# Parser: Syntactic Analysis

**Readings: EAC2 Chapter 3**

EECS4302 A:
Compilers and Interpreters
Summer 2025

CHEN-WEI WANG
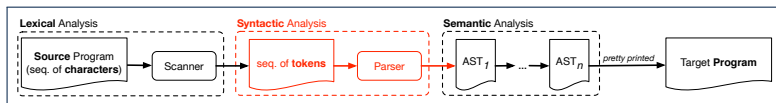
---

## Context-Free Languages: Introduction

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

---

## Parser in Context

- Recall:



- Treats the input programas 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 <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)* .

---

## CFG: Example (1.1)

- 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 <u>infinite</u> number of strings from

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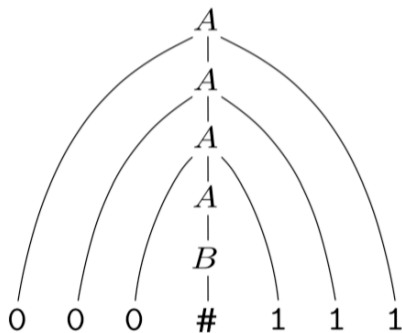A &\rightarrow B \\
B &\rightarrow \#
\end{aligned}
$$

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

## CFG: Example (1.2)

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.*

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

## CFG: Example (3)

Design a CFG for the following language:

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

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

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

## 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*.

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

---

## 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:
  1. All operations are <u>mixed</u> together.
  2. Relevant operations are <u>grouped</u> together.
     Try both!

---

## CFG: Example (5.2) Version 1

$$
\begin{array}{lll}
Expression & \rightarrow & IntegerConstant \\
 & | & -IntegerConstant \\
 & | & BooleanConstant \\
 & | & BinaryOp \\
 & | & UnaryOp \\
 & | & (\ Expression\ ) \\
\\
IntegerConstant & \rightarrow & Digit \\
 & | & Digit\ IntegerConstant \\
\\
Digit & \rightarrow & 0\,|\,1\,|\,2\,|\,3\,|\,4\,|\,5\,|\,6\,|\,7\,|\,8\,|\,9 \\
\\
BooleanConstant & \rightarrow & \text{TRUE} \\
 & | & \text{FALSE}
\end{array}
$$

---

## CFG: Example (5.3) Version 1

$$
\begin{array}{lll}
BinaryOp & \rightarrow & Expression + Expression \\
 & | & Expression - Expression \\
 & | & Expression * Expression \\
 & | & Expression\ /\ Expression \\
 & | & Expression\ \&\&\ Expression \\
 & | & Expression\ ||\ Expression \\
 & | & Expression => Expression \\
 & | & Expression == Expression \\
 & | & Expression\ /=\ Expression \\
 & | & Expression > Expression \\
 & | & Expression < Expression \\
\\
UnaryOp & \rightarrow & !\ Expression
\end{array}
$$

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.

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

  - $S \in V$ is is the **start variable**.
- Given strings $u, v, w \in (V \cup \Sigma)^*$, variable $A \in V$, a rule $A \rightarrow w$:
  - $\boxed{uAv \Rightarrow uwv}$ 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 $G$

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

## CFG: Formal Definition (2): Example

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

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

  Present your answer in a **formal** manner.
- $G = (\{S\}, \{(,)\}, R, S)$, where $R$ is

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

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

## CFG: Formal Definition (3): Example

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

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

  - $V = \{Expr, Term, Factor\}$
  - $\Sigma = \{\mathtt{a}, +, *, (, )\}$
  - $S = Expr$
- **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 $\varnothing$):

| | | |
|---|---|---|
| $L(\;\epsilon\;)$ | $=$ | $\{\epsilon\}$ |
| $L(\;a\;)$ | $=$ | $\{a\}$ |
| $L(\;E + F\;)$ | $=$ | $L(E) \cup L(F)$ |
| $L(\;EF\;)$ | $=$ | $L(E)L(F)$ |
| $L(\;E^*\;)$ | $=$ | $(L(E))^*$ |
| $L(\;(E)\;)$ | $=$ | $L(E)$ |

