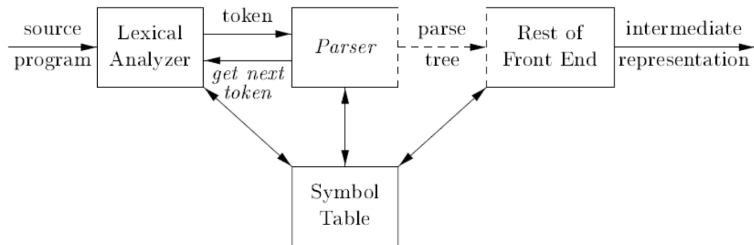# Front End: Syntax Analysis

# The Role of the Parser



Figure 4.1: Position of parser in compiler model

# The Role of the Parser

- Construct a parse tree
- Report and recover from errors
- Collect information into symbol tables

# Types of Parsers

- There are three general types of parsers for grammars:
  - Universal
  - Top-down
  - Bottom-up

- In compilers, the methods commonly used are either top-down or bottom-up.

- One input symbol at a time, from left to right.

- Efficiency is achieved by restricting to particular grammars: **LL** (manually) or **LR** (automated tools).

# Grammars for expressions

- **Universal** methods are suitable for general grammars, e.g.

$$E \quad \rightarrow \quad E + E \mid E * E \mid (E) \mid \mathbf{id}$$

  (no associativity, no precedence captured)

- **Bottom-up** methods: **LR** grammars, e.g.

$$E \quad \rightarrow \quad E + T \mid T$$
$$T \quad \rightarrow \quad T * F \mid F$$
$$F \quad \rightarrow \quad (E) \mid \mathbf{id}$$

  (associativity and precedence captured)

- **Top-down** methods: **LL** grammars, e.g.

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

# Context-free Grammars

A *Context-free grammar* (or *grammar*) systematically describes the syntax of programming language constructs.

$$
\begin{aligned}
expression &\rightarrow expression \text{ + } term \\
expression &\rightarrow expression \text{ - } term \\
expression &\rightarrow term \\
term &\rightarrow term \text{ * } factor \\
term &\rightarrow term \text{ / } factor \\
term &\rightarrow factor \\
factor &\rightarrow \text{( } expression \text{ )} \\
factor &\rightarrow \textbf{id}
\end{aligned}
$$

Figure 4.2: Grammar for simple arithmetic expressions

Terminal symbols: **id** + - * / ( )   Non-terminal: *expression*, *term*, *factor*.   Start symbol: *expression*

# CFG: Formal Definition

$$G = (T, N, P, S)$$

- T is a finite set of terminals
- N is a finite set of non-terminals
- P is a finite subset of production rules of the form
  - $A \rightarrow \alpha_1 \alpha_2 \ldots \alpha_k$ with $A \in N$, $\alpha_i \in T \cup N$
- S is the start symbol
  - $S \in N$

## Derivations

Using notational conventions the grammar in Fig.4.2 becomes

$$
\begin{aligned}
E &\rightarrow E + T \mid T \\
T &\rightarrow T * F \mid F \\
F &\rightarrow (E) \mid \mathbf{id}
\end{aligned}
$$

A derivation of a string of terminals in this grammar is a proof that the string is an expression.

*Leftmost* derivation: always choose the leftmost nonterminal

$$E \Rightarrow^{lm} E + T \Rightarrow^{lm} \mathbf{id} + T \Rightarrow^{lm} \mathbf{id} + F \Rightarrow^{lm} \mathbf{id} + \mathbf{id}$$

*Rightmost* derivation: always choose the righttmost nonterminal

$$E \Rightarrow^{rm} E+T \Rightarrow^{rm} E+F \Rightarrow^{rm} E+\mathbf{id} \Rightarrow^{rm} T+\mathbf{id} \Rightarrow^{rm} F+\mathbf{id} \Rightarrow^{rm} \mathbf{id}+\mathbf{id}$$

# Parse Trees

A parse tree is a graphical representation of a derivation: an interior node represents the head of a production; its children are labelled by the symbols in the body.

