# Parser: Syntactic Analysis

**Readings: EAC2 Chapter 3**
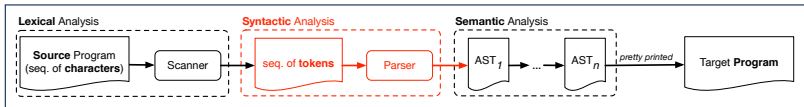
EECS4302 A:
Compilers and Interpreters
Fall 2022

CHEN-WEI WANG

# Parser in Context

○ Recall:



○ Treats the input programas as a ***a sequence of <u>classified</u> tokens/words***

○ Applies rules ==***parsing***== token sequences as

   ***abstract syntax trees (ASTs)***               [ ==***syntactic***== analysis ]

○ Upon termination:
   • Reports token sequences <u>not</u> derivable as ASTs
   • Produces an ***AST***

○ No longer considers ***every character*** in input program.

○ ==*Derivable*== token sequences constitute a

   ==*context-free language (CFL)*== .

# Context-Free Languages: Introduction

- We have seen **regular languages**:
  - Can be described using **finite automata** or **regular expressions**.
  - Satisfy the **pumping lemma**.
- Language with **recursive** structures 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.

- The following language that is **_non-regular_**

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

  can be described using a **_context-free grammar (CFG)_**:

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

- 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*).

## CFG: Example (1.2)

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

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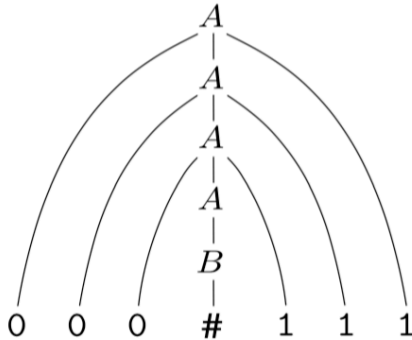$$

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

Given a CFG, a string's *derivation* can be shown as a *parse tree*.

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

## 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.*

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

Design a CFG for the following language:

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

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

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

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 $S$.

$$
\begin{array}{rcll}
S & \to & \epsilon & \{\textit{BC of } S\} \\
S & \to & AS & \{\textit{RC of } S\} \\
A & \to & \epsilon & \{\textit{BC of } A\} \\
A & \to & 01 & \{\textit{BC of } A\} \\
A & \to & 0A1 & \{\textit{RC of } A\text{: equal 0's and 1's}\} \\
A & \to & A1 & \{\textit{RC of } A\text{: more 1's}\}
\end{array}
$$

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

- Arithmetic operations: $+$, $-$, $\star$, $/$
- Relational operations: $>$, $<$, $>=$, $<=$, $==$, $/=$
- Logical operations: true, false, !, &&, ||, =>
  Start with the variable *Expression*.
- There are two possible versions:
  1. All operations are <u>mixed</u> together.
  2. Relevant operations are <u>grouped</u> together.
     Try both!

| *Expression* | → | *IntegerConstant* |
| | \| | −*IntegerConstant* |
| | \| | *BooleanConstant* |
| | \| | *BinaryOp* |
| | \| | *UnaryOp* |
| | \| | 〈 *Expression* 〉 |
| | | |
| *IntegerConstant* | → | *Digit* |
| | \| | *Digit IntegerConstant* |
| | | |
| *Digit* | → | 0 \| 1 \| 2 \| 3 \| 4 \| 5 \| 6 \| 7 \| 8 \| 9 |
| | | |
| *BooleanConstant* | → | TRUE |
| | \| | FALSE |

$$
\begin{aligned}
\textit{BinaryOp} \quad \rightarrow \quad & \textit{Expression} + \textit{Expression} \\
| \quad & \textit{Expression} - \textit{Expression} \\
| \quad & \textit{Expression} \star \textit{Expression} \\
| \quad & \textit{Expression} / \textit{Expression} \\
| \quad & \textit{Expression} \,\&\&\, \textit{Expression} \\
| \quad & \textit{Expression} \,||\, \textit{Expression} \\
| \quad & \textit{Expression} => \textit{Expression} \\
| \quad & \textit{Expression} == \textit{Expression} \\
| \quad & \textit{Expression} /= \textit{Expression} \\
| \quad & \textit{Expression} > \textit{Expression} \\
| \quad & \textit{Expression} < \textit{Expression}
\end{aligned}
$$

$$
\textit{UnaryOp} \quad \rightarrow \quad ! \; \textit{Expression}
$$

However, Version 1 of CFG:

- **Parses** string that requires further **semantic analysis** (e.g., type checking):
  e.g., 3 => 6
- Is **ambiguous**, meaning?
  - Some string may have <u>more than one</u> ways to interpreting it.
  - An interpretation is either visualized as a **parse tree**, or written as a sequence of **derivations**.

  e.g., Draw the parse tree(s) for 3 * 5 + 4

| *Expression* | → | *ArithmeticOp* |
| | | *RelationalOp* |
| | | *LogicalOp* |
| | | 〈 *Expression* 〉 |

| *IntegerConstant* | → | *Digit* |
| | | *Digit IntegerConstant* |

| *Digit* | → | 0 \| 1 \| 2 \| 3 \| 4 \| 5 \| 6 \| 7 \| 8 \| 9 |

| *BooleanConstant* | → | TRUE |
| | | FALSE |

| *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* |

However, Version 2 of CFG:

- Eliminates some cases for further semantic analysis:
  e.g., `(1 + 2) => (5 / 4)`                [ no parse tree ]
- Still *parses* strings that might require further *semantic analysis*:
  e.g., `(1 + 2) / (5 - (2 + 3))`
- Still is *ambiguous*.
  e.g., Draw the parse tree(s) for `3 * 5 + 4`

# 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*.           [$V \cap \Sigma = \varnothing$]
  - $R$ is a finite set of *rules* s.t.

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

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

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

- Design the **CFG** for strings of properly-nested parentheses.

  e.g., (), () (), ((( () ) ) (), *etc.*

  Present your answer in a **formal** manner.

- $G = (\{S\}, \{(,)\}, R, S)$, where $R$ is

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

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

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

$$
\begin{array}{rcl}
\textit{Expr} & \rightarrow & \textit{Expr} + \textit{Term} \\
& | & \textit{Term} \\
\textit{Term} & \rightarrow & \textit{Term} * \textit{Factor} \\
& | & \textit{Factor} \\
\textit{Factor} & \rightarrow & (\textit{Expr}) \\
& | & \texttt{a}
\end{array}
$$

  - $V = \{\textit{Expr}, \textit{Term}, \textit{Factor}\}$
  - $\Sigma = \{\texttt{a}, +, *, (, )\}$
  - $S = \textit{Expr}$
- **Precedence** of operators $+$, $*$ 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**.

## Regular Expressions to CFG's

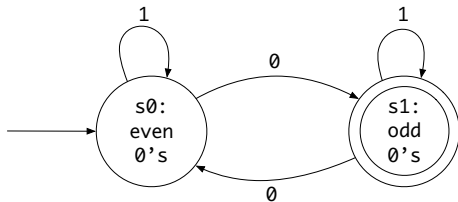- Recall the semantics of regular expressions (assuming that we do not consider $\varnothing$):

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

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

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

# 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 $R_i \to aR_j$ to the grammar if $\delta(q_i, a) = q_j$.
  - Add a rule $R_i \to \epsilon$ if $q_i \in F$.
- e.g., Grammar for



$$R_0 \quad \to \quad 1R_0 \mid 0R_1$$
$$R_1 \quad \to \quad 0R_0 \mid 1R_1 \mid \epsilon$$

# CFG: Leftmost Derivations (1)

$$
\begin{array}{rcl}
Expr & \rightarrow & Expr\ +\ Term \mid Term \\
Term & \rightarrow & Term\ *\ Factor \mid Factor \\
Factor & \rightarrow & (\,Expr\,) \mid a
\end{array}
$$

○ Given a string ($\in (V \cup \Sigma)^*$), a **left-most derivation (LMD)** keeps substituting the <u>leftmost</u> non-terminal ($\in V$).

