

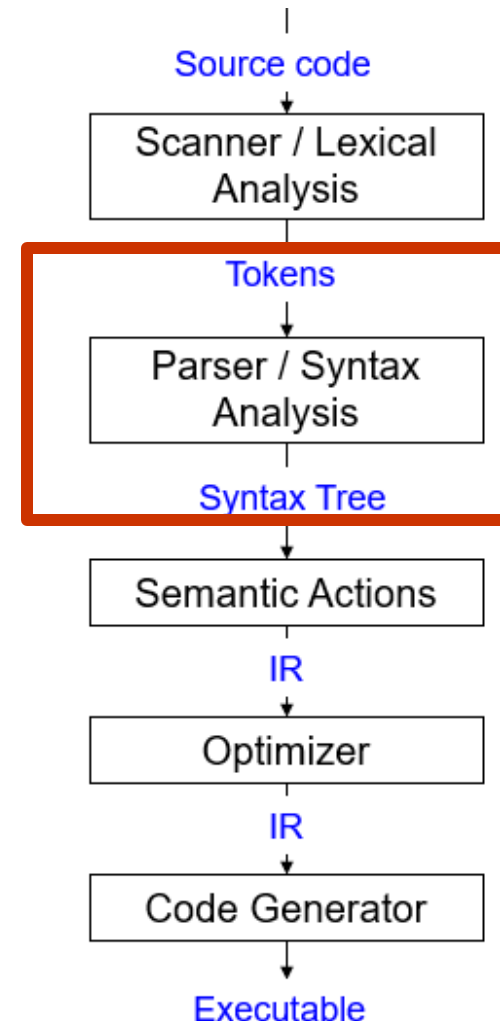
# CS406: Compilers

Spring 2021

## Week 3: Parsers

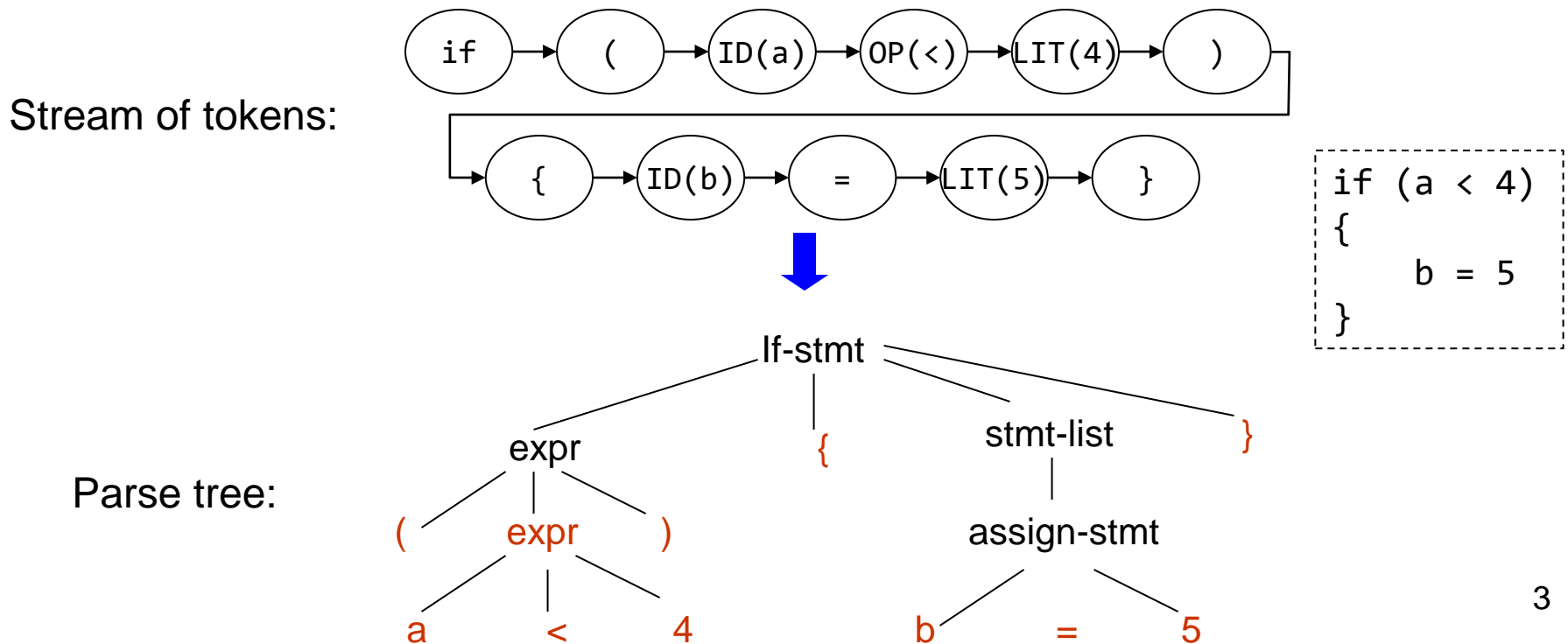
# Parsers - Overview

- Also called syntax analyzers
- Determine two things:
  1. Is a program syntactically valid?
    - Is an English sentence grammatically correct?
  2. What is the structure of programming language constructs? E.g. does the sequence `IF, ID(a), OP(<), ID(b), {, ID(a), ASSIGN, LIT(5), }, ;, }` refer to an if-statement?
    - Diagramming English sentences



# Parsers - Overview

- Input: stream of tokens
- Output: Parse tree
  - sometimes implicit



# Parsers – what do we need to know?

1. How do we define language constructs?
  - Context-free grammars
2. How do we determine: 1) valid strings in the language? 2) structure of program?
  - LL Parsers, LR Parsers
3. How do we write Parsers?
  - E.g. use a parser generator tool such as Bison

# Center Embeddings in English

The bird flew

The bird the boy saw flew

The bird the boy the dog chased saw flew

The bird the boy the dog the man owned chased saw flew

The bird the boy the dog the man the woman loved owned chased  
saw flew

...

*Exercise: write a regular expression that match the pattern. Note: the alphabets of your language are 'Noun', 'Verb' and 'the'*

# Languages

- A language is (possibly infinite) set of strings
- Regular expressions describe *regular languages*  
weakness: can't describe a string of the form:

$$\{ ({}^i )^i \mid i \geq 1 \}$$

E.g.  $((2+3)*5)$

Parenthesized expressions:  $(( ( \text{int } x; )) )$

*Programming language syntax is i.e. recursive*

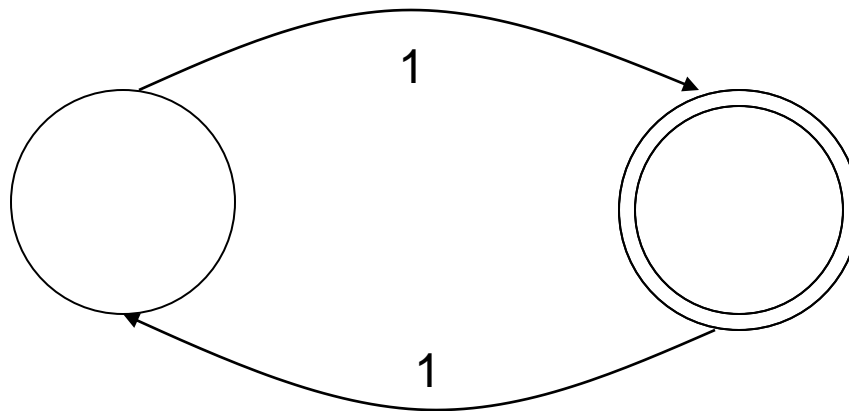
Nested structures:

IF  
  IF  
    IF  
      FI  
    IF  
  FI  
FI

is C regular?

# Trivia

- Regular expressions can describe strings:  
 $\{ \text{mod } k \mid k = \# \text{ states in FA} \}$



“accept all strings having odd number of 1s”

# Context Free Grammar (CFG)

- Natural notation for describing recursive structure definitions. Hence, suitable for specifying language constructs.
- Consist of:
  - A set of *Terminals* (T)
  - A set of *Non-terminals* (N)
  - A *Start Symbol* ( $S \in N$ )
  - A set of *Productions* ( $X \rightarrow Y_1 \dots Y_N$ )  
( aka. rules)  
$$P : X \longrightarrow Y_1 Y_2 Y_3 \dots Y_N \mid X \in N, Y_i \in N \cup T \cup \epsilon / \lambda$$



# Context Free Grammar (CFG)

- Grammar  $G = (T, N, S, P)$

E.g.  $G = (\{a, b\}, \{S, A, B\}, S, \{S \rightarrow AB, A \rightarrow Aa, A \rightarrow a, B \rightarrow Bb, B \rightarrow b\})$

- Implicit meanings
  - First rule listed in the set of productions contains start symbol (on the left-hand side)
  - In the set of productions, you can replace the symbol X (appearing on the right-hand side only) with the string of symbols that are on the right-hand side of a rule, which has X (on the left-hand side)