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

| | | |
|---|---|---|
| $S$ | $\rightarrow$ | $A \mid B$ |
| $A$ | $\rightarrow$ | $\epsilon \mid AC$ |
| $C$ | $\rightarrow$ | $00 \mid 1$ |
| $B$ | $\rightarrow$ | $\epsilon \mid BD$ |
| $D$ | $\rightarrow$ | $11 \mid 0$ |

## 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



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

---

## CFG: Rightmost Derivations (1)

| | | |
|---|---|---|
| *Expr* | $\rightarrow$ | *Expr* + *Term* \| *Term* |
| *Term* | $\rightarrow$ | *Term* $\star$ *Factor* \| *Factor* |
| *Factor* | $\rightarrow$ | ⟨*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 $\star$ a:

| | | |
|---|---|---|
| *Expr* | $\Rightarrow$ | *Expr* + *Term* |
| | $\Rightarrow$ | *Expr* + *Term* $\star$ *Factor* |
| | $\Rightarrow$ | *Expr* + *Term* $\star$ *a* |
| | $\Rightarrow$ | *Expr* + *Factor* $\star$ *a* |
| | $\Rightarrow$ | *Expr* + *a* $\star$ *a* |
| | $\Rightarrow$ | *Term* + *a* $\star$ *a* |
| | $\Rightarrow$ | *Factor* + *a* $\star$ *a* |
| | $\Rightarrow$ | *a* + *a* $\star$ *a* |

- This *RMD* suggests that a $\star$ a is the right operand of +.

---

## CFG: Leftmost Derivations (1)

| | | |
|---|---|---|
| *Expr* | $\rightarrow$ | *Expr* + *Term* \| *Term* |
| *Term* | $\rightarrow$ | *Term* $\star$ *Factor* \| *Factor* |
| *Factor* | $\rightarrow$ | ⟨*Expr*⟩ \| *a* |

- 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 $\star$ a:

| | | |
|---|---|---|
| *Expr* | $\Rightarrow$ | *Expr* + *Term* |
| | $\Rightarrow$ | *Term* + *Term* |
| | $\Rightarrow$ | *Factor* + *Term* |
| | $\Rightarrow$ | *a* + *Term* |
| | $\Rightarrow$ | *a* + *Term* $\star$ *Factor* |
| | $\Rightarrow$ | *a* + *Factor* $\star$ *Factor* |
| | $\Rightarrow$ | *a* + *a* $\star$ *Factor* |
| | $\Rightarrow$ | *a* + *a* $\star$ *a* |

- This *LMD* suggests that a $\star$ a is the right operand of +.

---

## CFG: Leftmost Derivations (2)

| | | |
|---|---|---|
| *Expr* | $\rightarrow$ | *Expr* + *Term* \| *Term* |
| *Term* | $\rightarrow$ | *Term* $\star$ *Factor* \| *Factor* |
| *Factor* | $\rightarrow$ | ⟨*Expr*⟩ \| *a* |

- *Unique LMD* for the string (a + a) $\star$ a:

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

- This *LMD* suggests that (a + a) is the left operand of $\star$.

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

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

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

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

● 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 ∵ Derivations with **distinct orders** of variable substitutions may still result in the **same parse tree**.

● For example:

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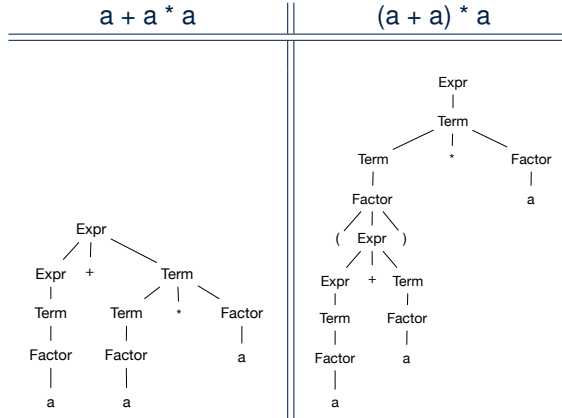Factor &\rightarrow (Expr) \mid a
\end{aligned}
$$

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

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



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

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**.

