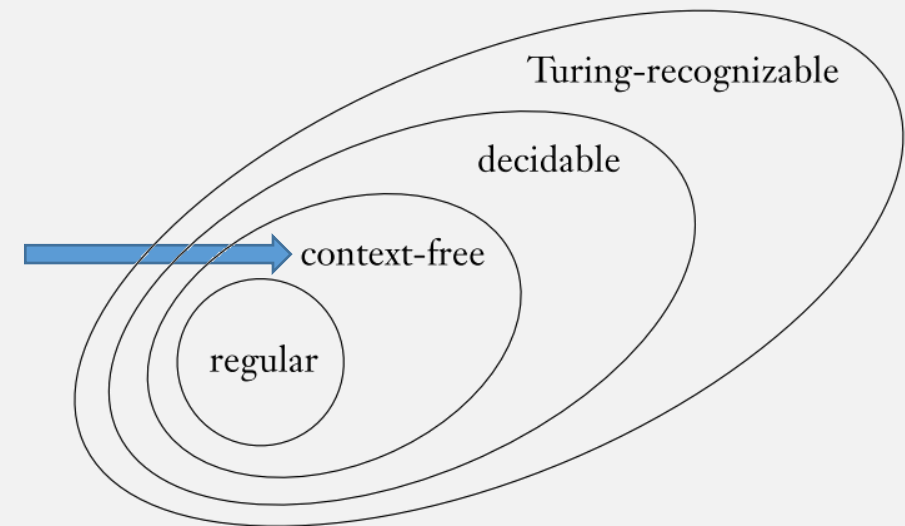


UMB CS 622

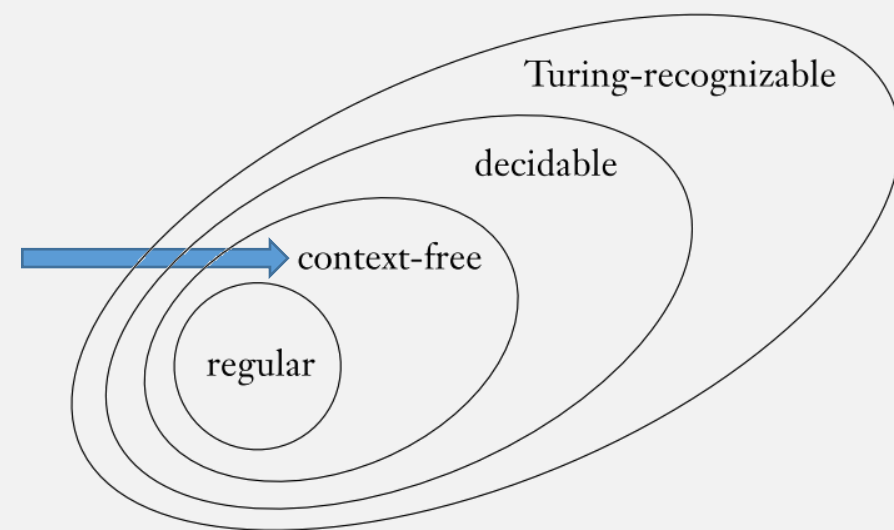
Context-Free Languages (CFLs)

Monday, March 18, 2024



Announcements

- HW 4 in
 - ~~due Mon 3/18 12pm noon~~
- HW 5 out
 - due Mon 3/25 12pm noon



Last Time:

Non-Regular Languages

Example:

