# CS 4300: Compiler Theory

# Chapter 4
# Syntax Analysis

*Dr. Xuejun Liang*

# Outlines (Sections)

# Quick Review of Last Lecture

- Writing a Grammar
  - Left Recursion Elimination Examples
  - Left Factoring
- Top-Down Parsing
  - FIRST Set, FOLLOW Set and examples
  - LL(1) Grammar and examples
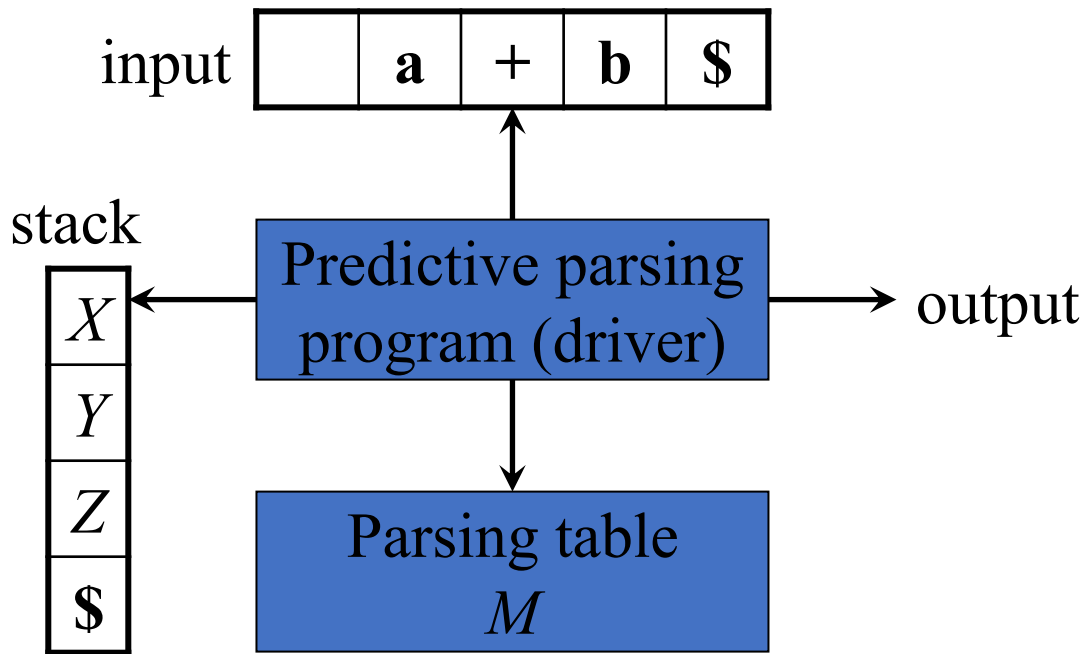
# Using FIRST and FOLLOW in a Recursive-Descent Parser

$expr \rightarrow term\ rest$

$rest \rightarrow + term\ rest$
$\quad\quad | \ - term\ rest$
$\quad\quad | \ \varepsilon$

$term \rightarrow \mathbf{id}$

**procedure** *rest*();
**begin**
  **if** *lookahead* in <u>FIRST(+ *term rest*)</u> **then**
    *match*( '+' ); *term*(); *rest*()
  **else if** *lookahead* in <u>FIRST(- *term rest*)</u> **then**
    *match*( '-' ); *term*(); *rest*()
  **else if** *lookahead* in <u>FOLLOW(*rest*)</u> **then**
    **return**
  **else** error**()**
**end**;

| where | FIRST(+ *term rest*) = { + } |
|---|---|
| | FIRST(- *term rest*) = { - } |
| | FOLLOW(*rest*) = { **$** } |