○ **Unique LMD** for the string `a + a * a`:

$$
\begin{array}{rcl}
Expr & \Rightarrow & Expr\ +\ Term \\
& \Rightarrow & Term\ +\ Term \\
& \Rightarrow & Factor\ +\ Term \\
& \Rightarrow & a\ +\ Term \\
& \Rightarrow & a\ +\ Term\ *\ Factor \\
& \Rightarrow & a\ +\ Factor\ *\ Factor \\
& \Rightarrow & a\ +\ a\ *\ Factor \\
& \Rightarrow & a\ +\ a\ *\ a
\end{array}
$$

○ This **LMD** suggests that `a * a` is the right operand of `+`.

22 of 96

| | | |
|---|---|---|
| *Expr* | → | *Expr* + *Term* \| *Term* |
| *Term* | → | *Term* * *Factor* \| *Factor* |
| *Factor* | → | ( *Expr* ) \| *a* |

○ Given a string ($\in (V \cup \Sigma)^*$), a *right-most derivation (RMD)* keeps substituting the <u>rightmost</u> non-terminal ($\in V$).

○ *Unique RMD* for the string a + a * a:

| *Expr* | ⇒ | *Expr* + *Term* |
|---|---|---|
| | ⇒ | *Expr* + *Term* * *Factor* |
| | ⇒ | *Expr* + *Term* * *a* |
| | ⇒ | *Expr* + *Factor* * *a* |
| | ⇒ | *Expr* + *a* * *a* |
| | ⇒ | *Term* + *a* * *a* |
| | ⇒ | *Factor* + *a* * *a* |
| | ⇒ | *a* + *a* * *a* |

○ This *RMD* suggests that a * a is the right operand of +.

| | | |
|---|---|---|
| *Expr* | → | *Expr* + *Term* \| *Term* |
| *Term* | → | *Term* * *Factor* \| *Factor* |
| *Factor* | → | ( *Expr* ) \| *a* |

○ **Unique LMD** for the string (a + a) * a:

| | | |
|---|---|---|
| *Expr* | ⇒ | *Term* |
| | ⇒ | *Term* * *Factor* |
| | ⇒ | *Factor* * *Factor* |
| | ⇒ | ( *Expr* ) * *Factor* |
| | ⇒ | ( *Expr* + *Term* ) * *Factor* |
| | ⇒ | ( *Term* + *Term* ) * *Factor* |
| | ⇒ | ( *Factor* + *Term* ) * *Factor* |
| | ⇒ | ( *a* + *Term* ) * *Factor* |
| | ⇒ | ( *a* + *Factor* ) * *Factor* |
| | ⇒ | ( *a* + *a* ) * *Factor* |
| | ⇒ | ( *a* + *a* ) * *a* |

○ This **LMD** suggests that (a + a) is the left operand of *.

| | | |
|---|---|---|
| *Expr* | → | *Expr* + *Term* \| *Term* |
| *Term* | → | *Term* * *Factor* \| *Factor* |
| *Factor* | → | ( *Expr* ) \| *a* |

○ **Unique RMD** for the string (a + a) * a:

| | | |
|---|---|---|
| *Expr* | ⇒ | *Term* |
| | ⇒ | *Term* * *Factor* |
| | ⇒ | *Term* * *a* |
| | ⇒ | *Factor* * *a* |
| | ⇒ | ( *Expr* ) * *a* |
| | ⇒ | ( *Expr* + *Term* ) * *a* |
| | ⇒ | ( *Expr* + *Factor* ) * *a* |
| | ⇒ | ( *Expr* + *a* ) * *a* |
| | ⇒ | ( *Term* + *a* ) * *a* |
| | ⇒ | ( *Factor* + *a* ) * *a* |
| | ⇒ | ( *a* + *a* ) * *a* |

○ This **RMD** suggests that (a + a) is the left operand of *.

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

| a + a * a | (a + a) * a |
|---|---|



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

- A string $w \in \Sigma^*$ may have <u>more than one</u> *derivations*.

  **Q**: distinct *derivations* for $w \in \Sigma^* \Rightarrow$ distinct *parse trees* for $w$?

  **A**: Not in general $\because$ Derivations with *distinct orders* of variable substitutions may still result in the *same parse tree*.

- For example:

$$
\begin{array}{rcl}
\textit{Expr} & \rightarrow & \textit{Expr} + \textit{Term} \,|\, \textit{Term} \\
\textit{Term} & \rightarrow & \textit{Term} * \textit{Factor} \,|\, \textit{Factor} \\
\textit{Factor} & \rightarrow & (\textit{Expr}) \,|\, a
\end{array}
$$

  For string $a + a * a$, the *LMD* and *RMD* have *distinct orders* of variable substitutions, but their corresponding *parse trees are the <u>same</u>*.

{}

# 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 LMDs*** or, equally,
two or more ***distinct RMDs***.

  We require that all such derivations are completed by following a
  <u>consisten</u> order (**leftmost** or **rightmost**) to avoid ***false positive***.

○ $G$ is  *ambiguous*  if it generates some string ambiguously.

{}

# CFG: Ambiguity: Exercise (1)

- Is the following grammar *ambiguous* ?

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

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



- *Distinct ASTs* (for the *same input*) imply *distinct semantic interpretations*: e.g., a pre-order traversal for evaluation
- **Exercise**: Show *LMDs* for the two parse trees.

- Is the following grammar *ambiguous* ?

$$Statement \quad \rightarrow \quad \texttt{if } Expr \texttt{ then } Statement$$
$$| \quad \texttt{if } Expr \texttt{ then } Statement \texttt{ else } Statement$$
$$| \quad Assignment$$
$$\ldots$$

- Yes ∵ it derives the following string *ambiguously* :

$\texttt{if } Expr_1 \texttt{ then if } Expr_2 \texttt{ then } Assignment_1 \texttt{ else } Assignment_2$



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

(*Meaning 1*) *Assignment$_2$* may be associated with the <u>inner</u> if:



(*Meaning 2*) *Assignment$_2$* may be associated with the <u>outer</u> if:

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

| | | |
|---|---|---|
| *Statement* | → | `if` *Expr* `then` *Statement* |
| | \| | `if` *Expr* `then` *WithElse* `else` *Statement* |
| | \| | *Assignment* |
| *WithElse* | → | `if` *Expr* `then` *WithElse* `else` *WithElse* |
| | \| | *Assignment* |

- 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.