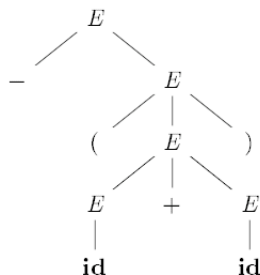$$E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid \textbf{id}$$



Figure 4.3: Parse tree for $-(\textbf{id} + \textbf{id})$

# Example

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E + E) \Rightarrow -(\mathbf{id} + E) \Rightarrow -(\mathbf{id} + \mathbf{id}) \qquad (4.8)$$
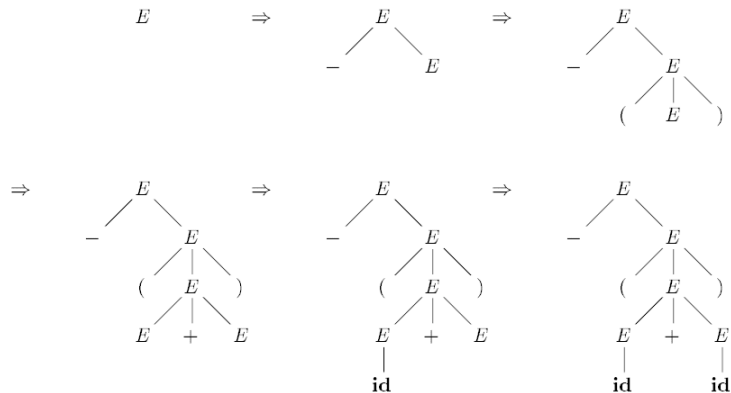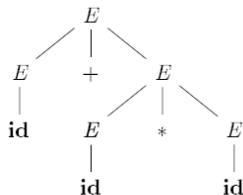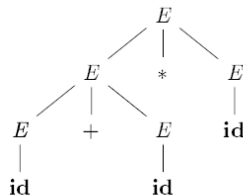


Figure 4.4: Sequence of parse trees for derivation (4.8)

# Ambiguity

A grammar that produces more than one parse tree for some sentence is called ambiguous.



Figure 4.5: Two parse trees for **id+id*id**

Problems: (1) Ambiguity can make parsing difficult; (2) Underlying structure is ill-defined.

# Language Generated by a Grammar

A grammar $G$ generates a language $L$ if we can show that:

- Every string generated by $G$ is in $L$, and
- Every string in $L$ can be generated by $G$.

**Example**: Show that the grammar

$$S \rightarrow (S)S \mid \varepsilon$$

generates all strings of balanced parentheses and only such strings.

# Grammars vs Regular Expressions

Every regular language is a context-free language but non vice-versa.

**Example:** The language generated by the regular expression

$$(a|b)^* abb$$

is equivalent to the grammar

$$
\begin{aligned}
A_0 &\rightarrow aA_0 \mid bA_0 \mid aA_1 \\
A_1 &\rightarrow bA_2 \\
A_2 &\rightarrow bA_3 \\
A_3 &\rightarrow \varepsilon
\end{aligned}
$$

# NFA-based Construction

From the NFA for the regular expression,

- For each state $i$ of the NFA, create a nonterminal $A_i$
- Add production $A_i \rightarrow aA_j$ for each transition from $i$ to $j$ on $a$
- If $i$ is accepting then add $A_i \rightarrow \varepsilon$
- If $i$ is the starting state, make $A_i$ the start symbol of the grammar.

# Grammar with no Corresponding Regular Expression

The language
$$L = \{a^n b^n \mid n \geq 1\}$$

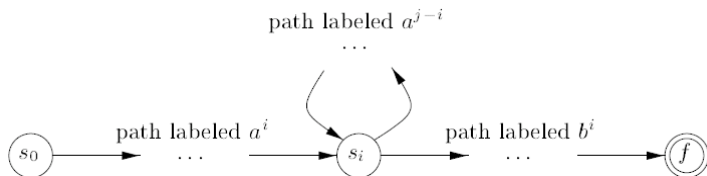can be described by a grammar but not by a regular expression. Why?



path labeled $a^{j-i}$

$\cdots$

path labeled $a^i$

$s_0$ $\quad\cdots\quad$ $s_i$

path labeled $b^i$

$\cdots$

$f$

Figure 4.6: DFA $D$ accepting both $a^i b^i$ and $a^j b^i$.

