# Syntactic Analysis (Top Down Parsing)

Dr. P K Singh

# Top-Down Parsing

- The parse tree is created top to bottom.

- Top-down parser

  - Recursive-Descent Parsing

    - Backtracking is needed (If a choice of a production rule does not work, we backtrack to try other alternatives.)

    - It is a general parsing technique, but not widely used.

    - Not efficient

  - Predictive Parsing

    - no backtracking

    - efficient

    - needs a special form of grammars (LL(1) grammars).

    - Recursive Predictive Parsing is a special form of Recursive Descent parsing without backtracking.

    - Non-Recursive (Table Driven) Predictive Parser is also known as LL(1) parser.

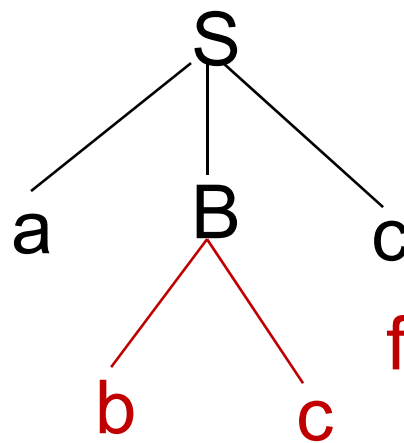# Recursive-Descent Parsing (uses Backtracking)

- Backtracking is needed.

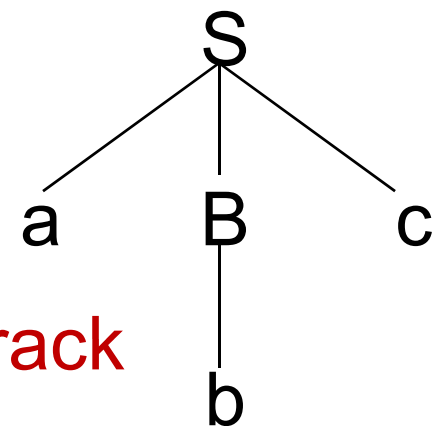- It tries to find the left-most derivation.

$S \rightarrow aBc$

$B \rightarrow bc \mid b$

Input: abc



fails, backtrack

# Predictive Parser

a grammar  ➔  ➔  a grammar suitable for

    eliminate    left    predictive  parsing (a LL(1)

    left recursion    factor    grammar)  no %100 guarantee.

- When re-writing a non-terminal in a derivation step, a predictive parser can uniquely choose a production rule by just looking the current symbol in the input string.

$A \rightarrow \alpha_1 \mid ... \mid \alpha_n$      input: ... a .......

                   current token

# Left Recursion

- A grammar is **left recursive** if it has a non-terminal A such that there is a derivation.

$$A \underset{+}{\Rightarrow} A\alpha \qquad \text{for some string } \alpha$$

- Top-down parsing techniques **cannot** handle left-recursive grammars.

- So, we have to convert our left-recursive grammar into an equivalent grammar which is not left-recursive.

- The left-recursion may appear in a single step of the derivation (*immediate left-recursion*), or may appear in more than one step of the derivation.

# Immediate Left-Recursion

$A \rightarrow A\,\alpha \mid \beta$       where $\beta$ does not start with A

$$\Downarrow$$       eliminate immediate left recursion

$A \rightarrow \beta\,A'$

$A' \rightarrow \alpha\,A' \mid \varepsilon$       an equivalent grammar

In general,

$A \rightarrow A\,\alpha_1 \mid ... \mid A\,\alpha_m \mid \beta_1 \mid ... \mid \beta_n$       where $\beta_1 ... \beta_n$ do not start with A

$$\Downarrow$$       eliminate immediate left recursion

$A \rightarrow \beta_1\,A' \mid ... \mid \beta_n\,A'$

$A' \rightarrow \alpha_1\,A' \mid ... \mid \alpha_m\,A' \mid \varepsilon$       an equivalent grammar

# Immediate Left-Recursion -- Example

$E \rightarrow E+T \mid T$

$T \rightarrow T*F \mid F$

$F \rightarrow id \mid (E)$

$\Downarrow$   eliminate immediate left recursion

$E \rightarrow T\,E'$

$E' \rightarrow +T\,E' \mid \varepsilon$

$T \rightarrow F\,T'$

$T' \rightarrow *F\,T' \mid \varepsilon$

$F \rightarrow id \mid (E)$

# Left-Recursion -- Problem

- A grammar cannot be immediately left-recursive, but it still can be left-recursive.

