

# Parser: Syntactic Analysis

Readings: EAC2 Chapter 3



EECS4302 M:  
Compilers and Interpreters  
Winter 2020

CHEN-WEI WANG



## Context-Free Languages: Introduction

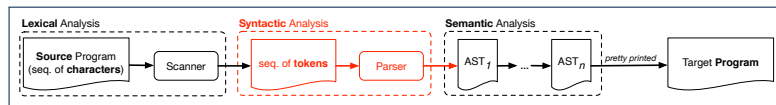
- We have seen **regular languages** :
  - Can be described using *finite automata* or *regular expressions*.
  - Satisfy the *pumping lemma*.
- Languages with a *recursive* structure are provably *non-regular*.  
e.g.,  $\{0^n 1^n \mid n \geq 0\}$
- **Context-free grammars (CFG's)** are used to describe strings that can be generated in a *recursive* fashion.
- **Context-free languages (CFL's)** are:
  - Languages that can be described using CFG's.
  - A proper superset of the set of regular languages.

3 of 96

## Parser in Context



- Recall:



- Treats the input programs as a **a sequence of classified tokens/words**
- Applies rules **parsing** token sequences as **abstract syntax trees (ASTs)** [ **syntactic** analysis ]
- Upon termination:
  - Reports token sequences not derivable as ASTs
  - Produces an **AST**
- No longer considers **every character** in input program.
- **Derivable** token sequences constitute a **context-free language (CFL)** .

2 of 96

## CFG: Example (1.1)



- The language that we previously proved as *non-regular*

$$\{0^n \# 1^n \mid n \geq 0\}$$

can be described using the following **grammar** :

$$\begin{aligned} A &\rightarrow 0A1 \\ A &\rightarrow B \\ B &\rightarrow \# \end{aligned}$$

- A grammar contains a collection of *substitution* or *production* rules, where:
  - A **terminal** is a word  $w \in \Sigma^*$  (e.g., 0, 1, etc.).
  - A **variable** or **non-terminal** is a word  $w \notin \Sigma^*$  (e.g., A, B, etc.).
  - A **start variable** occurs on the LHS of the topmost rule (e.g., A).

4 of 96

## CFG: Example (1.2)

- Given a grammar, generate a string by:
  - Write down the *start variable*.
  - Choose a production rule where the *start variable* appears on the LHS of the arrow, and *substitute* it by the RHS.
  - There are two cases of the re-written string:
    - It contains *no* variables, then you are done.
    - It contains *some* variables, then *substitute* each variable using the relevant *production rules*.
  - Repeat Step 3.
- e.g., We can generate an *infinite* number of strings from

$$\begin{aligned} A &\rightarrow 0A1 \\ A &\rightarrow B \\ B &\rightarrow \# \end{aligned}$$

- $A \Rightarrow B \Rightarrow \#$
- $A \Rightarrow 0A1 \Rightarrow 0B1 \Rightarrow 0\#1$
- $A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 00B11 \Rightarrow 00\#11$
- ...

5 of 96

## CFG: Example (2)

Design a CFG for the following language:

$$\{w \mid w \in \{0,1\}^* \wedge w \text{ is a palidrome}\}$$

e.g., 00, 11, 0110, 1001, etc.

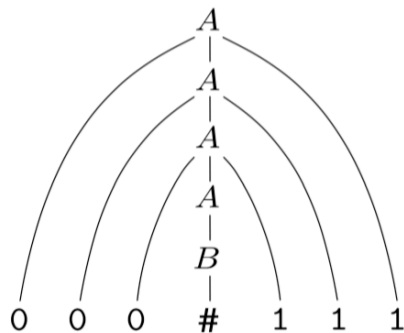
$$\begin{aligned} P &\rightarrow \epsilon \\ P &\rightarrow 0 \\ P &\rightarrow 1 \\ P &\rightarrow 0P0 \\ P &\rightarrow 1P1 \end{aligned}$$

7 of 96

## CFG: Example (1.2)

Given a CFG, the *derivation* of a string can be shown as a *parse tree*.

e.g., The derivation of 000#111 has the parse tree



6 of 96

## CFG: Example (3)

Design a CFG for the following language:

$$\{ww^R \mid w \in \{0,1\}^*\}$$

e.g., 00, 11, 0110, etc.

$$\begin{aligned} P &\rightarrow \epsilon \\ P &\rightarrow 0P0 \\ P &\rightarrow 1P1 \end{aligned}$$

8 of 96

## CFG: Example (4)



Design a CFG for the set of binary strings, where each block of 0's followed by at least as many 1's.

e.g., 000111, 0001111, etc.

- We use  $S$  to represent one such string, and  $A$  to represent each such block in  $S$ .

$S \rightarrow \epsilon$  {BC of  $S$ }  
 $S \rightarrow AS$  {RC of  $S$ }  
 $A \rightarrow \epsilon$  {BC of  $A$ }  
 $A \rightarrow 01$  {BC of  $A$ }  
 $A \rightarrow 0A1$  {RC of  $A$ : equal 0's and 1's}  
 $A \rightarrow A1$  {RC of  $A$ : more 1's}

9 of 96

## CFG: Example (5.2) Version 1



$Expression \rightarrow IntegerConstant$   
 $Expression \rightarrow -IntegerConstant$   
 $Expression \rightarrow BooleanConstant$   
 $Expression \rightarrow BinaryOp$   
 $Expression \rightarrow UnaryOp$   
 $Expression \rightarrow ( Expression )$

$IntegerConstant \rightarrow Digit$   
 $IntegerConstant \rightarrow Digit IntegerConstant$

$Digit \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

$BooleanConstant \rightarrow TRUE$   
 $BooleanConstant \rightarrow FALSE$

11 of 96

## CFG: Example (5.1) Version 1



Design the grammar for the following small expression language, which supports:

- Arithmetic operations: +, -, \*, /
- Relational operations: >, <, >=, <=, ==, /=
- Logical operations: true, false, !, &&, ||, =>

Start with the variable  $Expression$ .

- There are two possible versions:
  - All operations are mixed together. [e.g.,  $(1 + true)/false$ ]
  - Relevant operations are grouped together. Try both!

10 of 96

## CFG: Example (5.3) Version 1



$BinaryOp \rightarrow Expression + Expression$   
 $BinaryOp \rightarrow Expression - Expression$   
 $BinaryOp \rightarrow Expression * Expression$   
 $BinaryOp \rightarrow Expression / Expression$   
 $BinaryOp \rightarrow Expression \&\& Expression$   
 $BinaryOp \rightarrow Expression || Expression$   
 $BinaryOp \rightarrow Expression => Expression$   
 $BinaryOp \rightarrow Expression == Expression$   
 $BinaryOp \rightarrow Expression /= Expression$   
 $BinaryOp \rightarrow Expression > Expression$   
 $BinaryOp \rightarrow Expression < Expression$

$UnaryOp \rightarrow ! Expression$

12 of 96

## CFG: Example (5.4) Version 1



However, Version 1 of CFG:

- **Parses** string that requires further **semantic analysis** (e.g., type checking):  
e.g.,  $3 \Rightarrow 6$
- Is **ambiguous**, meaning that a string may have more than one ways to interpret it.  
e.g., Draw the **parse tree(s)** for  $3 * 5 + 4$

13 of 96

## CFG: Example (5.6) Version 2



```
ArithmeticOp → ArithmeticOp + ArithmeticOp
              | ArithmeticOp - ArithmeticOp
              | ArithmeticOp * ArithmeticOp
              | ArithmeticOp / ArithmeticOp
              | (ArithmeticOp)
              | IntegerConstant
              | -IntegerConstant
RelationalOp  → ArithmeticOp == ArithmeticOp
              | ArithmeticOp /= ArithmeticOp
              | ArithmeticOp > ArithmeticOp
              | ArithmeticOp < ArithmeticOp
LogicalOp     → LogicalOp && LogicalOp
              | LogicalOp || LogicalOp
              | LogicalOp => LogicalOp
              | ! LogicalOp
              | (LogicalOp)
              | RelationalOp
              | BooleanConstant
```

15 of 96

## CFG: Example (5.5) Version 2



```
Expression → ArithmeticOp
            | RelationalOp
            | LogicalOp
            | ( Expression )
IntegerConstant → Digit
                | Digit IntegerConstant
Digit           → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
BooleanConstant → TRUE
                | FALSE
```

14 of 96

## CFG: Example (5.7) Version 2



However, Version 2 of CFG:

- Eliminates some cases for further semantic analysis:  
e.g.,  $(1 + 2) \Rightarrow (5 / 4)$  [ no parse tree ]
- Still **Parses** string that might require further **semantic analysis** :  
e.g.,  $(1 + 2) / (5 - (2 + 3))$
- Is **ambiguous**, meaning that a string may have more than one ways to interpret it.  
e.g., Draw the **parse tree(s)** for  $3 * 5 + 4$

16 of 96



## CFG: Formal Definition (1)

- A **context-free grammar (CFG)** is a 4-tuple  $(V, \Sigma, R, S)$ :
  - $V$  is a finite set of **variables**.
  - $\Sigma$  is a finite set of **terminals**.  $[V \cap \Sigma = \emptyset]$
  - $R$  is a finite set of **rules** s.t.

$$R \subseteq \{v \rightarrow s \mid v \in V \wedge s \in (V \cup \Sigma)^*\}$$

- $S \in V$  is the **start variable**.
- Given strings  $u, v, w \in (V \cup \Sigma)^*$ , variable  $A \in V$ , and a rule  $A \rightarrow w$ :
  - $uAv \Rightarrow uwv$  means that  $uAv$  **yields**  $uwv$ .
  - $u \xRightarrow{*} v$  means that  $u$  **derives**  $v$ , if:
    - $u = v$ ; or
    - $u \Rightarrow u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_k \Rightarrow v$  [a yield sequence]
- Given a CFG  $G = (V, \Sigma, R, S)$ , the language of  $G$

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

17 of 96

## CFG: Formal Definition (3): Example

- Consider the grammar  $G = (V, \Sigma, R, S)$ :

- $R$  is

$$\begin{array}{lcl} \text{Expr} & \rightarrow & \text{Expr} + \text{Term} \\ & & | \text{Term} \\ \text{Term} & \rightarrow & \text{Term} * \text{Factor} \\ & & | \text{Factor} \\ \text{Factor} & \rightarrow & (\text{Expr}) \\ & & | a \end{array}$$

- $V = \{\text{Expr}, \text{Term}, \text{Factor}\}$
- $\Sigma = \{a, +, *, (, )\}$
- $S = \text{Expr}$
- Precedence** of operators  $+$  and  $*$  is embedded in the grammar.
  - "Plus" is specified at a **higher** level (**Expr**) than is "times" (**Term**).
  - Both operands of a multiplication (**Factor**) may be **parenthesized**.

19 of 96

## CFG: Formal Definition (2): Example

- Design the CFG for strings of properly-nested parentheses.  
e.g.,  $()$ ,  $()()$ ,  $((()()))()$ , etc.  
Present your answer in a formal manner.
- $G = (\{S\}, \{(), \epsilon\}, R, S)$ , where  $R$  is

$$S \rightarrow ( S ) \mid SS \mid \epsilon$$

- Draw **parse trees** for the above three strings that  $G$  generates.