# Non-Context-Free Grammars

Grammars alone can be not sufficient to specify some programming language construct.

This happens for constructs that are *context-dependent*.

The language

$$L_1 = \{wcw \mid w \text{ in } (\mathbf{a}|\mathbf{b})^*\}$$

is non-context-free. $L_1$ abstracts the requirements that identifiers are defined before their use (as in *C* and *Java*).

$$L_2 = \{a^n b^m c^n d^m \mid n \geq 0, m \geq 0\}$$

is non-context-free. $L_2$ abstracts the requirements that the number of formal parameters in a function declaration is the same as the number of actual parameters in a use of the function.

# Common Grammars Problems (CGP)

A grammar may have some 'bad' styles or ambiguity. Some CGP are:

- Ambiguity
- Left-recursion
- Left factors

We need to transform a grammar $G_1$ into a grammar $G_2$ with no CGP and such that $G_1$ and $G_2$ are equivalent, i.e. they define the same language.

# Eliminating Ambiguity

Consider the grammar:

$$stmt \quad \rightarrow \quad \textbf{if} \text{ expr } \textbf{then} \text{ stmt}$$
$$| \quad \textbf{if} \text{ expr } \textbf{then} \text{ stmt } \textbf{else} \text{ stmt}$$
$$| \quad \textbf{other}$$

The sentence

**if** E1 **then if** E2 **then** S1 **else** S2

is ambiguous (cf. Figure 4.9).

$$
\begin{aligned}
stmt \quad &\rightarrow \quad matched\_stmt \\
&| \quad open\_stmt \\
matched\_stmt \quad &\rightarrow \quad \textbf{if } expr \textbf{ then } matched\_stmt \textbf{ else } matched\_stmt \\
&| \quad \textbf{other} \\
open\_stmt \quad &\rightarrow \quad \textbf{if } expr \textbf{ then } stmt \\
&| \quad \textbf{if } expr \textbf{ then } matched\_stmt \textbf{ else } open\_stmt
\end{aligned}
$$

Figure 4.10: Unambiguous grammar for if-then-else statements

# Example



Figure 4.9: Two parse trees for an ambiguous sentence

# CGP: Left Recursion

Top-down parsing cannot handle left-recursive grammars.

We need to eliminate left recursion.

# Eliminating Left Recursion

Consider a grammar $G$ with a production
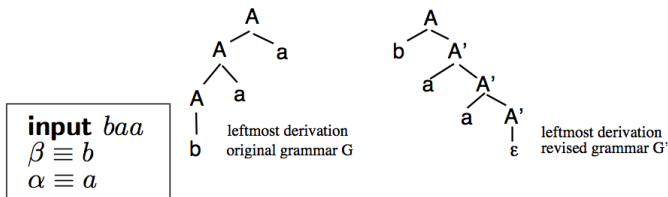
$$A \rightarrow A\alpha \mid \beta,$$

where $\beta$ does not start with $A$.

Transform $G$ in $G'$ by replacing it by

$$A \rightarrow \beta A'$$
$$A' \rightarrow \alpha A' \mid \varepsilon.$$

$G$ and $G'$ are equivalent: $L(G) = L(G')$.



**input** $baa$
$\beta \equiv b$
$\alpha \equiv a$

leftmost derivation
original grammar G

leftmost derivation
revised grammar G'

# The Grammar Expression Example

The non-left-recursive expression grammar

$$\begin{aligned}
E &\rightarrow TE' \\
E' &\rightarrow +TE' \mid \varepsilon \\
T &\rightarrow FT' \\
T' &\rightarrow *FT' \mid \varepsilon \\
F &\rightarrow (E) \mid \textbf{id}
\end{aligned}$$

is obtained by eliminating immediate left recursion from the expression grammar

$$\begin{aligned}
E &\rightarrow E + T \mid T \\
T &\rightarrow T * F \mid F \\
F &\rightarrow (E) \mid \textbf{id}
\end{aligned}$$

by applying the above transformation.

# Algorithm for Eliminating Left Recursion

**Input**: A grammar $G$ with no cycles and no $\varepsilon$-productions.
**Output**: An equivalent grammar with no left recursion..

