

# Ambiguity and Parsing Algorithms

CSCI 3136: Principles of Programming Languages

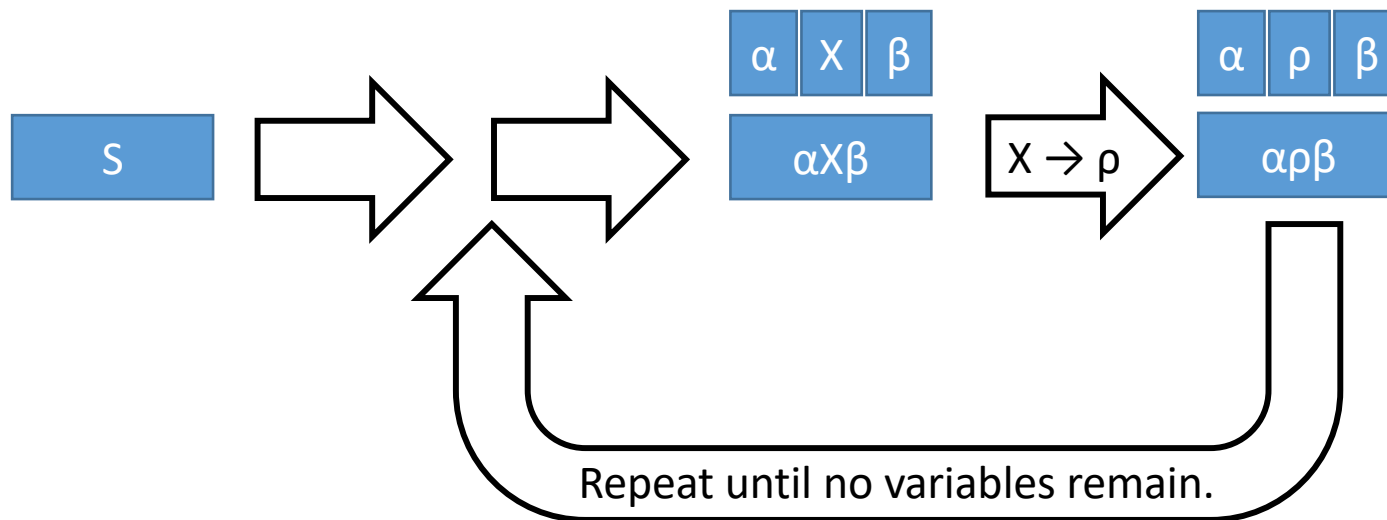
# Agenda

- Announcements
  - Assignment 3 is out and due June 7.
  - Extended office hours: WF 2:30 – 4pm (more coming if needed)
  - Midterm will be on June 19, 10am – 11:30am in **CHEB 170**
- Readings:
  - Today: 2.3.0, 2.3.1
  - Note: I recommend using alternative texts for this part of the course:
  - E..g, Hopcorft et al, “Introduction to Automata Theory”
- Lecture Contents
  - Ambiguous Grammars
  - Left and Right Parse Tree Derivations
  - LL(K) and LR(K) Parsing
  - S-Grammars
  - LL(1) Parsing
  - LR(1) Parsing

# Recall: A CFG for Expressions

- $V = \{E, Op\}$
- $\Sigma = \{\text{identifier, number, (, ), +, -, *, /}\}$
- $P = \{$ 
  - $E \rightarrow E Op E$
  - $E \rightarrow -E$
  - $E \rightarrow ( E )$
  - $E \rightarrow \text{number}$
  - $E \rightarrow \text{identifier}$
  - $Op \rightarrow +$
  - $Op \rightarrow -$
  - $Op \rightarrow *$
  - $Op \rightarrow /$ $\}$
- $S = E$