18 of 96

## Regular Expressions to CFG's

- Recall the semantics of regular expressions (assuming that we do not consider  $\emptyset$ ):

$$\begin{array}{lcl} L(\epsilon) & = & \{\epsilon\} \\ L(a) & = & \{a\} \\ L(E + F) & = & L(E) \cup L(F) \\ L(EF) & = & L(E)L(F) \\ L(E^*) & = & (L(E))^* \\ L(E) & = & L(E) \end{array}$$

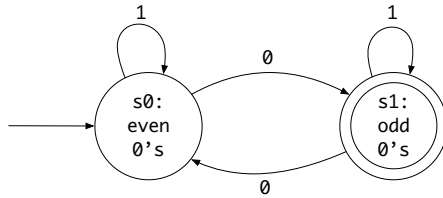
- e.g., Grammar for  $(00 + 1)^* + (11 + 0)^*$

$$\begin{array}{lcl} S & \rightarrow & A \mid B \\ A & \rightarrow & \epsilon \mid AC \\ C & \rightarrow & 00 \mid 1 \\ B & \rightarrow & \epsilon \mid BD \\ D & \rightarrow & 11 \mid 0 \end{array}$$

20 of 96

## DFA to CFG's

- Given a DFA  $M = (Q, \Sigma, \delta, q_0, F)$ :
  - Make a variable  $R_i$  for each state  $q_i \in Q$ .
  - Make  $R_0$  the start variable, where  $q_0$  is the start state of  $M$ .
  - Add a rule  $R_i \rightarrow aR_j$  to the grammar if  $\delta(q_i, a) = q_j$ .
  - Add a rule  $R_i \rightarrow \epsilon$  if  $q_i \in F$ .
- e.g., Grammar for



$$R_0 \rightarrow 1R_0 \mid 0R_1$$

$$R_1 \rightarrow 0R_0 \mid 1R_1 \mid \epsilon$$

21 of 96

## CFG: Rightmost Derivations (1)

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr} + \text{Term} \mid \text{Term} \\ \text{Term} &\rightarrow \text{Term} * \text{Factor} \mid \text{Factor} \\ \text{Factor} &\rightarrow (\text{Expr}) \mid a \end{aligned}$$

- Unique rightmost derivation** for the string  $a + a * a$ :

$$\begin{aligned} \text{Expr} &\Rightarrow \text{Expr} + \text{Term} \\ &\Rightarrow \text{Expr} + \text{Term} * \text{Factor} \\ &\Rightarrow \text{Expr} + \text{Term} * a \\ &\Rightarrow \text{Expr} + \text{Factor} * a \\ &\Rightarrow \text{Expr} + a * a \\ &\Rightarrow \text{Term} + a * a \\ &\Rightarrow \text{Factor} + a * a \\ &\Rightarrow a + a * a \end{aligned}$$

- This **rightmost derivation** suggests that  $a * a$  is the right operand of  $+$ .

23 of 96

## CFG: Leftmost Derivations (1)

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr} + \text{Term} \mid \text{Term} \\ \text{Term} &\rightarrow \text{Term} * \text{Factor} \mid \text{Factor} \\ \text{Factor} &\rightarrow (\text{Expr}) \mid a \end{aligned}$$

- Unique leftmost derivation** for the string  $a + a * a$ :

$$\begin{aligned} \text{Expr} &\Rightarrow \text{Expr} + \text{Term} \\ &\Rightarrow \text{Term} + \text{Term} \\ &\Rightarrow \text{Factor} + \text{Term} \\ &\Rightarrow a + \text{Term} \\ &\Rightarrow a + \text{Term} * \text{Factor} \\ &\Rightarrow a + \text{Factor} * \text{Factor} \\ &\Rightarrow a + a * \text{Factor} \\ &\Rightarrow a + a * a \end{aligned}$$

- This **leftmost derivation** suggests that  $a * a$  is the right operand of  $+$ .

22 of 96

## CFG: Leftmost Derivations (2)

$$\begin{aligned} \text{Expr} &\rightarrow \text{Expr} + \text{Term} \mid \text{Term} \\ \text{Term} &\rightarrow \text{Term} * \text{Factor} \mid \text{Factor} \\ \text{Factor} &\rightarrow (\text{Expr}) \mid a \end{aligned}$$

- Unique leftmost derivation** for the string  $(a + a) * a$ :

$$\begin{aligned} \text{Expr} &\Rightarrow \text{Term} \\ &\Rightarrow \text{Term} * \text{Factor} \\ &\Rightarrow \text{Factor} * \text{Factor} \\ &\Rightarrow (\text{Expr}) * \text{Factor} \\ &\Rightarrow (\text{Expr} + \text{Term}) * \text{Factor} \\ &\Rightarrow (\text{Term} + \text{Term}) * \text{Factor} \\ &\Rightarrow (\text{Factor} + \text{Term}) * \text{Factor} \\ &\Rightarrow (a + \text{Term}) * \text{Factor} \\ &\Rightarrow (a + \text{Factor}) * \text{Factor} \\ &\Rightarrow (a + a) * \text{Factor} \\ &\Rightarrow (a + a) * a \end{aligned}$$

- This **leftmost derivation** suggests that  $(a + a)$  is the left operand of  $*$ .

24 of 96

## CFG: Rightmost Derivations (2)



$Expr \rightarrow Expr + Term \mid Term$   
 $Term \rightarrow Term * Factor \mid Factor$   
 $Factor \rightarrow (Expr) \mid a$

- Unique rightmost derivation for the string  $(a + a) * a$ :

$Expr \Rightarrow Term$   
 $\Rightarrow Term * Factor$   
 $\Rightarrow Term * a$   
 $\Rightarrow Factor * a$   
 $\Rightarrow (Expr) * a$   
 $\Rightarrow (Expr + Term) * a$   
 $\Rightarrow (Expr + Factor) * a$   
 $\Rightarrow (Expr + a) * a$   
 $\Rightarrow (Term + a) * a$   
 $\Rightarrow (Factor + a) * a$   
 $\Rightarrow (a + a) * a$

This rightmost derivation suggests that  $(a + a)$  is the left operand of  $*$ .

25 of 96

## CFG: Parse Trees vs. Derivations (2)



- A string  $w \in \Sigma^*$  may have more than one derivations.
- Q: distinct derivations for  $w \in \Sigma^* \Rightarrow$  distinct parse trees for  $w$ ?
- A: Not in general  $\therefore$  Derivations with distinct orders of variable substitutions may still result in the same parse tree.
- For example:

$Expr \rightarrow Expr + Term \mid Term$   
 $Term \rightarrow Term * Factor \mid Factor$   
 $Factor \rightarrow (Expr) \mid a$

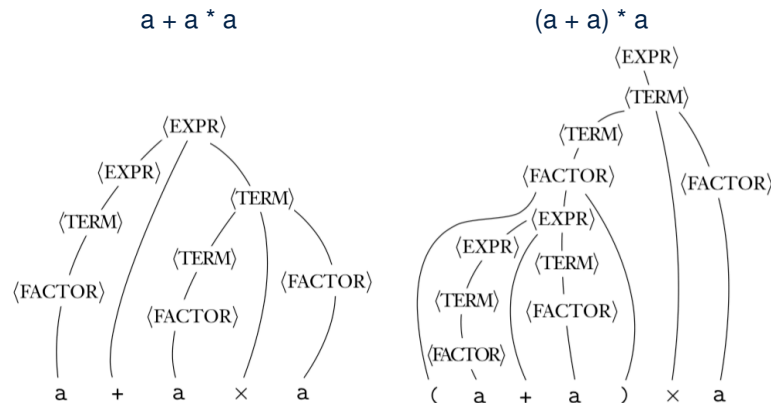
For string  $a + a * a$ , the leftmost and rightmost derivations have distinct orders of variable substitutions, but their corresponding parse trees are the same.

27 of 96

## CFG: Parse Trees vs. Derivations (1)



- Parse trees for (leftmost & rightmost) derivations of expressions:



- Orders in which derivations are performed are not reflected on parse trees.

26 of 96

## CFG: Ambiguity: Definition



Given a grammar  $G = (V, \Sigma, R, S)$ :

- A string  $w \in \Sigma^*$  is derived ambiguously in  $G$  if there exist two or more distinct parse trees or, equally, two or more distinct leftmost derivations or, equally, two or more distinct rightmost derivations.
- Here we require that all such derivations have been completed by following a particular order (leftmost or rightmost) to avoid false alarm.
- $G$  is ambiguous if it generates some string ambiguously.

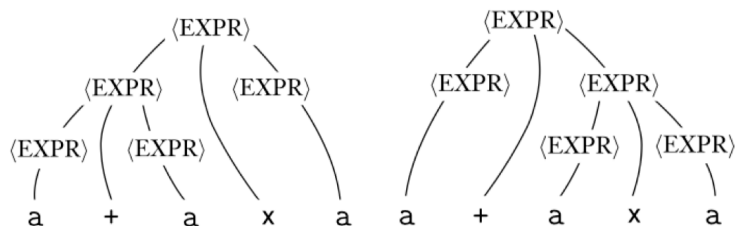
28 of 96

## CFG: Ambiguity: Exercise (1)

- Is the following grammar **ambiguous**?

$$\text{Expr} \rightarrow \text{Expr} + \text{Expr} \mid \text{Expr} * \text{Expr} \mid (\text{Expr}) \mid a$$

- Yes  $\therefore$  it generates the string  $a + a * a$  **ambiguously**:

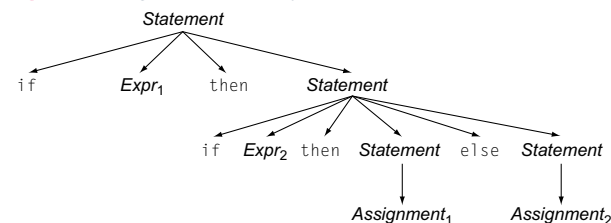


- Distinct ASTs** (for the **same input**) mean **distinct semantic interpretations**: e.g., when a post-order traversal is used to implement evaluation
- Exercise**: Show **leftmost** derivations for the two parse trees.

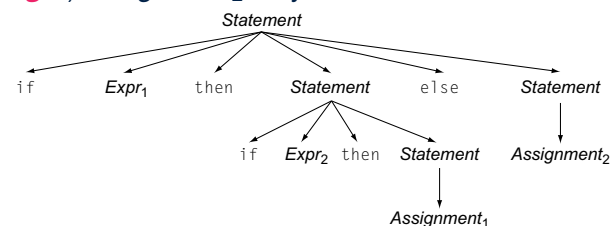
29 of 96

## CFG: Ambiguity: Exercise (2.2)

- (**Meaning 1**)  $\text{Assignment}_2$  may be associated with the inner **if**:



- (**Meaning 2**)  $\text{Assignment}_2$  may be associated with the outer **if**:



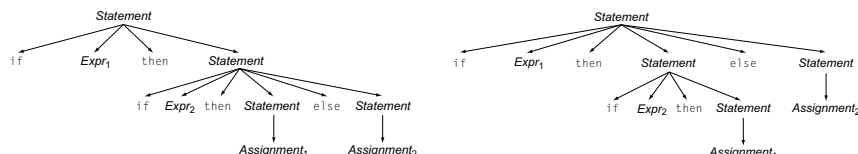
31 of 96

## CFG: Ambiguity: Exercise (2.1)

- Is the following grammar **ambiguous**?

$$\begin{aligned} \text{Statement} &\rightarrow \text{if } \text{Expr} \text{ then } \text{Statement} \\ &\quad \mid \text{if } \text{Expr} \text{ then } \text{Statement} \text{ else } \text{Statement} \\ &\quad \mid \text{Assignment} \end{aligned}$$

- Yes  $\therefore$  it generates the following string **ambiguously**:

$$\text{if } \text{Expr}_1 \text{ then if } \text{Expr}_2 \text{ then } \text{Assignment}_1 \text{ else } \text{Assignment}_2$$


- This is called the **dangling else** problem.
- Exercise**: Show **leftmost** derivations for the two parse trees.

