# CS 4300: Compiler Theory

# Chapter 4
# Syntax Analysis

*Dr. Xuejun Liang*

# Outlines (Sections)

1. Introduction
2. Context-Free Grammars
3. Writing a Grammar
4. Top-Down Parsing
5. Bottom-Up Parsing
6. Introduction to LR Parsing: Simple LR
7. More Powerful LR Parsers
8. Using Ambiguous Grammars
9. Parser Generators

# Quick Review of Last Lecture

- Introduction
  - The role of the Parser
  - Many levels of Programming Errors
  - Error Recovery Strategies
  - Representative Grammars
- Context-Free Grammars
  - Derivations and Languages
- Writing a Grammar
  - Lexical Versus Syntactic Analysis
  - Eliminating Ambiguity
  - Eliminating left recursion

# Left Recursion

- A grammar is **left recursive** if it has a nonterminal $A$ such that there is a derivation $A \xRightarrow{+} A\ \alpha$ for some string $\alpha$.

- When a grammar is left recursive then a predictive parser loops forever on certain inputs.

- **Immediate left recursion**, where there is a production of the form $A \rightarrow A\ \alpha$.

$$
\begin{aligned}
A &\rightarrow A\ \alpha \\
&\mid \beta \\
&\mid \gamma
\end{aligned}
\qquad \Longrightarrow \qquad
\begin{aligned}
A &\rightarrow \beta\ R \\
&\mid \gamma\ R \\
R &\rightarrow \alpha\ R \\
&\mid \varepsilon
\end{aligned}
$$

# Algorithm to eliminate left recursion

*Input: Grammar G with no cycles or ε-productions*

Arrange the nonterminals in some order $A_1, A_2, \ldots, A_n$
**for** $i = 1, \ldots, n$ **{**

    **for** $j = 1, \ldots, i$-1 **{**

        replace each

$$A_i \rightarrow A_j\, \gamma$$

        with

$$A_i \rightarrow \delta_1\, \gamma \mid \delta_2\, \gamma \mid \ldots \mid \delta_k\, \gamma$$

        where

$$A_j \rightarrow \delta_1 \mid \delta_2 \mid \ldots \mid \delta_k$$

    **}**

    eliminate the *immediate left recursion* in $A_i$

**}**

# Immediate Left-Recursion Elimination

Rewrite every left-recursive production

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \cdots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \mid \cdots \mid \beta_n$$

into a right-recursive production:

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \cdots \mid \beta_n A'$$
$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \cdots \mid \alpha_m A' \mid \epsilon$$

---

$$A \rightarrow A\,\alpha$$
$$\mid A\,\delta$$
$$\mid \beta$$
$$\mid \gamma$$

$\Longrightarrow$

$$A \rightarrow \beta\,A'$$
$$\mid \gamma\,A'$$
$$A' \rightarrow \alpha\,A'$$
$$\mid \delta\,A'$$
$$\mid \varepsilon$$

# Example Left Recursion Elim.

$$S \rightarrow A\,a \mid b$$
$$A \rightarrow A\,c \mid S\,d \mid \epsilon$$

$\Big\}$ Choose arrangement: $S, A$

$i = 1$: Nothing to do

$i = 2, j = 1$: *Replace S in A $\rightarrow$ S d with A a | b*

$$A \rightarrow A\,c \mid A\,a\,d \mid b\,d \mid \epsilon$$

Eliminate the *immediate left recursion* in $A$

$$S \rightarrow A\,a \mid b$$
$$A \rightarrow b\,d\,A' \mid A'$$
$$A' \rightarrow c\,A' \mid a\,d\,A' \mid \epsilon$$

# Example Left Recursion Elim.

$$A \rightarrow B\,C \mid \mathbf{a}$$
$$B \rightarrow C\,A \mid A\,\mathbf{b}$$
$$C \rightarrow A\,B \mid C\,C \mid \mathbf{a}$$

$\Bigg\}$ Choose arrangement: $A$, $B$, $C$

$i = 1$:      nothing to do

$i = 2, j = 1$:      $B \rightarrow C\,A \mid \underline{A}\,\mathbf{b}$

         $\Rightarrow$      $B \rightarrow C\,A \mid \underline{B\,C}\,\mathbf{b} \mid \underline{\mathbf{a}}\,\mathbf{b}$