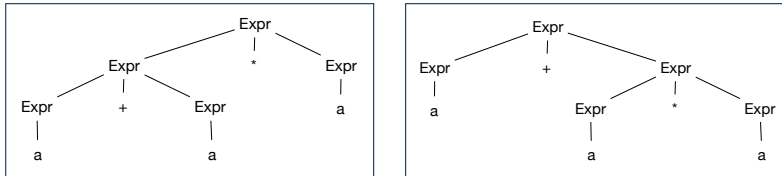○ $G$ is **ambiguous** if it generates some string ambiguously.

- 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
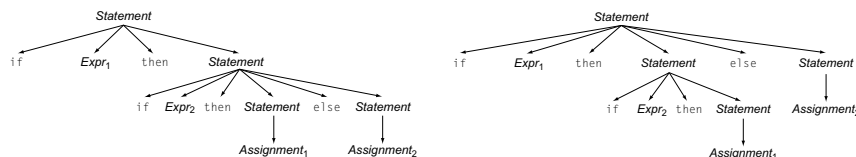- **Exercise**: Show **LMDs** for the two parse trees.

---

(**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`:

---

- Is the following grammar **ambiguous** ?

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

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

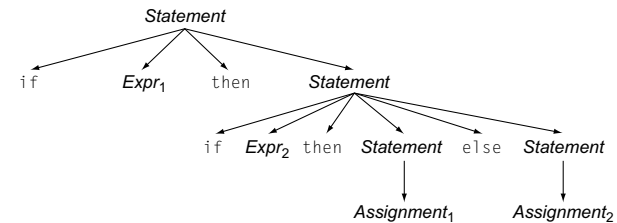`if` $Expr_1$ `then if` $Expr_2$ `then` $Assignment_1$ `else` $Assignment_2$



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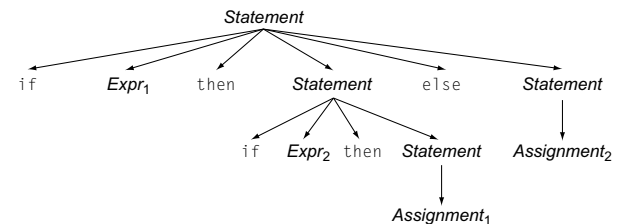- **Exercise**: Show **LMDs** for the two parse trees.

---

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

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

- 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.

## 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: $[\ w \in (V \cup \Sigma)^*, A \in V, \boxed{A \to w} \in R\ ]$
  - **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.

---

## TDP: Exercise (1)

- Given the following CFG **G**:

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

  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*.

---

## 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 (2)

- Given the following CFG **G**:

$$
\begin{aligned}
Expr &\to Term\ Expr' \\
Expr' &\to +\ Term\ Expr' \\
&\mid\ \epsilon \\
Term &\to Factor\ Term' \\
Term' &\to *\ Factor\ Term' \\
&\mid\ \epsilon \\
Factor &\to\ (\ Expr\ ) \\
&\mid\ a
\end{aligned}
$$

  **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., $Expr \to \epsilon$)?
  **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*.

## 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 $A \to A\alpha \in R$ 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 $A \to B\beta \in R$ for non-terminals $A, B \in V$, $B \overset{*}{\Rightarrow} A\gamma$

| | | |
|---|---|---|
| A | → | Br |
| B | → | Cd |
| C | → | At |

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

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

---

## 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 → βA'ᵢ, A'ᵢ → αA'ᵢ | 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: $\forall j : j < i \bullet \neg(\exists k \bullet k \le j \wedge A_j \to A_k \cdots \in R)$

---

## 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 $B_i$ ($1 \le i \le k$) 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 $A \to \epsilon$ ]

---

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

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

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

• Which are the **nullable** variables?       [S, A, B]

| | | | |
|---|---|---|---|
| S | → | A \| B \| AB | {S → ε not included} |
| A | → | aAA \| aA \| a | {A → aA duplicated} |
| B | → | bBB \| bB \| b | {B → bB duplicated} |

- 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 the **lookhead** of one symbol.
  - We also call such backtrack-free CFG a **predictive grammar**.