- By just eliminating the immediate left-recursion, we may not get a grammar which is not left-recursive.

$S \rightarrow Aa \mid b$
$A \rightarrow Sc \mid d$      This grammar is not immediately left-recursive, but it is still left-recursive.

$\underline{S} \Rightarrow Aa \Rightarrow \underline{S}ca$      or
$\underline{A} \Rightarrow Sc \Rightarrow \underline{A}ac$   causes to a left-recursion

- So, we have to eliminate all left-recursions from our grammar

# Eliminate Left-Recursion -- Algorithm

- Arrange non-terminals in some order: $A_1 \ldots A_n$

- **for** i **from** 1 **to** n **do** {

    - **for** j **from** 1 **to** i-1 **do** {

        replace each production

$$A_i \to A_j\, \gamma$$

           by

$$A_i \to \alpha_1\, \gamma \mid \ldots \mid \alpha_k\, \gamma$$

$$\text{where } A_j \to \alpha_1 \mid \ldots \mid \alpha_k$$

    }

    - eliminate immediate left-recursions among $A_i$ productions

}

# Eliminate Left-Recursion -- Example

S → Aa | b
A → Ac | Sd | f

- Order of non-terminals: S, A

for S:
- we do not enter the inner loop.
- there is no immediate left recursion in S.

for A:
- Replace A → Sd   with   A → Aad | bd
  So, we will have   A → Ac | Aad | bd | f
- Eliminate the immediate left-recursion in A

      A → bdA' | fA'
      A' → cA' |  adA' | ε

So, the resulting equivalent grammar which is not left-recursive is:
    S → Aa | b
    A → bdA' | fA'
    A' → cA' |  adA' | ε

# Eliminate Left-Recursion – Example2

S → Aa | b
A → Ac | Sd | f

- Order of non-terminals: A, S

for A:
    - we do not enter the inner loop.
    - Eliminate the immediate left-recursion in A
        A → SdA' | fA'
        A' → cA' | ε
for S:
    - Replace  S → Aa  with  S → SdA'a | fA'a
     So, we will have  S → SdA'a | fA'a | b
    - Eliminate the immediate left-recursion in S
        S → fA'aS' | bS'
        S' → dA'aS' | ε

So, the resulting equivalent grammar which is not left-recursive is:
    S → fA'aS' | bS'
    S' → dA'aS' | ε
    A → SdA' | fA'
    A' → cA' | ε

# Left-Factoring

- A predictive parser (a top-down parser without backtracking) insists that the grammar must be *left-factored*.

stmt → `if` expr `then` stmt `else` stmt  |

      `if` expr `then` stmt

- when we see `if`, we cannot now which production rule to choose to re-write *stmt* in the derivation.

# Left Factoring

- Rewriting productions to delay decisions

- Helpful for predictive parsing

- Not guaranteed to remove ambiguity

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$$

$$\Downarrow$$

$$A \rightarrow \alpha A'$$
$$A' \rightarrow \beta_1 \mid \beta_2$$

# Algorithm: Left Factoring

**Algorithm 4.2.** Left factoring a grammar.

*Input.* Grammar $G$.

*Output.* An equivalent left-factored grammar.

