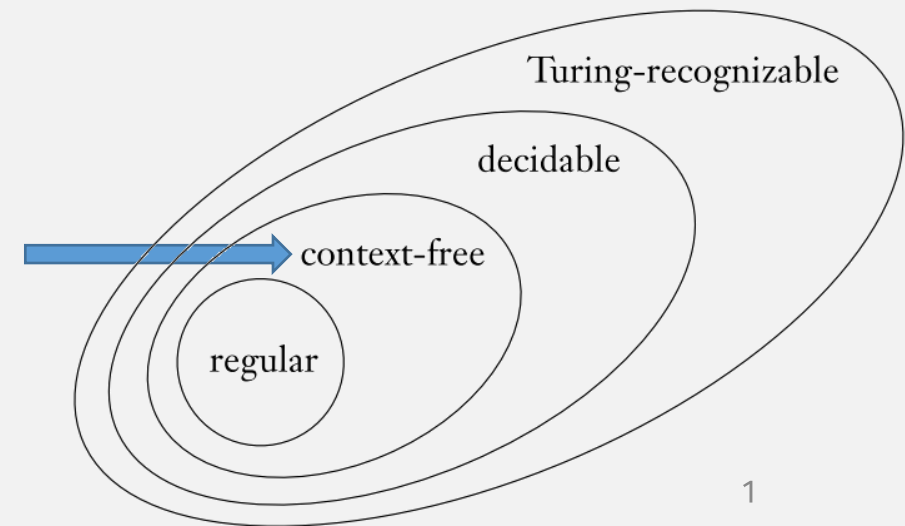


UMB CS 420

Context-Free Languages (CFLs)

Wednesday, February 16, 2022



Announcements

- HW3, Problem 3 (rev strings) has a new requirement!
 - See piazza post and hw3 page
- HW3 due Sun 2/20 11:59pm EST
- Reminder: **No class next Monday 2/21**

Last Time:

Pumping lemma If A is a regular language, then there is a number p (the pumping length) where if s is any string in A of length at least p , then s may be divided into three pieces, $s = xyz$, satisfying the following conditions:

1. for each $i \geq 0$, $xy^iz \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

Let B be the language $\{0^n 1^n \mid n \geq 0\}$. We use the pumping lemma to prove that B is not regular. The proof is by contradiction.

- **Assume** the language is regular
- So it must follow the Pumping Lemma:
 - All strings longer than length p ...
 - ... must be splittable into xyz ... where y is “pumpable”
- Find **counterexample** where Pumping Lemma does not hold: $0^p 1^p$
- Therefore, the language is not regular
 - This is the contrapositive of the Pumping Lemma
 - And also a **contradiction** of the assumption!

Last Time:

Pumping lemma If A is a regular language, then there is a number p (the pumping length) where if s is any string in A of length at least p , then s may be divided into three pieces, $s = xyz$, satisfying the following conditions:

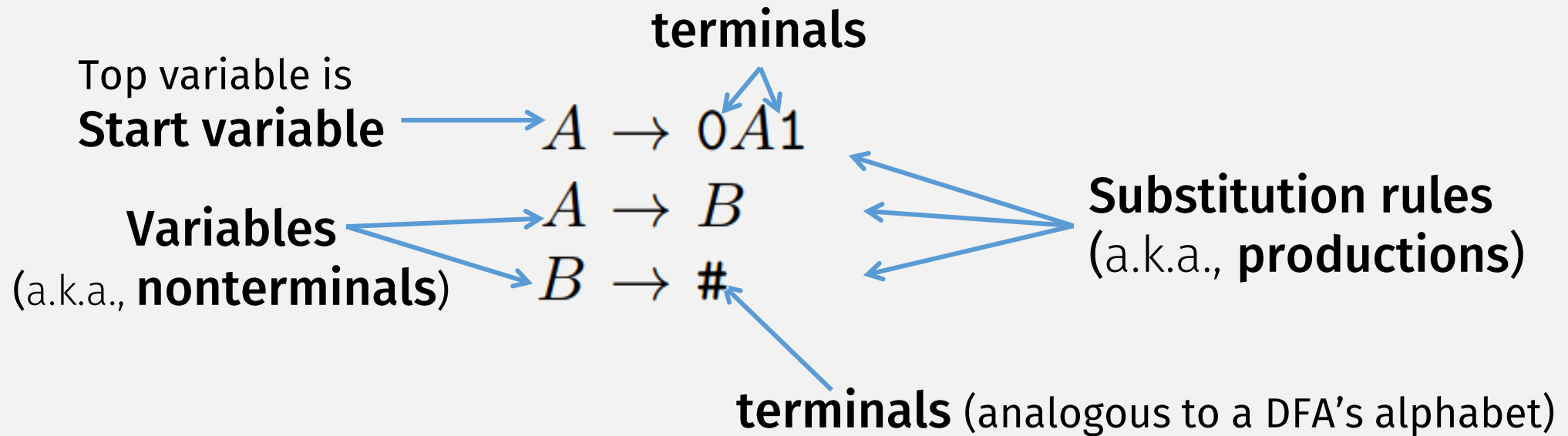
1. for each $i \geq 0$, $xy^iz \in A$,
2. $|y| > 0$, and
3. $|xy| \leq p$.