---

- 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 **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 |

---

- 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:

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

---

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

```
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
```

## Computing the FIRST Set: Extension

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

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

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

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

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

- More precisely:

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

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

---

## Extended 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 |

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**(+) ∧ *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*) ≜ set of symbols that can appear to the immediate right of a string derived from *v*.

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

---

## The FOLLOW 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 **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

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

```
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)
```

---

## 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

---

## 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:

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

---

## 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 |

- $\epsilon \notin \text{FIRST}(Factor) \Rightarrow \text{START}(Factor) = \text{FIRST}(Factor)$
- $\text{FIRST}(Factor \to \text{name})$     = {name}
- $\text{FIRST}(Factor \to \text{name} \ [ ArgList ])$     = {name}
- $\text{FIRST}(Factor \to \text{name} \ ( ArgList ))$     = {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.

53 of 96

## 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***:
    $$\exists cfg \mid L(cfg) = l \land isBacktrackFree(cfg)$$
- Given a CFG $g = (V, \Sigma, R, S)$, whether or not $g$ is ***backtrack-free*** is ***decidable***:
  For each $A \to \gamma_1 \mid \gamma_2 \mid \ldots \mid \gamma_n \in R$:
  $$\forall i, j : 1 \le i, j \le n \land i \ne j \bullet \textsc{Start}(\gamma_i) \cap \textsc{Start}(\gamma_j) = \varnothing$$

55 of 96

## Left-Factoring: Exercise

- 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 | | \| | $\epsilon$ |

- <u>Identify</u> common prefix `name` and <u>rewrite</u> rules 11, 12, and 13:
  $$\begin{aligned} Factor &\to \texttt{name}\ Arguments \\ Arguments &\to \texttt{[}\ ArgList\ \texttt{]} \\ &\mid \texttt{(}\ ArgList\ \texttt{)} \\ &\mid \epsilon \end{aligned}$$

  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 **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.

56 of 96

- 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.

---

- 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.

---

- **_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_ .