# Derivations in a Nutshell



# Parse Tree of a Derivation $(42+13)*11$

$\sigma = \mathbf{E}$

$\Rightarrow \mathbf{E} \text{ Op } \mathbf{E}$

$\Rightarrow (\mathbf{E}) \text{ Op } \mathbf{E}$

$\Rightarrow (\mathbf{E} \text{ Op } \mathbf{E}) \text{ Op } \mathbf{E}$

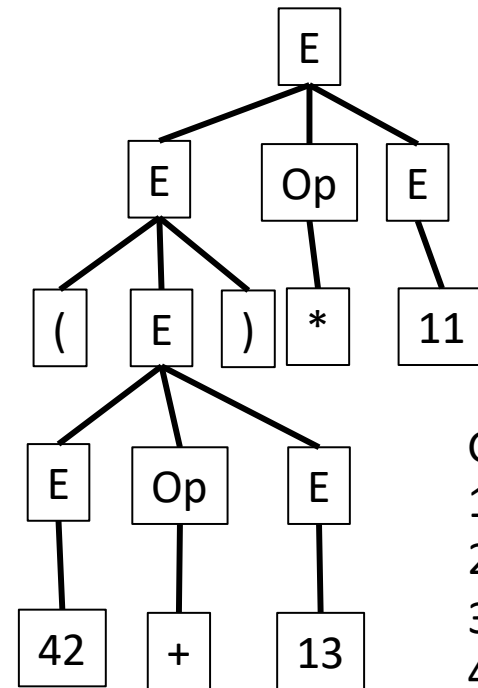
$\Rightarrow (42 \text{ **Op** } \mathbf{E}) \text{ Op } \mathbf{E}$

$\Rightarrow (42 + \mathbf{E}) \text{ Op } \mathbf{E}$

$\Rightarrow (42 + 13) \text{ **Op** } \mathbf{E}$

$\Rightarrow (42 + 13) * \mathbf{E}$

$\Rightarrow (42 + 13) * 11$



Grammar

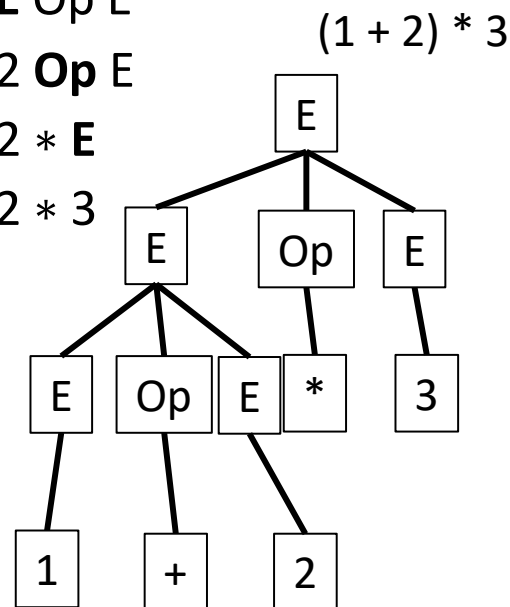
1.  $E \rightarrow E \text{ Op } E$
2.  $E \rightarrow -E$
3.  $E \rightarrow ( E )$
4.  $E \rightarrow \text{number}$
5.  $E \rightarrow \text{identifier}$
6.  $Op \rightarrow +$
7.  $Op \rightarrow -$
8.  $Op \rightarrow *$
9.  $Op \rightarrow /$

