# CS4212 – Compiler Design

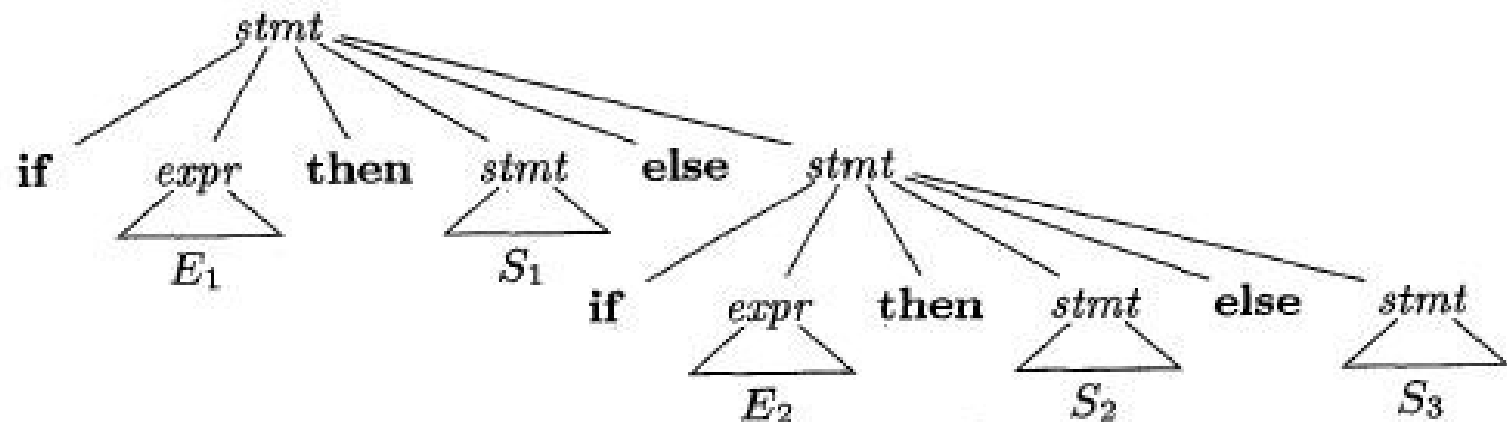# Syntactic Analysis 2

# Outline

- Grammar transformations

- Top-down parsing

  - Recursive-descent parsing

  - Predictive parsing

  - LL(1) grammars

- Bottom-up parsing

  - Shift-reduce parsing

  - LR-parsing

- Yacc

- Textbook: "Compilers: Principles, Techniques, and Practices", Aho, Lam, Sethi, Ullman

  - Chapter 4: 4.3-4.7,4.9

# Eliminating Ambiguity

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

**if** $E_1$ **then** $S_1$ **else if** $E_2$ **then** $S_2$ **else** $S_3$

# Eliminating Ambiguity

$$\textbf{if } E_1 \textbf{ then if } E_2 \textbf{ then } S_1 \textbf{ else } S_2$$

has the two parse trees

# Elimination of Ambiguity

Unamibiguous grammar

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

# Elimination of Left Recursion

Immediate left recursion can be eliminated by the following technique, which works for any number of $A$-productions. First, group the productions as

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

where no $\beta_i$ begins with an $A$. Then, replace the $A$-productions by

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

The nonterminal $A$ generates the same strings as before but is no longer left recursive.

# Elimination of Left Recursion

**Algorithm**        Eliminating left recursion.

**INPUT**: Grammar $G$ with no cycles or $\epsilon$-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 $A_i \to A_j\gamma$ by the productions $A_i \to \delta_1\gamma \mid \delta_2\gamma \mid \cdots \mid \delta_k\gamma$, where $A_j \to \delta_1 \mid \delta_2 \mid \cdots \mid \delta_k$ are all current $A_j$-productions
5)         }
6)     eliminate the immediate left recursion among the $A_i$-productions
7)   }

# Elimination of Left Recursion

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

Eliminating the immediate left recursion among these $A$-productions yields the following grammar.