*Method.* For each nonterminal $A$ find the longest prefix $\alpha$ common to two or more of its alternatives. If $\alpha \neq \epsilon$, i.e., there is a nontrivial common prefix, replace all the $A$ productions $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \cdots \mid \alpha\beta_n \mid \gamma$ where $\gamma$ represents all alternatives that do not begin with $\alpha$ by

$$A \rightarrow \alpha A' \mid \gamma$$
$$A' \rightarrow \beta_1 \mid \beta_2 \mid \cdots \mid \beta_n$$

Here $A'$ is a new nonterminal. Repeatedly apply this transformation until no two alternatives for a nonterminal have a common prefix. $\square$

# Left-Factoring – Example1

A → abB | aB | cdg | cdeB | cdfB

$$\Downarrow$$

A → aA' | cdg | cdeB | cdfB

A' → bB | B

$$\Downarrow$$

A → aA' | cdA''

A' → bB | B

A'' → g | eB | fB

# Left-Factoring – Example2

A → ad | a | ab | abc | b

$$\Downarrow$$

A → aA' | b

A' → d | ε | b | bc

$$\Downarrow$$

A → aA' | b

A' → d | ε | bA''

A'' → ε | c

# Top Down Parsing

- Can be viewed two ways:

  - Attempt to find leftmost derivation for input string

  - Attempt to create parse tree, starting from at root, creating nodes in preorder

- General form is recursive descent parsing

  - May require backtracking

  - Backtracking parsers not used frequently because not needed

# Predictive Parsing

- A special case of recursive-descent parsing that does not require backtracking

- Must always know which production to use based on current input symbol

- Can often create appropriate grammar:

    – removing left-recursion

    – left factoring the resulting grammar

# Predictive Parser (example)

stmt → `if` ......      |

       `while` ......      |

       `begin` ......      |

       `for` .....

- When we are trying to write the non-terminal *stmt*, if the current token is `if` we have to choose first production rule.

- When we are trying to write the non-terminal *stmt*, we can uniquely choose the production rule by just looking the current token.

- We eliminate the left recursion in the grammar, and left factor it. But it may not be suitable for predictive parsing (not LL(1) grammar).

# Recursive Predictive Parsing

- Each non-terminal corresponds to a procedure.

Ex:      A → aBb          (This is only the production rule for A)

proc A {

       - match the current token with a, and move to the next token;

       - call 'B';

       - match the current token with b, and move to the next token;

       }

# Recursive Predictive Parsing (cont.)

A → aBb | bAB

proc A {

    case of the current token {

        'a':  - match the current token with a, and move to the next token;

              - call 'B';

              - match the current token with b, and move to the next token;

        'b':  - match the current token with b, and move to the next token;

              - call 'A';

              - call 'B';

    }

}

# Recursive Predictive Parsing (cont.)

- When to apply $\varepsilon$-productions.


  $A \rightarrow aA \mid bB \mid \varepsilon$


- If all other productions fail, we should apply an $\varepsilon$-production. For example, if the current token is not a or b, we may apply the $\varepsilon$-production.

- Most correct choice: We should apply an $\varepsilon$-production for a non-terminal A when the current token is in the follow set of A (which terminals can follow A in the sentential forms).

# Transition Diagrams

- ## For parser:

  - One diagram for each nonterminal

  - Edge labels can be tokens or nonterminal

    - A transition on a token means we should take that transition if token is next input symbol

    - A transition on a nonterminal can be thought of as a call to a procedure for that nonterminal

- ## As opposed to lexical analyzers:

  - One (or more) diagrams for each token

  - Labels are symbols of input alphabet

# Creating Transition Diagrams

- First eliminate left recursion from grammar

- Then left factor grammar

- For each nonterminal A:

  - Create an initial and final state

  - For every production A $\rightarrow$ $X_1 X_2 \ldots X_n$, create a path from initial to final state with edges labeled $X_1$, $X_2$, …, $X_n$.
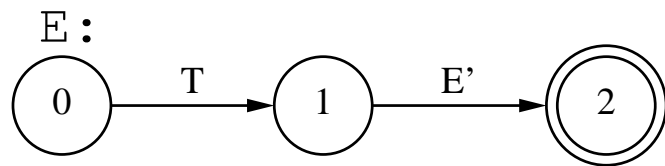
# Using Transition Diagrams