# Context Free Grammar (CFG)

1. Begin with only S as the initial string

2. Replace S

- S replaced with AB

3. Repeat 2 until the string contains only terminals

- AB replaced with aB
- aB replaced with bb

$G = (T, N, S, P)$   
 $P: \{ S \rightarrow AB, \\ A \rightarrow Aa, \\ A \rightarrow a, \\ B \rightarrow Bb, \\ B \rightarrow b \}$

**Summary:** we move from S to a string of terminals through a series of transformations:

$\alpha_0 \rightarrow \dots \rightarrow \alpha_n$  where  $\alpha_1 \dots \alpha_n$  are strings

Shorthand notation:  $\alpha_0 \xrightarrow{*} \alpha_n$

# Language of the Grammar

- Language  $L(G)$  of the context-free grammar  $G$ 
  - Set of strings that can be derived from  $S$
  - $\{a_1a_2a_3 \dots a_N \mid a_i \in T \ \forall i \text{ and } S \xrightarrow{*} a_1a_2a_3 \dots a_N \}$
  - Is called context-free language
    - All regular languages are context-free but not vice-versa.

# Context-Sensitive Grammar

- Can have context-sensitive grammar and languages (think:  $aB \rightarrow ab$ )
  - Cannot replace right-hand side with left-hand side irrespective of the context.
  - E.g.  $aB \rightarrow ab$  lays down a context: ‘a’ must be a prefix in order to transform the string “aB” to a string of terminals “ab”
    - ccaBb can be replaced by ccabb

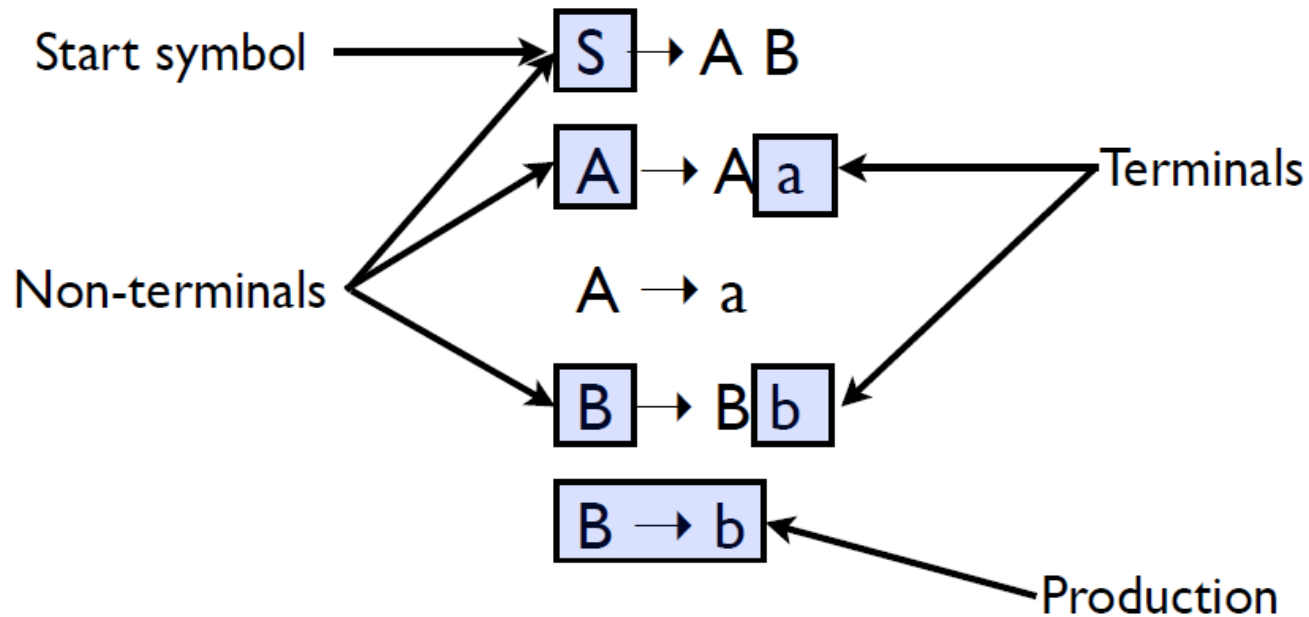
*Is grammar  $G$  context-free?*

$G = (T, N, S, P)$   
 $P: \{$   
     $S \rightarrow AB,$   
     $A \rightarrow Aa,$   
     $A \rightarrow a,$   
     $B \rightarrow Bb,$   
     $B \rightarrow b \}$