30 of 96

## CFG: Ambiguity: Exercise (2.3)

- We may remove the **ambiguity** by specifying that the **dangling else** is associated with the **nearest if**:

$$\begin{aligned} \text{Statement} &\rightarrow \text{if } \text{Expr} \text{ then } \text{Statement} \\ &\quad \mid \text{if } \text{Expr} \text{ then } \text{WithElse} \text{ else } \text{Statement} \\ &\quad \mid \text{Assignment} \\ \text{WithElse} &\rightarrow \text{if } \text{Expr} \text{ then } \text{WithElse} \text{ else } \text{WithElse} \\ &\quad \mid \text{Assignment} \end{aligned}$$

- When applying `if ... then WithElse else Statement`:
  - The **true** branch will be produced via **WithElse**.
  - The **false** branch will be produced via **Statement**.

There is **no circularity** between the two non-terminals.

32 of 96

## Discovering Derivations

- Given a CFG  $G = (V, \Sigma, R, S)$  and an input program  $p \in \Sigma^*$ :
  - So far we **manually** come up a valid derivation  $S \Rightarrow^* p$ .
  - A parser is supposed to **automate** this derivation process.  
Given an input sequence of  $(t, c)$  pairs, where **token**  $t$  (e.g., r241) belongs to some **syntactic category**  $c$  (e.g., register):  
Either output a **valid derivation** (as an **AST**), or signal an **error**.
- In the process of building an **AST** for the input program:
  - Root** of AST: **start symbol**  $S$  of  $G$
  - Internal nodes**: A subset of variables  $V$  of  $G$
  - Leaves** of AST: **token sequence** input by the scanner  
⇒ Discovering the **grammatical connections** (according to  $R$ ) between the **root**, **internal nodes**, and **leaves** is the hard part!
- Approaches to Parsing:  $[w \in (V \cup \Sigma)^*, A \in V, A \rightarrow w \in R]$ 
  - Top-down** parsing  
For a node representing  $A$ , extend it with a subtree representing  $w$ .
  - Bottom-up** parsing  
For a substring matching  $w$ , build a node representing  $A$  accordingly.

33 of 96

## TDP: Exercise (1)

- Given the following CFG  $G$ :

$$\begin{array}{l} \text{Expr} \rightarrow \text{Expr} + \text{Term} \\ \quad \quad | \quad \text{Term} \\ \text{Term} \rightarrow \text{Term} * \text{Factor} \\ \quad \quad | \quad \text{Factor} \\ \text{Factor} \rightarrow (\text{Expr}) \\ \quad \quad | \quad a \end{array}$$

Trace  $TDParse$  on how to build an AST for input  $a + a * a$ .

- Running  $TDParse$  with  $G$  results an **infinite loop** !!!
  - $TDParse$  focuses on the **leftmost** non-terminal.
  - The grammar  $G$  contains **left-recursions**.
- We must first convert left-recursions in  $G$  to **right-recursions**.

35 of 96

## TDP: Discovering Leftmost Derivation

```

ALGORITHM: TDParse
INPUT: CFG G = (V, Σ, R, S)
OUTPUT: Root of a Parse Tree or Syntax Error
PROCEDURE:
  root := a new node for the start symbol S
  focus := root
  initialize an empty stack trace
  trace.push(null)
  word := NextWord()
  while (true):
    if focus ∈ V then
      if ∃ unvisited rule focus → β1β2...βn ∈ R then
        create β1, β2...βn as children of focus
        trace.push(βnβn-1...β2)
        focus := β1
      else
        if focus = S then report syntax error
        else backtrack
    elseif word matches focus then
      word := NextWord()
      focus := trace.pop()
    elseif word = EOF ∧ focus = null then return root
    else backtrack
  
```

**backtrack** ≝ pop  $focus.siblings$ ;  $focus := focus.parent$ ;  $focus.resetChildren$

34 of 96

## TDP: Exercise (2)

- Given the following CFG  $G$ :

$$\begin{array}{l} \text{Expr} \rightarrow \text{Term Expr}' \\ \text{Expr}' \rightarrow + \text{Term Expr}' \\ \quad \quad | \quad \epsilon \\ \text{Term} \rightarrow \text{Factor Term}' \\ \text{Term}' \rightarrow * \text{Factor Term}' \\ \quad \quad | \quad \epsilon \\ \text{Factor} \rightarrow (\text{Expr}) \\ \quad \quad | \quad a \end{array}$$

**Exercise.** Trace  $TDParse$  on building AST for  $a + a * a$ .

**Exercise.** Trace  $TDParse$  on building AST for  $(a + a) * a$ .

**Q:** How to handle  $\epsilon$ -productions (e.g.,  $\text{Expr} \rightarrow \epsilon$ )?

**A:** Execute  $focus := trace.pop()$  to advance to next node.

- Running  $TDParse$  will **terminate** ∵  $G$  is **right-recursive**.
- We will learn about a systematic approach to converting left-recursions in a given grammar to **right-recursions**.

36 of 96

## Left-Recursions (LR): Direct vs. Indirect

Given CFG  $G = (V, \Sigma, R, S)$ ,  $\alpha, \beta, \gamma \in (V \cup \Sigma)^*$ ,  $G$  contains:

- A **cycle** if  $\exists A \in V \bullet A \xrightarrow{*} A$
- A **direct** LR if  $A \rightarrow A\alpha \in R$  for non-terminal  $A \in V$   
e.g.,

Expr	→	Expr + Term
		Term
Term	→	Term * Factor
		Factor
Factor	→	(Expr)
		a

e.g.,

Expr	→	Expr + Term
		Expr - Term
		Term
Term	→	Term * Factor
		Term / Factor
		Factor

- An **indirect** LR if  $A \rightarrow B\beta \in R$  for non-terminals  $A, B \in V$ ,  $B \xrightarrow{*} A\gamma$

A	→	Br
B	→	Cd
C	→	At

A	→	Ba		b
B	→	Cd		e
C	→	Df		g
D	→	f		Aa   Cg

$A \rightarrow Br, B \xrightarrow{*} Atd$

$A \rightarrow Ba, B \xrightarrow{*} Aafd$

37 of 96

## CFG: Eliminating $\epsilon$ -Productions (1)

- Motivations:
  - *TDParse* requires CFG with no  $\epsilon$ -productions.
  - *RemoveLR* produces CFG which may contain  $\epsilon$ -productions.
- $\epsilon \notin L \Rightarrow \exists$  CFG  $G = (V, \Sigma, R, S)$  s.t.  $G$  has no  $\epsilon$ -productions.  
An  **$\epsilon$ -production** has the form  $A \rightarrow \epsilon$ .
- A variable  $A$  is **nullable** if  $A \xrightarrow{*} \epsilon$ .
  - Each terminal symbol is **not nullable**.
  - Variable  $A$  is **nullable** if either:
    - $A \rightarrow \epsilon \in R$ ; or
    - $A \rightarrow B_1 B_2 \dots B_k \in R$ , where each variable  $B_i$  ( $1 \leq i \leq k$ ) is a **nullable**.
- Given a production  $B \rightarrow CAD$ , if only variable  $A$  is nullable, then there are 2 versions of  $B$ :  $B \rightarrow CAD \mid CD$
- In general, given a production  $A \rightarrow X_1 X_2 \dots X_k$  with  $k$  symbols, if  $m$  of the  $k$  symbols are **nullable**:
  - $m < k$ : There are  $2^m$  versions of  $A$ .
  - $m = k$ : There are  $2^m - 1$  versions of  $A$ . [excluding  $A \rightarrow \epsilon$ ]

39 of 96

## TDP: (Preventively) Eliminating LR

```

1 ALGORITHM: RemoveLR
2 INPUT: CFG G=(V, Σ, R, S)
3 ASSUME: G acyclic ∧ with no ε-productions
4 OUTPUT: G' s.t. G' ≡ G, G' has no
5         indirect & direct left-recursions
6 PROCEDURE:
7   impose an order on V: ⟨(A1, A2, ..., An)⟩
8   for i: 1 .. n:
9     for j: 1 .. i-1:
10      if ∃ Ai → Ajγ ∈ R ∧ Aj → δ1 | δ2 | ... | δm ∈ R then
11        replace Ai → Ajγ with Ai → δ1γ | δ2γ | ... | δmγ
12      end
13   for Ai → Aiα | β ∈ R:
14     replace it with: Ai → βA', A' → αA' | ε
    
```

L9 to L11: Remove **indirect** left-recursions from  $A_1$  to  $A_{i-1}$ .

L12 to L13: Remove **direct** left-recursions from  $A_1$  to  $A_{i-1}$ .

**Loop Invariant (outer for-loop)?** At the start of  $i^{th}$  iteration:

- No direct or indirect left-recursions for  $A_1, A_2, \dots, A_{i-1}$ .
- More precisely:  $\forall k: k < i \bullet \neg(\exists l \bullet l \leq k \wedge A_k \rightarrow A_l \dots \in R)$

38 of 96

## CFG: Eliminating $\epsilon$ -Productions (2)

- Eliminate  $\epsilon$ -productions from the following grammar:

$$\begin{aligned}
 S &\rightarrow AB \\
 A &\rightarrow aAA \mid \epsilon \\
 B &\rightarrow bBB \mid \epsilon
 \end{aligned}$$

- Which are the **nullable** variables?

[S, A, B]

$$\begin{aligned}
 S &\rightarrow A \mid B \mid AB \quad \{S \rightarrow \epsilon \text{ not included}\} \\
 A &\rightarrow aAA \mid aA \mid a \quad \{A \rightarrow aA \text{ duplicated}\} \\
 B &\rightarrow bBB \mid bB \mid b \quad \{B \rightarrow bB \text{ duplicated}\}
 \end{aligned}$$

40 of 96

## Backtrack-Free Parsing (1)

- TDParse automates the **top-down, leftmost** derivation process by consistently choosing production rules (e.g., in order of their appearance in CFG).
  - This **inflexibility** may lead to **inefficient** runtime performance due to the need to **backtrack**.
  - e.g., It may take the **construction of a giant subtree** to find out a **mismatch** with the input tokens, which end up requiring it to **backtrack** all the way back to the **root** (start symbol).
- We may avoid backtracking with a modification to the parser:
  - When deciding which production rule to choose, consider:
    - (1) the **current** input symbol
    - (2) the **consequential first** symbol if a rule was applied for **focus** [ **lookahead** symbol ]
  - Using a **one symbol lookahead**, w.r.t. a **right-recursive** CFG, each alternative for the **leftmost nonterminal** leads to a **unique terminal**, allowing the parser to decide on a choice that prevents **backtracking**.
  - Such CFG is **backtrack free** with a **lookahead** of one symbol.
  - We also call such backtrack-free CFG a **predictive grammar**.