An arbitrary count

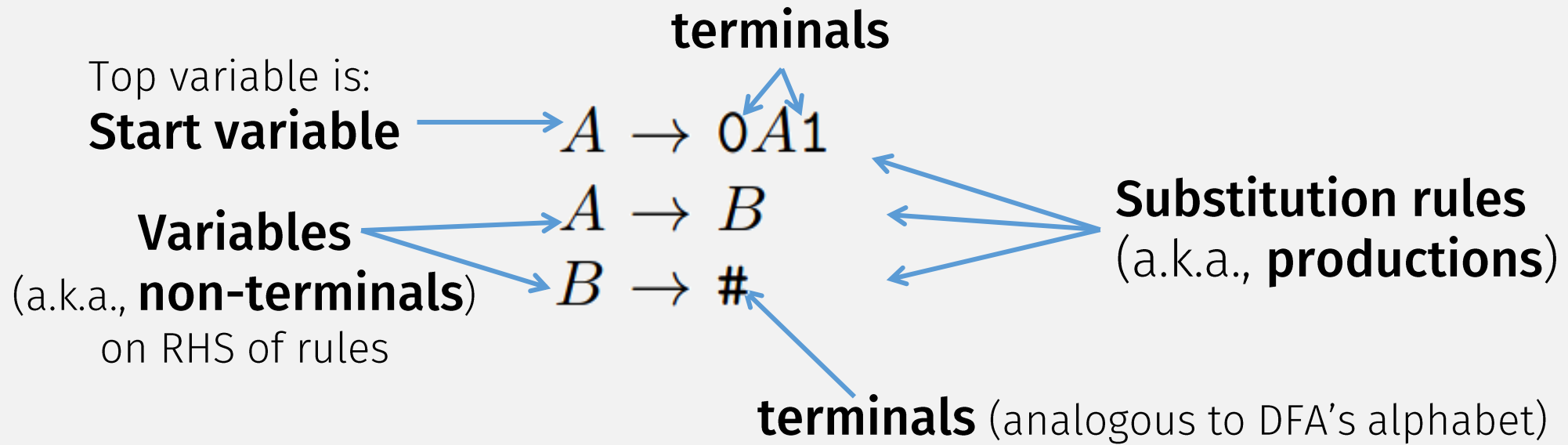
$$L = \{ 0^n 1^n \mid n \geq 0 \}$$

- A DFA recognizing L would require infinite states! (impossible)
 - States representing: zero 0s seen, one 0 seen, two 0s, ...
- This language is the same as many PLs, e.g., HTML!
 - To better see this replace:
 - "0" with "<tag>" or "("
 - "1" with "</tag>" or ")"
- The Problem: remembering nestedness
 - Need to count arbitrary nesting depths
 - E.g., `if { if { if { ... } } }`
 - Thus: most programming language syntax is not regular!

We can prove non-regularness ... with the **Pumping Lemma** (and proof by contradiction)

But ... what kind of language is it then?

A Context-Free Grammar (CFG)



Context-Free Grammar (CFG): Formally

Grammar $G_1 = (V, \Sigma, R, S)$

R is this set of rules (mappings): **terminals**

Top variable is:

Start variable $\rightarrow A \rightarrow 0A1$

Variables $\rightarrow A \rightarrow B$

(a.k.a., **non-terminals**) $\rightarrow B \rightarrow \#$

terminals (analogous to DFA's alphabet)

CFG Practical Application:
Used to describe
programming language
syntax!

Substitution rules
(a.k.a., **productions**)

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

[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

“productions” = rules

“nonterminal” = variable

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*.

“goal symbol” = Start variable

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.

A CFG specifies a language!

2.2. The Lexical Grammar

(definition of a **language**: sequence of symbols)

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).

Definition:
A CFG describes a context-free language!

Analogies

Regular Language	Context-Free Language (CFL)
Regular Expression	Context-Free Grammar (CFG)
thm	def
A Reg Expr <u>describes</u> a Regular lang	A CFG <u>describes</u> a CFL