# String Derivations

- How do we apply the grammar rules repeatedly to determine the validity of a string? (i.e. string belongs to the language specified by the context-free grammar)
  - Begin with  $S$
  - Replace  $S$
  - Repeat till string contains terminals only  
 *$L(G)$  must contain strings of terminals only*
- Notation:
  - We will use Greek letters to denote strings containing non-terminals and terminals

# Simple grammar



*Backus Naur Form (BNF)*

# Generating strings

$S \rightarrow A B$

$A \rightarrow A a$

$A \rightarrow a$

$B \rightarrow B b$

$B \rightarrow b$

- Given a start rule, productions tell us how to rewrite a non-terminal into a different set of symbols
- Some productions may rewrite to  $\lambda$ . That just removes the non-terminal

To derive the string “a a b b b” we can do the following rewrites:

$$\begin{aligned} S &\Rightarrow A B \Rightarrow A a B \Rightarrow a a B \Rightarrow a a B b \Rightarrow \\ &a a B b b \Rightarrow a a b b b \end{aligned}$$

# Exercise

Which of the below strings are accepted by the grammar:

$A \rightarrow aAa$

$A \rightarrow bBb$

$A \rightarrow \lambda$

$B \rightarrow cA$

$B \rightarrow \lambda$

1. abcba
2. abcbca
3. abba
4. abca



# Programming language syntax

- Programming language syntax is defined with CFGs
- Constructs in language become non-terminals
- May use auxiliary non-terminals to make it easier to define constructs

`if_stmt` → if ( `cond_expr` ) then `statement` `else_part`

`else_part` → else `statement`

`else_part` →  $\lambda$

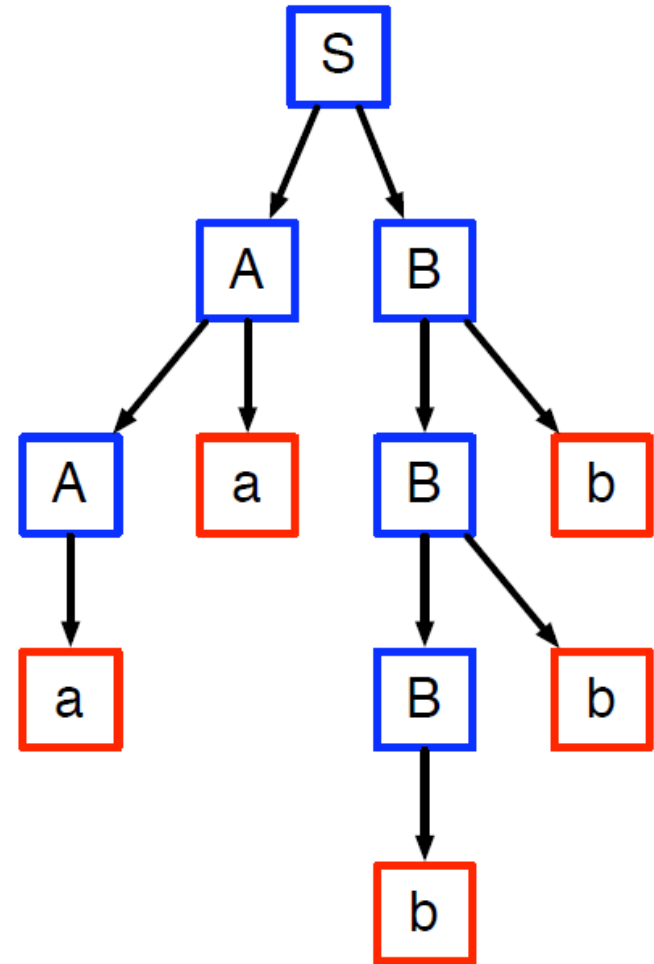
- Tokens in language become terminals

# CFG Contd..

- Is it enough if parsers answer “yes” or “no” to check if a string belongs to context-free language?
  - Also need a parse tree
- What if the answer is a “no”?
  - Handle errors
- How do we implement CFGs?
  - E.g. Bison

# Parse trees

- Tree which shows how a string was produced by a language
- Interior nodes of tree: non-terminals
  - Children: the terminals and non-terminals generated by applying a production rule
- Leaf nodes: terminals



# Parse Trees and String Derivations

- Recall: sequence of rules applied to produce a string is a derivation
- A derivation defines a parse tree
  - A parse tree may have many derivations

# Leftmost derivation

- Rewriting of a given string starts with the leftmost symbol
- Exercise: do a leftmost derivation of the input program

$F(V + V)$

using the following grammar:

E	→	Prefix (E)
E	→	V Tail
Prefix	→	F
Prefix	→	$\lambda$
Tail	→	+ E
Tail	→	$\lambda$

- What does the parse tree look like?