41 of 96

## The FIRST Set: Examples

- Consider this **right-recursive** CFG:

0	Goal	→	Expr	6	Term'	→	x Factor Term'
1	Expr	→	Term Expr'	7			+ Factor Term'
2	Expr'	→	+ Term Expr'	8			ε
3			- Term Expr'	9	Factor	→	( Expr )
4			ε	10			num
5	Term	→	Factor Term'	11			name

- Compute **FIRST** for each terminal (e.g., num, +, ( )):

	num	name	+	-	x	÷	(	)	eof	ε
FIRST	num	name	+	-	x	÷	(	)	eof	ε

- Compute **FIRST** for each non-terminal (e.g., Expr, Term'):

	Expr	Expr'	Term	Term'	Factor
FIRST	(, name, num	+, -, ε	(, name, num	x, ÷, ε	(, name, num

43 of 96

## The FIRST Set: Definition

- Say we write  $T \subset \mathbb{P}(\Sigma^*)$  to denote the set of valid tokens recognizable by the scanner.
- **FIRST**( $\alpha$ )  $\hat{=}$  set of symbols that can appear as the **first word** in some string derived from  $\alpha$ .
- More precisely:

$$\text{FIRST}(\alpha) = \begin{cases} \{\alpha\} & \text{if } \alpha \in T \\ \{w \mid w \in \Sigma^* \wedge \alpha \xRightarrow{*} w\beta \wedge \beta \in (V \cup \Sigma)^*\} & \text{if } \alpha \in V \end{cases}$$

42 of 96

## Computing the FIRST Set

$$\text{FIRST}(\alpha) = \begin{cases} \{\alpha\} & \text{if } \alpha \in T \\ \{w \mid w \in \Sigma^* \wedge \alpha \xRightarrow{*} w\beta \wedge \beta \in (V \cup \Sigma)^*\} & \text{if } \alpha \in V \end{cases}$$

**ALGORITHM: GetFirst**

**INPUT:** CFG  $G = (V, \Sigma, R, S)$   
 $T \subset \Sigma^*$  denotes valid terminals

**OUTPUT:** **FIRST** :  $V \cup T \cup \{\epsilon, eof\} \rightarrow \mathbb{P}(T \cup \{\epsilon, eof\})$

**PROCEDURE:**

```

for  $\alpha \in (T \cup \{eof, \epsilon\})$ : FIRST( $\alpha$ ) := { $\alpha$ }
for  $A \in V$ : FIRST( $A$ ) :=  $\emptyset$ 
lastFirst :=  $\emptyset$ 
while (lastFirst  $\neq$  FIRST):
    lastFirst := FIRST
    for  $A \rightarrow \beta_1\beta_2\dots\beta_k \in R$  s.t.  $\forall \beta_j: \beta_j \in (T \cup V)$ :
        rhs := FIRST( $\beta_1$ ) - { $\epsilon$ }
        for ( $i := 1$ ;  $\epsilon \in \text{FIRST}(\beta_i) \wedge i < k$ ;  $i++$ ):
            rhs := rhs  $\cup$  (FIRST( $\beta_{i+1}$ ) - { $\epsilon$ })
        if  $i = k \wedge \epsilon \in \text{FIRST}(\beta_k)$  then
            rhs := rhs  $\cup$  { $\epsilon$ }
        end
    FIRST( $A$ ) := FIRST( $A$ )  $\cup$  rhs
    
```

44 of 96



## Computing the FIRST Set: Extension



- Recall: **FIRST** takes as input a token or a variable.

$$\mathbf{FIRST} : V \cup T \cup \{\epsilon, eof\} \rightarrow \mathbb{P}(T \cup \{\epsilon, eof\})$$

- The computation of variable *rhs* in algorithm `GetFirst` actually suggests an extended, overloaded version:

$$\mathbf{FIRST} : (V \cup T \cup \{\epsilon, eof\})^* \rightarrow \mathbb{P}(T \cup \{\epsilon, eof\})$$

**FIRST** may also take as input a string  $\beta_1\beta_2\dots\beta_n$  (RHS of rules).

- More precisely:

$$\mathbf{FIRST}(\beta_1\beta_2\dots\beta_n) = \left\{ \begin{array}{l} \mathbf{FIRST}(\beta_1) \cup \mathbf{FIRST}(\beta_2) \cup \dots \cup \mathbf{FIRST}(\beta_k) \\ \wedge \\ \epsilon \notin \mathbf{FIRST}(\beta_k) \end{array} \middle| \begin{array}{l} \forall i: 1 \leq i < k \bullet \epsilon \in \mathbf{FIRST}(\beta_i) \end{array} \right\}$$

**Note.**  $\beta_k$  is the first symbol whose **FIRST** set does not contain  $\epsilon$ .

45 of 96

## Is the FIRST Set Sufficient



- Consider the following three productions:

$$\begin{array}{l} \text{Expr}' \rightarrow + \text{Term Term}' \quad (1) \\ \quad \quad \quad | - \text{Term Term}' \quad (2) \\ \quad \quad \quad | \epsilon \quad (3) \end{array}$$

In TDP, when the parser attempts to expand an *Expr'* node, it **looks ahead with one symbol** to decide on the choice of rule: **FIRST**(+) = {+}, **FIRST**(-) = {-}, and **FIRST**( $\epsilon$ ) = { $\epsilon$ }.

**Q.** When to choose rule (3) (causing *focus := trace.pop()*)?

**A?.** Choose rule (3) when *focus*  $\neq$  **FIRST**(+)  $\wedge$  *focus*  $\neq$  **FIRST**(-)?

- Correct** but **inefficient** in case of illegal input string: syntax error is only reported after possibly a long series of **backtrack**.
  - Useful if parser knows which words can appear, after an application of the  $\epsilon$ -production (rule (3)), as leading symbols.
- FOLLOW** ( $v : V$ )  $\triangleq$  set of symbols that can appear to the immediate right of a string derived from  $\alpha$ .

$$\mathbf{FOLLOW}(v) = \{w \mid w, x, y \in \Sigma^* \wedge v \xRightarrow{*} x \wedge S \xRightarrow{*} xwy\}$$

47 of 96

## Extended FIRST Set: Examples



Consider this **right**-recursive CFG:

0	Goal	$\rightarrow$	Expr	6	Term'	$\rightarrow$	$\times$ Factor Term'
1	Expr	$\rightarrow$	Term Expr'	7			$\div$ Factor Term'
2	Expr'	$\rightarrow$	+ Term Expr'	8			$\epsilon$
3			- Term Expr'	9	Factor	$\rightarrow$	( Expr )
4			$\epsilon$	10			num
5	Term	$\rightarrow$	Factor Term'	11			name

e.g., **FIRST**(Term Expr') = **FIRST**(Term) = { $\_$ , name, num}

e.g., **FIRST**(+ Term Expr') = **FIRST**(+) = {+}

e.g., **FIRST**(- Term Expr') = **FIRST**(-) = {-}

e.g., **FIRST**( $\epsilon$ ) = { $\epsilon$ }

46 of 96

## The FOLLOW Set: Examples



- Consider this **right**-recursive CFG:

0	Goal	$\rightarrow$	Expr	6	Term'	$\rightarrow$	$\times$ Factor Term'
1	Expr	$\rightarrow$	Term Expr'	7			$\div$ Factor Term'
2	Expr'	$\rightarrow$	+ Term Expr'	8			$\epsilon$
3			- Term Expr'	9	Factor	$\rightarrow$	( Expr )
4			$\epsilon$	10			num
5	Term	$\rightarrow$	Factor Term'	11			name

- Compute **FOLLOW** for each non-terminal (e.g., Expr, Term'):

	Expr	Expr'	Term	Term'	Factor
FOLLOW	eof, $\_$	eof, $\_$	eof, +, -, $\_$	eof, +, -, $\_$	eof, +, -, $\times$ , $\div$ , $\_$

48 of 96



## Computing the FOLLOW Set



$$\text{FOLLOW}(v) = \{w \mid w, x, y \in \Sigma^* \wedge v \xRightarrow{*} x \wedge S \xRightarrow{*} xwy\}$$

```

ALGORITHM: GetFollow
INPUT: CFG G = (V, Σ, R, S)
OUTPUT: FOLLOW: V → P(T ∪ {eof})
PROCEDURE:
for A ∈ V: FOLLOW(A) := ∅
FOLLOW(S) := {eof}
lastFollow := ∅
while (lastFollow ≠ FOLLOW):
    lastFollow := FOLLOW
    for A → β1β2...βk ∈ R:
        trailer := FOLLOW(A)
        for i: k .. 1:
            if βi ∈ V then
                FOLLOW(βi) := FOLLOW(βi) ∪ trailer
                if ε ∈ FIRST(βi)
                    then trailer := trailer ∪ (FIRST(βi) - ε)
                else trailer := FIRST(βi)
            else
                trailer := FIRST(βi)

```

49 of 96

## TDP: Lookahead with One Symbol



```

ALGORITHM: TDParse
INPUT: CFG G = (V, Σ, R, S)
OUTPUT: Root of a Parse Tree or Syntax Error
PROCEDURE:
root := a new node for the start symbol S
focus := root
initialize an empty stack trace
trace.push(null)
word := NextWord()
while (true):
    if focus ∈ V then % use FOLLOW set as well?
        if ∃ unvisited rule focus → β1β2...βn ∈ R ∧ word ∈ FIRST+(β) then
            create β1, β2...βn as children of focus
            trace.push(βnβn-1...β2)
            focus := β1
        else
            if focus = S then report syntax error
            else backtrack
    elseif word matches focus then
        word := NextWord()
        focus := trace.pop()
    elseif word = EOF ∧ focus = null then return root
    else backtrack

```

backtrack ≜ pop focus.siblings; focus := focus.parent; focus.resetChildren

51 of 96

## Backtrack-Free Grammar



- A **backtrack-free grammar** (for a **top-down parser**), when expanding the **focus internal node**, is always able to choose a **unique** rule with the **one-symbol lookahead** (or report a **syntax error** when no rule applies).
- To formulate this, we first define:

$$\text{FIRST}^+(A \rightarrow \beta) = \begin{cases} \text{FIRST}(\beta) & \text{if } \epsilon \notin \text{FIRST}(\beta) \\ \text{FIRST}(\beta) \cup \text{FOLLOW}(A) & \text{otherwise} \end{cases}$$

$\text{FIRST}(\beta)$  is the extended version where  $\beta$  may be  $\beta_1\beta_2\dots\beta_n$

- Now, a **backtrack-free grammar** has each of its productions  $A \rightarrow \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_n$  satisfying:

$$\forall i, j: 1 \leq i, j \leq n \wedge i \neq j \bullet \text{FIRST}^+(\gamma_i) \cap \text{FIRST}^+(\gamma_j) = \emptyset$$

50 of 96

## Backtrack-Free Grammar: Exercise



Is the following CFG **backtrack free**?