## 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(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 s_i'' then
      trace.push(word)
      trace.push(i)
      word := NextWord()
    else
      fail()
```

---

## 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:

|       | 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 #$j$

---

## 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 |

---

## 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 |

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* |

---

## 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.,

$$[A \to \beta \bullet \gamma,\ \texttt{a}]$$

where:

- A **production rule** $\boxed{A \to \beta\gamma}$ is currently being applied.
- A **terminal symbol** $\boxed{\texttt{a}}$ 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 $\texttt{a}$ 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 $\boxed{LR(1)\ item}$ can denote:

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

- ○ 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**       $[A \to \beta \bullet \gamma,\ \texttt{a}]$

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

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

- ○ Parser has progressed from $[A \to \bullet\beta\gamma,\ \texttt{a}]$ by recognizing $\beta\gamma$.
- ○ $\beta\gamma$ found in a context where an $A$ followed by $\texttt{a}$ would be valid.
- ○ If the current input <u>word</u> matches $\texttt{a}$, then:
  - • Current **complet item** is a $\boxed{handle}$ .
  - • 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 | $Goal \to List$ |
|---|---|
| 2 | $List\ \to List\ Pair$ |
| 3 | $\mid Pair$ |
| 4 | $Pair \to (\ Pair\ )$ |
| 5 | $\mid (\ )$ |

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

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

**Intermediate States:** Subset Construction

**Q.** Derive all $\boxed{LR(1)\ items}$ for the above grammar.

- ○ **FOLLOW**(*List*) = $\{\texttt{eof}, \texttt{(}\}$    **FOLLOW**(*Pair*) = $\{\texttt{eof}, \texttt{(},\texttt{)}\}$
- ○ For each production $A \to \beta$, given **FOLLOW**($A$), $\boxed{LR(1)\ items}$ are:

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

**Q.** Given production $A \to \beta$ (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, (, )}

○ All 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 $\beta$)?

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

---

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.

---

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

| | | |
|---|---|---|
| [*Goal* → • *List*, eof] | | |
| [*Goal* → *List* •, eof] | | |
| [*List* → • *List Pair*, eof] | [*List* → • *List Pair*, (] | |
| [*List* → *List* • *Pair*, eof] | [*List* → *List* • *Pair*, (] | |
| [*List* → *List Pair* •, eof] | [*List* → *List Pair* •, (] | |
| [*List* → • *Pair*, eof] | [*List* → • *Pair*, (] | |
| [*List* → *Pair* •, eof] | [*List* → *Pair* •, (] | |
| [*Pair* → • ( *Pair* ), eof] | [*Pair* → • ( *Pair* ), )] | [*Pair* → • ( *Pair* ), (] |
| [*Pair* → ( • *Pair* ), eof] | [*Pair* → ( • *Pair* ), )] | [*Pair* → ( • *Pair* ), (] |
| [*Pair* → ( *Pair* • ), eof] | [*Pair* → ( *Pair* • ), )] | [*Pair* → ( *Pair* • ), (] |
| [*Pair* → ( *Pair* ) •, eof] | [*Pair* → ( *Pair* ) •, )] | [*Pair* → ( *Pair* ) •, (] |
| [*Pair* → • ( ), eof] | [*Pair* → • ( ), (] | [*Pair* → • ( ), )] |
| [*Pair* → ( • ), eof] | [*Pair* → ( • ), (] | [*Pair* → ( • ), )] |
| [*Pair* → ( ) •, eof] | [*Pair* → ( ) •, (] | [*Pair* → ( ) •, )] |

---

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

Recall:

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

**Initial State**: [*Goal* → •*List*, eof]

**Desired Final State**: [*Goal* → *List*•, eof]

○ The *canonical collection*  [ Example of $\mathcal{CC}$ ]

$$\mathcal{CC} = \{cc_0, cc_1, cc_2, \ldots, cc_n\}$$

denotes the set of **valid subset states** of a **LR(1) parser**.
• Each $cc_i \in \mathcal{CC}$ ($0 \le i \le n$) is a set of **LR(1) items**.
• $\mathcal{CC} \subseteq \mathbb{P}(\textbf{LR(1) items})$    $|\mathcal{CC}|$?    [ $|\mathcal{CC}| \le 2^{|\textbf{LR(1) items}|}$ ]
○ 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*
✓ $\mathcal{CC}$ ≈ **subset** 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**: Given: matching $\gamma$ can reduce to $C$
- **Line 10**: Given: $b \in$ **FIRST**$(\delta a)$ is a valid lookahead symbol after reducing $\gamma$ to $C$
- **Line 11**: Add a new item $[\ C \to •\gamma,\ b]$ into $s$.
- **Line 6**: Termination is guaranteed.
  ∵ Each iteration adds ≥ 1 item to $s$ (otherwise $lastS \neq s$ is *false*).

```
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**: Given: item $[\alpha \to \beta \bullet x\delta,\ a]$ (where $x$ is the next to match)
**Line 8**: Add $[\alpha \to \beta x \bullet \delta,\ a]$ (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 → List
2   List  → List Pair
3         | Pair
4   Pair  → ( Pair )
5         | ( )
```

**Initial State:** $[Goal \to •List,\ \texttt{eof}]$

Calculate $cc_0 = $ ***closure***$(\{ [Goal \to •List,\ \texttt{eof}] \})$.

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

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

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