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

# Left Factoring

Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive, or top-down, parsing. When the choice between two alternative $A$-productions is not clear, we may be able to rewrite the productions to defer the decision until enough of the input has been seen that we can make the right choice.
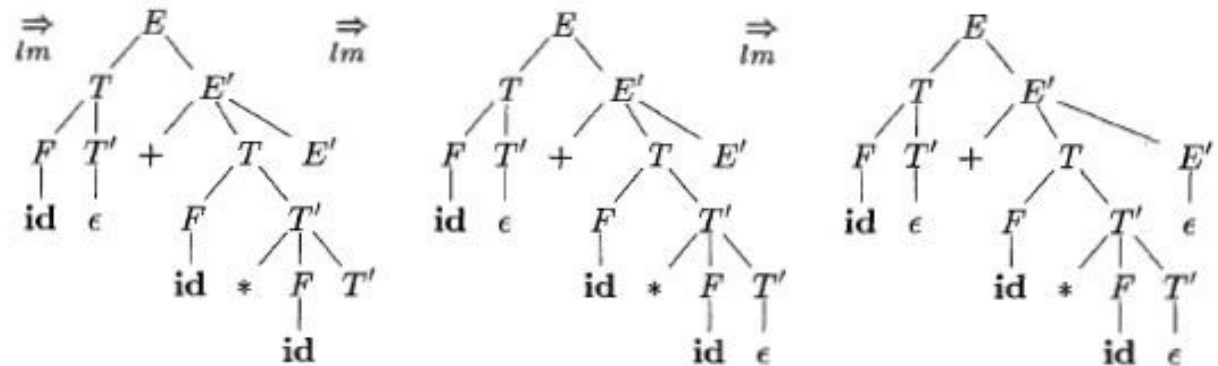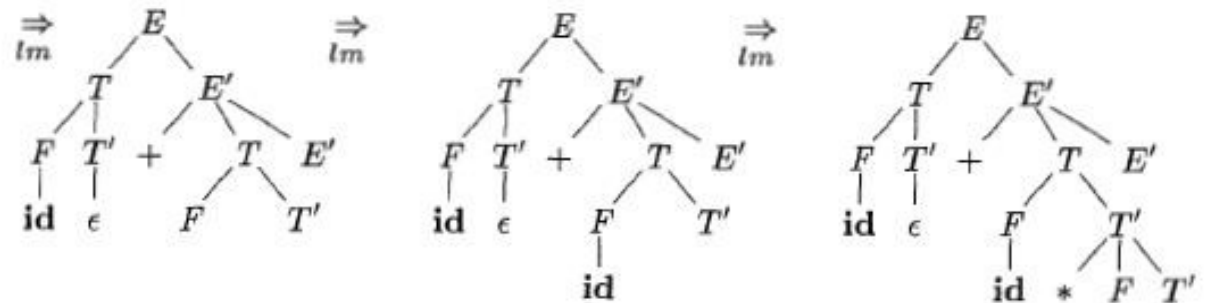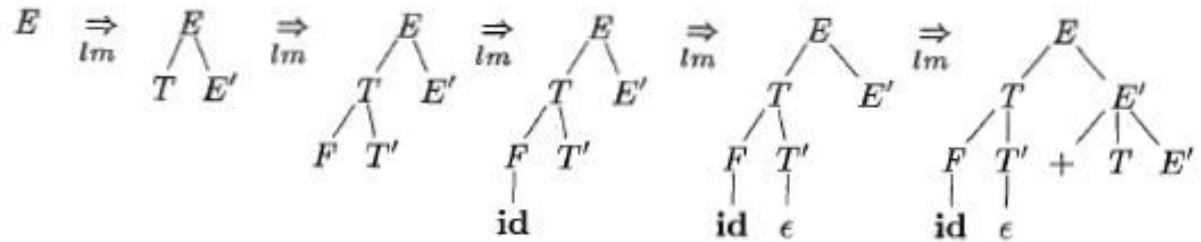
For example, if we have the two productions