11	Factor	→	name
12			name [ ArgList ]
13			name ( ArgList )
15	ArgList	→	Expr MoreArgs
16	MoreArgs	→	, Expr MoreArgs
17			ε

- $\epsilon \notin \text{FIRST}(\text{Factor}) \Rightarrow \text{FIRST}^+(\text{Factor}) = \text{FIRST}(\text{Factor})$
- $\text{FIRST}(\text{Factor} \rightarrow \text{name}) = \{\text{name}\}$
- $\text{FIRST}(\text{Factor} \rightarrow \text{name} [\text{ArgList}]) = \{\text{name}\}$
- $\text{FIRST}(\text{Factor} \rightarrow \text{name} (\text{ArgList})) = \{\text{name}\}$

∴ The above grammar is **not** backtrack free.

⇒ To expand an AST node of *Factor*, with a **lookahead** of name, the parser has no basis to choose among rules 11, 12, and 13.

52 of 96

## Backtrack-Free Grammar: Left-Factoring

- A CFG is not backtrack free if there exists a **common prefix** (name) among the RHS of **multiple** production rules.
- To make such a CFG **backtrack-free**, we may transform it using **left factoring**: a process of extracting and isolating **common prefixes** in a set of production rules.
  - Identify a common prefix  $\alpha$ :
 
$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_j$$

[ each of  $\gamma_1, \gamma_2, \dots, \gamma_j$  does not begin with  $\alpha$  ]
  - Rewrite that production rule as:
 
$$\begin{aligned} A &\rightarrow \alpha B \mid \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_j \\ B &\rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n \end{aligned}$$
  - New rule  $B \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$  may also contain **common prefixes**.
  - Rewriting **continues** until no common prefixes are identified.

53 of 96

## TDP: Terminating and Backtrack-Free

- Given an arbitrary CFG as input to a **top-down parser**:
  - Q. How do we avoid a **non-terminating** parsing process?
    - A. Convert left-recursions to right-recursion.
  - Q. How do we minimize the need of **backtracking**?
    - A. left-factoring & one-symbol lookahead using **FIRST<sup>+</sup>**
- Not** every context-free **language** has a corresponding **backtrack-free** context-free **grammar**.
  - Given a CFL  $L$ , the following is **undecidable**:
 
$$\exists \text{cfg} \mid L(\text{cfg}) = L \wedge \text{isBacktrackFree}(\text{cfg})$$
- Given a CFG  $g = (V, \Sigma, R, S)$ , whether or not  $g$  is **backtrack-free** is **decidable**:
  - For each  $A \rightarrow \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_n \in R$ :

$$\forall i, j: 1 \leq i, j \leq n \wedge i \neq j \bullet \text{FIRST}^+(\gamma_i) \cap \text{FIRST}^+(\gamma_j) = \emptyset$$

55 of 96

## Left-Factoring: Exercise

- Use **left-factoring** to remove all **common prefixes** from the following grammar.

11	<i>Factor</i>	$\rightarrow$	name
12			name [ <i>ArgList</i> ]
13			name ( <i>ArgList</i> )
15	<i>ArgList</i>	$\rightarrow$	<i>Expr</i> <i>MoreArgs</i>
16	<i>MoreArgs</i>	$\rightarrow$	, <i>Expr</i> <i>MoreArgs</i>
17			$\epsilon$

- Identify common prefix **name** and **rewrite** rules 11, 12, and 13:

<i>Factor</i>	$\rightarrow$	name	<i>Arguments</i>
<i>Arguments</i>	$\rightarrow$	[ <i>ArgList</i> ]	
			( <i>ArgList</i> )
			$\epsilon$

Any more **common prefixes**?

[ No ]

54 of 96

## Backtrack-Free Parsing (2.1)

- A **recursive-descent** parser is:
  - A top-down parser
  - Structured as a set of **mutually recursive** procedures
    - Each procedure corresponds to a **non-terminal** in the grammar.
    - See an **example**.
- Given a **backtrack-free** grammar, a tool (a.k.a. **parser generator**) can automatically generate:
  - FIRST**, **FOLLOW**, and **FIRST<sup>+</sup>** sets
  - An efficient **recursive-descent** parser
    - This generated parser is called an **LL(1) parser**, which:
      - Processes input from **L**eft to **r**ight
      - Constructs a **L**eftmost derivation
      - Uses a lookahead of **1** symbol
  - LL(1) grammars** are those working in an **LL(1)** scheme.
    - LL(1) grammars** are **backtrack-free** by definition.

56 of 96

## Backtrack-Free Parsing (2.2)

- Consider this CFG with  $FIRST^+$  sets of the RHSs:

	Production	$FIRST^+$
2	$Expr' \rightarrow + Term Expr'$	{+}
3	$  - Term Expr'$	{-}
4	$  \epsilon$	{ $\epsilon$ , eof, $\_$ }

- The corresponding *recursive-descent* parser is structured as:

```

ExprPrim()
if word = + ∨ word = - then /* Rules 2, 3 */
  word := NextWord()
  if (Term())
    then return ExprPrim()
    else return false
elseif word = ) ∨ word = eof then /* Rule 4 */
  return true
else
  report a syntax error
  return false
end

Term()
...

```

See: [parser generator](#)

57 of 96

## BUP: Discovering Rightmost Derivation

- In TDP, we build the start variable as the *root node*, and then work towards the *leaves*. [ **leftmost** derivation ]
  - In Bottom-Up Parsing (BUP):
    - Words (terminals) are still returned from **left to right** by the scanner.
    - As terminals, or a mix of terminals and variables, are identified as **reducible** to some variable  $A$  (i.e., matching the RHS of some production rule for  $A$ ), then a layer is added.
    - Eventually:
      - accept:** The *start variable* is reduced and **all** words have been consumed.
      - reject:** The next word is not eof, but no further **reduction** can be identified.
- Q.** Why can BUP find the *rightmost* derivation (RMD), if any?  
**A.** BUP discovers steps in a **RMD** in its **reverse** order.

59 of 96

## LL(1) Parser: Exercise

Consider the following grammar:

$L \rightarrow R a$	$R \rightarrow aba$	$Q \rightarrow bbc$
$  Q ba$	$  caba$	$  bc$
	$  R bc$	

**Q.** Is it suitable for a *top-down predictive* parser?

- If so, show that it satisfies the **LL(1)** condition.
- If not, identify the problem(s) and correct it (them). Also show that the revised grammar satisfies the **LL(1)** condition.

58 of 96

## BUP: Discovering Rightmost Derivation (1)

- table-driven LR(1)** parser: an implementation for BUP, which
  - Processes input from **Left** to right
  - Constructs a **Rightmost** derivation
  - Uses a lookahead of **1** symbol
- A language has the **LR(1)** property if it:
  - Can be parsed in a single **Left** to right scan,
  - To build a **reversed Rightmost** derivation,
  - Using a lookahead of **1** symbol to determine parsing actions.
- Critical step in a bottom-up parser is to find the **next handle**.

60 of 96

## BUP: Discovering Rightmost Derivation (2)



```

ALGORITHM: BUParse
INPUT: CFG G = (V, Σ, R, S), Action & Goto Tables
OUTPUT: Report Parse Success or Syntax Error
PROCEDURE:
  initialize an empty stack trace
  trace.push(S) /* start state */
  word := NextWord()
  while (true)
    state := trace.top()
    act := Action[state, word]
    if act = 'accept' then
      succeed()
    elseif act = 'reduce A → β' then
      trace.pop() 2 × |β| times /* word + state */
      state := trace.top()
      trace.push(A)
      next := Goto[state, A]
      trace.push(next)
    elseif act = 'shift si' then
      trace.push(word)
      trace.push(si)
      word := NextWord()
    else
      fail()
  
```

61 of 96

## BUP: Example Tracing (2.1)



Consider the steps of performing BUP on input  $()$ :

Iteration	State	word	Stack	Handle	Action
initial	—	(	\$ 0	— none —	—
1	0	(	\$ 0	— none —	shift 3
2	3	)	\$ 0 ( 3	— none —	shift 7
3	7	eof	\$ 0 ( 3 ) 7	( )	reduce 5
4	2	eof	\$ 0 Pair 2	Pair	reduce 3
5	1	eof	\$ 0 List 1	List	accept

63 of 96

## BUP: Example Tracing (1)



- Consider the following grammar for parentheses:

```

1 Goal → List
2 List → List Pair
3     | Pair
4 Pair → ( Pair )
5     | ( )
  
```

- Assume: tables **Action** and **Goto** constructed accordingly:

State	Action Table		Goto Table	
	eof	(	List	Pair
0		s 3	1	2
1	acc	s 3		4
2	r 3	r 3		
3		s 6	s 7	5
4	r 2	r 2		
5			s 8	
6		s 6	s 10	9
7	r 5	r 5		
8	r 4	r 4		
9			s 11	
10		r 5		
11		r 4		

In **Action** table:

- $s_j$ : shift to state  $i$
- $r_j$ : reduce to the LHS of production # $j$

62 of 96

## BUP: Example Tracing (2.2)



Consider the steps of performing BUP on input  $((())())$ :

Iteration	State	word	Stack	Handle	Action
initial	—	(	\$ 0	— none —	—
1	0	(	\$ 0	— none —	shift 3
2	3	(	\$ 0 ( 3	— none —	shift 6
3	6	)	\$ 0 ( 3 ( 6	— none —	shift 10
4	10	)	\$ 0 ( 3 ( 6 ) 10	( )	reduce 5
5	5	)	\$ 0 ( 3 Pair 5	— none —	shift 8
6	8	(	\$ 0 ( 3 Pair 5 ) 8	( Pair )	reduce 4
7	2	(	\$ 0 Pair 2	Pair	reduce 3
8	1	(	\$ 0 List 1	— none —	shift 3
9	3	)	\$ 0 List 1 ( 3	— none —	shift 7
10	7	eof	\$ 0 List 1 ( 3 ) 7	( )	reduce 5
11	4	eof	\$ 0 List 1 Pair 4	List Pair	reduce 2
12	1	eof	\$ 0 List 1	List	accept

64 of 96

## BUP: Example Tracing (2.3)



Consider the steps of performing BUP on input  $( ) )$ :

Iteration	State	word	Stack	Handle	Action
initial	—	(	\$ 0	— none —	—
1	0	(	\$ 0	— none —	shift 3
2	3	)	\$ 0 ( 3	— none —	shift 7
3	7	)	\$ 0 ( 3 ) 7	— none —	error

65 of 96

## LR(1) Items: Definition



- In **LR(1)** parsing, **Action** and **Goto** tables encode legitimate ways (w.r.t. a grammar) for finding **handles** (for **reductions**).
- In a **table-driven LR(1)** parser, the table-construction algorithm represents each potential **handle** (for a **reduction**) with an **LR(1) item** e.g.,

$$[A \rightarrow \beta \bullet \gamma, a]$$

where:

- A production rule  $A \rightarrow \beta\gamma$  is currently being applied.
- A placeholder,  $\bullet$ , indicates the position of the parser's **stack top**.
  - ✓ The parser's stack contains  $\beta$  ("left context").
  - ✓  $\gamma$  is yet to be matched.
- Remark.** Upon matching  $\beta\gamma$ , if  $a$  matches the current word, then we "replace"  $\beta\gamma$  (and their corresponding states) with  $A$  (and its corresponding state).
- A terminal symbol  $a$  serves as a **lookahead symbol**.

