Context-Free Languages Thoughts on Assignment 1

CS 331 Programming Languages Lecture Slides Wednesday, January 22, 2025

Glenn G. Chappell Department of Computer Science University of Alaska Fairbanks ggchappell@alaska.edu

© 2017–2025 Glenn G. Chappell

Unit Overview Formal Languages & Grammars

Topics

- ✓ Basic concepts
- Introduction to formal languages & grammars
- The Chomsky hierarchy
- Regular languages
- Regular languages & regular expressions
 - Context-free languages
 - Programming language syntax specification

Review

Review Regular Languages — Regular Grammars & Languages

A **regular grammar** is a grammar, each of whose productions looks like one of the following. We allow a production using the same

- $A \rightarrow bC$ $A \to \varepsilon \qquad A \to b$
- nonterminal twice: $A \rightarrow bA$
- A **regular language** is a language that is generated by some regular grammar.

This is a regular grammar: $S \rightarrow \varepsilon$ $S \rightarrow t$ $S \rightarrow xB$ $B \rightarrow yS$

Therefore, the language it generates is a regular language:

 $\{\varepsilon, xy, xyxy, xyxyxy, ..., t, xyt, xyxyt, xyxyxyt, ...\}$

2025-01-22

CS 331 Spring 2025

A deterministic finite automaton (Latin plural "automata"), or DFA, is a kind of recognizer for regular languages.



Rule. For each character in the alphabet, each state has *exactly* one transition leaving it that is associated with that character.

To use a DFA as a recognizer:

- Start in the start state; proceed in a series of steps.
- At each step, read a character from the input and follow the transition from at the current state, labeled with that character.
- If, when we reach the end of the input, we are in an accepting state, then we accept the input.
- The set of all inputs that are accepted is the language **recognized** by the DFA.

Exercise

4. What language is recognized by the DFA whose diagram is shown?



To use a DFA as a recognizer:

- Start in the start state; proceed in a series of steps.
- At each step, read a character from the input and follow the transition from at the current state, labeled with that character.
- If, when we reach the end of the input, we are in an accepting state, then we accept the input.
- The set of all inputs that are accepted is the language **recognized** by the DFA.

Answer

4. What language is recognized by the DFA whose diagram is shown?

The set of strings consisting either of *b* then zero or more *c*'s, or of one or more *c*'s.

{*b*, *bc*, *bcc*, *bccc*, *...*, *c*, *cc*, *ccc*, *...*}



A **regular expression** is a generator for a regular language.

We specified the syntax of regular expressions by showing how to build them up from small pieces.

- A *single character* is a regular expression: *a*.
- The *empty string* is a regular expression: *ε*.

If A and B are regular expressions, then so are the following.

- A* Kleene star
- AB (say KLAY-nee)
- A|B

The above are listed from high to low precedence. All are leftassociative. Override precedence using parentheses.

• (A)

Here is a regular expression: $(a|x)^*cb$

A regular expression is said to **match** certain strings.

Semantics of regular expressions:

- A single character matches itself, and nothing else.
- The empty string matches itself, and nothing else.
- A* matches the concatenation of zero of more strings, each of which is matched by A.
 - Note that A* matches the empty string, no matter what A is.
- *AB* matches the concatenation of any string matched by *A* and any string matched by *B*.
- A|B matches all strings matched by A and also all strings matched by B.
- (*A*) matches the same strings that are matched by *A*.

The language **generated** by a regular expression consists of all strings that it matches.

Regular expression libraries typically include shortcuts in their syntax. The ones below do *not* change which languages can be generated. They may be used this semester in answers to assignments/quizzes/exams.

Matches any single character

- [abcd] Matches a single a, b, c, or d, like a|b|c|d
- [a-d] Same as above
- [^a-d] Matches anything except a, b, c, or d
- x+ Matches one or more x characters: x, xx, xxx, xxx, etc.
- x? Matches zero or one x characters—an optional x
- Matches a dot: "."
 Matches a backslash: "\"
 Using a backslash in this way is called escaping.

Slashes might be delimiters: / [0-9] / matches any ASCII digit.

CS 331 Spring 2025

The regular languages are precisely:

- The languages generated by regular grammars.
- The languages recognized by DFAs.
- The languages generated by regular expressions (in the strict sense, with no features beyond those covered).

That is, these three classes of languages are identical.

We will use ideas about regular languages when we do **lexical analysis** (**lexing**): breaking up a program into **lexemes** (words, roughly).

Context-Free Languages

We now turn to the second smallest class of languages in the Chomsky hierarchy: the *context-free languages*.

- Context-free languages are important because, for most programming languages, the set of all syntactically correct programs forms a context-free language.
- And for those PLs that do not have this property, it is still common for the *techniques* used in dealing with context-free languages to be useful.
- Context-free languages, and the associated grammars, are thus important in **parsing**: determining whether input (for example, a program) is syntactically correct, and, if so, finding its structure.