# Another Example: $1 + 2 * 3$

## This is ambiguous!

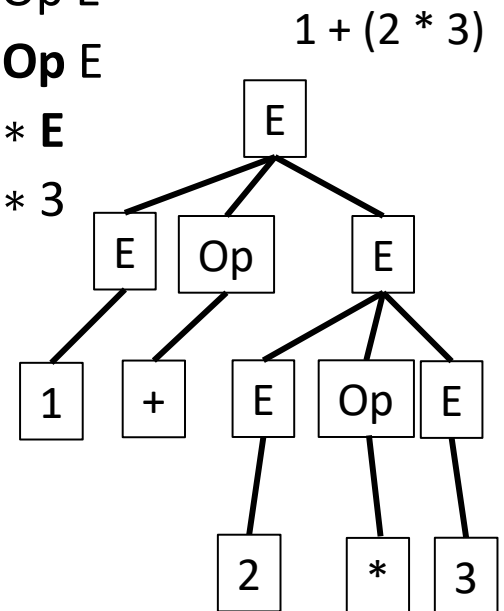
$\sigma = E$

$\Rightarrow E \text{ Op } E$   
 $\Rightarrow E \text{ Op } E \text{ Op } E$   
 $\Rightarrow 1 \text{ Op } E \text{ Op } E$   
 $\Rightarrow 1 + E \text{ Op } E$   
 $\Rightarrow 1 + 2 \text{ Op } E$   
 $\Rightarrow 1 + 2 * E$   
 $\Rightarrow 1 + 2 * 3$



$\sigma = E$

$\Rightarrow E \text{ Op } E$   
 $\Rightarrow 1 \text{ Op } E$   
 $\Rightarrow 1 + E$   
 $\Rightarrow 1 + E \text{ Op } E$   
 $\Rightarrow 1 + 2 \text{ Op } E$   
 $\Rightarrow 1 + 2 * E$   
 $\Rightarrow 1 + 2 * 3$



Grammar

1.  $E \rightarrow E \text{ Op } E$
2.  $E \rightarrow -E$
3.  $E \rightarrow ( E )$
4.  $E \rightarrow \text{number}$
5.  $E \rightarrow \text{identifier}$
6.  $\text{Op} \rightarrow +$
7.  $\text{Op} \rightarrow -$
8.  $\text{Op} \rightarrow *$
9.  $\text{Op} \rightarrow /$

# Ambiguity

- Observations:
  - There are infinitely many grammars to specify the same language
  - There may be multiple parse trees for the same sentence!
- Definition: If multiple parse trees can be generated by  $G$  for the same sentence, then  $G$  is *ambiguous*.
- Definition: If  $L$  does not have an unambiguous grammar, then  $L$  is *inherently ambiguous*
  - Usually not the case for programming languages!

# An Unambiguous Expression Grammar

Grammar

1.  $E \rightarrow T$
2.  $E \rightarrow E + T$
3.  $E \rightarrow E - T$
4.  $T \rightarrow F$
5.  $T \rightarrow T * F$
6.  $T \rightarrow T / F$
7.  $F \rightarrow \text{number}$
8.  $F \rightarrow \text{identifier}$
9.  $F \rightarrow (E)$

- Try deriving  $1 + 2 * 3$



# Derivation Order

- Derivation orders refer to the order in which variables are replaced in the current partial derivation.
- The two most common ones are:
  - **Leftmost derivation** replaces the leftmost variable in each step
  - **Rightmost derivation** replaces the rightmost variable in each step