Let B be the language $\{0^n 1^n \mid n \geq 0\}$. We use the pumping lemma to prove that B is not regular. The proof is by contradiction.

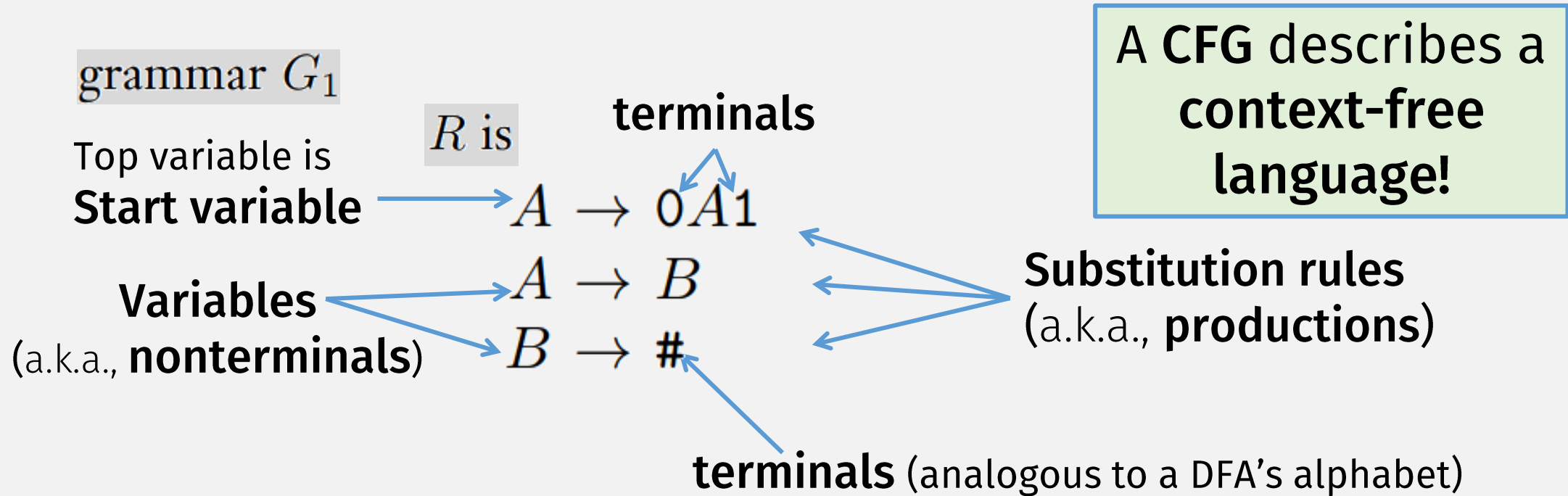
If this language is not regular, then what is it???

Maybe? ... a **context-free language (CFL)**?

A Context-Free Grammar (CFG)



A Context-Free Grammar (CFG)



A *context-free grammar* is a 4-tuple (V, Σ, R, S) , where

1. V is a finite set called the *variables*,
2. Σ is a finite set, disjoint from V , called the *terminals*,
3. R is a finite set of *rules*, with each rule being a variable and a string of variables and terminals, and
4. $S \in V$ is the start variable.

$$V = \{A, B\},$$

$$\Sigma = \{0, 1, \#\},$$

$$S = A,$$

Java Syntax: Described with CFGs

ORACLE

[Java SE](#) > [Java SE Specifications](#) > [Java Language Specification](#)

Chapter 2. Grammars

[Prev](#)

Chapter 2. Grammars

This chapter describes the **context-free grammars** used in this specification to define the lexical and syntactic structure of a program.