# Non-Recursive Predictive Parsing: Table-Driven Parsing

- Given an LL(1) grammar $G = (N, T, P, S)$ construct a table $M[A, a]$ for $A \in N, a \in T$ and use a *driver program* with a *stack*
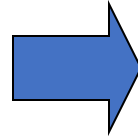
input | | **a** | **+** | **b** | **$** |

stack

| $X$ |
| $Y$ |
| $Z$ |
| **$** |

Predictive parsing program (driver)

output

Parsing table
$M$

# Constructing an LL(1) Predictive Parsing Table

**for** each production $A \rightarrow \alpha$ {
    **for** each $a \in \text{FIRST}(\alpha)$ {
        add $A \rightarrow \alpha$ to $M[A, a]$
    **}**
    **if** $\varepsilon \in \text{FIRST}(\alpha)$ {
        **for** each $b \in \text{FOLLOW}(A)$ {
            add $A \rightarrow \alpha$ to $M[A, b]$
        **}**
    **}**
**}**
Mark each undefined entry in $M$ error

# Example Table

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

| $A \rightarrow \alpha$ | FIRST($\alpha$) | FOLLOW($A$) |
|---|---|---|
| $E \rightarrow T\,E'$ | ( id | $ ) |
| $E' \rightarrow +\,T\,E'$ | + | $ ) |
| $E' \rightarrow \varepsilon$ | $\varepsilon$ | |
| $T \rightarrow F\,T'$ | ( id | + $ ) |
| $T' \rightarrow *\,F\,T'$ | * | + $ ) |
| $T' \rightarrow \varepsilon$ | $\varepsilon$ | |
| $F \rightarrow (\,E\,)$ | ( | * + $ ) |
| $F \rightarrow \mathbf{id}$ | id | |

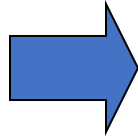| | **id** | + | * | ( | ) | $ |
|---|---|---|---|---|---|---|
| $E$ | $E \rightarrow T\,E'$ | | | $E \rightarrow T\,E'$ | | |
| $E'$ | | $E' \rightarrow +\,T\,E'$ | | | $E' \rightarrow \varepsilon$ | $E' \rightarrow \varepsilon$ |
| $T$ | $T \rightarrow F\,T'$ | | | $T \rightarrow F\,T'$ | | |
| $T'$ | | $T' \rightarrow \varepsilon$ | $T' \rightarrow *\,F\,T'$ | | $T' \rightarrow \varepsilon$ | $T' \rightarrow \varepsilon$ |
| $F$ | $F \rightarrow \mathbf{id}$ | | | $F \rightarrow (\,E\,)$ | | |

# LL(1) Grammars are Unambiguous

Ambiguous grammar

$S \rightarrow \mathbf{i}\, E\, \mathbf{t}\, S\, S' \mid \mathbf{a}$

$S' \rightarrow \mathbf{e}\, S \mid \varepsilon$

$E \rightarrow \mathbf{b}$

| $A \rightarrow \alpha$ | FIRST($\alpha$) | FOLLOW($A$) |
|---|---|---|
| $S \rightarrow \mathbf{i}\, E\, \mathbf{t}\, S\, S'$ | $\mathbf{i}$ | **e $** |
| $S \rightarrow \mathbf{a}$ | $\mathbf{a}$ | |
| $S' \rightarrow \mathbf{e}\, S$ | $\mathbf{e}$ | **e $** |
| $S' \rightarrow \varepsilon$ | $\varepsilon$ | |
| $E \rightarrow \mathbf{b}$ | $\mathbf{b}$ | $\mathbf{t}$ |

Error: duplicate table entry

| | **a** | **b** | **e** | **i** | **t** | **$** |
|---|---|---|---|---|---|---|
| $S$ | $S \rightarrow \mathbf{a}$ | | | $S \rightarrow \mathbf{i}\, E\, \mathbf{t}\, S\, S'$ | | |
| $S'$ | | | $S' \rightarrow \varepsilon$<br>$S' \rightarrow \mathbf{e}\, S$ | | | $S' \rightarrow \varepsilon$ |
| $E$ | | $E \rightarrow \mathbf{b}$ | | | | |