# 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.t. $S \overset{*}{\Rightarrow} p$.
  - A *parser* is supposed to **automate** this *derivation* process.
    - Input : **A sequence of** $(t, c)$ **pairs**, where each *token* $t$ (e.g., r241) belongs to a *syntactic category* $c$ (e.g., register); and a *CFG G*.
    - Output : A *valid derivation* (as an *AST*); or A *parse error*.
- In the process of constructing an *AST* for the input program:
  - *Root* of AST: The *start symbol S* of *G*
  - *Internal nodes*: A *subset of variables V* of *G*
  - *Leaves* of AST: A *token/terminal* sequence
    $\Rightarrow$ Discovering the *grammatical connections* (w.r.t. *R* of *G*) between the *root*, *internal nodes*, and *leaves* is the hard part!
- Approaches to Parsing:      [ $w \in (V \cup \Sigma)^*$, $A \in V$, $A \to w \in R$ ]
  - **Top-down** parsing
    For a node representing *A*, <u>extend it with a subtree</u> representing *w*.
  - **Bottom-up** parsing
    For a substring matching *w*, <u>build a node</u> representing *A* accordingly.

# 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

# TDP: Exercise (1)

- Given the following CFG **G**:

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

  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*.

- Given the following CFG **G**:

$$
\begin{array}{rcl}
Expr & \rightarrow & Term \ \ Expr' \\
Expr' & \rightarrow & + \ Term \ \ Expr' \\
& | & \epsilon \\
Term & \rightarrow & Factor \ \ Term' \\
Term' & \rightarrow & * \ Factor \ \ Term' \\
& | & \epsilon \\
Factor & \rightarrow & (Expr) \\
& | & a
\end{array}
$$

  **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., *Expr* $\rightarrow \epsilon$)?
  **A**: Execute *focus* := *trace*.pop() to advance to next node.

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

# **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 $\exists A \in V \bullet A \overset{*}{\Rightarrow} A$

○ A **direct** LR if $A \rightarrow A\alpha \in R$ for non-terminal $A \in V$

e.g.,

| Expr | $\rightarrow$ | Expr + Term |
|---|---|---|
| | \| | Term |
| Term | $\rightarrow$ | Term * Factor |
| | \| | Factor |
| Factor | $\rightarrow$ | ( Expr ) |
| | \| | a |

e.g.,

| Expr | $\rightarrow$ | Expr + Term |
|---|---|---|
| | \| | Expr − Term |
| | \| | Term |
| Term | $\rightarrow$ | Term * Factor |
| | \| | Term / Factor |
| | \| | Factor |

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

| A | $\rightarrow$ | Br |
|---|---|---|
| B | $\rightarrow$ | Cd |
| C | $\rightarrow$ | At |

| A | $\rightarrow$ | Ba | \| | b |
|---|---|---|---|---|
| B | $\rightarrow$ | Cd | \| | e |
| C | $\rightarrow$ | Df | \| | g |
| D | $\rightarrow$ | f | \| | Aa | \| | Cg |

$A \rightarrow Br, B \overset{*}{\Rightarrow} Atd$        $A \rightarrow Ba, B \overset{*}{\Rightarrow} Aafd$

## TDP: (Preventively) Eliminating LRs

```
1    ALGORITHM: RemoveLR
2      INPUT: CFG G = (V, Σ, R, S)
3      ASSUME: G has 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 → βAᵢ', Aᵢ' → αAᵢ' | ε
```

- **L9** to **L12**: Remove *indirect* left-recursions from $A_1$ to $A_{i-1}$.
- **L13** to **L14**: 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 $A_1, A_2, \ldots, A_{i-1}$.
  - More precisely: $\forall j : j < i \bullet \neg(\exists k \bullet k \leq j \wedge A_j \to A_k \cdots \in R)$

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

- Motivations:
  - **TDParse** handles each $\epsilon$-production as a special case.
  - **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 $A \to \epsilon$.

- A variable $A$ is *nullable* if $A \overset{*}{\Rightarrow} \epsilon$.
  - Each terminal symbol is *not nullable*.
  - Variable $A$ is *nullable* if either:
    - $A \to \epsilon \in R$; or
    - $A \to B_1 B_2 \ldots B_k \in R$, where each variable $B_i$ $(1 \leq i \leq k)$ is a *nullable*.
- Given a production $B \to CAD$, if only variable $A$ is *nullable*,
  then there are 2 versions of $B$: $B \to CAD \mid CD$
- In general, given a production $A \to X_1 X_2 \ldots X_k$ 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 $A \to \epsilon$ ]

- 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]

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

# 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 <u>consequential</u> ***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 the ***lookahead*** of one symbol.

  - We also call such backtrack-free CFG a ***predictive grammar***.

# The FIRST Set: Definition

- Say we write $T \subset \mathbb{P}(\Sigma^*)$ to denote the set of valid tokens recognizable by the scanner.

- **FIRST** $(\alpha) \triangleq$ set of symbols that can appear as the **_first word_** in some string derived from $\alpha$.

- More precisely:

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

## 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 | | | $\epsilon$ |
| 3 | | | - *Term Expr′* | 9 | *Factor* | → | ( *Expr* ) |
| 4 | | | $\epsilon$ | 10 | | | num |
| 5 | *Term* | → | *Factor Term′* | 11 | | | name |

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

| | num | name | + | - | × | ÷ | ( | ) | eof | $\epsilon$ |
|---|---|---|---|---|---|---|---|---|---|---|
| FIRST | num | name | + | - | x | ÷ | ( | ) | eof | $\epsilon$ |

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

| | *Expr* | *Expr'* | *Term* | *Term'* | *Factor* |
|---|---|---|---|---|---|
| FIRST | ( , name , num | + , - , $\epsilon$ | ( , name , num | x , ÷ , $\epsilon$ | ( , name , num |

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

```
ALGORITHM: GetFirst
  INPUT: CFG G = (V, Σ, R, S)
    T ⊂ Σ* denotes valid terminals
  OUTPUT: FIRST : V ∪ T ∪ {ε, eof} ⟶ ℙ(T ∪ {ε, eof})
PROCEDURE:
  for α ∈ (T ∪ {eof, ε}): FIRST(α) := {α}
  for A ∈ V: FIRST(A) := ∅
  lastFirst := ∅
  while (lastFirst ≠ FIRST):
     lastFirst := FIRST
     for A → β₁β₂...βₖ ∈ R s.t. ∀βⱼ : βⱼ ∈ (T ∪ V):
        rhs := FIRST(β₁) − {ε}
        for (i := 1; ε ∈ FIRST(βᵢ) ∧ i < k; i++):
           rhs := rhs ∪ (FIRST(βᵢ₊₁) − {ε})
        if i = k ∧ ε ∈ FIRST(βₖ) then
           rhs := rhs ∪ {ε}
        end
        FIRST(A) := FIRST(A) ∪ rhs
```

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

$$\text{FIRST} : V \cup T \cup \{\epsilon, \textit{eof}\} \longrightarrow \mathbb{P}(T \cup \{\epsilon, \textit{eof}\})$$

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

$$\text{FIRST} : (V \cup T \cup \{\epsilon, \textit{eof}\})^* \longrightarrow \mathbb{P}(T \cup \{\epsilon, \textit{eof}\})$$

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

- More precisely:

$\text{FIRST}(\beta_1 \beta_2 \ldots \beta_n) =$
$$\left\{ \text{FIRST}(\beta_1) \cup \text{FIRST}(\beta_2) \cup \cdots \cup \text{FIRST}(\beta_{k-1}) \cup \text{FIRST}(\beta_k) \left| \begin{array}{l} \forall i : 1 \leq i < k \bullet \epsilon \in \text{FIRST}(\beta_i) \\ \wedge \\ \epsilon \notin \text{FIRST}(\beta_k) \end{array} \right. \right.$$

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

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 | | \| | $\epsilon$ |
| 3 | | \| | - *Term Expr'* | | 9 | *Factor* | → | ( *Expr* ) |
| 4 | | \| | $\epsilon$ | | 10 | | \| | num |
| 5 | *Term* | → | *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**($\epsilon$) = {$\epsilon$}

# Is the FIRST Set Sufficient

- Consider the following three productions:

| Expr′ | → | + | Term | Term′ | (1) |
|-------|---|---|------|-------|-----|
|       | \| | − | Term | Term′ | (2) |
|       | \| | $\epsilon$ |      |       | (3) |

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* ≠ **FIRST**$(+) \wedge$ *focus* ≠ **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 leadling symbols.