66 of 96

## LR(1) Items: Scenarios



An **LR(1) item** can be:

- POSSIBILITY**  $[A \rightarrow \bullet\beta\gamma, a]$ 
  - In the current parsing context, an  $A$  would be valid.
  - $\bullet$  represents the position of the parser's **stack top**.
  - Recognizing a  $\beta$  next would be one step towards discovering an  $A$ .
- PARTIALLY COMPLETION**  $[A \rightarrow \beta \bullet \gamma, a]$ 
  - The parser has progressed from  $[A \rightarrow \bullet\beta\gamma, a]$  by recognizing  $\beta$ .
  - Recognizing a  $\gamma$  next would be one step towards discovering an  $A$ .
- COMPLETION**  $[A \rightarrow \beta\gamma \bullet, a]$ 
  - Parser has progressed from  $[A \rightarrow \bullet\beta\gamma, a]$  by recognizing  $\beta\gamma$ .
  - $\beta\gamma$  found in a context where an  $A$  followed by  $a$  would be valid.
  - If the current input word matches  $a$ , then:
    - Current **complet item** is a **handle**.
    - Parser can **reduce**  $\beta\gamma$  to  $A$  (and replace  $\beta\gamma$  with  $A$  in its stack).

67 of 96

## LR(1) Items: Example (1.1)



Consider the following grammar for parentheses:

1	$Goal \rightarrow List$
2	$List \rightarrow List Pair$
3	$  Pair$
4	$Pair \rightarrow ( Pair )$
5	$  ( )$

**Initial State:**  $[Goal \rightarrow \bullet List, eof]$

**Desired Final State:**  $[Goal \rightarrow List \bullet, eof]$

**Intermediate States:** Subset Construction

Q. Derive all **LR(1) items** for the above grammar.

- FOLLOW(List)** =  $\{eof, (, )\}$     **FOLLOW(Pair)** =  $\{eof, (, )\}$
- For each production  $A \rightarrow \beta$ , given **FOLLOW(A)**, **LR(1) items** are:
 
$$\begin{aligned} & \{ [A \rightarrow \bullet\beta\gamma, a] \mid a \in \text{FOLLOW}(A) \} \\ & \cup \\ & \{ [A \rightarrow \beta \bullet \gamma, a] \mid a \in \text{FOLLOW}(A) \} \\ & \cup \\ & \{ [A \rightarrow \beta\gamma \bullet, a] \mid a \in \text{FOLLOW}(A) \} \end{aligned}$$

68 of 96

## LR(1) Items: Example (1.2)

Q. Given production  $A \rightarrow \beta$  (e.g.,  $Pair \rightarrow ( Pair )$ ), how many

**LR(1) items** can be generated?

- The current parsing progress (on matching the RHS) can be:
  - $\bullet( Pair )$
  - $( \bullet Pair )$
  - $( Pair \bullet )$
  - $( Pair ) \bullet$
- Lookahead symbol following  $Pair$ ?  $FOLLOW(Pair) = \{eof, (, )\}$
- All possible **LR(1) items** related to  $Pair \rightarrow ( Pair )$ ?
 

$\checkmark [ \bullet( Pair ), eof ]$	$\checkmark [ \bullet( Pair ), ( ]$	$\checkmark [ \bullet( Pair ), ) ]$
$\checkmark [ ( \bullet Pair ), eof ]$	$\checkmark [ ( \bullet Pair ), ( ]$	$\checkmark [ ( \bullet Pair ), ) ]$
$\checkmark [ ( Pair \bullet ), eof ]$	$\checkmark [ ( Pair \bullet ), ( ]$	$\checkmark [ ( Pair \bullet ), ) ]$
$\checkmark [ ( Pair ) \bullet, eof ]$	$\checkmark [ ( Pair ) \bullet, ( ]$	$\checkmark [ ( Pair ) \bullet, ) ]$

A. How many in general (in terms of  $A$  and  $\beta$ )?

$$\underbrace{|\beta| + 1}_{\text{possible positions of } \bullet} \times \underbrace{|FOLLOW(A)|}_{\text{possible lookahead symbols}}$$

possible positions of  $\bullet$  possible lookahead symbols

## LR(1) Items: Example (1.3)

A. There are 33 **LR(1) items** in the parentheses grammar.

$[Goal \rightarrow \bullet List, eof]$		
$[Goal \rightarrow List \bullet, eof]$		
$[List \rightarrow \bullet List Pair, eof]$	$[List \rightarrow \bullet List Pair, (]$	
$[List \rightarrow List \bullet Pair, eof]$	$[List \rightarrow List \bullet Pair, (]$	
$[List \rightarrow List Pair \bullet, eof]$	$[List \rightarrow List Pair \bullet, (]$	
$[List \rightarrow \bullet Pair, eof]$	$[List \rightarrow \bullet Pair, (]$	
$[List \rightarrow Pair \bullet, eof]$	$[List \rightarrow Pair \bullet, (]$	
$[Pair \rightarrow \bullet ( Pair ), eof]$	$[Pair \rightarrow \bullet ( Pair ), )]$	$[Pair \rightarrow \bullet ( Pair ), (]$
$[Pair \rightarrow ( \bullet Pair ), eof]$	$[Pair \rightarrow ( \bullet Pair ), )]$	$[Pair \rightarrow ( \bullet Pair ), (]$
$[Pair \rightarrow ( Pair \bullet ), eof]$	$[Pair \rightarrow ( Pair \bullet ), )]$	$[Pair \rightarrow ( Pair \bullet ), (]$
$[Pair \rightarrow ( Pair ) \bullet, eof]$	$[Pair \rightarrow ( Pair ) \bullet, )]$	$[Pair \rightarrow ( Pair ) \bullet, (]$
$[Pair \rightarrow \bullet ( ), eof]$	$[Pair \rightarrow \bullet ( ), (]$	$[Pair \rightarrow \bullet ( ), )]$
$[Pair \rightarrow ( \bullet ), eof]$	$[Pair \rightarrow ( \bullet ), (]$	$[Pair \rightarrow ( \bullet ), )]$
$[Pair \rightarrow ( ) \bullet, eof]$	$[Pair \rightarrow ( ) \bullet, (]$	$[Pair \rightarrow ( ) \bullet, )]$

## LR(1) Items: Example (2)

Consider the following grammar for expressions:

0	$Goal \rightarrow Expr$	6	$Term' \rightarrow \times Factor Term'$
1	$Expr \rightarrow Term Expr'$	7	$  \div Factor Term'$
2	$Expr' \rightarrow + Term Expr'$	8	$  \epsilon$
3	$  - Term Expr'$	9	$Factor \rightarrow ( Expr )$
4	$  \epsilon$	10	$  num$
5	$Term \rightarrow Factor Term'$	11	$  name$

Q. Derive all **LR(1) items** for the above grammar.

Hints. First compute **FOLLOW** for each non-terminal:

	<b>Expr</b>	<b>Expr'</b>	<b>Term</b>	<b>Term'</b>	<b>Factor</b>
FOLLOW	$eof, )$	$eof, )$	$eof, +, -, )$	$eof, +, -, )$	$eof, +, -, \times, \div, )$

Tips. Ignore  $\epsilon$  production such as  $Expr' \rightarrow \epsilon$

since the **FOLLOW** sets already take them into consideration.

## Canonical Collection (CC) vs. LR(1) items

1	$Goal \rightarrow List$
2	$List \rightarrow List Pair$
3	$  Pair$
4	$Pair \rightarrow ( Pair )$
5	$  ( )$

Recall:

**LR(1) Items:** 33 items

**Initial State:**  $[Goal \rightarrow \bullet List, eof]$

**Desired Final State:**  $[Goal \rightarrow List \bullet, eof]$

o The **canonical collection**

$$CC = \{CC_0, CC_1, CC_2, \dots, CC_n\}$$

denotes the set of **valid states** of a **LR(1) parser**.

• Each  $CC_i \in CC$  ( $0 \leq i \leq n$ ) is a set of **LR(1) items**.

•  $CC \subseteq \mathbb{P}(\text{LR(1) items})$   $|CC|?$   $[|CC| \leq 2^{|\text{LR(1) items}|}]$

o To model a **LR(1) parser**, we use techniques similar to how we construct a DFA from an NFA (subset construction and  $\epsilon$ -closure).

o **Analogies.**

$\checkmark$  **LR(1) items**  $\approx$  states of source **NFA**

$\checkmark$  **CC**  $\approx$  states of target **DFA**

## Constructing $CC$ : The *closure* Procedure (1)



```

1 ALGORITHM: closure
2 INPUT: CFG  $G = (V, \Sigma, R, S)$ , a set  $s$  of LR(1) items
3 OUTPUT: a set of LR(1) items
4 PROCEDURE:
5   lastS :=  $\emptyset$ 
6   while (lastS  $\neq$  s):
7     lastS := s
8     for  $[A \rightarrow \dots \bullet C \delta, a] \in s$ :
9       for  $C \rightarrow \gamma \in R$ :
10        for  $b \in \text{FIRST}(\delta a)$ :
11          s :=  $s \cup \{ [C \rightarrow \bullet \gamma, b] \}$ 
12   return s
    
```

- Line 8:  $[A \rightarrow \dots \bullet C \delta, a] \in s$  indicates that the parser's next task is to match  $C \delta$  with a lookahead symbol  $a$ .
- Line 9: Given: matching  $\gamma$  can reduce to  $C$
- Line 10: Given:  $b \in \text{FIRST}(\delta a)$  is a valid lookahead symbol after reducing  $\gamma$  to  $C$
- Line 11: Add a new item  $[C \rightarrow \bullet \gamma, b]$  into  $s$ .
- Line 6: Termination is guaranteed.  
 $\therefore$  Each iteration adds  $\geq 1$  item to  $s$  (otherwise  $\text{lastS} \neq s$  is *false*).

73 of 96

## Constructing $CC$ : The *goto* Procedure (1)



```

1 ALGORITHM: goto
2 INPUT: a set  $s$  of LR(1) items, a symbol  $x$ 
3 OUTPUT: a set of LR(1) items
4 PROCEDURE:
5   moved :=  $\emptyset$ 
6   for item  $\in s$ :
7     if item =  $[\alpha \rightarrow \beta \bullet x \delta, a]$  then
8       moved :=  $\text{moved} \cup \{ [\alpha \rightarrow \beta x \bullet \delta, a] \}$ 
9   end
10  return closure(moved)
    
```

- Line 7: Given: item  $[\alpha \rightarrow \beta \bullet x \delta, a]$  (where  $x$  is the next to match)
- Line 8: Add  $[\alpha \rightarrow \beta x \bullet \delta, a]$  (indicating  $x$  is matched) to *moved*
- Line 10: Calculate and return *closure(moved)* as the "next state" from  $s$  with a "transition"  $x$ .

75 of 96

## Constructing $CC$ : The *closure* Procedure (2.1)



```

1 Goal  $\rightarrow$  List
2 List  $\rightarrow$  List Pair
3   | Pair
4 Pair  $\rightarrow$  ( Pair )
5   | ( )
    
```

**Initial State:**  $[Goal \rightarrow \bullet List, eof]$

Calculate  $cc_0 = \text{closure}([Goal \rightarrow \bullet List, eof])$ .