# Predictive Parsing Program (Driver)

read w$ into the input buffer; // w is the input
push($); push(*S*);
a = lookahead;               // a is the first symbol of w
X = pop();
while ( X ≠ $ ) {
    if ( X = a ) {a = lookahead;}        // a is next symbol;
    else if ( X is a terminal ) error();
    else if ( M [X, a] is an error entry ) error();
    else if ( M[X, a] = X → $Y_1 Y_2 \ldots Y_k$ ) {
        output the production X → $Y_1 Y_2 \ldots Y_k$ ;
        push ($Y_k$); push($Y_{k-l}$) , ... , push($Y_1$);
    }
    X = pop();
}

Example: Moves of table-driven parsing on input id + id * id

|  | id | + | * | $ |
|---|---|---|---|---|
| $E$ | $E \to T E'$ |  |  |  |
| $E'$ |  | $E' \to + T E'$ |  | $E' \to \varepsilon$ |
| $T$ | $T \to F T'$ |  |  |  |
| $T'$ |  | $T' \to \varepsilon$ | $T' \to * F T'$ | $T' \to \varepsilon$ |
| $F$ | $F \to \mathbf{id}$ |  |  |  |

| MATCHED | STACK | INPUT | Action |
|---|---|---|---|
|  | $E\$$ | $\mathbf{id + id * id}\$$ |  |
|  | $T E'\$$ | $\mathbf{id + id * id}\$$ | output $E \to T E'$ |
|  | $F T' E'\$$ | $\mathbf{id + id * id}\$$ | output $T \to F T'$ |
|  | $\mathbf{id}\, T' E'\$$ | $\mathbf{id + id * id}\$$ | output $F \to \mathbf{id}$ |
| $\mathbf{id}$ | $T' E'\$$ | $+ \mathbf{id} * \mathbf{id}\$$ | match $\mathbf{id}$ |
| $\mathbf{id}$ | $E'\$$ | $+ \mathbf{id} * \mathbf{id}\$$ | output $T' \to \epsilon$ |
| $\mathbf{id}$ | $+ T E'\$$ | $+ \mathbf{id} * \mathbf{id}\$$ | output $E' \to + T E'$ |
| $\mathbf{id} +$ | $T E'\$$ | $\mathbf{id} * \mathbf{id}\$$ | match $+$ |
| $\mathbf{id} +$ | $F T' E'\$$ | $\mathbf{id} * \mathbf{id}\$$ | output $T \to F T'$ |
| $\mathbf{id} +$ | $\mathbf{id}\, T' E'\$$ | $\mathbf{id} * \mathbf{id}\$$ | output $F \to \mathbf{id}$ |
| $\mathbf{id} + \mathbf{id}$ | $T' E'\$$ | $* \mathbf{id}\$$ | match $\mathbf{id}$ |
| $\mathbf{id} + \mathbf{id}$ | $* F T' E'\$$ | $* \mathbf{id}\$$ | output $T' \to * F T'$ |
| $\mathbf{id} + \mathbf{id} *$ | $F T' E'\$$ | $\mathbf{id}\$$ | match $*$ |
| $\mathbf{id} + \mathbf{id} *$ | $\mathbf{id}\, T' E'\$$ | $\mathbf{id}\$$ | output $F \to \mathbf{id}$ |
| $\mathbf{id} + \mathbf{id} * \mathbf{id}$ | $T' E'\$$ | $\$$ | match $\mathbf{id}$ |
| $\mathbf{id} + \mathbf{id} * \mathbf{id}$ | $E'\$$ | $\$$ | output $T' \to \epsilon$ |
| $\mathbf{id} + \mathbf{id} * \mathbf{id}$ | $\$$ | $\$$ | output $E' \to \epsilon$ |