```
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)
```

Resulting transition table:

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

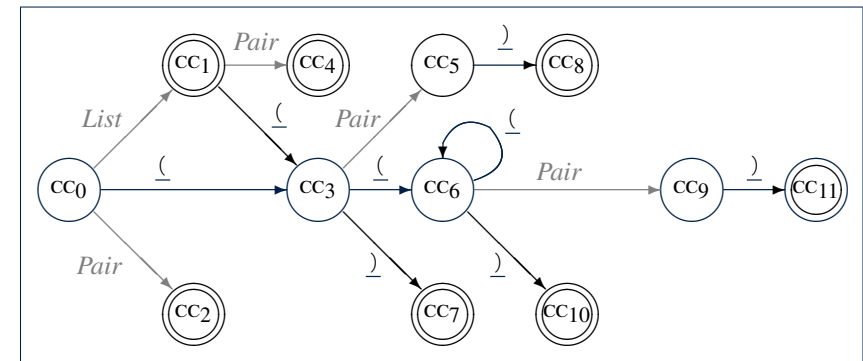| 1 | *Goal* → *List* |
|---|---|
| 2 | *List* → *List Pair* |
| 3 | \| *Pair* |
| 4 | *Pair* → ( *Pair* ) |
| 5 | \| ( ) |

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

Resulting DFA for the parser:

# Constructing $\mathcal{CC}$: The Algorithm (2.4.1)

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

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

$$cc_1 = \left\{ \begin{array}{lll} [Goal \to List\, \bullet, \texttt{eof}] & [List \to List\, \bullet\, Pair, \texttt{eof}] & [List \to List\, \bullet\, Pair, \underline{(}] \\ [Pair \to \bullet\, \underline{(}\ Pair\ \underline{)}, \texttt{eof}] & [Pair \to \bullet\, \underline{(}\ Pair\ \underline{)}, \underline{(}] & [Pair \to \bullet\, \underline{(}\ \underline{)}, \texttt{eof}] \\ & [Pair \to \bullet\, \underline{(}\ \underline{)}, \underline{(}] & \end{array} \right\}$$

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

$$cc_3 = \left\{ \begin{array}{lll} [Pair \to \underline{(}\, \bullet\, Pair\ \underline{)}, \underline{)}] & [Pair \to \underline{(}\, \bullet\, Pair\ \underline{)}, \texttt{eof}] & [Pair \to \underline{(}\, \bullet\, Pair\ \underline{)}, \underline{(}] \\ [Pair \to \underline{(}\, \bullet\, \underline{)}, \underline{)}] & [Pair \to \underline{(}\, \bullet\, \underline{)}, \texttt{eof}] & [Pair \to \underline{(}\, \bullet\, \underline{)}, \underline{(}] \end{array} \right\}$$

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

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

$$cc_6 = \left\{ \begin{array}{ll} [Pair \to \underline{(}\, \bullet\, Pair\ \underline{)}, \underline{)}] & [Pair \to \underline{(}\, \bullet\, Pair\ \underline{)}, \underline{)}] \\ [Pair \to \underline{(}\, \bullet\, \underline{)}, \underline{)}] & [Pair \to \underline{(}\, \bullet\, \underline{)}, \underline{)}] \end{array} \right\}$$

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

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

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

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

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

---

# 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*

---

# Constructing *Action* and *Goto* Tables (2)

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 | | |

---

# BUP: Discovering Ambiguity (1)

| 1 | *Goal* | → | *Stmt* |
|---|---|---|---|
| 2 | *Stmt* | → | 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 $\delta : \mathcal{CC} \times \Sigma \to \mathcal{CC}$

Resulting transition table:

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

Resulting canonical collection $\mathcal{CC}$:

$$cc_8 = \{[Stmt \rightarrow \texttt{assign} \bullet, \{\texttt{eof}, \texttt{else}\}]\}$$

$$cc_9 = \begin{cases} [Stmt \rightarrow \texttt{if expr then } Stmt \texttt{ else} \bullet Stmt, \texttt{eof}], \\ [Stmt \rightarrow \bullet \texttt{ if expr then } Stmt, \texttt{eof}], \\ [Stmt \rightarrow \bullet \texttt{ if expr then } Stmt \texttt{ else } Stmt, \texttt{eof}], \\ [Stmt \rightarrow \bullet \texttt{ assign}, \texttt{eof}] \end{cases}$$