(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 not expressible with 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
```

<https://docs.python.org/3/reference/grammar.html>

Java Syntax: Described with CFGs

ORACLE

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

[Prev](#)

Definition:

A **CFG** describes a **context-free language!**
but what strings are in the language?

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).

Generating Strings with a CFG

In-class exercise:

Write 3 more strings that can be generated by this grammar

Definition:

A **CFG** describes a **context-free language!** but what strings are in the language?

1st rule $\longrightarrow A \rightarrow 0A1$

2nd rule $\longrightarrow A \rightarrow B$

Last rule $\longrightarrow B \rightarrow \#$

“Applying a rule”
= replace LHS variable
with RHS sequence

At each step, *arbitrarily*
choose any variable to
replace, and any rule to apply

Stop when: string is all terminals

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

Example:

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

Start with:
Start variable

Apply 1st rule

1st rule again

1st rule again

Apply 2nd rule

Apply last rule

Generating Strings with a CFG

Definition:

A **CFG** describes a **context-free language!**
but what strings are in the language?

$G_1 =$

$A \rightarrow 0A1$

$A \rightarrow B$

$B \rightarrow \#$

Strings in CFG's language
= all possible **generated** / **derived** strings

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

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

Example:

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

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:

$\alpha, \beta \in (V \cup \Sigma)^*$ ← sequence of terminals or variables

$A \in V$ ← Variable

$A \rightarrow \gamma \in R$ ← Rule

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:

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

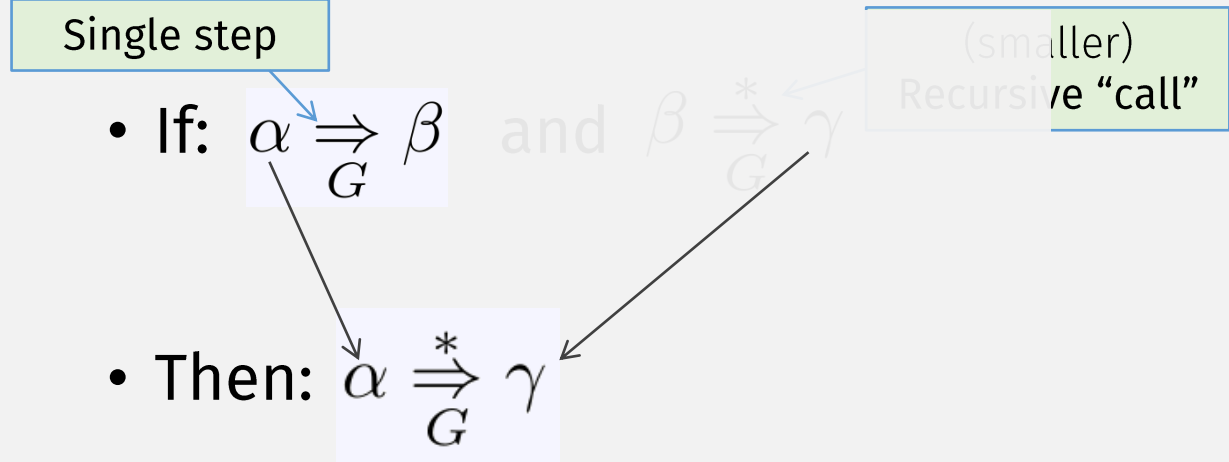