2.1. Context-Free Grammars

A *context-free grammar* consists of a number of **productions**. Each production has an abstract symbol called a **nonterminal** as its *left-hand side*, and a sequence of one or more nonterminal and **terminal symbols** as its *right-hand side*. For each grammar, the terminal symbols are drawn from a specified *alphabet*.

Starting from a sentence consisting of a single distinguished nonterminal, called the *goal symbol*, a given context-free grammar specifies a language, namely, the set of possible sequences of terminal symbols that can result from repeatedly replacing any nonterminal in the sequence with a right-hand side of a production for which the nonterminal is the left-hand side.

2.2. The Lexical Grammar

A *lexical grammar* for the Java programming language is given in §3. This grammar has as its terminal symbols the characters of the Unicode character set. It defines a set of productions, starting from the goal symbol *Input* (§3.5), that describe how sequences of Unicode characters (§2.1) are translated into a sequence of input elements (§2.5).

(partially)

Python Syntax: Described with a CFG

10. Full Grammar specification

This is the full Python grammar, as it is read by the parser generator and used to parse Python source files:

```
# Grammar for Python

# NOTE WELL: You should also follow all the steps listed at
# https://devguide.python.org/grammar/

# Start symbols for the grammar:
#     single_input is a single interactive statement;
#     file_input is a module or sequence of commands read from an input file;
#     eval_input is the input for the eval() functions.
#     func_type_input is a PEP 484 Python 2 function type comment
# NB: compound_stmt in single_input is followed by extra NEWLINE!
# NB: due to the way TYPE_COMMENT is tokenized it will always be followed by a NEWLINE
single_input: NEWLINE | simple_stmt | compound_stmt NEWLINE
file_input: (NEWLINE | stmt)* ENDMARKER
eval_input: testlist NEWLINE* ENDMARKER
```

(indentation checking
probably not
describable with a CFG)

Many Other Language (partially)

~~Python~~ Syntax: Described with a CFG

10. Full Grammar specification

This is the full Python grammar, as it is read by the parser generator and used to parse Python source files:

```
# Grammar for Python

# NOTE WELL: You should also follow all the steps listed at
# https://devguide.python.org/grammar/

# Start symbols for the grammar:
#     single_input is a single interactive statement;
#     file_input is a module or sequence of commands read from an input file;
#     eval_input is the input for the eval() functions.
#     func_type_input is a PEP 484 Python 2 function type comment
# NB: compound_stmt in single_input is followed by extra NEWLINE!
# NB: due to the way TYPE_COMMENT is tokenized it will always be followed by a NEWLINE
single_input: NEWLINE | simple_stmt | compound_stmt NEWLINE
file_input: (NEWLINE | stmt)* ENDMARKER
eval_input: testlist NEWLINE* ENDMARKER
```

Generating Strings with a CFG

A **CFG** represents a **context free language!**

Strings in CFG's language
= all possible generated strings

$L(G_1)$ is $\{0^n \# 1^n \mid n \geq 0\}$

$G_1 =$

1st rule $\rightarrow A \rightarrow 0A1$
 $A \rightarrow B$
 $B \rightarrow \#$

Stop when string is all terminals

A CFG **generates** a string, by repeatedly applying substitution rules:

$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$

Start variable

After applying 1st rule

1st rule again

1st rule again

Use 2nd rule

Use last rule ⁸

Derivations: Formally

A *context-free grammar* is a 4-tuple (V, Σ, R, S) , where

1. V is a finite set called the *variables*,
2. Σ is a finite set, disjoint from V , called the *terminals*,
3. R is a finite set of *rules*, with each rule being a variable and a string of variables and terminals, and
4. $S \in V$ is the start variable.

Let $G = (V, \Sigma, R, S)$

Single-step

$$\alpha A \beta \xRightarrow{G} \alpha \gamma \beta$$

Where:

$$\alpha, \beta \in (V \cup \Sigma)^* \leftarrow \begin{array}{|l} \text{Strings of terminals} \\ \text{and variables} \end{array}$$

$$A \in V \leftarrow \begin{array}{|l} \text{Variable} \end{array}$$

$$A \rightarrow \gamma \in R \leftarrow \begin{array}{|l} \text{Rule} \end{array}$$

Extended Derivation

Base case: $\alpha \xRightarrow{G}^* \alpha$ (0 steps)

Recursive case: (multistep)