10

# Panic Mode Recovery

Add synchronizing actions to undefined entries based on FOLLOW

FOLLOW($E$) = { **)** **$** }
FOLLOW($E^{'}$) = { **)** **$** }
FOLLOW($T$) = { **+** **)** **$** }
FOLLOW($T^{'}$) = { **+** **)** **$** }
FOLLOW($F$) = { **+** * **)** **$** }

Example: As $ ∈ Follow(E), M(E, $) = ***synch***

|   | **id** | **+** | ***** | **(** | **)** | **$** |
|---|---|---|---|---|---|---|
| $E$ | $E \rightarrow T\,E^{'}$ |   |   | $E \rightarrow T\,E^{'}$ | ***synch*** | ***synch*** |
| $E^{'}$ |   | $E^{'} \rightarrow + T\,E^{'}$ |   |   | $E^{'} \rightarrow \varepsilon$ | $E^{'} \rightarrow \varepsilon$ |
| $T$ | $T \rightarrow F\,T^{'}$ | ***synch*** |   | $T \rightarrow F\,T^{'}$ | ***synch*** | ***synch*** |
| $T^{'}$ |   | $T^{'} \rightarrow \varepsilon$ | $T^{'} \rightarrow * F\,T^{'}$ |   | $T^{'} \rightarrow \varepsilon$ | $T^{'} \rightarrow \varepsilon$ |
| $F$ | $F \rightarrow \textbf{id}$ | ***synch*** | ***synch*** | $F \rightarrow (\,E\,)$ | ***synch*** | ***synch*** |

The driver pops current nonterminal $A$ if M(A, a) in the above table is ***synch***, or skips input (lookahead) token a if M(A, a) is blank

11

# Example: Moves of parsing and error recovery on the erroneous input + *id* * + *id*

| STACK | INPUT | REMARK |
|---|---|---|
| $E\ \$$ | $+\ \mathbf{id} * + \mathbf{id}\ \$$ | error, skip $+$ |
| $E\ \$$ | $\mathbf{id} * + \mathbf{id}\ \$$ | $\mathbf{id}$ is in FIRST$(E)$ |
| $TE'\ \$$ | $\mathbf{id} * + \mathbf{id}\ \$$ | |
| $FT'E'\ \$$ | $\mathbf{id} * + \mathbf{id}\ \$$ | |
| $\mathbf{id}\ T'E'\$$ | $\mathbf{id} * + \mathbf{id}\ \$$ | |
| $T'E'\ \$$ | $* + \mathbf{id}\ \$$ | |
| $* FT'E'\ \$$ | $* + \mathbf{id}\ \$$ | |
| $FT'E'\ \$$ | $+ \mathbf{id}\ \$$ | error, $M[F, +] = $ synch |
| $T'E'\ \$$ | $+ \mathbf{id}\ \$$ | $F$ has been popped |
| $E'\ \$$ | $+ \mathbf{id}\ \$$ | |
| $+ TE'\ \$$ | $+ \mathbf{id}\ \$$ | |
| $TE'\ \$$ | $\mathbf{id}\ \$$ | |
| $FT'E'\ \$$ | $\mathbf{id}\ \$$ | |
| $\mathbf{id}\ T'E'\ \$$ | $\mathbf{id}\ \$$ | |
| $T'E'\ \$$ | $\$$ | |
| $E'\ \$$ | $\$$ | |
| $\$$ | $\$$ | |

| | **id** | **+** | **\*** | **$** |
|---|---|---|---|---|
| $E$ | $E \rightarrow T E'$ | | | *synch* |
| $E'$ | | $E' \rightarrow + T E'$ | | $E' \rightarrow \varepsilon$ |
| $T$ | $T \rightarrow F T'$ | *synch* | | *synch* |
| $T'$ | | $T' \rightarrow \varepsilon$ | $T' \rightarrow * F T'$ | $T' \rightarrow \varepsilon$ |
| $F$ | $F \rightarrow \mathbf{id}$ | *synch* | *synch* | *synch* |