- **FOLLOW** $(v : V) \triangleq$ set of symbols that can appear to the
  underline{immediate right} of a string derived from *v*.

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

# The FOLLOW Set: Examples

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

| 0 | *Goal* | $\rightarrow$ | *Expr* | 6 | *Term'* | $\rightarrow$ | x *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, +, -, x, ÷, ) |

# Computing the FOLLOW Set

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

```
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)
```

# Backtrack-Free Grammar

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

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

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

- A **backtrack-free grammar** has each of its productions $A \to \gamma_1 \mid \gamma_2 \mid \ldots \mid \gamma_n$ satisfying:

$$\forall i, j : 1 \le i, j \le n \land i \ne j \bullet \text{START}(\gamma_i) \cap \text{START}(\gamma_j) = \varnothing$$

# 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

     if ∃ unvisited rule focus → β₁β₂...βₙ ∈ R ∧   word ∈ START(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

Is the following CFG *backtrack free*?

| 11 | *Factor* | → | name |
|----|----------|---|------|
| 12 | | | name [ *ArgList* ] |
| 13 | | | name ( *ArgList* ) |
| 15 | *ArgList* | → | *Expr MoreArgs* |
| 16 | *MoreArgs* | → | , *Expr MoreArgs* |
| 17 | | | $\epsilon$ |

○ $\epsilon \notin$ **FIRST**(*Factor*) ⇒ **START**(*Factor*) = **FIRST**(*Factor*)

○ **FIRST**(*Factor* → name) = {name}

○ **FIRST**(*Factor* → name [*ArgList*]) = {name}

○ **FIRST**(*Factor* → name (*ArgList*)) = {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.

# Backtrack-Free Grammar: Left-Factoring

- A CFG is <u>not</u> 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 $\alpha$:

    $$A \to \alpha\beta_1 \mid \alpha\beta_2 \mid \ldots \mid \alpha\beta_n \mid \gamma_1 \mid \gamma_2 \mid \ldots \mid \gamma_j$$

    [ each of $\gamma_1, \gamma_2, \ldots, \gamma_j$ does not begin with $\alpha$ ]

  - Rewrite  that production rule as:

    $$\begin{aligned} A &\to \alpha B \mid \gamma_1 \mid \gamma_2 \mid \ldots \mid \gamma_j \\ B &\to \beta_1 \mid \beta_2 \mid \ldots \mid \beta_n \end{aligned}$$

  - New rule $B \to \beta_1 \mid \beta_2 \mid \ldots \mid \beta_n$ may <u>also</u> contain ***common prefixes***.
  - Rewriting  continues  until no common prefixes are identified.

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

| 11 | *Factor* | → | name |
|----|----------|---|------|
| 12 | | \| | name [ *ArgList* ] |
| 13 | | \| | name ( *ArgList* ) |
| 15 | *ArgList* | → | *Expr MoreArgs* |
| 16 | *MoreArgs* | → | , *Expr MoreArgs* |
| 17 | | \| | $\epsilon$ |

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

| *Factor* | → | name *Arguments* |
|----------|---|------------------|
| *Arguments* | → | [ *ArgList* ] |
| | \| | ( *ArgList* ) |
| | \| | $\epsilon$ |

Any more **common prefixes**?                                        [ No ]

# TDP: Terminating and Backtrack-Free

- Given an <u>arbitrary</u> 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 <u>minimize</u> the need of **backtracking**?
    **A.** left-factoring & one-symbol lookahead using **START**
- **Not** every context-free <u>language</u> has a corresponding **backtrack-free** context-free <u>grammar</u>.

  Given a CFL $l$, the following is **undecidable**:

  $$\exists cfg \mid L(cfg) = l \wedge isBacktrackFree(cfg)$$

- Given a CFG $g = (V, \Sigma, R, S)$, whether or not $g$ is **backtrack-free** is **decidable**:

  For each $A \rightarrow \gamma_1 \mid \gamma_2 \mid \ldots \mid \gamma_n \in R$:

  $$\forall i, j : 1 \leq i, j \leq n \wedge i \neq j \bullet \textbf{START}(\gamma_i) \cap \textbf{START}(\gamma_j) = \varnothing$$

# 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 **START** sets
  - An efficient ***recursive-descent*** parser
    This generated parser is called an *LL(1) parser*, which:
    - Processes input from **L**eft to right
    - 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.

## Backtrack-Free Parsing (2.2)

- Consider this CFG with **START** sets of the RHSs:

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

- 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

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.

# BUP: Discovering Rightmost Derivation

- In TDP, we build the <u>start variable</u> 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 <mark>*reducible*</mark> 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 <mark>*reduction*</mark> can be identified.

  **Q.** Why can BUP find the ***rightmost*** derivation (RMD), if any?

  **A.** BUP discovers steps in a ***RMD*** in its <mark>*reverse*</mark> order.

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

```
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(0) /* start state */
 word := NextWord()
 while(true)
   state := trace.top()
   act := Action[state, word]
   if act = ``accept'' then
     succeed()
   elseif act = ``reduce based on 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 to sᵢ'' then
     trace.push(word)
     trace.push(i)
     word := NextWord()
   else
     fail()
```

- Consider the following grammar for parentheses:

| | | |
|---|---|---|
| 1 | $Goal \rightarrow$ | $List$ |
| 2 | $List \rightarrow$ | $List\ Pair$ |
| 3 | | $\mid Pair$ |
| 4 | $Pair \rightarrow$ | $\underline{(\ Pair\ )}$ |
| 5 | | $\mid \underline{(\ )}$ |

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

| | Action **Table** | | | Goto **Table** | |
|---|---|---|---|---|---|
| **State** | 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:

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

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* |

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* |

Consider the steps of performing BUP on input $\boxed{\text{( ) )}}$ :

| Iteration | State | *word* | Stack | Handle | Action |
|-----------|-------|--------|-------|--------|--------|
| *initial* | — | <u>(</u> | $ 0 | *— none —* | — |
| 1 | 0 | <u>(</u> | $ 0 | *— none —* | *shift 3* |
| 2 | 3 | <u>)</u> | $ 0 <u>(</u> 3 | *— none —* | *shift 7* |
| 3 | 7 | <u>)</u> | $ 0 <u>(</u> 3 <u>)</u> 7 | *— none —* | *error* |

# LR(1) Items: Definition

- In **LR(1)** parsing, *Action* and *Goto* tabeles encode legitimate ways (w.r.t. a CFG) 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.,

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

where:

- A *production rule* $\boxed{A \to \beta\gamma}$ is currently being applied.
- A *terminal symbol* $\boxed{a}$ servers as a *lookahead symbol*.
- A *placeholder* $\boxed{\bullet}$ indicates the parser's *stack top*.
  - ✓ The parser's *stack* contains $\beta$ ("left context").
  - ✓ $\gamma$ is yet to be matched.
  - Upon matching $\beta\gamma$, if a matches the current <u>word</u>, then we "replace" $\beta\gamma$ (and their associated <u>states</u>) with *A* (and its associated <u>state</u>).