74 of 96

## Constructing $CC$ : The *goto* Procedure (2)



```

1 Goal  $\rightarrow$  List
2 List  $\rightarrow$  List Pair
3   | Pair
4 Pair  $\rightarrow$  ( Pair )
5   | ( )
    
```

$$cc_0 = \left\{ \begin{array}{lll} [Goal \rightarrow \bullet List, eof] & [List \rightarrow \bullet List Pair, eof] & [List \rightarrow \bullet List Pair, \_] \\ [List \rightarrow \bullet Pair, eof] & [List \rightarrow \bullet Pair, \_] & [Pair \rightarrow \bullet ( Pair ), eof] \\ [Pair \rightarrow \bullet ( Pair ), \_] & [Pair \rightarrow \bullet ( ), eof] & [Pair \rightarrow \bullet ( ), \_] \end{array} \right\}$$

Calculate  $\text{goto}(cc_0, ($ .

["next state" from  $cc_0$  taking  $($

76 of 96

## Constructing $CC$ : The Algorithm (1)



```

1  ALGORITHM: BuildCC
2  INPUT: a grammar  $G = (V, \Sigma, R, S)$ , goal production  $S \rightarrow S'$ 
3  OUTPUT:
4  (1) a set  $CC = \{cc_0, cc_1, \dots, cc_n\}$  where  $cc_i \in G'$ 's LR(1) items
5  (2) a transition function
6  PROCEDURE:
7   $cc_0 := \text{closure}(\{[S' \rightarrow \bullet S, \text{eof}]\})$ 
8   $CC := \{cc_0\}$ 
9   $processed := \{cc_0\}$ 
10  $lastCC := \emptyset$ 
11 while ( $lastCC \neq CC$ ):
12    $lastCC := CC$ 
13   for  $cc_j$  s.t.  $cc_j \in CC \wedge cc_j \notin processed$ :
14      $processed := processed \cup \{cc_j\}$ 
15     for  $x$  s.t.  $[\dots \rightarrow \dots \bullet x \dots] \in cc_j$ :
16        $temp := \text{goto}(cc_j, x)$ 
17       if  $temp \notin CC$  then
18          $CC := CC \cup \{temp\}$ 
19       end
20    $\delta := \delta \cup (cc_j, x, temp)$ 

```

77 of 96

## Constructing $CC$ : The Algorithm (2.2)



Resulting transition table:

Iteration	Item	Goal	List	Pair	(	)	eof
0	$CC_0$	$\emptyset$	$CC_1$	$CC_2$	$CC_3$	$\emptyset$	$\emptyset$
1	$CC_1$	$\emptyset$	$\emptyset$	$CC_4$	$CC_3$	$\emptyset$	$\emptyset$
	$CC_2$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
	$CC_3$	$\emptyset$	$\emptyset$	$CC_5$	$CC_6$	$CC_7$	$\emptyset$
2	$CC_4$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
	$CC_5$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$CC_8$	$\emptyset$
	$CC_6$	$\emptyset$	$\emptyset$	$CC_9$	$CC_6$	$CC_{10}$	$\emptyset$
	$CC_7$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
3	$CC_8$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
	$CC_9$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$CC_{11}$	$\emptyset$
	$CC_{10}$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
4	$CC_{11}$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$

79 of 96

## Constructing $CC$ : The Algorithm (2.1)



```

1  Goal  $\rightarrow$  List
2  List  $\rightarrow$  List Pair
3      | Pair
4  Pair  $\rightarrow$  ( Pair )
5      | ( )

```

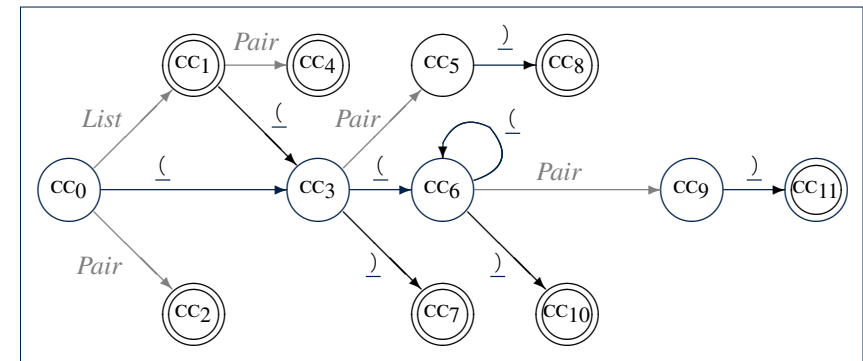
- Calculate  $CC = \{cc_0, cc_1, \dots, cc_{11}\}$
- Calculate the transition function  $\delta : CC \times \Sigma \rightarrow CC$

78 of 96

## Constructing $CC$ : The Algorithm (2.3)



Resulting DFA for the parser:



80 of 96



## Constructing $CC$ : The Algorithm (2.4.1)



Resulting canonical collection  $CC$ :

$$\begin{aligned}
 CC_0 &= \left\{ \begin{array}{l} [Goal \rightarrow \bullet List, eof] \quad [List \rightarrow \bullet List Pair, eof] \quad [List \rightarrow \bullet List Pair, \_] \\ [List \rightarrow \bullet Pair, eof] \quad [List \rightarrow \bullet Pair, \_] \quad [Pair \rightarrow \bullet \_ Pair \_, eof] \\ [Pair \rightarrow \bullet \_ Pair \_, \_] \quad [Pair \rightarrow \bullet \_ \_, eof] \quad [Pair \rightarrow \bullet \_ \_, \_] \end{array} \right\} & CC_1 &= \left\{ \begin{array}{l} [Goal \rightarrow List \bullet, eof] \quad [List \rightarrow List \bullet Pair, eof] \quad [List \rightarrow List \bullet Pair, \_] \\ [Pair \rightarrow \bullet \_ Pair \_, eof] \quad [Pair \rightarrow \bullet \_ Pair \_, \_] \quad [Pair \rightarrow \bullet \_ \_, eof] \\ [Pair \rightarrow \bullet \_ \_, \_] \end{array} \right\} \\
 CC_2 &= \left\{ [List \rightarrow Pair \bullet, eof] \quad [List \rightarrow Pair \bullet, \_] \right\} & CC_3 &= \left\{ \begin{array}{l} [Pair \rightarrow \bullet \_ Pair \_, \_] \quad [Pair \rightarrow \_ \bullet Pair \_, eof] \quad [Pair \rightarrow \_ \bullet Pair \_, \_] \\ [Pair \rightarrow \bullet \_ \_, \_] \quad [Pair \rightarrow \_ \bullet \_, eof] \quad [Pair \rightarrow \_ \bullet \_, \_] \end{array} \right\} \\
 CC_4 &= \left\{ [List \rightarrow List Pair \bullet, eof] \quad [List \rightarrow List Pair \bullet, \_] \right\} & CC_5 &= \left\{ [Pair \rightarrow \_ Pair \bullet \_, eof] \quad [Pair \rightarrow \_ Pair \bullet \_, \_] \right\} \\
 CC_6 &= \left\{ \begin{array}{l} [Pair \rightarrow \bullet \_ Pair \_, \_] \quad [Pair \rightarrow \_ \bullet Pair \_, \_] \\ [Pair \rightarrow \bullet \_ \_, \_] \quad [Pair \rightarrow \_ \bullet \_, \_] \end{array} \right\} & CC_7 &= \left\{ [Pair \rightarrow \_ \_ \bullet, eof] \quad [Pair \rightarrow \_ \_ \bullet, \_] \right\} \\
 CC_8 &= \left\{ [Pair \rightarrow \_ Pair \_ \bullet, eof] \quad [Pair \rightarrow \_ Pair \_ \bullet, \_] \right\} & CC_9 &= \left\{ [Pair \rightarrow \_ Pair \bullet \_, \_] \right\} \\
 CC_{10} &= \left\{ [Pair \rightarrow \_ \_ \bullet, \_] \right\} & CC_{11} &= \left\{ [Pair \rightarrow \_ Pair \_ \bullet, \_] \right\}
 \end{aligned}$$

81 of 96

## Constructing Action and Goto Tables (2)



Resulting Action and Goto tables:

State	Action Table			Goto Table	
	eof	(	)	List	Pair
0		s 3		1	2
1	acc	s 3			4
2	r 3	r 3			
3		s 6	s 7		5
4	r 2	r 2			
5			s 8		
6		s 6	s 10		9
7	r 5	r 5			
8	r 4	r 4			
9			s 11		
10			r 5		
11			r 4		

83 of 96

## Constructing Action and Goto Tables (1)



```

1  ALGORITHM: BuildActionGotoTables
2  INPUT:
3  (1) a grammar  $G = (V, \Sigma, R, S)$ 
4  (2) goal production  $S \rightarrow S'$ 
5  (3) a canonical collection  $CC = \{cc_0, cc_1, \dots, cc_n\}$ 
6  (4) a transition function  $\delta: CC \times \Sigma \rightarrow CC$ 
7  OUTPUT: Action Table & Goto Table
8  PROCEDURE:
9  for  $cc_i \in CC$ :
10 for item  $\in cc_i$ :
11 if item =  $[A \rightarrow \beta \bullet x \gamma, a]$  \pause  $\delta(cc_i, x) = cc_j$  then
12   Action[i, x] := shift j
13 elseif item =  $[A \rightarrow \beta \bullet, a]$  then
14   Action[i, a] := reduce  $A \rightarrow \beta$ 
15 elseif item =  $[S \rightarrow S' \bullet, eof]$  then
16   Action[i, eof] := accept
17 end
18 for  $v \in V$ :
19 if  $\delta(cc_i, v) = cc_j$  then
20   Goto[i, v] = j
21 end
    
```

- L12, 13: Next valid step in discovering  $A$  is to match terminal symbol  $x$ .
- L14, 15: Having recognized  $\beta$ , if current word matches lookahead  $a$ , reduce  $\beta$  to  $A$ .
- L16, 17: Accept if input exhausted and what's recognized reducible to start var.  $S$ .
- L20, 21: Record consequence of a reduction to non-terminal  $v$  from state  $i$

82 of 96

## BUP: Discovering Ambiguity (1)



1	Goal	→	Stmt
2	Stmt	→	if expr then Stmt
3			if expr then Stmt else Stmt
4			assign

- Calculate  $CC = \{cc_0, cc_1, \dots\}$
- Calculate the transition function  $\delta: CC \times \Sigma \rightarrow CC$

84 of 96

## BUP: Discovering Ambiguity (2.1)



Resulting transition table:

Item	Goal	Stmt	if	expr	then	else	assign	eof
0	CC <sub>0</sub>	∅	CC <sub>1</sub>	CC <sub>2</sub>	∅	∅	∅	CC <sub>3</sub>
1	CC <sub>1</sub>	∅	∅	∅	∅	∅	∅	∅
	CC <sub>2</sub>	∅	∅	∅	CC <sub>4</sub>	∅	∅	∅
	CC <sub>3</sub>	∅	∅	∅	∅	∅	∅	∅
2	CC <sub>4</sub>	∅	∅	∅	∅	CC <sub>5</sub>	∅	∅
3	CC <sub>5</sub>	∅	CC <sub>6</sub>	CC <sub>7</sub>	∅	∅	∅	CC <sub>8</sub>
4	CC <sub>6</sub>	∅	∅	∅	∅	∅	CC <sub>9</sub>	∅
	CC <sub>7</sub>	∅	∅	∅	∅	∅	∅	∅
	CC <sub>8</sub>	∅	∅	∅	∅	∅	∅	∅
5	CC <sub>9</sub>	∅	CC <sub>11</sub>	CC <sub>2</sub>	∅	∅	∅	CC <sub>3</sub>
	CC <sub>10</sub>	∅	∅	∅	∅	CC <sub>12</sub>	∅	∅
6	CC <sub>11</sub>	∅	∅	∅	∅	∅	∅	∅
	CC <sub>12</sub>	∅	CC <sub>13</sub>	CC <sub>7</sub>	∅	∅	∅	CC <sub>8</sub>
7	CC <sub>13</sub>	∅	∅	∅	∅	∅	∅	∅
8	CC <sub>14</sub>	∅	CC <sub>15</sub>	CC <sub>7</sub>	∅	∅	∅	CC <sub>8</sub>
9	CC <sub>15</sub>	∅	∅	∅	∅	∅	∅	∅