$$A \in V \leftarrow \begin{array}{l} \text{Variable} \end{array}$$

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

Multi-step (recursively defined)

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

Recursive case: (> 0 steps)



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.

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

“the language of a grammar G is ...”

“all possible sequences of terminal symbols ...”

... “that can be generated with rules of grammar G ”

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

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

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

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

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

Statements and Justifications?

Proof: $L = \{0^n 1^n \mid n \geq 0\}$ is a CFL

Statements

1. If a CFG describes a language, then it is a CFL

2. CFG G_1 describes L

$$\begin{array}{l} A \rightarrow 0A1 \\ A \rightarrow B \\ B \rightarrow \varepsilon \end{array}$$

3. $L = \{0^n 1^n \mid n \geq 0\}$ is a CFL

Justifications

1. Definition of CFL

2. (Did you come up with examples???)

3. By Statements #1 and #2

"|" symbol = Shorthand for multiple rules with same LHS variable

A String Can Have Multiple Derivations

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

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

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

RIGHTMOST DERIVATION

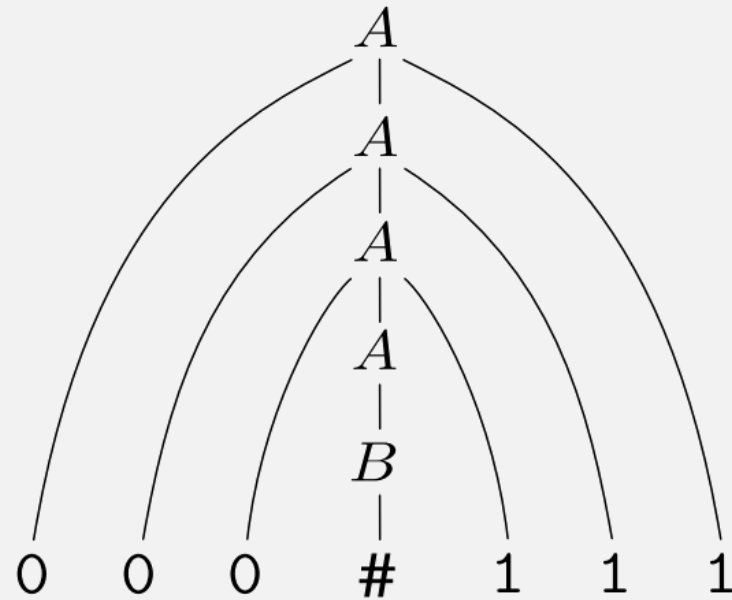
- EXPR ⇒
- EXPR + TERM ⇒
- TERM + TERM ⇒
- FACTOR + TERM ⇒
- **a + TERM**
- ...

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**



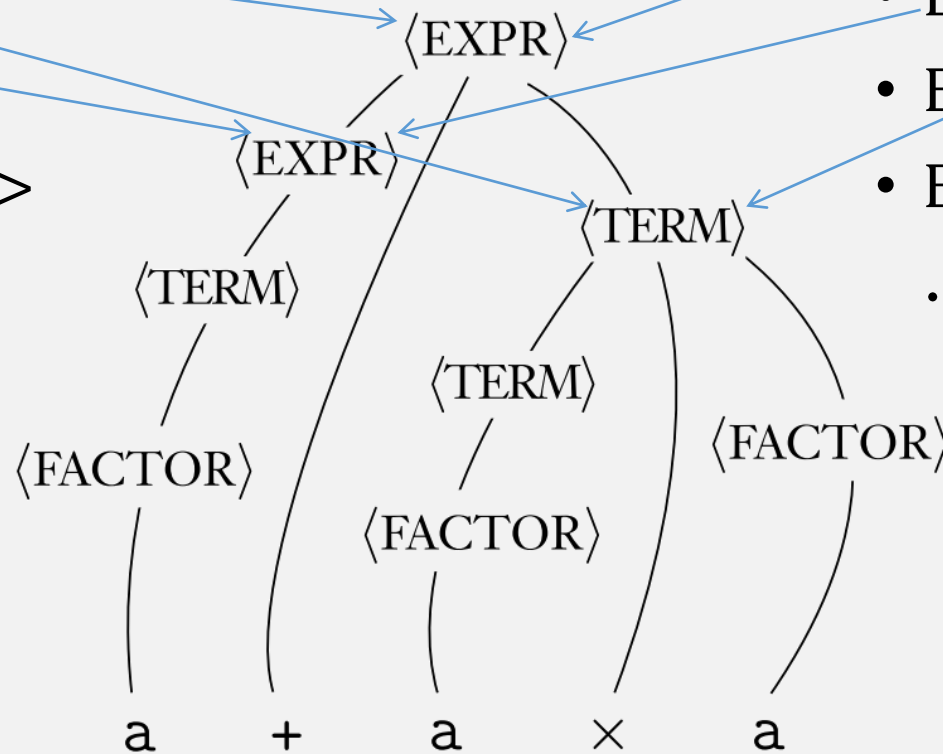
Multiple Derivations, Single Parse Tree

Leftmost derivation

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

Rightmost derivation

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



Same parse tree

A parse tree represents a CFG computation ... like a sequence of states represents a DFA computation

A Parse Tree gives "meaning" to a string