- Predictive parsers:
  - Start at start symbol of grammar
  - From state s with edge to state t labeled with token a, if next input token is a:
    - State changes to t
    - Input cursor moves one position right
  - If edge labeled by nonterminal A:
    - State changes to start state for A
    - Input cursor is not moved
    - If final state of A reached, then state changes to t
  - If edge labeled by ε, state changes to t

- Can be recursive or non-recursive using stack
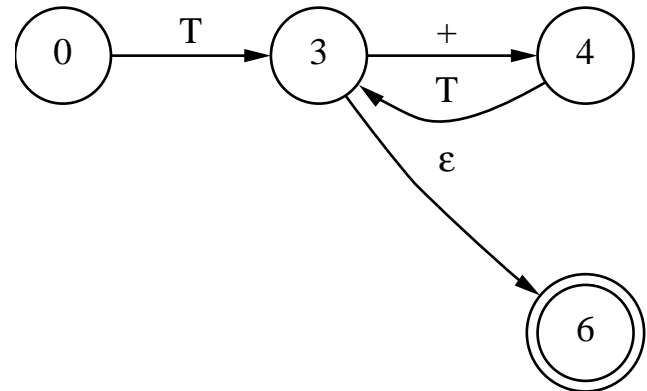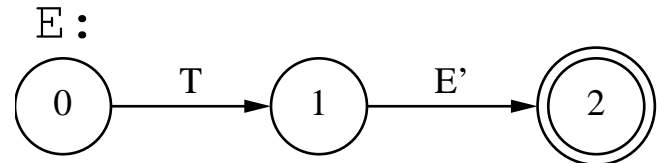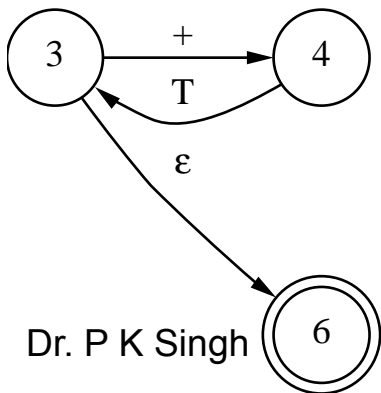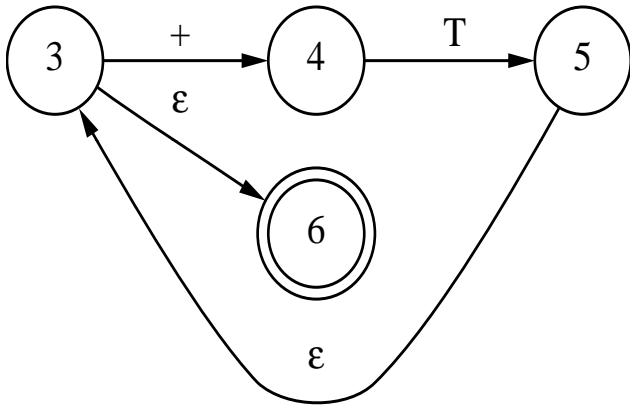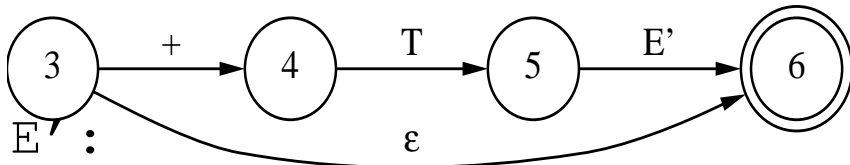
# Transition Diagram Example