$$
\begin{aligned}
stmt \quad &\rightarrow \quad \textbf{if } expr \textbf{ then } stmt \textbf{ else } stmt \\
&\mid \quad \textbf{if } expr \textbf{ then } stmt
\end{aligned}
$$

on seeing the input **if**, we cannot immediately tell which production to choose to expand $stmt$. In general, if $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ are two $A$-productions, and the input begins with a nonempty string derived from $\alpha$, we do not know whether to expand $A$ to $\alpha\beta_1$ or $\alpha\beta_2$. However, we may defer the decision by expanding $A$ to $\alpha A'$. Then, after seeing the input derived from $\alpha$, we expand $A'$ to $\beta_1$ or to $\beta_2$. That is, left-factored, the original productions become

$$
\begin{aligned}
A &\rightarrow \alpha A' \\
A' &\rightarrow \beta_1 \mid \beta_2
\end{aligned}
$$

# Top-Down Parsing

# The Idea of Predictive Parsing

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

id+id*id

# Recursive Descent Parsing

- One procedure for each non-terminal of the grammar.

- In the procedure, a production is chosen first.

  - Subsequent choices are performed by backtracking.

- For each symbol on rhs of production:

  - If symbol is terminal, match with input tape.

    - If match fails, backtrack.
    - If match succeeds, advance input tape

  - If symbol is non-terminal B, call B()

- Inefficient

  - Production is „guessed"

  - Wrong guess detected too late

  - Backtracking rewinds input tape, leading to multiple leads of same symbol

  - Worst case: exponential

- Does not work on left-recursive grammars

# Recursive Descent Parsing

```
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 */;
        }
}
```

# FIRST and FOLLOW

Define $FIRST(\alpha)$, where $\alpha$ is any string of grammar symbols, to be the set of terminals that begin strings derived from $\alpha$. If $\alpha \overset{*}{\Rightarrow} \epsilon$, then $\epsilon$ is also in $FIRST(\alpha)$. For example, in Fig. 4.15, $A \overset{*}{\Rightarrow} c\gamma$, so $c$ is in $FIRST(A)$.

Define $FOLLOW(A)$, for nonterminal $A$, to be the set of terminals $a$ that can appear immediately to the right of $A$ in some sentential form; that is, the set of terminals $a$ such that there exists a derivation of the form $S \overset{*}{\Rightarrow} \alpha A a\beta$, for some $\alpha$ and $\beta$, as in Fig. 4.15. Note that there may have been symbols between $A$ and $a$, at some time during the derivation, but if so, they derived $\epsilon$ and disappeared. In addition, if $A$ can be the rightmost symbol in some sentential form, then $\$$ is in $FOLLOW(A)$; recall that $\$$ is a special "endmarker" symbol that is assumed not to be a symbol of any grammar.

# FIRST

To compute $\mathrm{FIRST}(X)$ for all grammar symbols $X$, apply the following rules until no more terminals or $\epsilon$ can be added to any $\mathrm{FIRST}$ set.

1. If $X$ is a terminal, then $\mathrm{FIRST}(X) = \{X\}$.

2. If $X$ is a nonterminal and $X \to Y_1 Y_2 \cdots Y_k$ is a production for some $k \geq 1$, then place $a$ in $\mathrm{FIRST}(X)$ if for some $i$, $a$ is in $\mathrm{FIRST}(Y_i)$, and $\epsilon$ is in all of $\mathrm{FIRST}(Y_1), \ldots, \mathrm{FIRST}(Y_{i-1})$; that is, $Y_1 \cdots Y_{i-1} \overset{*}{\Rightarrow} \epsilon$. If $\epsilon$ is in $\mathrm{FIRST}(Y_j)$ for all $j = 1, 2, \ldots, k$, then add $\epsilon$ to $\mathrm{FIRST}(X)$. For example, everything in $\mathrm{FIRST}(Y_1)$ is surely in $\mathrm{FIRST}(X)$. If $Y_1$ does not derive $\epsilon$, then we add nothing more to $\mathrm{FIRST}(X)$, but if $Y_1 \overset{*}{\Rightarrow} \epsilon$, then we add $\mathrm{FIRST}(Y_2)$, and so on.