1)     arrange the nonterminals in some order $A_1, A_2, \ldots, A_n$.
2)     **for** ( each $i$ from 1 to $n$ ) {
3)            **for** ( each $j$ from 1 to $i - 1$ ) {
4)                 replace each production of the form $A_i \rightarrow A_j \gamma$ by the
                         productions $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \cdots \mid \delta_k \gamma$, where
                         $A_j \rightarrow \delta_1 \mid \delta_2 \mid \cdots \mid \delta_k$ are all current $A_j$-productions
5)            }
6)            eliminate the immediate left recursion among the $A_i$-productions
7)     }

Figure 4.11: Algorithm to eliminate left recursion from a grammar

# Applying the Algorithm

**for** $i = 1$ **to** $n$ **do**
- **for** $j = 1$ **to** $i - 1$ **do**
  - ▷ *replace* $A_i \to A_j \gamma$
    *with* $A_i \to \delta_1 \gamma \mid \cdots \mid \delta_k \gamma$
    *where* $A_j \to \delta_1 \mid \cdots \mid \delta_k$ *are all the current* $A_j$*-productions.*
- **Eliminate immediate left-recursion for** $A_i$
  - ▷ *New nonterminals generated above are numbered* $A_{i+n}$

- **Original Grammar:**
  - **(1)** $S \to Aa \mid b$
  - **(2)** $A \to Ac \mid Sd \mid e$
- **Ordering of nonterminals:** $S \equiv A_1$ **and** $A \equiv A_2$.
- $i = 1$
  - **do nothing as there is no immediate left-recursion for** $S$
- $i = 2$
  - **replace** $A \to Sd$ **by** $A \to Aad \mid bd$
  - **hence (2) becomes** $A \to Ac \mid Aad \mid bd \mid e$
  - **after removing immediate left-recursion:**
    - ▷ $A \to bdA' \mid eA'$
    - ▷ $A' \to cA' \mid adA' \mid \epsilon$
- **Resulting grammar:**
  - ▷ $S \to Aa \mid b$
  - ▷ $A \to bdA' \mid eA'$
  - ▷ $A' \to cA' \mid adA' \mid \epsilon$

# CGP: Left Factor

The *left factor* problem occurs when for some nonterminal $A$ there are $A$- productions whose bodies have a common prefix.
**Example**

$$stmt \rightarrow \textbf{if } expr \textbf{ then } stmt \textbf{ else } stmt$$
$$\mid \quad \textbf{if } expr \textbf{ then } stmt$$

On input **if**, we have no way to decide which production to choose.

Idea: Expand with the full common factor!

# Eliminating Left Factors

The algorithm below produces on input $G$ an equivalent left-factored $G'$.

**Input: context free grammar $G$**

**Output: equivalent** left-factored **context-free grammar $G'$**

**for each nonterminal $A$ do**
- find the longest non-$\epsilon$ prefix $\alpha$ that is common to right-hand sides of two or more productions;
- replace
  - ▷ $A \rightarrow \alpha\beta_1 \mid \cdots \mid \alpha\beta_n \mid \gamma_1 \mid \cdots \mid \gamma_m$

  **with**
  - ▷ $A \rightarrow \alpha A' \mid \gamma_1 \mid \cdots \mid \gamma_m$
  - ▷ $A' \rightarrow \beta_1 \mid \cdots \mid \beta_n$
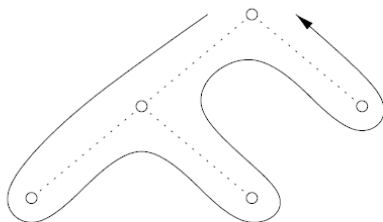- repeat the above step until the grammar has no two productions with a common prefix;

# Top-down Parsing

Constructing a parse tree for the input string starting from the root in a depth-first manner (leftmost derivation).