# LR(1) Items: Scenarios

An LR(1) item can denote:

**1. POSSIBILITY** $[A \rightarrow \bullet\beta\gamma, \ \mathtt{a}]$

- 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*.

**2. PARTIAL COMPLETION** $[A \rightarrow \beta \bullet \gamma, \ \mathtt{a}]$

- The parser has progressed from $[A \rightarrow \bullet\beta\gamma, \ \mathtt{a}]$ by recognizing $\beta$.
- Recognizing a $\gamma$ next would be one step towards discovering an *A*.

**3. COMPLETION** $[A \rightarrow \beta\gamma\bullet, \ \mathtt{a}]$

- Parser has progressed from $[A \rightarrow \bullet\beta\gamma, \ \mathtt{a}]$ by recognizing $\beta\gamma$.
- $\beta\gamma$ found in a context where an *A* followed by $\mathtt{a}$ would be valid.
- If the current input <u>word</u> matches $\mathtt{a}$, then:
  - Current ***complet item*** is a *handle* .
  - Parser can ***reduce*** $\beta\gamma$ to *A*
  - Accordingly, in the ***stack***, $\beta\gamma$ (and their associated <u>states</u>) are replaced with *A* (and its associated <u>state</u>).

# LR(1) Items: Example (1.1)

Consider the following grammar for parentheses:

| 1 | $Goal \rightarrow List$ |
|---|---|
| 2 | $List \rightarrow List\ Pair$ |
| 3 | $\mid Pair$ |
| 4 | $Pair \rightarrow \underline{(}\ Pair\ \underline{)}$ |
| 5 | $\mid \underline{(}\ \underline{)}$ |

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

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

**Intermediate States:** Subset Construction

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

○ **FOLLOW**($List$) = $\{\texttt{eof}, \texttt{(}\}$   **FOLLOW**($Pair$) = $\{\texttt{eof}, \texttt{(}, \texttt{)}\}$

○ For each production $A \rightarrow \beta$, given **FOLLOW**($A$), **LR(1) items** are:

$$\{\ [A \rightarrow \bullet\beta\gamma,\ a] \mid a \in \textbf{FOLLOW}(A)\ \}$$
$$\cup$$
$$\{\ [A \rightarrow \beta \bullet \gamma,\ a] \mid a \in \textbf{FOLLOW}(A)\ \}$$
$$\cup$$
$$\{\ [A \rightarrow \beta\gamma\bullet,\ a] \mid a \in \textbf{FOLLOW}(A)\ \}$$

# LR(1) Items: Example (1.2)

**Q.** Given production $A \rightarrow \beta$ (e.g., *Pair* → ( *Pair* ) ), how many LR(1) items can be generated?

○ The current parsing progress (on matching the RHS) can be:
  **1.** •( *Pair* )
  **2.** ( •*Pair* )
  **3.** ( *Pair*• )
  **4.** ( *Pair* ) •

○ Lookahead symbol following *Pair*?    **FOLLOW**(*Pair*) = {eof, (,)}

○ <u>All</u> possible LR(1) items related to *Pair* → ( *Pair* ) ?
  ✓ [•( *Pair* ), eof]   [•( *Pair* ), (]   [•( *Pair* ), )]
  ✓ [( •*Pair* ), eof]   [( •*Pair* ), (]   [( •*Pair* ), )]
  ✓ [( *Pair*• ), eof]   [( *Pair*• ), (]   [( *Pair*• ), )]
  ✓ [( *Pair* )•, eof]   [( *Pair* )•, (]   [( *Pair* )•, )]

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

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

# LR(1) Items: Example (1.3)

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

$[Goal \rightarrow \bullet\ List, \text{eof}]$

$[Goal \rightarrow List\ \bullet, \text{eof}]$

| | |
|---|---|
| $[List \rightarrow \bullet\ List\ Pair, \text{eof}]$ | $[List \rightarrow \bullet\ List\ Pair, \underline{(}\ ]$ |
| $[List \rightarrow List\ \bullet\ Pair, \text{eof}]$ | $[List \rightarrow List\ \bullet\ Pair, \underline{(}\ ]$ |
| $[List \rightarrow List\ Pair\ \bullet, \text{eof}]$ | $[List \rightarrow List\ Pair\ \bullet, \underline{(}\ ]$ |

| | |
|---|---|
| $[List \rightarrow \bullet\ Pair, \text{eof}\ ]$ | $[List \rightarrow \bullet\ Pair, \underline{(}\ ]$ |
| $[List \rightarrow Pair\ \bullet, \text{eof}\ ]$ | $[List \rightarrow Pair\ \bullet, \underline{(}\ ]$ |

| | | |
|---|---|---|
| $[Pair \rightarrow \bullet\ \underline{(}\ Pair\ \underline{)}, \text{eof}]$ | $[Pair \rightarrow \bullet\ \underline{(}\ Pair\ \underline{)}, \underline{)}]$ | $[Pair \rightarrow \bullet\ \underline{(}\ Pair\ \underline{)}, \underline{(}]$ |
| $[Pair \rightarrow \underline{(}\ \bullet\ Pair\ \underline{)}, \text{eof}]$ | $[Pair \rightarrow \underline{(}\ \bullet\ Pair\ \underline{)}, \underline{)}]$ | $[Pair \rightarrow \underline{(}\ \bullet\ Pair\ \underline{)}, \underline{(}]$ |
| $[Pair \rightarrow \underline{(}\ Pair\ \bullet\ \underline{)}, \text{eof}]$ | $[Pair \rightarrow \underline{(}\ Pair\ \bullet\ \underline{)}, \underline{)}]$ | $[Pair \rightarrow \underline{(}\ Pair\ \bullet\ \underline{)}, \underline{(}]$ |
| $[Pair \rightarrow \underline{(}\ Pair\ \underline{)}\ \bullet, \text{eof}]$ | $[Pair \rightarrow \underline{(}\ Pair\ \underline{)}\ \bullet, \underline{)}]$ | $[Pair \rightarrow \underline{(}\ Pair\ \underline{)}\ \bullet, \underline{(}]$ |

| | | |
|---|---|---|
| $[Pair \rightarrow \bullet\ \underline{(}\ \underline{)}, \text{eof}]$ | $[Pair \rightarrow \bullet\ \underline{(}\ \underline{)}, \underline{(}]$ | $[Pair \rightarrow \bullet\ \underline{(}\ \underline{)}, \underline{)}]$ |
| $[Pair \rightarrow \underline{(}\ \bullet\ \underline{)}, \text{eof}]$ | $[Pair \rightarrow \underline{(}\ \bullet\ \underline{)}, \underline{(}]$ | $[Pair \rightarrow \underline{(}\ \bullet\ \underline{)}, \underline{)}]$ |
| $[Pair \rightarrow \underline{(}\ \underline{)}\ \bullet, \text{eof}]$ | $[Pair \rightarrow \underline{(}\ \underline{)}\ \bullet, \underline{(}]$ | $[Pair \rightarrow \underline{(}\ \underline{)}\ \bullet, \underline{)}]$ |

# LR(1) Items: Example (2)

Consider the following grammar for expressions:

| | | | | | | |
|---|---|---|---|---|---|---|
| 0 | *Goal* | → | *Expr* | 6 | *Term′* | → | x *Factor Term′* |
| 1 | *Expr* | → | *Term Expr′* | 7 | | | ÷ *Factor Term′* |
| 2 | *Expr′* | → | + *Term Expr′* | 8 | | | $\epsilon$ |
| 3 | | | - *Term Expr′* | 9 | *Factor* | → | ( *Expr* ) |
| 4 | | | $\epsilon$ | 10 | | | num |
| 5 | *Term* | → | *Factor Term′* | 11 | | | name |

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

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

| | *Expr* | *Expr′* | *Term* | *Term′* | *Factor* |
|---|---|---|---|---|---|
| FOLLOW | eof, ) | eof, ) | eof, +, -, ) | eof, +, -, ) | eof, +, -, x, ÷, ) |

**Tips.** Ignore $\epsilon$ **production** such as *Expr′* → $\epsilon$
since the **FOLLOW** sets already take them into consideration.