3. If $X \to \epsilon$ is a production, then add $\epsilon$ to $\mathrm{FIRST}(X)$.

Now, we can compute $\mathrm{FIRST}$ for any string $X_1 X_2 \cdots X_n$ as follows. Add to $\mathrm{FIRST}(X_1 X_2 \cdots X_n)$ all non-$\epsilon$ symbols of $\mathrm{FIRST}(X_1)$. Also add the non-$\epsilon$ symbols of $\mathrm{FIRST}(X_2)$, if $\epsilon$ is in $\mathrm{FIRST}(X_1)$; the non-$\epsilon$ symbols of $\mathrm{FIRST}(X_3)$, if $\epsilon$ is in $\mathrm{FIRST}(X_1)$ and $\mathrm{FIRST}(X_2)$; and so on. Finally, add $\epsilon$ to $\mathrm{FIRST}(X_1 X_2 \cdots X_n)$ if, for all $i$, $\epsilon$ is in $\mathrm{FIRST}(X_i)$.

# FOLLOW

To compute $\text{FOLLOW}(A)$ for all nonterminals $A$, apply the following rules until nothing can be added to any FOLLOW set.

1. Place $\$$ in $\text{FOLLOW}(S)$, where $S$ is the start symbol, and $\$$ is the input right endmarker.

2. If there is a production $A \rightarrow \alpha B \beta$, then everything in $\text{FIRST}(\beta)$ except $\epsilon$ is in $\text{FOLLOW}(B)$.

3. If there is a production $A \rightarrow \alpha B$, or a production $A \rightarrow \alpha B \beta$, where $\text{FIRST}(\beta)$ contains $\epsilon$, then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$.

# Example

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

1. $\text{FIRST}(F) = \text{FIRST}(T) = \text{FIRST}(E) = \{(, \textbf{id}\}$. To see why, note that the two productions for $F$ have bodies that start with these two terminal symbols, $\textbf{id}$ and the left parenthesis. $T$ has only one production, and its body starts with $F$. Since $F$ does not derive $\epsilon$, $\text{FIRST}(T)$ must be the same as $\text{FIRST}(F)$. The same argument covers $\text{FIRST}(E)$.

2. $\text{FIRST}(E') = \{+, \epsilon\}$. The reason is that one of the two productions for $E'$ has a body that begins with terminal $+$, and the other's body is $\epsilon$. Whenever a nonterminal derives $\epsilon$, we place $\epsilon$ in $\text{FIRST}$ for that nonterminal.

3. $\text{FIRST}(T') = \{*, \epsilon\}$. The reasoning is analogous to that for $\text{FIRST}(E')$.

# Example

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

4. $\mathrm{FOLLOW}(E) = \mathrm{FOLLOW}(E') = \{),\$\}$. Since $E$ is the start symbol, $\mathrm{FOLLOW}(E)$ must contain $\$$. The production body $(\,E\,)$ explains why the right parenthesis is in $\mathrm{FOLLOW}(E)$. For $E'$, note that this nonterminal appears only at the ends of bodies of $E$-productions. Thus, $\mathrm{FOLLOW}(E')$ must be the same as $\mathrm{FOLLOW}(E)$.

5. $\mathrm{FOLLOW}(T) = \mathrm{FOLLOW}(T') = \{+,),\$\}$. Notice that $T$ appears in bodies only followed by $E'$. Thus, everything except $\epsilon$ that is in $\mathrm{FIRST}(E')$ must be in $\mathrm{FOLLOW}(T)$; that explains the symbol $+$. However, since $\mathrm{FIRST}(E')$ contains $\epsilon$ (i.e., $E' \overset{*}{\Rightarrow} \epsilon$), and $E'$ is the entire string following $T$ in the bodies of the $E$-productions, everything in $\mathrm{FOLLOW}(E)$ must also be in $\mathrm{FOLLOW}(T)$. That explains the symbols $\$$ and the right parenthesis. As for $T'$, since it appears only at the ends of the $T$-productions, it must be that $\mathrm{FOLLOW}(T') = \mathrm{FOLLOW}(T)$.