         $\Rightarrow_{(\text{imm})}$ $B \rightarrow C\,A\,B_R \mid \mathbf{a}\,\mathbf{b}\,B_R$

               $B_R \rightarrow C\,\mathbf{b}\,B_R \mid \varepsilon$

$i = 3, j = 1$:      $C \rightarrow \underline{A}\,B \mid C\,C \mid \mathbf{a}$

         $\Rightarrow$      $C \rightarrow \underline{B\,C}\,B \mid \underline{\mathbf{a}}\,B \mid C\,C \mid \mathbf{a}$

$i = 3, j = 2$:      $C \rightarrow \underline{B}\,C\,B \mid \mathbf{a}\,B \mid C\,C \mid \mathbf{a}$

         $\Rightarrow$      $C \rightarrow \underline{C\,A\,B_R}\,C\,B \mid \underline{\mathbf{a}\,\mathbf{b}\,B_R}\,C\,B \mid \mathbf{a}\,B \mid C\,C \mid \mathbf{a}$

         $\Rightarrow_{(\text{imm})}$ $C \rightarrow \mathbf{a}\,\mathbf{b}\,B_R\,C\,B\,C_R \mid \mathbf{a}\,B\,C_R \mid \mathbf{a}\,C_R$

               $C_R \rightarrow A\,B_R\,C\,B\,C_R \mid C\,C_R \mid \varepsilon$

# Example Left Recursion Elim.

$$A \rightarrow B\,C \mid \mathbf{a}$$
$$B \rightarrow C\,A \mid A\,\mathbf{b}$$
$$C \rightarrow A\,B \mid C\,C \mid \mathbf{a}$$

$\bigg\}$ Choose arrangement: $A,\,B,\,C$

$i = 1$:      nothing to do

$i = 2, j = 1$:      $B \rightarrow C\,A \mid \underline{A}\,\mathbf{b}$

$\Rightarrow$      $B \rightarrow C\,A \mid \underline{B\,C}\,\mathbf{b} \mid \underline{\mathbf{a}}\,\mathbf{b}$

$\Rightarrow_{\text{(imm)}}$ $B \rightarrow C\,A\,B_R \mid \mathbf{a}\,\mathbf{b}\,B_R$

$B_R \rightarrow C\,\mathbf{b}\,B_R \mid \varepsilon$

$i = 3, j = 1$:      $C \rightarrow \underline{A}\,B \mid C\,C \mid \mathbf{a}$

$\Rightarrow$      $C \rightarrow \underline{B\,C}\,B \mid \underline{\mathbf{a}}\,B \mid C\,C \mid \mathbf{a}$

$i = 3, j = 2$:      $C \rightarrow \underline{B}\,C\,B \mid \mathbf{a}\,B \mid C\,C \mid \mathbf{a}$

$\Rightarrow$      $C \rightarrow \underline{C\,A\,B_R}\,C\,B \mid \underline{\mathbf{a}\,\mathbf{b}\,B_R}\,C\,B \mid \mathbf{a}\,B \mid C\,C \mid \mathbf{a}$

$\Rightarrow_{\text{(imm)}}$ $C \rightarrow \mathbf{a}\,\mathbf{b}\,B_R\,C\,B\,C_R \mid \mathbf{a}\,B\,C_R \mid \mathbf{a}\,C_R$

$C_R \rightarrow A\,B_R\,C\,B\,C_R \mid C\,C_R \mid \varepsilon$

# Left Factoring

- When a nonterminal has two or more productions whose right-hand sides start with the same grammar symbols, the grammar is not LL(1) and cannot be used for predictive parsing

- Replace productions

$$A \rightarrow \alpha\, \beta_1 \mid \alpha\, \beta_2 \mid \dots \mid \alpha\, \beta_n \mid \gamma$$

with

$$A \rightarrow \alpha\, A_R \mid \gamma$$
$$A_R \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

- Example:

$$S \rightarrow i\, E\, t\, S \mid i\, E\, t\, S\, e\, S \mid a \quad \Longrightarrow \quad \begin{aligned} S &\rightarrow i\, E\, t\, S\, S' \mid a \\ S' &\rightarrow e\, S \mid \epsilon \end{aligned}$$