# Rightmost derivation

- Rewrite using the rightmost non-terminal, instead of the left
- What is the rightmost derivation of this string?

$F(V + V)$

E	→	Prefix (E)
E	→	V Tail
Prefix	→	F
Prefix	→	$\lambda$
Tail	→	+ E
Tail	→	$\lambda$

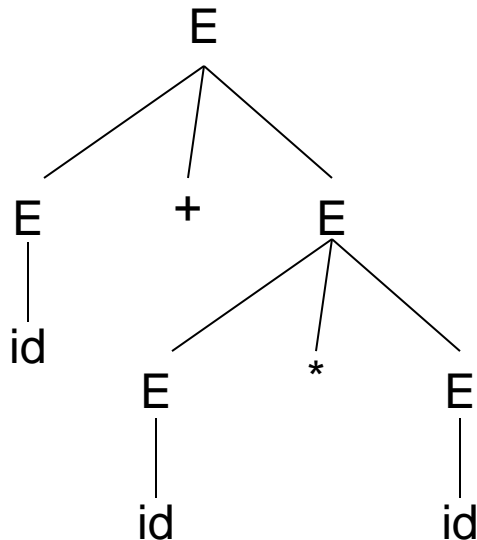
# Ambiguity

- Grammar that produces more than one parse tree for some string

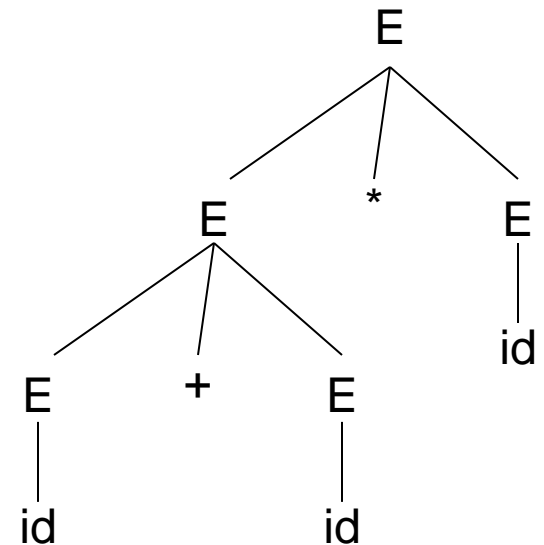
E.g.  $E \rightarrow E + E \mid E * E \mid id$

String: **id+id\*id**

$E \rightarrow E + E$   
 $E \rightarrow id + E$   
 $E \rightarrow id + E * E$   
 $E \rightarrow id + id * E$   
 $E \rightarrow id + id * id$



$E \rightarrow E * E$   
 $E \rightarrow E + E * E$   
 $E \rightarrow id + E * E$   
 $E \rightarrow id + id * E$   
 $E \rightarrow id + id * id$



# Ambiguity – what to do?

- Ignore it
  - Give hints to other components of the compiler on how to resolve it
- Fix it
  - Manually
  - May make the grammar complicated and difficult to maintain

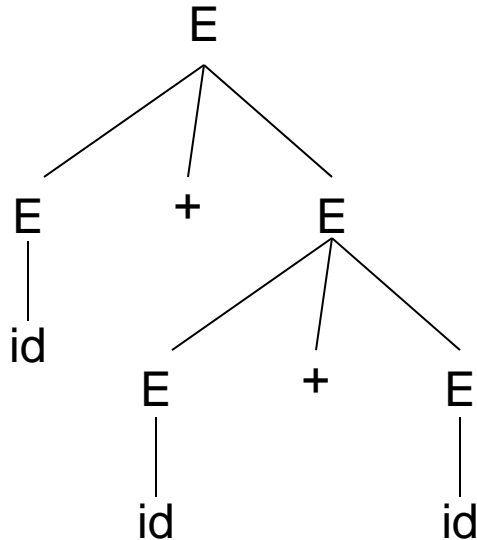


# Ambiguity – ignore

- $E \rightarrow E + E \mid id$

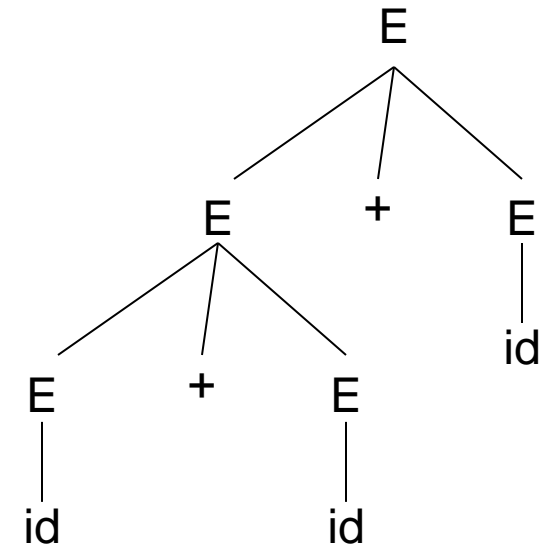
$E \rightarrow E + E$   
 $E \rightarrow id + E$   
 $E \rightarrow id + E + E$   
 $E \rightarrow id + id + E$   
 $E \rightarrow id + id + id$