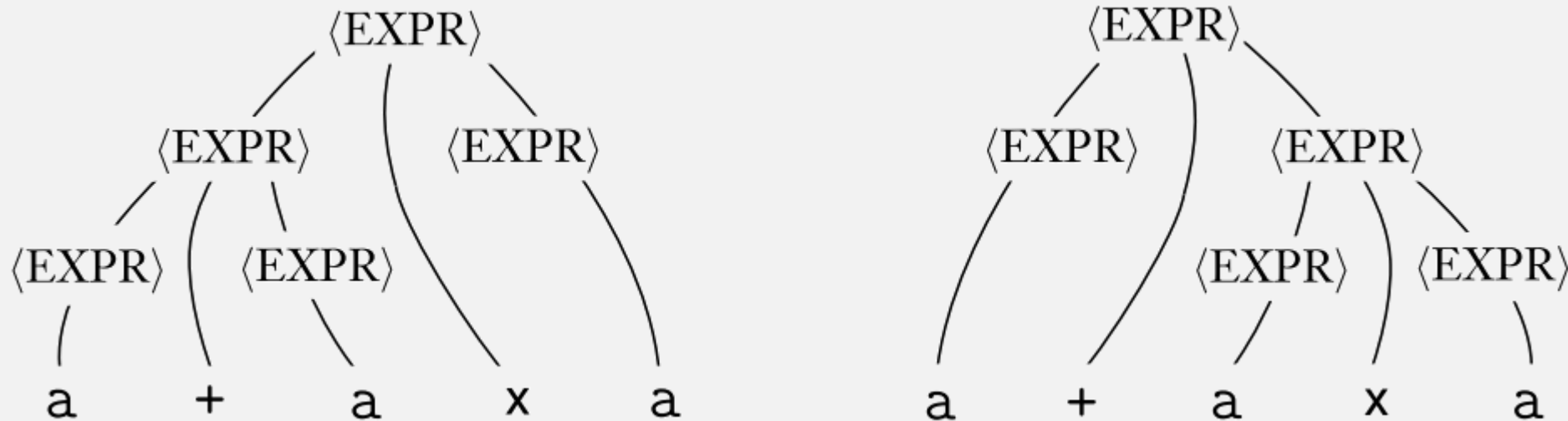
Ambiguity

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!**

So this string has
two meanings!



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, ie represent two different computations!
(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. A computation on a string should ideally have only one result.

Thus in practice, we typically focus on the **unambiguous subset** of CFGs (CFLs) (more on this later)

Problem is, there's no easy way to create an **unambiguous** grammar (it's 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 and rule combines two smaller grammars
 - “|” = “or” = union (combines 2 rules with same left side)

(Closed) Operations for 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”: $S' \rightarrow S' S_1 \mid \epsilon$

Could you write out
the full proof?

In-class Example: Designing grammars

alphabet Σ is $\{0,1\}$

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

1) come up with examples: In the language: **010, 101, 11011** **1, 0 ?**
Not in the language: **10, 01, 110** $\epsilon ?$

2) Create CFG:

$S \rightarrow 0C'0 \mid 1C'1 \mid 0 \mid 1$ “string starts/ends with same symbol, middle can be anything”

$C' \rightarrow C'C \mid \epsilon$ “middle: all possible terminals, repeated (ie, all possible strings)”

$C \rightarrow 0 \mid 1$ “all possible terminals”

3) Check CFG: generates examples in the language; **does not generate examples** not in language

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	A ??? <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

Next Time:

	Regular Languages	Context-Free Languages (CFLs)
thm	Regular Expression	Context-Free Grammar (CFG)
	A Reg Expr <u>describes</u> a Regular Lang	A CFG <u>describes</u> a CFL
def	Finite Automaton (FSM)	Push-down Automaton (PDA)
	An FSM <u>recognizes</u> a Regular Lang	A PDA <u>recognizes</u> a CFL
	<u>DIFFERENCE:</u>	<u>DIFFERENCE:</u>
	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

Submit in-class work 3/18

On gradescope