# 4. Top-Down Parsing

- Constructing a parse tree for the input string, starting from the root and creating the nodes of the parse tree in preorder

- Equivalently, finding the leftmost derivation for the input string

Grammar:

$$E \rightarrow T + T$$
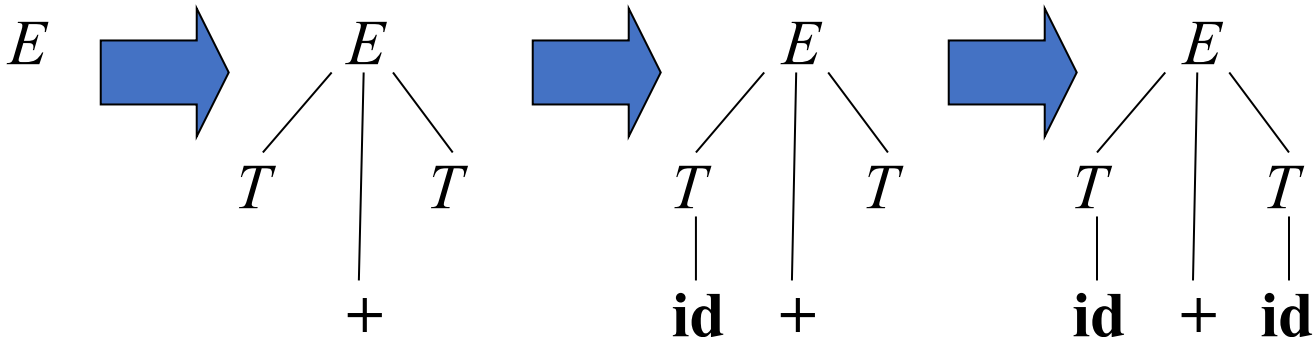$$T \rightarrow ( E )$$
$$T \rightarrow \textbf{-} E$$
$$T \rightarrow \textbf{id}$$

Leftmost derivation:

$$E \Rightarrow_{lm} T + T$$
$$\Rightarrow_{lm} \textbf{id} + T$$
$$\Rightarrow_{lm} \textbf{id} + \textbf{id}$$

# Parsing Methods

- *Universal* (any C-F grammar)
  - Cocke-Younger-Kasimi
  - Earley

- *Top-down* (C-F grammar with restrictions)
  - Recursive descent (predictive parsing)
  - LL (Left-to-right, Leftmost derivation) methods

- *Bottom-up* (C-F grammar with restrictions)
  - Operator precedence parsing
  - LR (Left-to-right, Rightmost derivation) methods
    - SLR, canonical LR, LALR

# Predictive Parsing

- Eliminate left recursion from grammar

- Left factor the grammar

- Compute FIRST and FOLLOW

- Two variants:
  - Recursive (recursive-descent parsing)
  - Non-recursive (table-driven parsing)

- LL(k) class of grammars
  - It can be used to construct predictive parsers looking k symbols ahead in the input.

# FIRST Set

- FIRST($\alpha$) = { *terminals that begin strings*
  *derived from* $\alpha$ }

  FIRST($a$) = {$a$}      if $a \in T$
  FIRST($\varepsilon$) = {$\varepsilon$}
  FIRST($A$) = $\cup_{A \to \alpha}$ FIRST($\alpha$) for $A \to \alpha \in P$
  FIRST($X_1 X_2 ... X_k$) =
      **if** for all $j$ = 1, ..., $i$-1 : $\varepsilon \in$ FIRST($X_j$) **then**
        add non-$\varepsilon$ in FIRST($X_i$) to FIRST($X_1 X_2 ... X_k$)
      **if** for all $j$ = 1, ..., $k$ : $\varepsilon \in$ FIRST($X_j$) **then**
        add $\varepsilon$ to FIRST($X_1 X_2 ... X_k$)

# FOLLOW Set

- FOLLOW(*A*) = { *the set of terminals that can immediately follow nonterminal A* }

FOLLOW(*A*) =
 **for** all (*B* → α *A* β) ∈ *P* **do**
   add FIRST(β)\{ε} to FOLLOW(*A*)
 **for** all (*B* → α *A* β) ∈ *P* and ε ∈ FIRST(β) **do**
   add FOLLOW(*B*) to FOLLOW(*A*)
 **for** all (*B* → α *A*) ∈ *P* **do**
   add FOLLOW(*B*) to FOLLOW(*A*)
 **if** *A* is the start symbol *S* **then**
   add **$** to FOLLOW(*A*)

