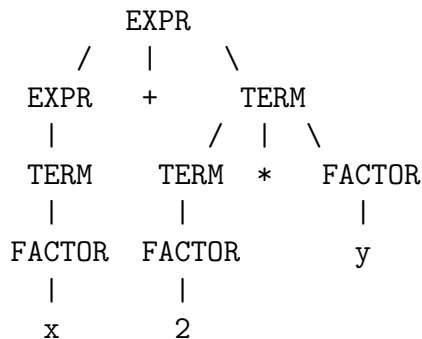


How would we modify the grammar to make it respect precedence rules? One way of doing it is to give different names to parts of a sum vs. parts of a product and treat them differently.

**Example 7.** Consider the following grammar  $G_1$ , with  $EXPR$  the start symbol.

$$\begin{aligned}
 EXPR &\rightarrow EXPR + TERM | TERM \\
 TERM &\rightarrow TERM * FACTOR | FACTOR \\
 FACTOR &\rightarrow (EXPR) | x | y | z | 0 | 1 | 2 | 3
 \end{aligned}$$

Note that in this case for the arithmetic expression there is only one possible parse tree:



### 3 Pumping lemma for CFLs

Just like we proved that certain languages are not regular by defining a property (pumping lemma) which all regular languages satisfy, and then showing that some languages don't, we will define a same kind of property for context-free languages.

Consider a simple grammar  $A \rightarrow aAb, A \rightarrow c$ . It is clear that it generates a non-regular language: a pumping lemma proof would say that either a's or b's would have to be pumped separately, or c would be repeated multiple times, resulting in strings not in the language. However, if we could pump both a's and b's at the same time, this would work.

With this intuition, we state the pumping lemma for context-free languages.

**Lemma 7.** *Let  $L$  be a context-free language. Then there exists a natural number  $p$  such that  $\forall s \in L, |s| \geq p, s = uvxyz$  where*

- 1)  $\forall i \geq 0, uv^i xy^i z \in L$
- 2)  $|vy| > 0$
- 3)  $|vxy| \leq p$ .

That is, if regular languages could be split into three parts so that the middle part could be iterated without creating strings not in the language, for CFLs the split is into 5 parts so that the middle and the sides stay, and the 2nd and 4th parts are iterated simultaneously. For the grammar we just saw  $u$  can consist of several  $a$ 's, and  $y$  of the same number of  $b$ 's; the middle part  $x$  contains the  $c$ .

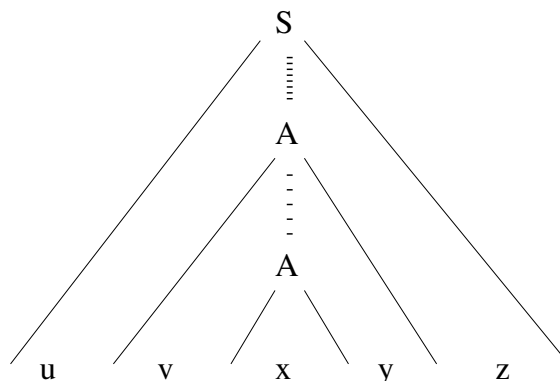
*Proof.* In the proof for regular languages, we found a repeating part by looking at the states of the automaton and finding a repeating states. Here, we will do the same, but with parse trees.

Here is an idea. Consider a parse tree which is so high that its height is more than the number of variables in the grammar. Why should such a tree exist? It is because if the language is infinite, then there are infinitely many parse trees, at least one distinct tree for a different string in the language (remember that a string is a sequence of leaves of the parse tree). There are finitely many rules, and each rule is of finite length, so the only way to have infinitely many trees is to allow them to grow to an arbitrary height.

Let's look more carefully at the possible parse trees. Suppose that the number of symbols (variables and terminals) in the body of the longest rule is  $d$ : then our tree would be  $d$ -ary (every node having at most  $d$  children). We need to know what is the shortest string guaranteed to have a parse tree of height at least  $|V| + 1$ . Recall that the longest possible

string produced by a  $d$ -ary tree of height  $h$  has length  $d^h$ : it corresponds to a sequence of leaves of a complete  $d$ -ary tree of height  $h$ . Therefore, if we take a string of length at least  $d^{|V|} + 1$  it is guaranteed to have all parse trees of height at least  $|V| + 1$ . For convenience, we can choose even larger  $p$ ; let  $p = d^{|V|+1}$ .

Now, why would we be interested in trees of this form? Consider the largest path from the root to a leaf in this tree. Since the height of the tree is greater than the number of variables, there would be a repeating variable, say  $A$ , on this path. A variable occurring in a parse tree means that a rule with this variable at the head was applied. But now notice that there is no way in the grammar to differentiate these occurrences and say how many there should be; thus, a subtree rooted at the first occurrence of  $A$  can be replaced by subtree in the second occurrence, and vice versa. The first one gives us  $i = 0$  from the pumping lemma, and iterating the second one gives any  $i$  we desire (see the figure).



To check the second condition, suppose that we start with the smallest possible parse tree for the string. Then doing the substitution of the first occurrence by the second would give us an even smaller tree, which is a contradiction.

Finally, let's prove the third condition. Starting from the bottom of the tree, choose the first time from the bottom a variable repeats; call that variable  $A$ . A subtree under  $A$  is of height at most  $|V| + 1$ , so it generates a string of length at most  $p$ . Here, the second occurrence of  $A$  generates  $xy$  where the bottom occurrence generates  $x$ .  $\square$

Now, let's see how this lemma can be used to show that some languages are not context-free.

**Example 8.** The language  $\{a^n b^n c^n\}$  is not context-free.

Suppose it is. Then let  $p$  be the pumping length, and take  $s = a^p b^p c^p$ . Now, just like in the proof that  $\{0^n 1^n\}$  is not regular, no matter how we split  $s = uvxyz$ , either  $v$  and  $y$  are monochromatic and thus repeating them disrupts the equality of the numbers, or they are not monochromatic and repeating them disrupts the order.

Now, recall that the language  $\{a^n b^n c^m\}$  is context-free: it is a special case of  $\{a^i b^j c^k \mid i = j \text{ or } i = k\}$ . Also,  $\{a^n b^m c^n\}$  is context-free, for the same reason. But their intersection is exactly  $\{a^n b^n c^n\}$ , which we just shown to be not context-free! This leads us to a surprising corollary:

**Corollary 8.** *The class of context-free languages are not closed under intersection.*

**Example 9.** The language  $\{ww|w \in \{0,1\}^*\}$  is not context-free.

Take  $s = 0^p 1^p 0^p 1^p$ . By the 3rd condition of the pumping lemma, that  $|vxy| \leq p$ , the sequence  $vxy$  overlaps at most two out of four blocks in our string. Either it is part of  $0^p 1^p$ ; then repeating any part of it will disrupt the equality with the second half of  $s$ . Or it overlaps the  $1^p 0^p$  piece in the middle; but then either  $v$  or  $y$  are monochromatic: say  $y$  is, and it consists of just 0s. But now repeating  $y$  will increase the number of 0s in the second copy of  $w$ , but not the first, again disrupting the equality.