- A **context-free grammar** (**CFG**) is a grammar, each of whose productions has a left-hand side consisting of a single nonterminal.
- All of the grammars we have looked at have been CFGs. In particular, every regular grammar is a CFG.
- A **context-free language** (**CFL**) is a language that is generated by some context-free grammar.
- Every regular language is a CFL. But there are context-free languages that are not regular.

Context-Free Languages Context-Free Grammars & Languages — Examples [1/2]

Here is a CFG.

- $S \rightarrow aSa$
- $S \rightarrow b$

"Context-free" refers to the fact that a nonterminal can be expanded at any time. Chomsky defined a larger class of grammars, *context-sensitive grammars*, in which productions can sometimes only be applied if a nonterminal has certain characters around it, that is, only in a certain *context*. We will not study context-sensitive grammars this semester.

This grammar generates the following language.

{b, aba, aabaa, aaabaaa, aaaabaaaa, ...}

We can also write this language as follows.

 $\{ a^k b a^k \mid k \ge 0 \}$

As we have noted, this is not a regular language. But since it is generated by a CFG, it is a CFL.

Regular grammars are not powerful enough to handle things like matching parentheses. But CFGs are powerful enough.

Consider the following grammar—where "(" and ")" are terminal symbols.

$$S \rightarrow SS$$

 $S \rightarrow (S)$
 $S \rightarrow \varepsilon$

The language generated by the above grammar consists of all sequences of properly matched parentheses. For example, here is one string in this language.

It is common for a CFG to have multiple productions with the same left-hand side. As a shortcut, we allow writing the left-hand side and the arrow only once, with the various right-hand sides separated by vertical bars ("|").

For example, our first CFG can be rewritten as follows.

 $S \to aSa$ $\longrightarrow S \to aSa \mid b$

We might place the right-hand sides on separate lines.

$$S \rightarrow aSa$$

| b

Parsing involves finding the structure of a program. One way to represent this structure is to use a **parse tree**.

We introduce parse trees using *Grammar A*, below. To the right is a derivation for the string *ppy* based on this CFG.

Grammar A	Derivation of ppy
$S \rightarrow AB$	<u>S</u>
$A \rightarrow pA \mid \varepsilon$	A <u>B</u>
$B \rightarrow x \mid y$	<u>A</u> y
	р <u>А</u> у
There are no parse trees on this slide! See the next slide for a drawing of a parse tree.	рр <u>А</u> у рру

Context-Free Languages Parse Trees — Definition [2/3]

Grammar A	Derivation of ppy
$S \rightarrow AB$	<u>S</u>
$A \rightarrow pA \mid \varepsilon$	A <u>B</u>
$B \rightarrow x \mid y$	<u>A</u> y
	р <u>А</u> у
A parse tree is a rooted tree with one symbol	рр <u>А</u> у
in each node, based on a derivation.The root node holds the start symbol.	рру
 The symbols a nonterminal is expanded into left to right, one symbol per tree node. 	become its children— Parse
Here is a parse tree based on the above derivation	ation. Tree S
Every terminal symbol is in a leaf of the parse We can read off the final string by looking a leaves that contain terminal symbols.	tree. A B It the P A y p A

Context-Free Languages Parse Trees — Definition [3/3]

Another grammar, derivation, and associated parse tree. Here, "+" is a terminal symbol.



Again, we can read off the final string by looking at the leaves that contain terminal symbols.

Context-Free Languages Parse Trees — TRY IT (Exercises)

Grammar C	Derivation of ac
$S \rightarrow XY$	<u>S</u>
$X \rightarrow a \mid b$	$\underline{X}Y$
$Y \rightarrow CY \mid C$	a <u>Y</u>
	ас

Exercises

- 1. Based on Grammar C and the given derivation, draw a parse tree for the string *ac*.
- 2. Based on Grammar C, draw a parse tree for the string *bcc*.

Context-Free Languages Parse Trees — TRY IT (Answers)

Gr	ammar C	Derivatio	n of <i>ac</i>
<i>S</i> -	$\rightarrow XY$	<u>S</u>	
Х -	→a b	<u>X</u> Y	
Y –	→ cY c	а <u>Ү</u>	
		ас	
An	swers		S
1.	Based on Grammar C and the g derivation, draw a parse tree for	iven or the string <i>ac</i> . <i>S</i>	X Y I I a C
2.	Based on Grammar C, draw a parse tree for the string <i>bcc</i> .	X Y b c Y l	In both exercises, the answer is unique.

If you drew something like this for Exercise 1 ...



... then be aware that the above is not a parse tree. A parse tree has *one symbol in each node*.