6. $\mathrm{FOLLOW}(F) = \{+,*,),\$\}$. The reasoning is analogous to that for $T$ in point (5).

# LL(1) Grammars

A grammar $G$ is LL(1) if and only if whenever $A \to \alpha \mid \beta$ are two distinct productions of $G$, the following conditions hold:

1. For no terminal $a$ do both $\alpha$ and $\beta$ derive strings beginning with $a$.

2. At most one of $\alpha$ and $\beta$ can derive the empty string.

3. If $\beta \overset{*}{\Rightarrow} \epsilon$, then $\alpha$ does not derive any string beginning with a terminal in FOLLOW($A$). Likewise, if $\alpha \overset{*}{\Rightarrow} \epsilon$, then $\beta$ does not derive any string beginning with a terminal in FOLLOW($A$).

The first two conditions are equivalent to the statement that FIRST($\alpha$) and FIRST($\beta$) are disjoint sets. The third condition is equivalent to stating that if $\epsilon$ is in FIRST($\beta$), then FIRST($\alpha$) and FOLLOW($A$) are disjoint sets, and likewise if $\epsilon$ is in FIRST($\alpha$).

# LL(1) Parsing Table

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

Blanks are error entries; nonblanks indicate a production with which to expand a non-terminal.

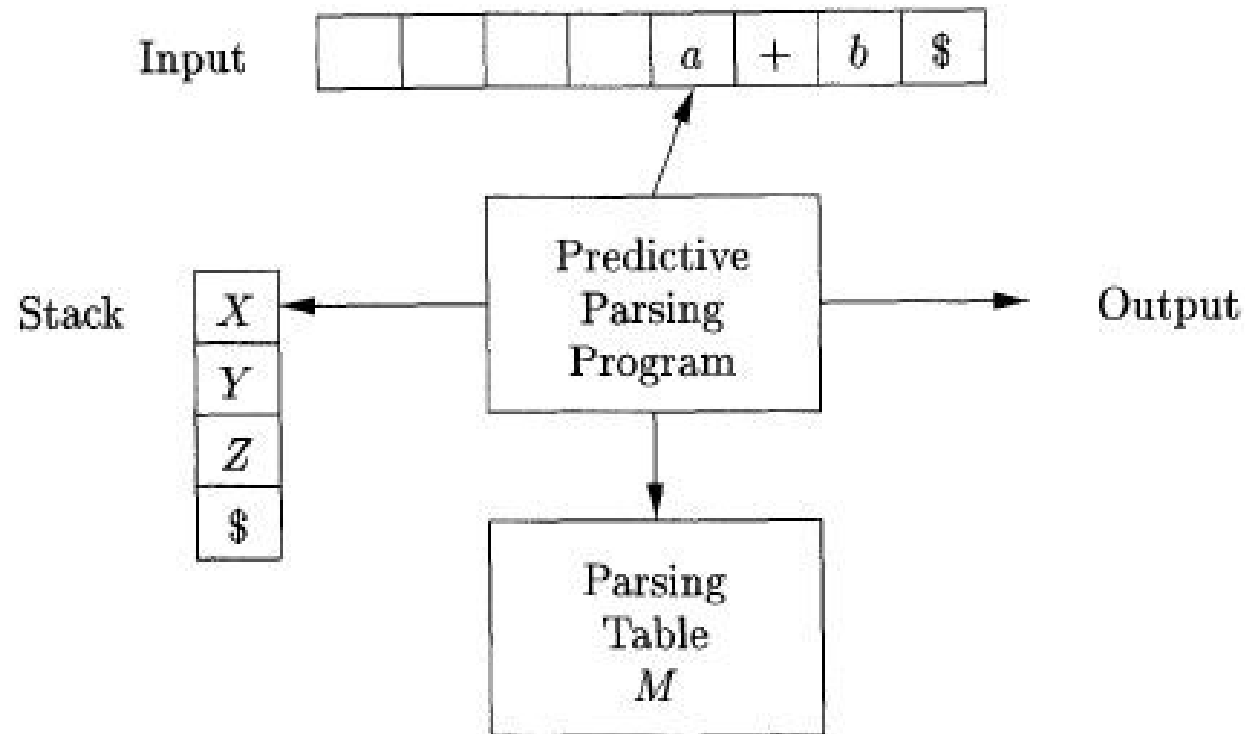| NON-TERMINAL | INPUT SYMBOL | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | **id** | + | * | ( | ) | $ |
| $E$ | $E \rightarrow TE'$ | | | $E \rightarrow TE'$ | | |
| $E'$ | | $E' \rightarrow +TE'$ | | | $E' \rightarrow \epsilon$ | $E' \rightarrow \epsilon$ |
| $T$ | $T \rightarrow FT'$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

