

# CSCI 301 - Lecture 24: CFG Application: Parsing

Given a grammar  $G$  and a string  $w$ , is  $w \in L(G)$  ?

Of practical interest if you ever want to:

- write and read data to a file or network
- write an interpreter or compiler for a programming language
- extract meaning from structured data

This is called parsing, and it's (in general) hard !

Why? There can be many ways to derive the same string.

Example: consider a grammar for arithmetic expressions.

$$V = \{E\}$$

$$\Sigma = \{0, 1, \dots, 9, +, -, *, /, (, )\}$$

$$S = E$$

$$R = \boxed{\begin{array}{l} E \rightarrow E + E \\ | \quad E - E \\ | \quad E * E \\ | \quad E / E \\ | \quad ( E ) \\ | \quad 0 | 1 | 2 | \dots | 9 \end{array}}$$

$L(G)$  includes:

$$\begin{aligned} &| + | \\ &(4 + 7) * (6 / 2) \end{aligned}$$

$L(G)$  does not include:

$$\begin{aligned} &| + \\ &2 * 4 \end{aligned}$$

$E \rightarrow E + E$
$E - E$
$E * E$
$E / E$
$( E )$
$0   1   2   \dots   9$

String:  $w = 1 + 1 * 4$

Derivations:

①

② left-most

③ right-most

$$E \Rightarrow \underline{E} + E$$

$$\underline{E} + \underline{E * E}$$

$$| + \underline{E} * E$$

$$\underline{|} + \underline{1} * \underline{E}$$

$$| + 1 * \underline{4}$$

$$\underline{E} \Rightarrow \underline{E} + E$$

$$| + \underline{E}$$

$$| + \underline{E} * E$$

$$| + \underline{1} * \underline{E}$$

$$| + 1 * 4$$

$$E \Rightarrow E * \underline{E}$$

$$\underline{E} * 4$$

$$E + \underline{E} * 4$$

$$E + \underline{1} * 4$$

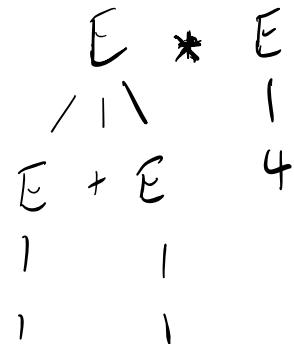
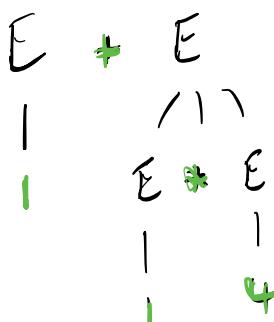
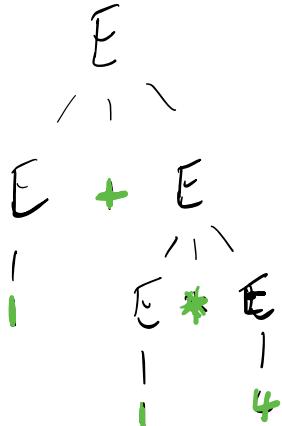
$$\underline{1} + 1 * 4$$

①

Parse trees:

②

③



## Definitions

left-most

: each step applies a rule to the leftmost nonterminal.

right-most

derivation

: each step applies a rule to the rightmost nonterminal.

Def: A grammar is ambiguous if some string  $w$  has more than one parse tree.

Equivalently: A grammar is ambiguous if there is more than one left-most derivation for some string.

Ex. PFA

# Parsing Techniques - Overview

Given a grammar  $G = (V, \Sigma, R, S)$  and a string  $w \in \Sigma^*$ , determine whether  $w \in L(G)$ , i.e., whether  $S \Rightarrow^* w$ .

Two broad categories of approaches:

- - Top-down: Start with  $S$ , apply rules until you've derived  $w$ .
- Bottom-up: Start with  $w$ , apply substitutions "backwards" to get  $S$ .

In both cases, the crux is: which rule should we apply?

Which nonterminal? leftmost or rightmost

Which rule? ??

General, brute force solution: backtracking

$E \rightarrow E + E$
$E \rightarrow E - E$
$E \rightarrow E * E$
$E \rightarrow E / E$
$E \rightarrow ( E )$
$E \rightarrow 0   1   2   \dots   9$

$$E \Rightarrow E + E \quad | * | + | - | \dots$$

~~$E \rightarrow E - E$~~

Can we do better?

Top-down:

- LL( $k$ )
    - left-to-right
    - ↓ leftmost
    - lookahead
- Lab 8 → LL(1)

Bottom-up:

- LR( $k$ )
  - ↑ rightmost
- LALR

Left Recursion

$w = aabb$

$$S \rightarrow aS$$

$$\begin{array}{l} | \\ bS \\ | \\ \epsilon \end{array}$$

$$S \rightarrow Sa$$

$$\begin{array}{l} | \\ Sb \\ | \\ \epsilon \end{array}$$