Produces:  
 $id + (id + id)$



$E \rightarrow E + E$   
 $E \rightarrow E + E + E$   
 $E \rightarrow id + E + E$   
 $E \rightarrow id + id + E$   
 $E \rightarrow id + id + id$

Produces:  
 $(id + id) + id$



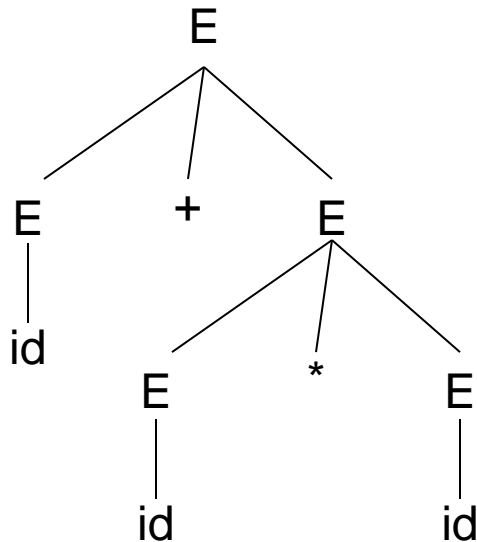
- Associativity declaration in Bison:  
`%left +`

Picks the parse tree on the right

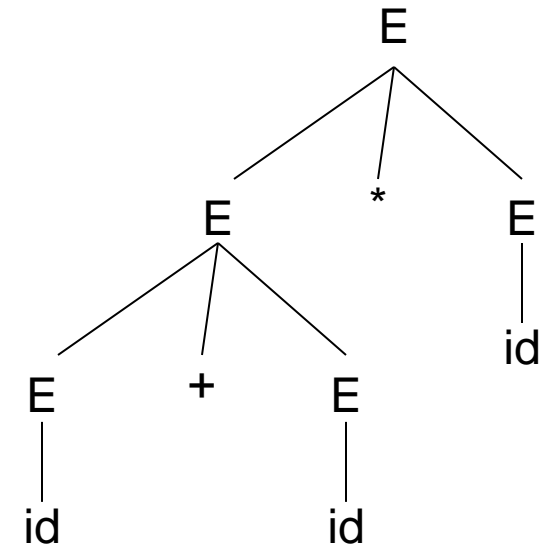
# Ambiguity - ignore

- $E \rightarrow E + E \mid E * E \mid id$

$E \rightarrow E + E$   
 $E \rightarrow id + E$   
 $E \rightarrow id + E * E$   
 $E \rightarrow id + id * E$   
 $E \rightarrow id + id * id$



$E \rightarrow E * E$   
 $E \rightarrow E + E * E$   
 $E \rightarrow id + E * E$   
 $E \rightarrow id + id * E$   
 $E \rightarrow id + id * id$



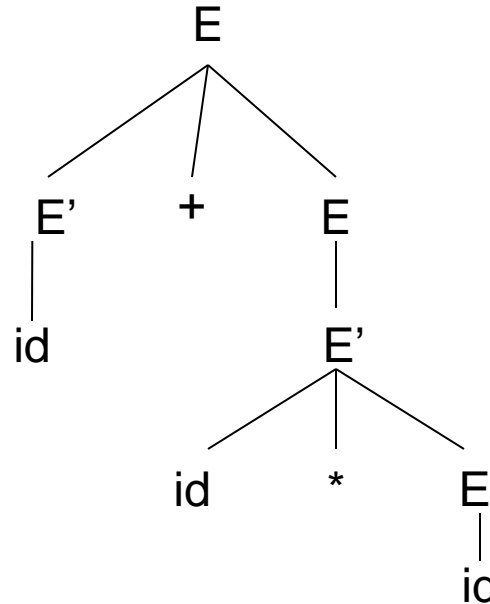
**%left +**  
**%left \***

*Tells that \* has higher precedence over + and both are left associative. So we get the tree on left.*

# Ambiguity – fixing

- Rewrite  $E \rightarrow E + E \mid E * E \mid id$  as  
 $E \rightarrow E' + E \mid E'$   
 $E' \rightarrow id * E' \mid id \mid (E) * E' \mid (E)$

$E \rightarrow E' + E$   
 $E \rightarrow id + E$   
 $E \rightarrow id + E'$   
 $E \rightarrow id + id * E'$   
 $E \rightarrow id + id * id$



`E` controls generation of `+`

`E'` controls generation of `*`. `*`'s are nested deeper in the parse tree.