**procedure** *visit*(node $N$) {
    **for** ( each child $C$ of $N$, from left to right ) {
        *visit*($C$);
    }
    evaluate semantic rules at node $N$;
}

Figure 2.11: A depth-first traversal of a tree

# Example

Given the grammar

$$\begin{aligned}
E &\rightarrow TE' \\
E' &\rightarrow +TE' \mid \varepsilon \\
T &\rightarrow FT' \\
T' &\rightarrow *FT' \mid \varepsilon \\
F &\rightarrow (E) \mid \textbf{id}
\end{aligned}$$

the sequence of trees given in the next slide corresponds to a
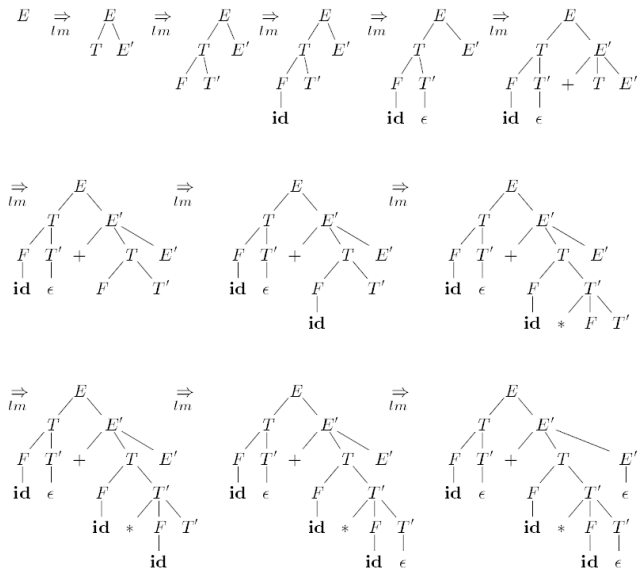leftmost derivation of the input string $\textbf{id} + \textbf{id} * \textbf{id}$.

# Example (ctdn.)



Figure 4.12: Top-down parse for **id** + **id** * **id**

# Recursive-descent Parsing

A recursive-descent parsing program is a set of procedures, one for each nonterminal, of the form:

```
      void A() {
1)          Choose an A-production, A → X₁X₂···Xₖ;
2)          for ( i = 1 to k ) {
3)                  if ( Xᵢ is a nonterminal )
4)                          call procedure Xᵢ();
5)                  else if ( Xᵢ equals the current input symbol a )
6)                          advance the input to the next symbol;
7)                  else /* an error has occurred */;
            }
      }
```

Figure 4.13: A typical procedure for a nonterminal in a top-down parser
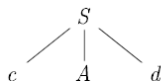
# Backtracking

Top-down parsing may require repeated scans over the input: if an *A*-production leads to a failure, we must *backtrack* and try with another one.
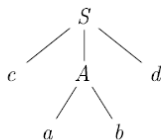
**Example**

$$S \rightarrow cAd$$
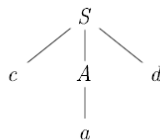$$A \rightarrow ab \mid a$$

On input $w = cad$ we apply recursive-descent parsing. Since the choice of the first production leads to failure, we backtrack and try the second.



(a)   (b)   (c)

# Predictive Parsing

The previous approach may be very inefficient due to backtracking.
A predictive parser is a recursive-descent parser needing no backtracking.

A predictive parser can choose one of the available productions for a nonterminal $A$ by looking at the next input symbol(s).

The class of **LL(1)** grammars [Lewis&Stearns 1968] can parsed by a predictive parsers in $O(n)$ time.

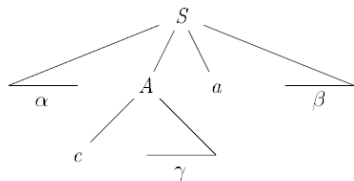We first need to introduce two important functions:
FIRST and FOLLOW.



Figure 4.15: Terminal $c$ is in FIRST($A$) and $a$ is in FOLLOW($A$)

# Computing FIRST

To compute FIRST($X$) for any symbol $X$, apply the rules:

1. If X is a terminal, then FIRST(X) = {X}.
2. if X → ε is a production then place ε in FIRST(X)
3. If X is a nonterminal and X → $Y_1 Y_2$ ... $Y_k$ is a production for some k ≥ 1, then place a in FIRST(X) if for some i, a is in FIRST($Y_i$), and ε is in all of FIRST($Y_i$), ... ,FIRST($Y_{i-1}$); that is, $Y_1$ ...$Y_{i-1}$ ⇒* ε. If ε is in FIRST($Y_j$) for all j = 1,2, ... ,k, then add ε to FIRST(X).