85 of 96

## BUP: Discovering Ambiguity (2.2.2)



Resulting canonical collection  $CC$ :

$$CC_8 = \{ [Stmt \rightarrow assign \bullet, \{eof, else\}] \}$$

$$CC_9 = \left\{ \begin{array}{l} [Stmt \rightarrow if \ expr \ then \ Stmt \ else \ \bullet \ Stmt, eof], \\ [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt, eof], \\ [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt \ else \ Stmt, eof], \\ [Stmt \rightarrow \bullet \ assign, eof] \end{array} \right\}$$

$$CC_{10} = \left\{ \begin{array}{l} [Stmt \rightarrow if \ expr \ \bullet \ then \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow if \ expr \ \bullet \ then \ Stmt \ else \ Stmt, \{eof, else\}] \end{array} \right\}$$

$$CC_{11} = \{ [Stmt \rightarrow if \ expr \ then \ Stmt \ else \ Stmt \ \bullet, eof] \}$$

$$CC_{12} = \left\{ \begin{array}{l} [Stmt \rightarrow if \ expr \ then \ \bullet \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow if \ expr \ then \ \bullet \ Stmt \ else \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt \ else \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow \bullet \ assign, \{eof, else\}] \end{array} \right\}$$

$$CC_{13} = \left\{ \begin{array}{l} [Stmt \rightarrow if \ expr \ then \ Stmt \ \bullet, \{eof, else\}], \\ [Stmt \rightarrow if \ expr \ then \ Stmt \ \bullet \ else \ Stmt, \{eof, else\}] \end{array} \right\}$$

$$CC_{14} = \left\{ \begin{array}{l} [Stmt \rightarrow if \ expr \ then \ Stmt \ else \ \bullet \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt \ else \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow \bullet \ assign, \{eof, else\}] \end{array} \right\}$$

87 of 96

## BUP: Discovering Ambiguity (2.2.1)



Resulting canonical collection  $CC$ :

$$CC_0 = \left\{ \begin{array}{l} [Goal \rightarrow \bullet \ Stmt, eof] \\ [Stmt \rightarrow \bullet \ assign, eof] \end{array} \right\} \quad \left\{ \begin{array}{l} [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt, eof] \\ [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt \ else \ Stmt, eof] \end{array} \right\}$$

$$CC_1 = \{ [Goal \rightarrow Stmt \ \bullet, eof] \}$$

$$CC_2 = \left\{ \begin{array}{l} [Stmt \rightarrow if \ \bullet \ expr \ then \ Stmt, eof], \\ [Stmt \rightarrow if \ \bullet \ expr \ then \ Stmt \ else \ Stmt, eof] \end{array} \right\}$$

$$CC_3 = \{ [Stmt \rightarrow assign \ \bullet, eof] \}$$

$$CC_4 = \left\{ \begin{array}{l} [Stmt \rightarrow if \ expr \ \bullet \ then \ Stmt, eof], \\ [Stmt \rightarrow if \ expr \ \bullet \ then \ Stmt \ else \ Stmt, eof] \end{array} \right\}$$

$$CC_5 = \left\{ \begin{array}{l} [Stmt \rightarrow if \ expr \ then \ \bullet \ Stmt, eof], \\ [Stmt \rightarrow if \ expr \ then \ \bullet \ Stmt \ else \ Stmt, eof], \\ [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow \bullet \ assign, \{eof, else\}], \\ [Stmt \rightarrow \bullet \ if \ expr \ then \ Stmt \ else \ Stmt, \{eof, else\}] \end{array} \right\}$$

$$CC_6 = \left\{ \begin{array}{l} [Stmt \rightarrow if \ expr \ then \ Stmt \ \bullet, eof], \\ [Stmt \rightarrow if \ expr \ then \ Stmt \ \bullet \ else \ Stmt, eof] \end{array} \right\}$$

$$CC_7 = \left\{ \begin{array}{l} [Stmt \rightarrow if \ \bullet \ expr \ then \ Stmt, \{eof, else\}], \\ [Stmt \rightarrow if \ \bullet \ expr \ then \ Stmt \ else \ Stmt, \{eof, else\}] \end{array} \right\}$$

86 of 96

## BUP: Discovering Ambiguity (3)



- Consider  $cc_{13}$

$$cc_{13} = \left\{ \begin{array}{l} [Stmt \rightarrow if \ expr \ then \ Stmt \ \bullet, \{eof, else\}], \\ [Stmt \rightarrow if \ expr \ then \ Stmt \ \bullet \ else \ Stmt, \{eof, else\}] \end{array} \right\}$$

**Q.** What does it mean if the current word to consume is `else`?  
**A.** We can either **shift** (then expecting to match another *Stmt*) or reduce to a *Stmt*.

A single *Action* table entry cannot hold these two alternatives.

This is known as the **shift-reduce conflict**.

- Consider another scenario, say:

$$\begin{array}{l} [A \rightarrow \gamma \delta \bullet, a] \\ [B \rightarrow \gamma \delta \bullet, a] \end{array}$$

**Q.** What does it mean if the current word to consume is `a`?

**A.** We can either **reduce** to *A* or **reduce** to *B*.

A single *Action* table entry cannot hold these two alternatives.

This is known as the **reduce-reduce conflict**.

88 of 96

## Index (1)



[Parser in Context](#)

[Context-Free Languages: Introduction](#)

[CFG: Example \(1.1\)](#)

[CFG: Example \(1.2\)](#)

[CFG: Example \(1.2\)](#)

[CFG: Example \(2\)](#)

[CFG: Example \(3\)](#)

[CFG: Example \(4\)](#)

[CFG: Example \(5.1\) Version 1](#)

[CFG: Example \(5.2\) Version 1](#)

[CFG: Example \(5.3\) Version 1](#)

89 of 96

## Index (2)



[CFG: Example \(5.4\) Version 1](#)

[CFG: Example \(5.5\) Version 2](#)

[CFG: Example \(5.6\) Version 2](#)

[CFG: Example \(5.7\) Version 2](#)

[CFG: Formal Definition \(1\)](#)

[CFG: Formal Definition \(2\): Example](#)

[CFG: Formal Definition \(3\): Example](#)

[Regular Expressions to CFG's](#)

[DFA to CFG's](#)

[CFG: Leftmost Derivations \(1\)](#)

[CFG: Rightmost Derivations \(1\)](#)

90 of 96

## Index (3)



[CFG: Leftmost Derivations \(2\)](#)

[CFG: Rightmost Derivations \(2\)](#)

[CFG: Parse Trees vs. Derivations \(1\)](#)

[CFG: Parse Trees vs. Derivations \(2\)](#)

[CFG: Ambiguity: Definition](#)

[CFG: Ambiguity: Exercise \(1\)](#)

[CFG: Ambiguity: Exercise \(2.1\)](#)

[CFG: Ambiguity: Exercise \(2.2\)](#)

[CFG: Ambiguity: Exercise \(2.3\)](#)

[Discovering Derivations](#)

[TDP: Discovering Leftmost Derivation](#)

91 of 96

## Index (4)



[TDP: Exercise \(1\)](#)

[TDP: Exercise \(2\)](#)

[Left-Recursions \(LF\): Direct vs. Indirect](#)

[TDP: \(Preventively\) Eliminating LRs](#)

[CFG: Eliminating  \$\epsilon\$ -Productions \(1\)](#)

[CFG: Eliminating  \$\epsilon\$ -Productions \(2\)](#)

[Backtrack-Free Parsing \(1\)](#)

[The first Set: Definition](#)

[The first Set: Examples](#)

[Computing the first Set](#)

[Computing the first Set: Extension](#)

92 of 96

## Index (5)



[Extended first Set: Examples](#)

[Is the first Set Sufficient?](#)

[The follow Set: Examples](#)

[Computing the follow Set](#)

[Backtrack-Free Grammar](#)

[TDP: Lookahead with One Symbol](#)

[Backtrack-Free Grammar: Exercise](#)

[Backtrack-Free Grammar: Left-Factoring](#)

[Left-Factoring: Exercise](#)

[TDP: Terminating and Backtrack-Free](#)

[Backtrack-Free Parsing \(2.1\)](#)

93 of 96

## Index (6)



[Backtrack-Free Parsing \(2.2\)](#)

[LL\(1\) Parser: Exercise](#)

[BUP: Discovering Rightmost Derivation](#)

[BUP: Discovering Rightmost Derivation \(1\)](#)

[BUP: Discovering Rightmost Derivation \(2\)](#)

[BUP: Example Tracing \(1\)](#)

[BUP: Example Tracing \(2.1\)](#)

[BUP: Example Tracing \(2.2\)](#)

[BUP: Example Tracing \(2.3\)](#)

[LR\(1\) Items: Definition](#)

[LR\(1\) Items: Scenarios](#)

94 of 96

## Index (7)



[LR\(1\) Items: Example \(1.1\)](#)

[LR\(1\) Items: Example \(1.2\)](#)

[LR\(1\) Items: Example \(1.3\)](#)

[LR\(1\) Items: Example \(2\)](#)

[Canonical Collection \( \$\mathcal{CC}\$ \) vs. LR\(1\) items](#)

[Constructing  \$\mathcal{CC}\$ : The \*closure\* Procedure \(1\)](#)

[Constructing  \$\mathcal{CC}\$ : The \*closure\* Procedure \(2.1\)](#)

[Constructing  \$\mathcal{CC}\$ : The \*goto\* Procedure \(1\)](#)

[Constructing  \$\mathcal{CC}\$ : The \*goto\* Procedure \(2\)](#)

[Constructing  \$\mathcal{CC}\$ : The Algorithm \(1\)](#)

[Constructing  \$\mathcal{CC}\$ : The Algorithm \(2.1\)](#)

95 of 96

## Index (8)



[Constructing  \$\mathcal{CC}\$ : The Algorithm \(2.2\)](#)

[Constructing  \$\mathcal{CC}\$ : The Algorithm \(2.3\)](#)

[Constructing  \$\mathcal{CC}\$ : The Algorithm \(2.4\)](#)

[Constructing \*Action\* and \*Goto\* Tables \(1\)](#)

[Constructing \*Action\* and \*Goto\* Tables \(2\)](#)

[BUP: Discovering Ambiguity \(1\)](#)

[BUP: Discovering Ambiguity \(2.1\)](#)

[BUP: Discovering Ambiguity \(2.2.1\)](#)

[BUP: Discovering Ambiguity \(2.2.2\)](#)

[BUP: Discovering Ambiguity \(3\)](#)

96 of 96