# Ambiguity - fixing

```
stmt -> if expr then stmt |  
      if expr then stmt else stmt |  
      other
```

**String:** if E1 then if E2 then S1 else S2

**Exercise:** *verify if the above grammar is ambiguous. If so, rewrite the grammar to make it unambiguous.*

```
stmt -> matched | open  
matched -> if expr then matched else matched |  
         other  
open -> if expr then stmt |  
       if expr then matched else open
```

# Error Handling

- Objective: detect invalid programs and provide meaningful feedback to programmer
  - Report errors accurately
  - Recover from errors quickly
  - Don't slow down compilation

# Error Types

- Many types of errors:
  - Lexical – use `Size` instead of `size`
  - Syntactic – extra brace
  - Semantic – `float sqr; sqr(2);`
  - Logical – use `=` instead of `==`

# Error Handling - Types

1. Panic mode
2. Error production
3. Automatic local or global correction

# Panic Mode Error Handling

- Simplest, most popular
- Discards tokens until one from a set of *synchronizing tokens* is found
- Synchronizing tokens have a clear role  
e.g. semicolons, braces
- E.g. `i=i++j`

*policy:* while parsing an expression, discard all tokens until an integer is found. *This policy skips the additional +*

- Specifying policy in bison: **error** keyword

```
E -> E + E | (E) | id | error int | error
```



# Error Productions

- Anticipate common errors
  - 2x instead of 2 \*
- Augment the grammar
  - $E \rightarrow EE \mid \dots$
- Disadvantages:
  - Complicates the grammar

# Error Corrections

- Rewrite the program – find a “nearby” correct program
  - Local corrections – insert a semicolon, replace a comma with semicolon etc.
  - Global corrections – modify the parse tree with “edit distance” metric in mind
- Disadvantages?
  - Implementation difficulty
  - Slows down compilation
  - Not sure if “nearby” program is intended

# Top-down Parsing

- Also called recursive-descent parsing
- Equivalent to finding the left-derivation for an input string
  - Recall: expand the leftmost non-terminal in a parse tree
  - Expand the parse tree in pre-order i.e. identify parent nodes before children

# Top-down Parsing

$S \rightarrow cAd$

$A \rightarrow ab \mid a$

String: cad

↑: next symbol to  
be read

*We need to backtrack  
after step 3 and reset  
input pointer*

*Can we do better ?*

Step	Input string	Parse tree
1	cad ↑	S
2	cad ↑ ↑	<pre>graph TD     S --&gt; c     S --&gt; A     S --&gt; d</pre>
3	cad ↑	<pre>graph TD     S --&gt; c     S --&gt; A     S --&gt; d     A --&gt; a     A --&gt; b</pre>
4	cad ↑	<pre>graph TD     S --&gt; c     S --&gt; A     S --&gt; d     A --&gt; a</pre>

# Top-down Parsing

- 1)  $S \rightarrow F$
- 2)  $S \rightarrow (S + F)$
- 3)  $F \rightarrow a$

string: (a+a)

string': (a+a)\$

	(	)	a	+	\$
S	2	-	1	-	-
F	-	-	3	-	-

*Assume that the table is given.*

- Table-driven (Parse Table) approach doesn't require backtracking

*But how do we construct such a table?*

# First and follow sets

- $\text{First}(\alpha)$ : the set of terminals (and/or  $\lambda$ ) that begin all strings that can be derived from  $\alpha$

- $\text{First}(A) = \{x, y, \lambda\}$

- $\text{First}(xA) = \{x\}$

- $\text{First}(AB) = \{x, y, b\}$

- $\text{Follow}(A)$ : the set of terminals (and/or \$, but no  $\lambda$ s) that can appear immediately after A in some partial derivation