# Leftmost Derivation Example $1+2*3$

$E \Rightarrow E + T$

$\Rightarrow T + T$

$\Rightarrow F + T$

$\Rightarrow 1 + T$

$\Rightarrow 1 + T * F$

$\Rightarrow 1 + F * F$

$\Rightarrow 1 + 2 * F$

$\Rightarrow 1 + 2 * 3$

Grammar

1.  $E \rightarrow T$

2.  $E \rightarrow E + T$

3.  $E \rightarrow E - T$

4.  $T \rightarrow F$

5.  $T \rightarrow T * F$

6.  $T \rightarrow T / F$

7.  $F \rightarrow \text{number}$

8.  $F \rightarrow \text{identifier}$

9.  $F \rightarrow (E)$

# Rightmost Derivation Example $1+2*3$

$E \Rightarrow E + T$

$\Rightarrow E + T * F$

$\Rightarrow E + T * 3$

$\Rightarrow E + F * 3$

$\Rightarrow E + 2 * 3$

$\Rightarrow T + 2 * 3$

$\Rightarrow F + 2 * 3$

$\Rightarrow 1 + 2 * 3$

Grammar

1.  $E \rightarrow T$

2.  $E \rightarrow E + T$

3.  $E \rightarrow E - T$

4.  $T \rightarrow F$

5.  $T \rightarrow T * F$

6.  $T \rightarrow T / F$

7.  $F \rightarrow \text{number}$

8.  $F \rightarrow \text{identifier}$

9.  $F \rightarrow (E)$

# Where Are We?

- CFGs are used to specify programming language syntax
- Parsing finds the parse tree of the program (token stream)
- CFGs for programming languages must unambiguously capture the program structure.
- Parsers must be efficient:
  - A parser can be generated from a CFG that runs in  $O(n^3)$  time
  - We prefer (require) linear time.
- How do we get this?

# Regular Grammars: A Brief Aside

- A CFG is *right-linear* if all productions are of the form
  - $A \rightarrow \sigma B, \sigma \in \Sigma^*, B \in V$
  - $A \rightarrow \sigma, \sigma \in \Sigma^*$
- A CFG is *left-linear* if all productions are of the form
  - $A \rightarrow B\sigma, \sigma \in \Sigma^*, B \in V$
  - $A \rightarrow \sigma, \sigma \in \Sigma^*$
- A CFG is regular if it is right-linear or left-linear
- Regular grammars specify exactly the set of regular languages
- Regular grammars are too weak to specify most programming languages
- But, parsers generated from them run in linear time!
  - Why?
  - Are there more complex grammars for which linear time parsers exist?

# LL and LR Grammars

Two kinds of unambiguous grammars that can be parsed efficiently

- LL(k) grammars
  - Are scanned Left-to-right and generate a Leftmost derivation
  - If the first letter in the current sentential form is a variable, k tokens look-ahead in the input suffice to decide which production to use to expand it.
- LR(k) grammars
  - Are scanned Left-to-right and generate a Rightmost derivation
  - The next k tokens in the input suffice to choose the next step the parser should perform.
- The syntax of almost every programming language can be described by LL(1) or LR(1) grammars!
  - How? Why?

# S-Grammars

- First let's consider a very simple grammar
- An *S-grammar* or *simple grammar* is a special case of an LL(1)-grammar
- A CFG is an S-grammar if
  - Every production starts with a terminal
  - Productions for the same LHS start with different terminals

E.g., If  $G$  contains  $A \rightarrow aA$  and  $A \rightarrow a$  then  $G$  is not simple!
- Idea: When using S-Grammars, selecting which rule to apply is easy.

# Example: LL(1) Parsing (top-down)

## S-Grammar for Polish Notation

1.  $S \rightarrow + SS$
2.  $S \rightarrow - SS$
3.  $S \rightarrow * SS$
4.  $S \rightarrow / SS$
5.  $S \rightarrow \text{neg } S$
6.  $S \rightarrow \text{integer}$

Expression:

$- * + 1 2 3 4$

Is interpreted as:

$(1 + 2) * 3 - 4$

- How do we derive

$- * + 1 2 3 4$

$S \Rightarrow - S S$

$\Rightarrow - * S S S$

$\Rightarrow - * + S S S S$

$\Rightarrow - * + 1 S S S$

$\Rightarrow - * + 1 2 S S$

$\Rightarrow - * + 1 2 3 S$

$\Rightarrow - * + 1 2 3 4$

- This is an example of LL(1) parsing
- How does a parser do this?