# Predictive Parsing Machine



Figure 4.19: Model of a table-driven predictive parser

# Predictive Parsing Algorithm

**INPUT**: A string $w$ and a parsing table $M$ for grammar $G$.

**OUTPUT**: If $w$ is in $L(G)$, a leftmost derivation of $w$; otherwise, an error indication.

**METHOD**: Initially, the parser is in a configuration with $w\$$ in the input buffer and the start symbol $S$ of $G$ on top of the stack, above $\$$. The program in Fig. 4.20 uses the predictive parsing table $M$ to produce a predictive parse for the input. □

```
set ip to point to the first symbol of w;
set X to the top stack symbol;
while ( X ≠ $ ) { /* stack is not empty */
        if ( X is a ) pop the stack and advance ip;
        else if ( X is a terminal ) error();
        else if ( M[X, a] is an error entry ) error();
        else if ( M[X, a] = X → Y₁Y₂···Yₖ ) {
                output the production X → Y₁Y₂···Yₖ;
                pop the stack;
                push Yₖ, Yₖ₋₁, ... , Y₁ onto the stack, with Y₁ on top;
        }
        set X to the top stack symbol;
}
```

# LL(1) Parsing - Discussion

- Caveats
  - can't work on left recursive grammars
    - left associativity is difficult to implement
  - difficult to add semantic actions
    - top-down parsing does not favor post-order traversal
  - more difficult to provide meaningful errors compared to bottom-up
- Tool: ANTLR (Java)
  - semantic actions encoded as virtual non-terminals – somewhat non-intuitive

# Bottom-up Parsing

# Bottom-up = Leftmost

$$E \Rightarrow T \Rightarrow T * F \Rightarrow T * \mathbf{id} \Rightarrow F * \mathbf{id} \Rightarrow \mathbf{id} * \mathbf{id}$$

# Shift-Reduce Parsing

# Shift-Reduce Parsing

| STACK | INPUT | ACTION |
|---|---|---|
| $\$$ | $\mathbf{id_1} * \mathbf{id_2} \$$ | shift |
| $\$\ \mathbf{id_1}$ | $* \mathbf{id_2} \$$ | reduce by $F \to \mathbf{id}$ |
| $\$\ F$ | $* \mathbf{id_2} \$$ | reduce by $T \to F$ |
| $\$\ T$ | $* \mathbf{id_2} \$$ | shift |
| $\$\ T *$ | $\mathbf{id_2} \$$ | shift |
| $\$\ T * \mathbf{id_2}$ | $\$$ | reduce by $F \to \mathbf{id}$ |
| $\$\ T * F$ | $\$$ | reduce by $T \to T * F$ |
| $\$\ T$ | $\$$ | reduce by $E \to T$ |
| $\$\ E$ | $\$$ | accept |

# LR(0) Parsing

# Yacc
# (separate set of slides)