$$cc_{10} = \begin{cases} [Stmt \rightarrow \texttt{if expr} \bullet \texttt{ then } Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \texttt{if expr} \bullet \texttt{ then } Stmt \texttt{ else } Stmt, \{\texttt{eof}, \texttt{else}\}] \end{cases}$$

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

$$cc_{12} = \begin{cases} [Stmt \rightarrow \texttt{if expr then} \bullet Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \texttt{if expr then} \bullet Stmt \texttt{ else } Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \bullet \texttt{if expr then } Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \bullet \texttt{if expr then } Stmt \texttt{ else } Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \bullet \texttt{assign}, \{\texttt{eof}, \texttt{else}\}] \end{cases}$$

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

$$cc_{14} = \begin{cases} [Stmt \rightarrow \texttt{if expr then } Stmt \texttt{ else} \bullet Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \bullet \texttt{ if expr then } Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \bullet \texttt{ if expr then } Stmt \texttt{ else } Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \bullet \texttt{ assign}, \{\texttt{eof}, \texttt{else}\}] \end{cases}$$

Resulting canonical collection $\mathcal{CC}$:

$$cc_0 = \begin{cases} [Goal \rightarrow \bullet Stmt, \texttt{eof}] \quad [Stmt \rightarrow \texttt{if expr then } Stmt, \texttt{eof}] \\ [Stmt \rightarrow \bullet \texttt{assign}, \texttt{eof}] \quad [Stmt \rightarrow \bullet \texttt{if expr then } Stmt \texttt{ else } Stmt, \texttt{eof}] \end{cases}$$

$$cc_1 = \{[Goal \rightarrow Stmt \bullet, \texttt{eof}]\}$$

$$cc_2 = \begin{cases} [Stmt \rightarrow \texttt{if} \bullet \texttt{ expr then } Stmt, \texttt{eof}], \\ [Stmt \rightarrow \texttt{if} \bullet \texttt{ expr then } Stmt \texttt{ else } Stmt, \texttt{eof}] \end{cases}$$

$$cc_3 = \{[Stmt \rightarrow \texttt{assign} \bullet, \texttt{eof}]\}$$

$$cc_4 = \begin{cases} [Stmt \rightarrow \texttt{if expr} \bullet \texttt{ then } Stmt, \texttt{eof}], \\ [Stmt \rightarrow \texttt{if expr} \bullet \texttt{ then } Stmt \texttt{ else } Stmt, \texttt{eof}] \end{cases}$$

$$cc_5 = \begin{cases} [Stmt \rightarrow \texttt{if expr then} \bullet Stmt, \texttt{eof}], \\ [Stmt \rightarrow \texttt{if expr then} \bullet Stmt \texttt{ else } Stmt, \texttt{eof}], \\ [Stmt \rightarrow \bullet \texttt{ if expr then } Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \bullet \texttt{ assign}, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \bullet \texttt{ if expr then } Stmt \texttt{ else } Stmt, \{\texttt{eof}, \texttt{else}\}] \end{cases}$$

$$cc_6 = \begin{cases} [Stmt \rightarrow \texttt{if expr then } Stmt \bullet, \texttt{eof}], \\ [Stmt \rightarrow \texttt{if expr then } Stmt \bullet \texttt{ else } Stmt, \texttt{eof}] \end{cases}$$

$$cc_7 = \begin{cases} [Stmt \rightarrow \texttt{if} \bullet \texttt{ expr then } Stmt, \{\texttt{eof}, \texttt{else}\}], \\ [Stmt \rightarrow \texttt{if} \bullet \texttt{ expr then } Stmt \texttt{ else } Stmt, \{\texttt{eof}, \texttt{else}\}] \end{cases}$$

- Consider $cc_{13}$

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

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

- Consider another scenario:

$$cc_i = \begin{cases} [A \rightarrow \gamma\delta\bullet, \texttt{ a}], \\ [B \rightarrow \gamma\delta\bullet, \texttt{ a}] \end{cases}$$

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