A CFG in which a single string has more than one parse tree, is said to be **ambiguous**.

So Grammar B is ambiguous.

2025-01-22

CS 331 Spring 2025

Context-Free Languages Ambiguity — Eliminating Ambiguity [1/3]



Ambiguity is a property of *grammars*, not of languages. And it is generally a property that we do not like.

Grammar B is ambiguous; however, in this case we can actually find a non-ambiguous CFG that generates the same language. Before finding such a grammar, we first note that, assuming "+" represents addition, we prefer parse tree #1, since it expresses the left associativity that we usually want addition to have: n+n+n = (n+n)+n.

Context-Free Languages Ambiguity — Eliminating Ambiguity [2/3]



Below is a non-ambiguous grammar that generates the same language and expresses the left-associativity of "+". Also shown: a derivation of n+n+n and the unique parse tree.

Derivation	Parse Tree S
<u>S</u>	This parse tree $S + n$
<u>S</u> +n	is unique!
<u>S</u> +n+n	S + n
n+n+n	n
	Derivation <u>S</u> <u>S</u> +n <u>S</u> +n+n n+n+n

Sometimes ambiguity cannot be eliminated. There are CFLs that are only generated by ambiguous CFGs. Such a CFL is **inherently ambiguous**.

Here is a standard example of an inherently ambiguous CFL.

$$\{ a^{m}b^{m}c^{n}d^{n} \mid m, n \geq 0 \} \cup \{ a^{p}b^{r}c^{r}d^{p} \mid p, r \geq 0 \}$$
Same

It can be demonstrated that, no matter how we write a CFG for this language, there will be some string that has two different parse trees.

Remember:

- Ambiguity is a property of grammars (CFGs).
- Inherent ambiguity is a property of *languages* (CFLs).

The CFG below generates only xyz. There are multiple derivations.

Grammar D	Derivation #1	Derivation #2	Derivation #3
$S \rightarrow ABC$	<u>S</u>	<u>S</u>	<u>S</u>
$A \rightarrow x$	<u>A</u> BC	AB <u>C</u>	А <u>В</u> С
$B \rightarrow y$	<u>х</u> <u>В</u> С	A <u>B</u> z	Ау <u>С</u>
$C \rightarrow Z$	<i>ху<u>С</u></i>	<u>A</u> yz	<u>A</u> yz
	xyz	XYZ	xyz

But there is only one parse tree. Grammar D is *not* ambiguous.



CS 331 Spring 2025

Even though they correspond to the same parse tree, these three derivations of *xyz* differ in a noteworthy way.

- In derivation #1, the leftmost nonterminal it expanded at each step. We call this a leftmost derivation.
- In derivation #2, the rightmost nonterminal it expanded at each step. We call this a rightmost derivation.
- Derivation #3 is neither leftmost nor rightmost.

#1: Leftmost derivation	#2: Rightmost derivation	#3: Neither
<u>S</u>	<u>S</u>	<u>S</u>
<u>A</u> BC	<i>AB<u>C</u></i>	А <u>В</u> С
х <u>В</u> С	A <u>B</u> z	Ау <u>С</u>
<i>ху<u>С</u></i>	<u>A</u> yz	<u>A</u> yz
ХУZ	XYZ	хуz

#1: Leftmost derivation	#2: Rightmost derivation	#3: Neither
<u>S</u>	<u>S</u>	<u>S</u>
<u>A</u> BC	<i>AB<u>C</u></i>	А <u>В</u> С
<u>х</u> <u>В</u> С	A <u>B</u> z	Ау <u>С</u>
<i>ху<u>С</u></i>	<u>A</u> yz	<u>A</u> yz
XYZ	XYZ	ХУZ

These concepts will come up later, in our study of parsing.

- A parser goes through the steps required to find a derivation. Some parsers go through the derivation in forward order, **expanding** the leftmost nonterminal first, producing a leftmost derivation.
- Other parsers go though the derivation in reverse, repeatedly contracting a substring to a nonterminal. Typically, the left part of the input is contracted first. Viewed in forward order, the rightmost nonterminal is expanded first, producing a rightmost derivation.

Do not confuse parse trees with derivations!

- A parse tree is a *rooted tree*.
- A derivation is a *list of strings*.

For every parse tree, there is a corresponding leftmost derivation and rightmost derivation.

Ambiguity is about multiple parse trees, not multiple derivations.



Thoughts on Assignment 1

Three quick notes:

- Turn in your work as a **PDF** file.
- The point of the first exercise is making sure you are able to execute Lua code. You will be writing a lot of Lua this semester. Now is the time to be sure you can execute it.
- We have covered the required material for all parts of the assignment, except the last exercise. This concerns something called *BNF grammars*, which we will discuss next time.