• If $\alpha \xRightarrow{G} \beta$ and $\beta \xRightarrow{G}^* \gamma$

Single step

Recursive call

• Then: $\alpha \xRightarrow{G}^* \gamma$

Formal Definition of a CFL

A *context-free grammar* is a 4-tuple (V, Σ, R, S) , where

1. V is a finite set called the *variables*,
2. Σ is a finite set, disjoint from V , called the *terminals*,
3. R is a finite set of *rules*, with each rule being a variable and a string of variables and terminals, and
4. $S \in V$ is the start variable.

$$G = (V, \Sigma, R, S)$$

$$L(G) = \left\{ w \in \Sigma^* \mid S \xRightarrow[G]{*} w \right\}$$

Any language that can be generated by some context-free grammar is called a *context-free language*

Flashback: $\{0^n 1^n \mid n \geq 0\}$

- Pumping Lemma says it's not a regular language
- It's a context-free language!
 - Proof?
 - Come up with CFG describing it ...
 - Hint: It's similar to:

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \cancel{\#} \epsilon$$

$$L(G_1) \text{ is } \{0^n \cancel{\#} 1^n \mid n \geq 0\}$$

A String Can Have Multiple Derivations

$$\begin{aligned}\langle \text{EXPR} \rangle &\rightarrow \langle \text{EXPR} \rangle + \langle \text{TERM} \rangle \mid \langle \text{TERM} \rangle \\ \langle \text{TERM} \rangle &\rightarrow \langle \text{TERM} \rangle \times \langle \text{FACTOR} \rangle \mid \langle \text{FACTOR} \rangle \\ \langle \text{FACTOR} \rangle &\rightarrow (\langle \text{EXPR} \rangle) \mid a\end{aligned}$$

Want to generate this string: **a + a × a**

- EXPR ⇒
- EXPR + TERM ⇒
- EXPR + TERM × FACTOR ⇒
- EXPR + TERM × a ⇒
- ...

RIGHTMOST DERIVATION

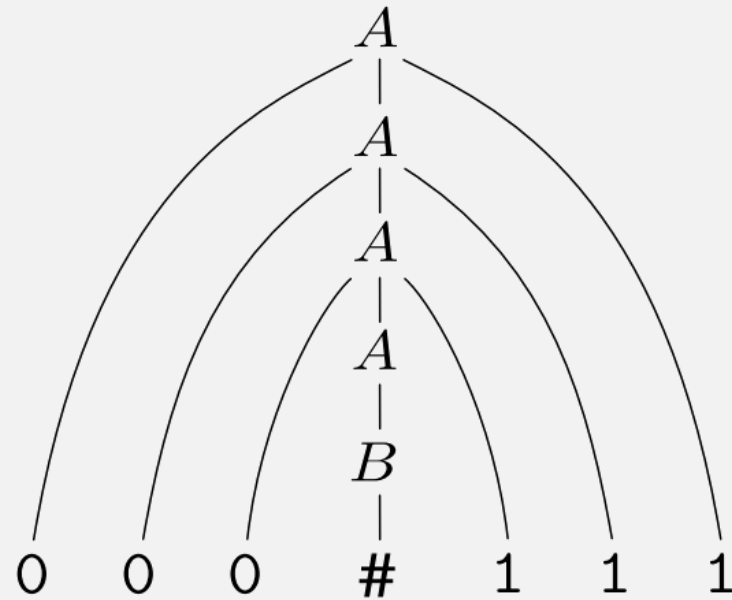
- EXPR ⇒
- EXPR + TERM ⇒
- TERM + TERM ⇒
- FACTOR + TERM ⇒
- a + TERM
- ...