# Computing FIRST (ctd.)

**Let $\alpha = X_1 X_2 \cdots X_n$. Perform the following steps in sequence:**
- **FIRST$(\alpha) \Leftarrow$ FIRST$(X_1) - \{\epsilon\}$;**
- **if $\epsilon \in$ FIRST$(X_1)$, then**
  - ▷ *put FIRST$(X_2) - \{\epsilon\}$ into FIRST$(\alpha)$;*
- **if $\epsilon \in$ FIRST$(X_1) \cap$ FIRST$(X_2)$, then**
  - ▷ *put FIRST$(X_3) - \{\epsilon\}$ into FIRST$(\alpha)$;*
- **$\cdots$**
- **if $\epsilon \in \cap_{i=1}^{n-1}$FIRST$(X_i)$, then**
  - ▷ *put FIRST$(X_n) - \{\epsilon\}$ into FIRST$(\alpha)$;*
- **if $\epsilon \in \cap_{i=1}^{n}$FIRST$(X_i)$, then**
  - ▷ *put $\{\epsilon\}$ into FIRST$(\alpha)$.*

# Computing FOLLOW

To compute Follow($X$) for all nonterminal $X$, apply the following rules until nothing can be added to any FOLLOW set.

1. Place $ in FOLLOW(S), (S start symbol, $ the input right endmarker).

2. If there is a production A → αB or a production A → αBβ where FIRST(β) contains ε then everything in FOLLOW(A) is in FOLLOW(B).

3. If there is a production A → αBβ then everything in in FIRST(β) except ε is in FOLLOW(B).

# *FIRST and FOLLOW Example*

```
E  → T E'
E' → + T E'  |  ε
T  → F T'
T' → * F T'  |  ε
F  → ( E )  |  id
```

## *Computing FOLLOW(A)*

- Place $ into FOLLOW(S)
- Repeat until nothing changes:
  - if A → αBβ then add FIRST(β)\{ε} to FOLLOW(B)
  - if A → αB then add FOLLOW(A) to FOLLOW(B)
  - if A → αBβ and ε is in FIRST(β) then add FOLLOW(A) to FOLLOW(B)

- FIRST(F) = FIRST(T) = FIRST(E) = {(, id }
- FIRST(E') = {+, ε}
- FIRST(T') = {*, ε}
- FOLLOW(E) = FOLLOW(E') = {), $}
- FOLLOW(T) = FOLLOW(T') = {+,),$}
- FOLLOW(F) = {+, *, ), $}

# **LL(1)** Grammars

**L**eft to right parsers producing a **L**eftmost derivation *looking* ahead by at most **1** input symbol.

## Definition

A grammar $G$ is **LL(1)** if and only if whenever $A \rightarrow \alpha \mid \beta$ are two distinct productions in $G$, then

- FIRST($\alpha$) and FIRST($\beta$) are disjoint sets
- If $\varepsilon$ is in FIRST($\beta$) then FIRST($\alpha$) and FOLLOW(A) are disjoint sets
- If $\varepsilon$ is in FIRST($\alpha$) then FIRST($\beta$) and FOLLOW(A) are disjoint sets.

Most programming language constructs are **LL(1)** but careful grammar writing is required.

If a grammar is **LL(1)** then it does not have CGP, but the vice-versa does not hold.

# Predictive Parsing Table

To construct a parsing table $M$ for a grammar $G$, for each production $A \to \alpha$ in G:

- If $a$ is in FIRST($a$), add $A \to \alpha$ in $M[A, a]$.
- If $\varepsilon$ is in FIRST($\alpha$), add $A \to \alpha$ in $M[A, b]$ for each $b$ in FOLLOW(A).
- If $\varepsilon$ is in FIRST($\alpha$) and \$ is in FOLLOW(A), add $A \to \alpha$ in $M[A, \$]$.