- $\text{Follow}(A) = \{b\}$

$$S \rightarrow A B \$$$

$$A \rightarrow x a A$$

$$A \rightarrow y a A$$

$$A \rightarrow \lambda$$

$$B \rightarrow b$$

# First and follow sets

- $\text{First}(\alpha) = \{a \in V_t \mid \alpha \Rightarrow^* a\beta\} \cup \{\lambda \mid \text{if } \alpha \Rightarrow^* \lambda\}$
- $\text{Follow}(A) = \{a \in V_t \mid S \Rightarrow^+ \dots Aa \dots\} \cup \{\$ \mid \text{if } S \Rightarrow^+ \dots A \$\}$

S: start symbol

a: a terminal symbol

A: a non-terminal symbol

$\alpha, \beta$ : a string composed of terminals and non-terminals (typically,  $\alpha$  is the RHS of a production

$\Rightarrow$ : derived in 1 step

$\Rightarrow^*$ : derived in 0 or more steps

$\Rightarrow^+$ : derived in 1 or more steps

# Computing first sets

- Terminal:  $\text{First}(a) = \{a\}$
- Non-terminal:  $\text{First}(A)$ 
  - Look at all productions for  $A$   
$$A \rightarrow X_1 X_2 \dots X_k$$
  - $\text{First}(A) \supseteq (\text{First}(X_1) - \lambda)$
  - If  $\lambda \in \text{First}(X_1)$ ,  $\text{First}(A) \supseteq (\text{First}(X_2) - \lambda)$
  - If  $\lambda$  is in  $\text{First}(X_i)$  for all  $i$ , then  $\lambda \in \text{First}(A)$
- Computing  $\text{First}(\alpha)$ : similar procedure to computing  $\text{First}(A)$



# Top-down Parsing – predictive parsers

- Idea: we know sentence has to start with initial symbol
- Build up partial derivations by *predicting* what rules are used to expand non-terminals
  - Often called *predictive parsers*
- If partial derivation has terminal characters, *match* them from the input stream

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

- A sentence in the grammar:

$x a c c \$$

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

special “end of input” symbol

- A sentence in the grammar:

$x a c c \$$

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

• A sentence in the grammar:

$x a c c \$$

Current derivation:  $S$

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

• A sentence in the grammar:

$x a c c \$$

Current derivation:  $A B c \$$

Predict rule

# A simple example

$S \rightarrow A B c \$$

Choose based on  
*first set* of rules

$A \rightarrow x a A$   
 $A \rightarrow y a A$   
 $A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

• A sentence in the grammar:

$x a c c \$$

Current derivation:  $x a A B c \$$

Predict rule *based on next token*

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

• A sentence in the grammar:

$x a c c \$$

Current derivation:  $x a A B c \$$

Match token

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

• A sentence in the grammar:

$x a c c \$$

Current derivation:  $x a A B c \$$

Match token



# A simple example

$$S \rightarrow A B c \$$$

Choose based on  
*first set* of rules

$$\begin{array}{l} A \rightarrow x a A \\ A \rightarrow y a A \\ A \rightarrow c \end{array}$$

$$B \rightarrow b$$

• A sentence in the grammar:

$$B \rightarrow \lambda$$

$x a c c \$$

Current derivation:  $x a c B c \$$

Predict rule *based on next token*

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

• A sentence in the grammar:

$x a c c \$$

Current derivation:  $x a c B c \$$

Match token

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

Choose based on  
*follow set*

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

- A sentence in the grammar:

$x a c c \$$

Current derivation:  $x a c \lambda c \$$

Predict rule *based on next token*

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

• A sentence in the grammar:

$x a c c \$$

Current derivation:  $x a c c \$$

Match token

# A simple example

$S \rightarrow A B c \$$

$A \rightarrow x a A$

$A \rightarrow y a A$

$A \rightarrow c$

$B \rightarrow b$

$B \rightarrow \lambda$

- A sentence in the grammar:

$x a c c \$$

Current derivation:  $x a c c \$$

Match token

# Top-down vs. Bottom-up parsers

- Top-down parsers expand the parse tree in *pre-order*
  - Identify parent nodes before the children
- Bottom-up parsers expand the parse tree in *post-order*
  - Identify children before the parents
- Notation:
  - LL(1): Top-down derivation with 1 symbol lookahead
  - LL(k): Top-down derivation with k symbols lookahead
  - LR(1): Bottom-up derivation with 1 symbol lookahead

# Suggested Reading

- Alfred V. Aho, Monica S. Lam, Ravi Sethi and Jeffrey D. Ullman: Compilers: Principles, Techniques, and Tools, 2/E, AddisonWesley 2007
  - Chapter 4 (Sections: 4.1 to 4.4)
- Fisher and LeBlanc: Crafting a Compiler with C
  - Chapter 4, Chapter 5 (Sections 5.1 to 5.5, 5.9)