LEFTMOST DERIVATION

Derivations and Parse Trees

$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$

A derivation may also be represented as a **parse tree**



Multiple Derivations, Single Parse Tree

Leftmost derivation

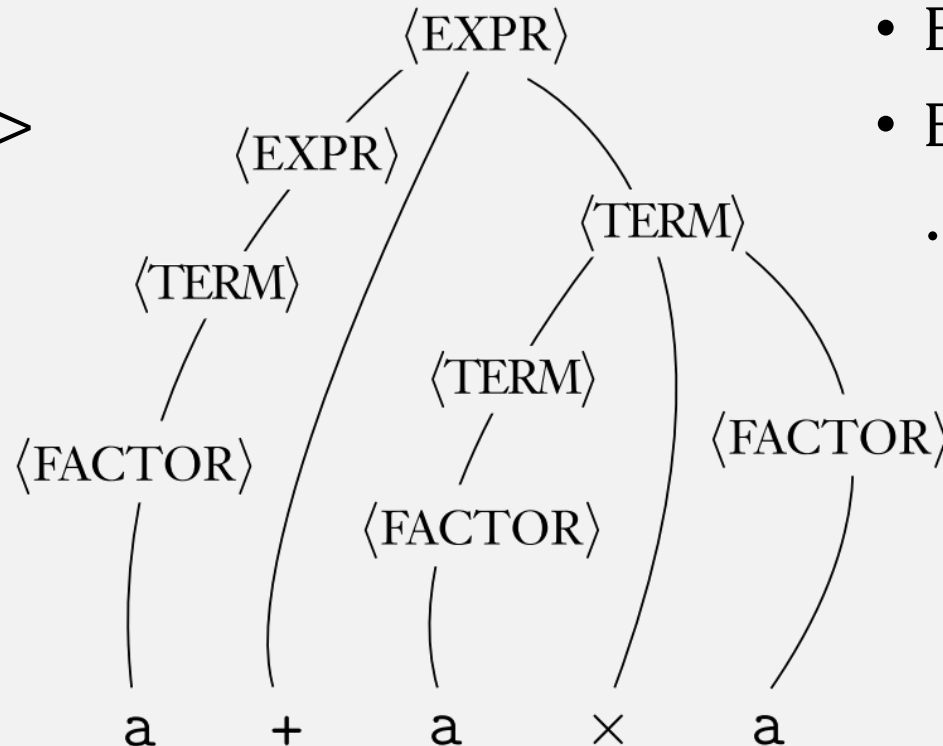
- EXPR =>
- EXPR + TERM =>
- TERM + TERM =>
- FACTOR + TERM =>
- a + TERM
- ...

Since the “meaning” (i.e., parse tree) is same, by convention we just use **leftmost** derivation

Rightmost derivation

- EXPR =>
- EXPR + TERM =>
- EXPR + TERM x FACTOR =>
- EXPR + TERM x a =>
- ...

Same parse tree



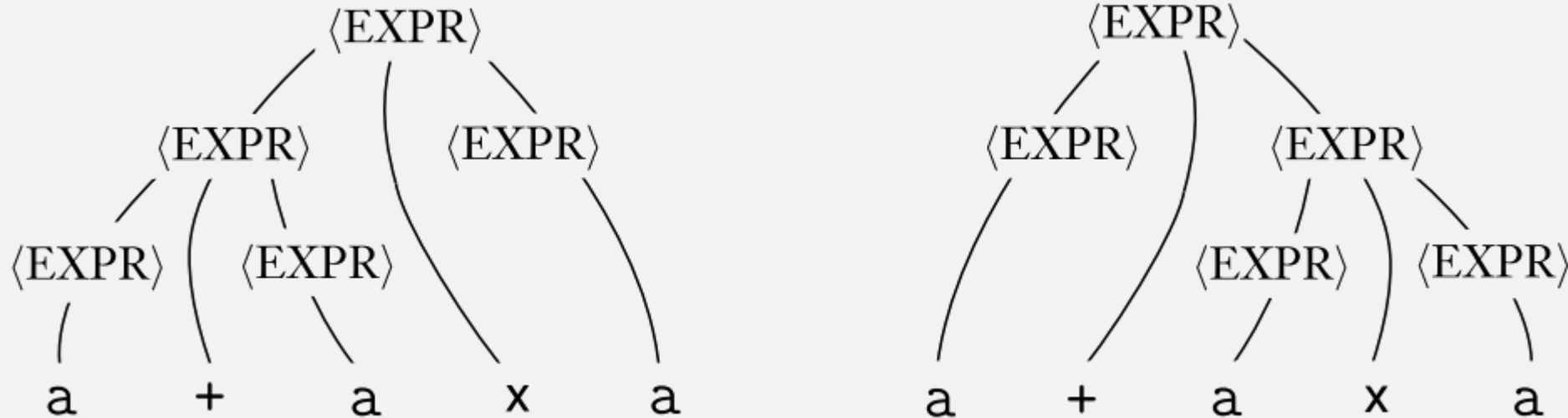
A Parse Tree gives “meaning” to a string

Ambiguity

grammar G_5 :

$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle \mid \langle \text{EXPR} \rangle \times \langle \text{EXPR} \rangle \mid (\langle \text{EXPR} \rangle) \mid a$

Same string,
different derivation,
and different parse tree!