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

| | |
|---|---|
| 1 | $Goal \rightarrow List$ |
| 2 | $List \rightarrow List\ Pair$ |
| 3 | $\mid Pair$ |
| 4 | $Pair \rightarrow \underline{(}\ Pair\ \underline{)}$ |
| 5 | $\mid \underline{(}\ \underline{)}$ |

<u>Recall</u>:

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

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

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

○ The *canonical collection*          [ Example of $\mathcal{CC}$ ]
$$\mathcal{CC} = \{cc_0, cc_1, cc_2, \ldots, cc_n\}$$
denotes the set of **valid <u>subset</u> states** of a **LR(1) parser**.
  - Each $cc_i \in \mathcal{CC}$ ($0 \leq i \leq n$) is a set of **LR(1) items**.
  - $\mathcal{CC} \subseteq \mathbb{P}(\textbf{LR(1) items})$      $|\mathcal{CC}|$?      [ $|\mathcal{CC}| \leq 2^{|LR(1)\ items|}$ ]

○ To model a **LR(1) parser**, we use techniques analogous to how an $\epsilon$-NFA is converted into a DFA (subset construction and $\epsilon$-closure).

○ **Analogies.**
  - ✓ **LR(1) items** ≈ states of source *NFA*
  - ✓ $\mathcal{CC}$ ≈ <u>subset</u> states of target *DFA*

```
1   ALGORITHM: closure
2    INPUT: CFG G = (V, Σ, R, S), a set s of LR(1) items
3    OUTPUT: a set of LR(1) items
4   PROCEDURE:
5    lastS := ∅
6    while (lastS ≠ s) :
7      lastS := s
8      for [A → ⋯ • C δ, a] ∈ s:
9        for C → γ ∈ R:
10         for b ∈ FIRST(δa):
11           s := s ∪ { [ C → •γ, b] }
12   return s
```

- **Line 8**: $[A \to \cdots \bullet\; C\; \delta,\; a] \in s$ indicates that the parser's next task is to match $C\; \delta$ with a lookahead symbol $a$.
- **Line 9**: <u>Given</u>: matching $\gamma$ can reduce to $C$
- **Line 10**: <u>Given</u>: $b \in$ **FIRST**$(\delta a)$ is a valid lookahead symbol after reducing $\gamma$ to $C$
- **Line 11**: Add a new item $[\; C \to \bullet\gamma,\; b]$ into $s$.
- **Line 6**: Termination is guaranteed.

  ∵ Each iteration adds ≥ 1 item to $s$ (otherwise *lastS* ≠ *s* is *false*).

| | |
|---|---|
| 1 | $Goal \rightarrow List$ |
| 2 | $List \;\; \rightarrow List \; Pair$ |
| 3 | $\qquad | \;\; Pair$ |
| 4 | $Pair \;\; \rightarrow \underline{(} \;\; Pair \;\; \underline{)}$ |
| 5 | $\qquad | \;\; \underline{(} \;\; \underline{)}$ |

*Initial State***:** $[Goal \rightarrow \bullet List, \; \texttt{eof}]$

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

# **Constructing $\mathcal{CC}$: 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 := ∅
6    for item ∈ s:
7      if item = [α → β • xδ, a] then
8        moved := moved ∪ { [α → βx • δ, a] }
9      end
10   return closure(moved)
```

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

| 1 | $Goal \rightarrow List$ |
| 2 | $List \rightarrow List\ Pair$ |
| 3 | $\mid Pair$ |
| 4 | $Pair \rightarrow \underline{(}\ Pair\ \underline{)}$ |
| 5 | $\mid \underline{(}\ \underline{)}$ |

$$cc_0 = \begin{Bmatrix} [Goal \rightarrow \bullet List, \texttt{eof}] & [List \rightarrow \bullet List\ Pair, \texttt{eof}] & [List \rightarrow \bullet List\ Pair, \underline{(}] \\ [List \rightarrow \bullet Pair, \texttt{eof}] & [List \rightarrow \bullet Pair, \underline{(}] & [Pair \rightarrow \bullet \underline{(}\ Pair\ \underline{)}, \texttt{eof}] \\ [Pair \rightarrow \bullet \underline{(}\ Pair\ \underline{)}, \underline{(}] & [Pair \rightarrow \bullet \underline{(}\ \underline{)}, \texttt{eof}] & [Pair \rightarrow \bullet \underline{(}\ \underline{)}, \underline{(}] \end{Bmatrix}$$

Calculate $goto(cc_0,\ \underline{(}\ )$.          ["next state" from $cc_0$ taking $\underline{(}$]