# Example

For the expression grammar the algorithm produces the following table.

| NON - TERMINAL | INPUT SYMBOL | | | | | |
|---|---|---|---|---|---|---|
| | **id** | + | * | ( | ) | $ |
| $E$ | $E \to TE'$ | | | $E \to TE'$ | | |
| $E'$ | | $E' \to +TE'$ | | | $E' \to \epsilon$ | $E' \to \epsilon$ |
| $T$ | $T \to FT'$ | | | $T \to FT'$ | | |
| $T'$ | | $T' \to \epsilon$ | $T' \to *FT'$ | | $T' \to \epsilon$ | $T' \to \epsilon$ |
| $F$ | $F \to \mathbf{id}$ | | | $F \to (E)$ | | |

Figure 4.17: Parsing table $M$ for Example 4.32

# Stack-based Predictive Parser



Figure 4.19: Model of a table-driven predictive parser

**let** $a$ be the first symbol of $w$;
**let** $X$ be the top stack symbol;
**while** ( $X \neq \$ $ ) { /* stack is not empty */
    **if** ( $X = a$ ) pop the stack and **let** $a$ be the next symbol of $w$;
    **else if** ( $X$ is a terminal ) *error*();
    **else if** ( $M[X, a]$ is an error entry ) *error*();
    **else if** ( $M[X, a] = X \rightarrow Y_1 Y_2 \cdots Y_k$ ) {
        output the production $X \rightarrow Y_1 Y_2 \cdots Y_k$;
        pop the stack;
        push $Y_k, Y_{k-1}, \ldots, Y_1$ onto the stack, with $Y_1$ on top;
    }
    **let** $X$ be the top stack symbol;
}

Figure 4.20: Predictive parsing algorithm

# Example

| MATCHED | STACK | INPUT | ACTION |
|---------|-------|-------|--------|
| | $E\$$ | $\mathbf{id} + \mathbf{id} * \mathbf{id}\$$ | |
| | $TE'\$$ | $\mathbf{id} + \mathbf{id} * \mathbf{id}\$$ | output $E \rightarrow TE'$ |
| | $FT'E'\$$ | $\mathbf{id} + \mathbf{id} * \mathbf{id}\$$ | output $T \rightarrow FT'$ |
| | $\mathbf{id}\,T'E'\$$ | $\mathbf{id} + \mathbf{id} * \mathbf{id}\$$ | output $F \rightarrow \mathbf{id}$ |
| $\mathbf{id}$ | $T'E'\$$ | $+ \mathbf{id} * \mathbf{id}\$$ | match $\mathbf{id}$ |
| $\mathbf{id}$ | $E'\$$ | $+ \mathbf{id} * \mathbf{id}\$$ | output $T' \rightarrow \epsilon$ |
| $\mathbf{id}$ | $+ TE'\$$ | $+ \mathbf{id} * \mathbf{id}\$$ | output $E' \rightarrow + TE'$ |
| $\mathbf{id} +$ | $TE'\$$ | $\mathbf{id} * \mathbf{id}\$$ | match $+$ |
| $\mathbf{id} +$ | $FT'E'\$$ | $\mathbf{id} * \mathbf{id}\$$ | output $T \rightarrow FT'$ |
| $\mathbf{id} +$ | $\mathbf{id}\,T'E'\$$ | $\mathbf{id} * \mathbf{id}\$$ | output $F \rightarrow \mathbf{id}$ |
| $\mathbf{id} + \mathbf{id}$ | $T'E'\$$ | $* \mathbf{id}\$$ | match $\mathbf{id}$ |
| $\mathbf{id} + \mathbf{id}$ | $* FT'E'\$$ | $* \mathbf{id}\$$ | output $T' \rightarrow * FT'$ |
| $\mathbf{id} + \mathbf{id} *$ | $FT'E'\$$ | $\mathbf{id}\$$ | match $*$ |
| $\mathbf{id} + \mathbf{id} *$ | $\mathbf{id}\,T'E'\$$ | $\mathbf{id}\$$ | output $F \rightarrow \mathbf{id}$ |
| $\mathbf{id} + \mathbf{id} * \mathbf{id}$ | $T'E'\$$ | $\$$ | match $\mathbf{id}$ |
| $\mathbf{id} + \mathbf{id} * \mathbf{id}$ | $E'\$$ | $\$$ | output $T' \rightarrow \epsilon$ |
| $\mathbf{id} + \mathbf{id} * \mathbf{id}$ | $\$$ | $\$$ | output $E' \rightarrow \epsilon$ |

Figure 4.21: Moves made by a predictive parser on input $\mathbf{id} + \mathbf{id} * \mathbf{id}$