Ambiguity

A string w is derived *ambiguously* in context-free grammar G if it has two or more different leftmost derivations. Grammar G is *ambiguous* if it generates some string ambiguously.

An ambiguous grammar can give a string multiple meanings!
(why is this **bad**?)

Real-life Ambiguity (“Dangling” else)

- What is the result of this C program?

```
if (1) if (0) printf("a"); else printf("2");
```



```
if (1)
  if (0)
    printf("a");
else
  printf("2");
```

VS

```
if (1)
  if (0)
    printf("a");
else
  printf("2");
```

This string has 2 parsings, and thus 2 meanings!

Ambiguous grammars are confusing. In a (programming) language, a string (program) should have only **one meaning** (result).

Problem is, there's no guaranteed way to create an unambiguous grammar (up to language designers to “be careful”)

Designing Grammars : Basics

1. Think about what you want to “link” together

- E.g., 0^n1^n
 - $A \rightarrow 0A1$
 - # 0s and # 1s are “linked”
- E.g., XML
 - $\text{ELEMENT} \rightarrow \langle \text{TAG} \rangle \text{CONTENT} \langle / \text{TAG} \rangle$
 - Start and end tags are “linked”

2. Start with small grammars and then combine (just like FSMs)

Designing Grammars: Building Up

- Start with small grammars and then combine (just like FSMs)
 - To create a grammar for the language $\{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\}$

- First create grammar for lang $\{0^n 1^n \mid n \geq 0\}$:

$$S_1 \rightarrow 0S_1 1 \mid \epsilon$$

- Then create grammar for lang $\{1^n 0^n \mid n \geq 0\}$:

$$S_2 \rightarrow 1S_2 0 \mid \epsilon$$

- Then combine: $S \rightarrow S_1 \mid S_2$

$$S_1 \rightarrow 0S_1 1 \mid \epsilon$$

$$S_2 \rightarrow 1S_2 0 \mid \epsilon$$

New start variable & rule combines two smaller grammars

"|" = "or" = union
(combines 2 rules with same left side)

Closed Operations on CFLs

- Start with small grammars and then combine (just like FSMs)

- “Or”: $S \rightarrow S_1 \mid S_2$

- “Concatenate”: $S \rightarrow S_1 S_2$

- “Repetition”: $S' \rightarrow S' S_1 \mid \epsilon$

In-class Example: Designing grammars

alphabet Σ is $\{0,1\}$

$\{w \mid w \text{ starts and ends with the same symbol}\}$

- $S \rightarrow 0C'0 \mid 1C'1 \mid \varepsilon$ “string starts/ends with same symbol, middle can be anything”
- $C' \rightarrow C'C \mid \varepsilon$ “middle: all possible terminals, repeated (ie, all possible strings)”
- $C \rightarrow 0 \mid 1$ “all possible terminals”

Next Time:

Regular Languages	Context-Free Languages (CFLs)
Regular Expression	Context-Free Grammar (CFG)
A Reg expr <u>describes</u> a Regular lang	A CFG <u>describes</u> a CFL
Finite automaton (FSM)	???
An FSM <u>recognizes</u> a Regular lang	???

Next Time:

Regular Languages	Context-Free Languages (CFLs)
Regular Expression	Context-Free Grammar (CFG)
A Reg expr <u>describes</u> a Regular lang	A CFG <u>describes</u> a CFL
Finite automaton (FSM)	Push-down automaton (PDA)
An FSM <u>recognizes</u> a Regular lang	A PDA <u>recognizes</u> a CFL

Next Time:

Regular Languages	Context-Free Languages (CFLs)
Regular Expression	Context-Free Grammar (CFG)
A Reg expr <u>describes</u> a Regular lang	A CFG <u>describes</u> a CFL
Finite automaton (FSM)	Push-down automaton (PDA)
An FSM <u>recognizes</u> a Regular lang	A PDA <u>recognizes</u> a CFL
DIFFERENCE:	DIFFERENCE:
A Regular lang is <u>defined</u> with a FSM	A CFL is <u>defined</u> with a CFG
<i>Proved:</i> Reg expr \Leftrightarrow Reg lang	<i>Must prove:</i> PDA \Leftrightarrow CFL

Check-in Quiz 2/16

On gradescope