# Phrase-Level Recovery

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

Change input stream by inserting missing tokens
For example: **id id** is changed into **id * id**

Pro:    Can be fully automated
Cons:  Recovery not always intuitive

Can then continue here

| | **id** | **+** | ***** | **(** | **)** | **$** |
|---|---|---|---|---|---|---|
| $E$ | $E \rightarrow T\,E'$ | | | $E \rightarrow T\,E'$ | *synch* | *synch* |
| $E'$ | | $E' \rightarrow +\,T\,E'$ | | | $E' \rightarrow \varepsilon$ | $E' \rightarrow \varepsilon$ |
| $T$ | $T \rightarrow F\,T'$ | *synch* | | $T \rightarrow F\,T'$ | *synch* | *synch* |
| $T'$ | *insert **** | $T' \rightarrow \varepsilon$ | $T' \rightarrow *\,F\,T'$ | | $T' \rightarrow \varepsilon$ | $T' \rightarrow \varepsilon$ |
| $F$ | $F \rightarrow \mathbf{id}$ | *synch* | *synch* | $F \rightarrow (\,E\,)$ | *synch* | *synch* |

***insert*** *****: driver inserts missing * and retries the production

13

# 5. Bottom-Up Parsing

- LR methods (Left-to-right, Rightmost derivation)
  - SLR, Canonical LR, LALR

- Other special cases:
  - Shift-reduce parsing
  - Operator-precedence parsing

# Shift-Reduce Parsing

Grammar:

$S \rightarrow \mathbf{a}\ A\ B\ \mathbf{e}$

$A \rightarrow A\ \mathbf{b}\ \mathbf{c} \mid \mathbf{b}$

$B \rightarrow \mathbf{d}$

These match production's right-hand sides

Reducing a sentence:

$\mathbf{a}\ \underline{\mathbf{b}}\ \mathbf{b}\ \mathbf{c}\ \mathbf{d}\ \mathbf{e}$

$\mathbf{a}\ \underline{A\ \mathbf{b}\ \mathbf{c}}\ \mathbf{d}\ \mathbf{e}$

$\mathbf{a}\ A\ \underline{\mathbf{d}}\ \mathbf{e}$

$\underline{\mathbf{a}\ A\ B\ \mathbf{e}}$

$S$

Shift-reduce corresponds to the reverse of a rightmost derivation:

$S \Rightarrow_{rm} \mathbf{a}\ A\ B\ \mathbf{e}$

$\Rightarrow_{rm} \mathbf{a}\ A\ \mathbf{d}\ \mathbf{e}$

$\Rightarrow_{rm} \mathbf{a}\ A\ \mathbf{b}\ \mathbf{c}\ \mathbf{d}\ \mathbf{e}$

$\Rightarrow_{rm} \mathbf{a}\ \mathbf{b}\ \mathbf{b}\ \mathbf{c}\ \mathbf{d}\ \mathbf{e}$



15

# Handles

A handle is a substring that matches the body of a production, and whose reduction represents one step along the reverse of a rightmost derivation

Grammar:

$S \rightarrow \mathbf{a}\ A\ B\ \mathbf{e}$

$A \rightarrow A\ \mathbf{b}\ \mathbf{c}\ |\ \mathbf{b}$

$B \rightarrow \mathbf{d}$

**a** $\underline{\mathbf{b}}$ **b c d e**

**a** $\underline{A}$ **b c d e**

**a** $A$ $\underline{\mathbf{d}}$ **e**

$\underline{\mathbf{a}\ A\ B\ \mathbf{e}}$

$S$

*Handle*

**a** $\underline{\mathbf{b}}$ **b c d e**

**a** $A$ $\underline{\mathbf{b}}$ **c d e**

**a** $A\ A$ **e**

… ?

NOT a handle, because further reductions will fail (result is not a sentential form)