```
1   ALGORITHM: BuildCC
2    INPUT: a grammar G = (V, Σ, R, S), goal production S → S′
3    OUTPUT:
4       (1) a set CC = {cc_0, cc_1, ..., cc_n} where cc_i ⊆ G′s LR(1) items
5       (2) a transition function
6   PROCEDURE:
7    cc_0 := closure({[S → •S′, eof]})
8    CC := {cc_0}
9    processed := {cc_0}
10   lastCC := ∅
11   while (lastCC ≠ CC):
12     lastCC := CC
13     for cc_i s.t. cc_i ∈ CC ∧ cc_i ∉ processed:
14       processed := processed ∪ {cc_i}
15       for x s.t. [···→···•x...] ∈ cc_i
16         temp := goto(cc_i, x)
17         if temp ∉ CC then
18           CC := CC ∪ {temp}
19         end
20         δ := δ ∪ (cc_i, x, temp)
```

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

- Calculate $\mathcal{CC} = \{cc_0, cc_1, \ldots, cc_{11}\}$
- Calculate the transition function $\delta : \mathcal{CC} \times (\Sigma \cup V) \rightarrow \mathcal{CC}$
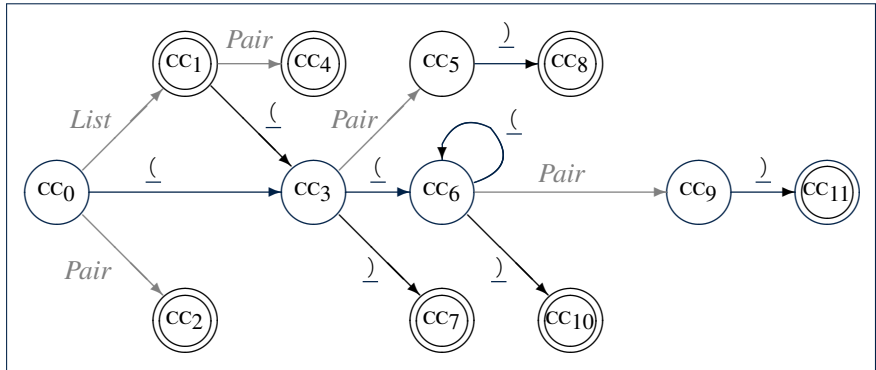
Resulting transition table:

| 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 | $cc_{11}$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ | $\emptyset$ |

Resulting DFA for the parser:

Resulting canonical collection $\mathcal{CC}$:         [ Def. of $\mathcal{CC}$ ]

$$cc_0 = \begin{cases} [Goal \to \bullet \, List, \text{eof}] & [List \to \bullet \, List \, Pair, \text{eof}] & [List \to \bullet \, List \, Pair, \underline{(}] \\ [List \to \bullet \, Pair, \text{eof}] & [List \to \bullet \, Pair, \underline{(}] & [Pair \to \bullet \, \underline{(} \, Pair \, \underline{)}, \text{eof}] \\ [Pair \to \bullet \, \underline{(} \, Pair \, \underline{)}, \underline{(}] & [Pair \to \bullet \, \underline{(} \, \underline{)}, \text{eof}] & [Pair \to \bullet \, \underline{(} \, \underline{)}, \underline{(}] \end{cases}$$

$$cc_1 = \begin{cases} [Goal \to List \, \bullet, \text{eof}] & [List \to List \, \bullet \, Pair, \text{eof}] & [List \to List \, \bullet \, Pair, \underline{(}] \\ [Pair \to \bullet \, \underline{(} \, Pair \, \underline{)}, \text{eof}] & [Pair \to \bullet \, \underline{(} \, Pair \, \underline{)}, \underline{(}] & [Pair \to \bullet \, \underline{(} \, \underline{)}, \text{eof}] \\ & [Pair \to \bullet \, \underline{(} \, \underline{)}, \underline{(}] \end{cases}$$

$$cc_2 = \left\{ [List \to Pair \, \bullet, \text{eof}] \quad [List \to Pair \, \bullet, \underline{(}] \right\}$$

$$cc_3 = \begin{cases} [Pair \to \bullet \, \underline{(} \, Pair \, \underline{)}, \underline{)}] & [Pair \to \bullet \, \underline{(} \, Pair \, \underline{)}, \text{eof}] & [Pair \to \underline{(} \, \bullet \, Pair \, \underline{)}, \underline{(}] \\ [Pair \to \bullet \, \underline{(} \, \underline{)}, \underline{)}] & [Pair \to \underline{(} \, \bullet \, \underline{)}, \text{eof}] & [Pair \to \underline{(} \, \bullet \, \underline{)}, \underline{(}] \end{cases}$$

$$cc_4 = \left\{ [List \to List \, Pair \, \bullet, \text{eof}] \quad [List \to List \, Pair \, \bullet, \underline{(}] \right\}$$

$$cc_5 = \left\{ [Pair \to \underline{(} \, Pair \, \bullet \, \underline{)}, \text{eof}] \quad [Pair \to \underline{(} \, Pair \, \bullet \, \underline{)}, \underline{(}] \right\}$$

$$cc_6 = \begin{cases} [Pair \to \bullet \, \underline{(} \, Pair \, \underline{)}, \underline{)}] & [Pair \to \underline{(} \, \bullet \, Pair \, \underline{)}, \underline{)}] \\ [Pair \to \bullet \, \underline{(} \, \underline{)}, \underline{)}] & [Pair \to \underline{(} \, \bullet \, \underline{)}, \underline{)}] \end{cases}$$

$$cc_7 = \left\{ [Pair \to \underline{(} \, \underline{)} \, \bullet, \text{eof}] \quad [Pair \to \underline{(} \, \underline{)} \, \bullet, \underline{(}] \right\}$$

$$cc_8 = \left\{ [Pair \to \underline{(} \, Pair \, \underline{)} \, \bullet, \text{eof}] \quad [Pair \to \underline{(} \, Pair \, \underline{)} \, \bullet, \underline{(}] \right\}$$

$$cc_9 = \left\{ [Pair \to \underline{(} \, Pair \, \bullet \, \underline{)}, \underline{)}] \right\}$$

$$cc_{10} = \left\{ [Pair \to \underline{(} \, \underline{)} \, \bullet, \underline{)}] \right\}$$

$$cc_{11} = \left\{ [Pair \to \underline{(} \, Pair \, \underline{)} \, \bullet, \underline{)}] \right\}$$

# Constructing *Action* and *Goto* Tables (1)

```
1   ALGORITHM: BuildActionGotoTables
2     INPUT:
3       (1) a grammar G = (V, Σ, R, S)
4       (2) goal production S → S′
5       (3) a canonical collection CC = {cc₀, cc₁, ..., ccₙ}
6       (4) a transition function δ : CC × Σ → CC
7     OUTPUT: Action Table & Goto Table
8     PROCEDURE:
9       for ccᵢ ∈ CC:
10        for item ∈ ccᵢ:
11          if item = [A → β • xγ, a] ∧ δ(ccᵢ, x) = ccⱼ then
12            Action[i, x] := shift j
13          elseif item = [A → β•, a] then
14            Action[i, a] := reduce A → β
15          elseif item = [S → S′•, eof] then
16            Action[i, eof] := accept
17          end
18        for v ∈ V:
19          if δ(ccᵢ, v) = ccⱼ 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*

Resulting **Action** and **Goto** tables:

|  | *Action* **Table** | | | *Goto* **Table** | |
|---|---|---|---|---|---|
| **State** | 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 |  |  |

| 1 | *Goal* | $\rightarrow$ | *Stmt* |
|---|---|---|---|
| 2 | *Stmt* | $\rightarrow$ | if expr then *Stmt* |
| 3 | | \| | if expr then *Stmt* else *Stmt* |
| 4 | | \| | assign |

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

Resulting transition table:

|   | Item | *Goal* | *Stmt* | if | expr | then | else | assign | eof |
|---|------|--------|--------|-----|------|------|------|--------|-----|
| 0 | $cc_0$ | Ø | $cc_1$ | $cc_2$ | Ø | Ø | Ø | $cc_3$ | Ø |
| 1 | $cc_1$ | Ø | Ø | Ø | Ø | Ø | Ø | Ø | Ø |
|   | $cc_2$ | Ø | Ø | Ø | $cc_4$ | Ø | Ø | Ø | Ø |
|   | $cc_3$ | Ø | Ø | Ø | Ø | Ø | Ø | Ø | Ø |
| 2 | $cc_4$ | Ø | Ø | Ø | Ø | $cc_5$ | Ø | Ø | Ø |
| 3 | $cc_5$ | Ø | $cc_6$ | $cc_7$ | Ø | Ø | Ø | $cc_8$ | Ø |
| 4 | $cc_6$ | Ø | Ø | Ø | Ø | Ø | $cc_9$ | Ø | Ø |
|   | $cc_7$ | Ø | Ø | Ø | $cc_{10}$ | Ø | Ø | Ø | Ø |
|   | $cc_8$ | Ø | Ø | Ø | Ø | Ø | Ø | Ø | Ø |
| 5 | $cc_9$ | Ø | $cc_{11}$ | $cc_2$ | Ø | Ø | Ø | $cc_3$ | Ø |
|   | $cc_{10}$ | Ø | Ø | Ø | Ø | $cc_{12}$ | Ø | Ø | Ø |
| 6 | $cc_{11}$ | Ø | Ø | Ø | Ø | Ø | Ø | Ø | Ø |
|   | $cc_{12}$ | Ø | $cc_{13}$ | $cc_7$ | Ø | Ø | Ø | $cc_8$ | Ø |
| 7 | $cc_{13}$ | Ø | Ø | Ø | Ø | Ø | $cc_{14}$ | Ø | Ø |
| 8 | $cc_{14}$ | Ø | $cc_{15}$ | $cc_7$ | Ø | Ø | Ø | $cc_8$ | Ø |
| 9 | $cc_{15}$ | Ø | Ø | Ø | Ø | Ø | Ø | Ø | Ø |