# Example

$$E \rightarrow T\ E'$$
$$E' \rightarrow + T\ E' \mid \epsilon$$
$$T \rightarrow F\ T'$$
$$T' \rightarrow * F\ T' \mid \epsilon$$
$$F \rightarrow (\ E\ ) \mid \mathbf{id}$$

FIRST(F)
= FIRST( ( E ) ) $\cup$ FIRST(**id**)
= FIRST( ( ) $\cup$ {**id**}
= {**(**} $\cup$ {**id**} = {**(, id**}

FIRST(T)
= FIRST(F) = {**(,** i**d**}

FIRST(E)
= FIRST(T) = {**(, id**}

FIRST(E')
FIRST(+TE') $\cup$ FIRST($\epsilon$)
= {+, **ε**}

FIRST(T')
FIRST(*FT') $\cup$ FIRST($\epsilon$)
= {***, ε**}

# Example

$$
\begin{aligned}
E &\rightarrow T\ E' \\
E' &\rightarrow +\ T\ E' \mid \epsilon \\
T &\rightarrow F\ T' \\
T' &\rightarrow *\ F\ T' \mid \epsilon \\
F &\rightarrow (\ E\ ) \mid \mathbf{id}
\end{aligned}
$$

FIRST(E') = {+, $\boldsymbol{\varepsilon}$}

FIRST(T') = {*, $\boldsymbol{\varepsilon}$}

FOLLOW(E) = {**)**, **\$**}

FOLLOW(E') = FOLLOW(E) = {**)**, **\$**}

FOLLOW(T) = (FIRST(E')\\{$\boldsymbol{\varepsilon}$}) $\cup$ FOLLOW(E) = {+, **)**, **\$**}

FOLLOW(T') = FOLLOW(T) = {+, **)**, **\$**}

FOLLOW(F) = (FIRST(T')\\{$\boldsymbol{\varepsilon}$}) $\cup$ FOLLOW(T) = {+, *, **)**, **\$**}

# LL(1) Grammar

- Predictive parsers, that is, recursive-descent parsers needing no backtracking, can be constructed for a class of grammars called LL(1)

- A grammar *G* is LL(1) if it is not left recursive and for each collection of productions
  $$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$
  for nonterminal *A* the following holds:

  1. $\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \varnothing$ for all $i \neq j$

  2. if $\alpha_i \Rightarrow^* \varepsilon$ then
     2.a.  $\alpha_j \not\Rightarrow^* \varepsilon$ for all $j \neq i$
     2.b.  $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \varnothing$ for all $j \neq i$

# Non-LL(1) Examples

| *Grammar* | *Not LL(1) because:* |
|---|---|
| $S \rightarrow S\,\mathbf{a} \mid \mathbf{a}$ | Left recursive |
| $S \rightarrow \mathbf{a}\,S \mid \mathbf{a}$ | $\mathrm{FIRST}(\mathbf{a}\,S) \cap \mathrm{FIRST}(\mathbf{a}) \neq \varnothing$ |
| $S \rightarrow \mathbf{a}\,R \mid \varepsilon$ <br> $R \rightarrow S \mid \varepsilon$ | For $R$: <br> $S \Rightarrow^* \varepsilon$ and $\varepsilon \Rightarrow^* \varepsilon$ |
| $S \rightarrow \mathbf{a}\,R\,\mathbf{a}$ <br> $R \rightarrow S \mid \varepsilon$ | For $R$: <br> $\mathrm{FIRST}(S) \cap \mathrm{FOLLOW}(R) \neq \varnothing$ |
| $S \rightarrow \mathbf{i}\,E\,\mathbf{t}\,S\,S' \mid \mathbf{a}$ <br> $S' \rightarrow \mathbf{e}\,S \mid \varepsilon$ <br> $E \rightarrow \mathbf{b}$ | For $S'$: <br> $\mathrm{FIRST}(\mathbf{e}\,S) \cap \mathrm{FOLLOW}(S') \neq \varnothing$ |