# LL(1) Parsing of S-Grammars

```
# Use a stack to store the
# current sentential form
push(S) # push start variable
t = next_token()
loop until no more tokens:
    x = pop()
    if x == t:
        t = next_token()
        continue
    elseif x ∈ V:
        select production x → t α
        add children t α to node x
        push(tα)
    else:
        error
```

Grammar

1.  $S \rightarrow + SS$
2.  $S \rightarrow - SS$
3.  $S \rightarrow * SS$
4.  $S \rightarrow / SS$
5.  $S \rightarrow \text{neg } S$
6.  $S \rightarrow \text{integer}$

Parse Expression:

- \* + 1 2 3 4

**This takes linear time!**

# Example: LR(1) Parsing (bottom-up)

## S-Grammar for Polish Notation

1.  $S \rightarrow + SS$
2.  $S \rightarrow - SS$
3.  $S \rightarrow * SS$
4.  $S \rightarrow / SS$
5.  $S \rightarrow \text{neg } S$
6.  $S \rightarrow \text{integer}$

Expression:

$- * + 1 2 3 4$

Is interpreted as:

$(1 + 2) * 3 - 4$

- How do we derive

$- * + 1 2 3 4$

$\Leftarrow - * + S 2 3 4$  [use rule 6]

$\Leftarrow - * + S S 3 4$  [use rule 6]

$\Leftarrow - * S 3 4$  [use rule 1]

$\Leftarrow - * S S 4$  [use rule 6]

$\Leftarrow - S 4$  [use rule 3]

$\Leftarrow - S S$  [use rule 6]

$\Leftarrow S$  [use rule 2]

- This is an example of LR(1) parsing

- How does a parser do this?

# LR(1) Parsing of S-Grammars

```
# Use a stack to store the
# what has been seen so far
push(next_token()) # init stack
loop until no more tokens:
    if  $\exists (P \rightarrow \alpha)$  such that  $\alpha \in \text{Stack}$ 
        # reduce operation
        pop( $\alpha$ )
        push(P)
        add children  $\alpha$  to node P
    else:
        # shift operation
        push(next_token())
```

Grammar

1.  $S \rightarrow + SS$
2.  $S \rightarrow - SS$
3.  $S \rightarrow * SS$
4.  $S \rightarrow / SS$
5.  $S \rightarrow \text{neg } S$
6.  $S \rightarrow \text{integer}$

Parse Expression:

$- * + 1 2 3 4$

**This takes linear time!**

# Building Parsers

- We now have some intuition about parsing algorithms
- But ...
  - The above algorithms are for S-Grammars (too simple)
  - Want to generate parser given a grammar
- So ...
  - Assume that we will be using more complex grammars
  - How do we generate the parser?

# Building an LL(1) Parser

- Basic Challenge: Given current token, which production does the parser select if next item in sentential form is a nonterminal

E.g., if  $S$  is on the stack and input is  $+$ , then parser must select production  $S \rightarrow +SS$

- In general: for input  $\mathbf{a}$ , sentential form  $A \dots$ , either
  - $A \Rightarrow \alpha \Rightarrow^* \mathbf{a}\beta$
  - $A \Rightarrow \alpha \Rightarrow^* \epsilon$  and derivation of  $A$  is succeeded by  $\mathbf{a}$ .
- Intuitively,  $\mathbf{a}$  is in the *predictor set* of  $A \rightarrow \alpha$   
if  $A\beta \Rightarrow \alpha\beta \Rightarrow^* \mathbf{a}\gamma$ , for  $\beta, \gamma \in \Sigma^*$   
I.e., the parser selects  $A \rightarrow \alpha$  if  $\mathbf{a}$  is the input and in the *predictor set* of  $A \rightarrow \alpha$

# LL(1) Grammars

- **Definition:** A grammar is LL(1) if the predictor sets of all productions with the same LHS are disjoint.
- E.g. S-Grammars are LL(1)  
Grammar
  1.  $S \rightarrow + SS$
  2.  $S \rightarrow - SS$
  3.  $S \rightarrow * SS$
  4.  $S \rightarrow / SS$
  5.  $S \rightarrow \text{neg } S$
  6.  $S \rightarrow \text{integer}$

Production	Predictor Set
$S \rightarrow + SS$	{+}
$S \rightarrow - SS$	{-}
$S \rightarrow * SS$	{*}
$S \rightarrow / SS$	{/}
$S \rightarrow \text{neg } S$	{neg}
$S \rightarrow \text{integer}$	{integer}