Resulting canonical collection $\mathcal{CC}$:

$$cc_0 = \begin{Bmatrix} [Goal \to \bullet \, Stmt, \text{eof}] & [Stmt \to \bullet \, \text{if expr then } Stmt, \text{eof}] \\ [Stmt \to \bullet \, \text{assign}, \text{eof}] & [Stmt \to \bullet \, \text{if expr then } Stmt \text{ else } Stmt, \text{eof}] \end{Bmatrix}$$

$$cc_1 = \Big\{ [Goal \to Stmt \, \bullet, \text{eof}] \Big\}$$

$$cc_2 = \begin{Bmatrix} [Stmt \to \text{if} \bullet \text{ expr then } Stmt, \text{eof}], \\ [Stmt \to \text{if} \bullet \text{ expr then } Stmt \text{ else } Stmt, \text{eof}] \end{Bmatrix}$$

$$cc_3 = \Big\{ [Stmt \to \text{assign} \, \bullet, \text{eof}] \Big\}$$

$$cc_4 = \begin{Bmatrix} [Stmt \to \text{if expr} \bullet \text{ then } Stmt, \text{eof}], \\ [Stmt \to \text{if expr} \bullet \text{ then } Stmt \text{ else } Stmt, \text{eof}] \end{Bmatrix}$$

$$cc_5 = \begin{Bmatrix} [Stmt \to \text{if expr then} \bullet \, Stmt, \text{eof}], \\ [Stmt \to \text{if expr then} \bullet \, Stmt \text{ else } Stmt, \text{eof}], \\ [Stmt \to \bullet \, \text{if expr then } Stmt, \{\text{eof}, \text{else}\}], \\ [Stmt \to \bullet \, \text{assign}, \{\text{eof}, \text{else}\}], \\ [Stmt \to \bullet \, \text{if expr then } Stmt \text{ else } Stmt, \{\text{eof}, \text{else}\}] \end{Bmatrix}$$

$$cc_6 = \begin{Bmatrix} [Stmt \to \text{if expr then } Stmt \, \bullet, \text{eof}], \\ [Stmt \to \text{if expr then } Stmt \, \bullet \text{ else } Stmt, \text{eof}] \end{Bmatrix}$$

$$cc_7 = \begin{Bmatrix} [Stmt \to \text{if} \bullet \text{ expr then } Stmt, \{\text{eof}, \text{else}\}], \\ [Stmt \to \text{if} \bullet \text{ expr then } Stmt \text{ else } Stmt, \{\text{eof}, \text{else}\}] \end{Bmatrix}$$

Resulting canonical collection $\mathcal{CC}$:

$$\text{CC}_8 = \{[\textit{Stmt} \rightarrow \text{assign} \bullet, \{\text{eof}, \text{else}\}]\}$$

$$\text{CC}_{10} = \begin{cases} [\textit{Stmt} \rightarrow \text{if expr} \bullet \text{then } \textit{Stmt}, \{\text{eof}, \text{else}\}], \\ [\textit{Stmt} \rightarrow \text{if expr} \bullet \text{then } \textit{Stmt} \text{ else } \textit{Stmt}, \{\text{eof}, \text{else}\}] \end{cases}$$

$$\text{CC}_{12} = \begin{cases} [\textit{Stmt} \rightarrow \text{if expr then} \bullet \textit{Stmt}, \{\text{eof}, \text{else}\}], \\ [\textit{Stmt} \rightarrow \text{if expr then} \bullet \textit{Stmt} \text{ else } \textit{Stmt}, \{\text{eof}, \text{else}\}], \\ [\textit{Stmt} \rightarrow \bullet \text{if expr then } \textit{Stmt}, \{\text{eof}, \text{else}\}], \\ [\textit{Stmt} \rightarrow \bullet \text{if expr then } \textit{Stmt} \text{ else } \textit{Stmt}, \{\text{eof}, \text{else}\}], \\ [\textit{Stmt} \rightarrow \bullet \text{assign}, \{\text{eof}, \text{else}\}] \end{cases}$$

$$\text{CC}_{14} = \begin{cases} [\textit{Stmt} \rightarrow \text{if expr then } \textit{Stmt} \text{ else} \bullet \textit{Stmt}, \{\text{eof}, \text{else}\}], \\ [\textit{Stmt} \rightarrow \bullet \text{if expr then } \textit{Stmt}, \{\text{eof}, \text{else}\}], \\ [\textit{Stmt} \rightarrow \bullet \text{if expr then } \textit{Stmt} \text{ else } \textit{Stmt}, \{\text{eof}, \text{else}\}], \\ [\textit{Stmt} \rightarrow \bullet \text{assign}, \{\text{eof}, \text{else}\}] \end{cases}$$

$$\text{CC}_9 = \begin{cases} [\textit{Stmt} \rightarrow \text{if expr then } \textit{Stmt} \text{ else} \bullet \textit{Stmt}, \text{eof}], \\ [\textit{Stmt} \rightarrow \bullet \text{if expr then } \textit{Stmt}, \text{eof}], \\ [\textit{Stmt} \rightarrow \bullet \text{if expr then } \textit{Stmt} \text{ else } \textit{Stmt}, \text{eof}], \\ [\textit{Stmt} \rightarrow \bullet \text{assign}, \text{eof}] \end{cases}$$

$$\text{CC}_{11} = \{[\textit{Stmt} \rightarrow \text{if expr then } \textit{Stmt} \text{ else } \textit{Stmt} \bullet, \text{eof}]\}$$

$$\text{CC}_{13} = \begin{cases} [\textit{Stmt} \rightarrow \text{if expr then } \textit{Stmt} \bullet, \{\text{eof}, \text{else}\}], \\ [\textit{Stmt} \rightarrow \text{if expr then } \textit{Stmt} \bullet \text{ else } \textit{Stmt}, \{\text{eof}, \text{else}\}] \end{cases}$$

# BUP: Discovering Ambiguity (3)

- Consider $cc_{13}$

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

  **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*.
  **Action**[13, `else`] cannot hold **shift** and **reduce** simultaneously.
  $\Rightarrow$ This is known as the  *shift-reduce conflict* .

- Consider another scenario:

$$cc_i = \left\{ \begin{array}{l} [A \rightarrow \gamma\delta\bullet, \ \texttt{a}], \\ [B \rightarrow \gamma\delta\bullet, \ \texttt{a}] \end{array} \right\}$$

  **Q**. What does it mean if the current word to consume is `a`?
  **A**. We can either **reduce** to *A* or **reduce** to *B*.
  **Action**[$i$, $a$] cannot hold *A* and *B* simultaneously.
  $\Rightarrow$ This is known as the  *reduce-reduce conflict* .