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

$$\Longrightarrow$$

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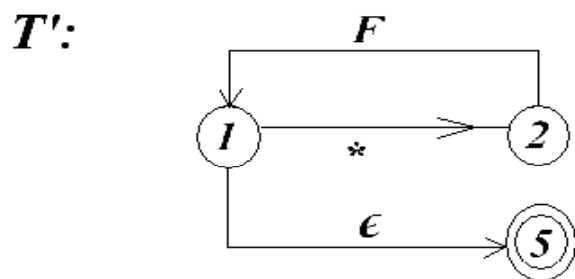$$T' \rightarrow *FT' \mid \varepsilon$$
$$F \rightarrow (E) \mid \textbf{id}$$



E:  0 →T→ 1 →E'→ 2

T':  10 →*→ 11 →F→ 12 →T'→ 13,  10 →ε→ 13

E':  3 →+→ 4 →T→ 5 →E'→ 6,  3 →ε→ 6

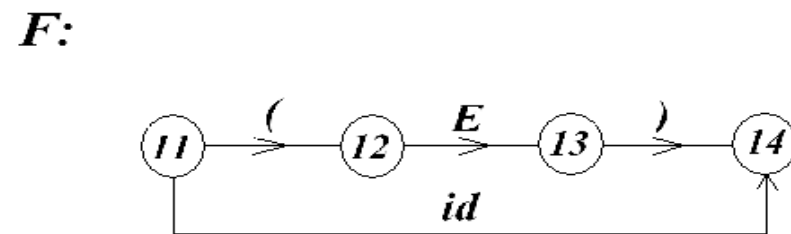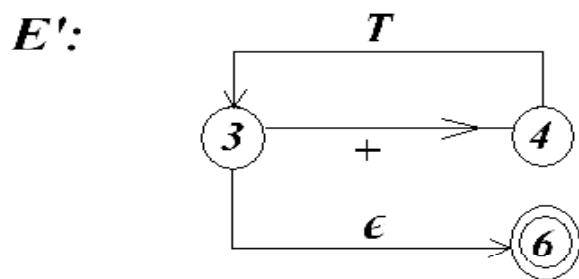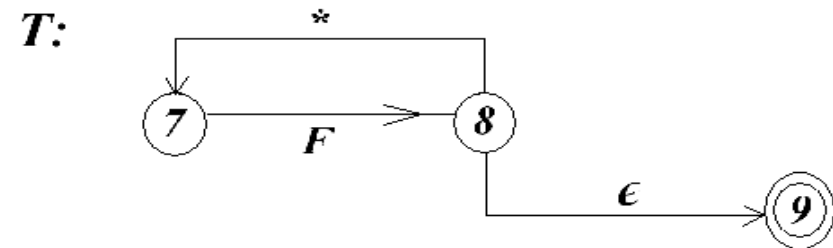T:  7 →F→ 8 →T'→ 9

F:  14 →(→ 15 →E→ 16 →)→ 17,  14 →id→ 17

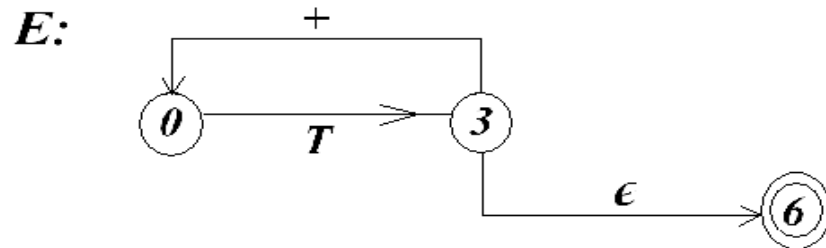# Simplifying Transition Diagrams

# Simplified Transition Diagrams

# Recursive Predictive Parsing (Example)

A → aBe | cBd | C

B → bB | ε

C → f

proc A {
    case of the current token {
        a:  - match the current token with a,
           and move to the next token;
           - call B;
           - match the current token with e,
           and move to the next token;
        c:  - match the current token with c,
           and move to the next token;
           - call B;
           - match the current token with d,
           and move to the next token;
      f:  - call C
    }
}

**first set of C**

proc C {    match the current token with f,
           and move to the next token; }

proc B {
    case of the current token {
        b:- match the current token with b,
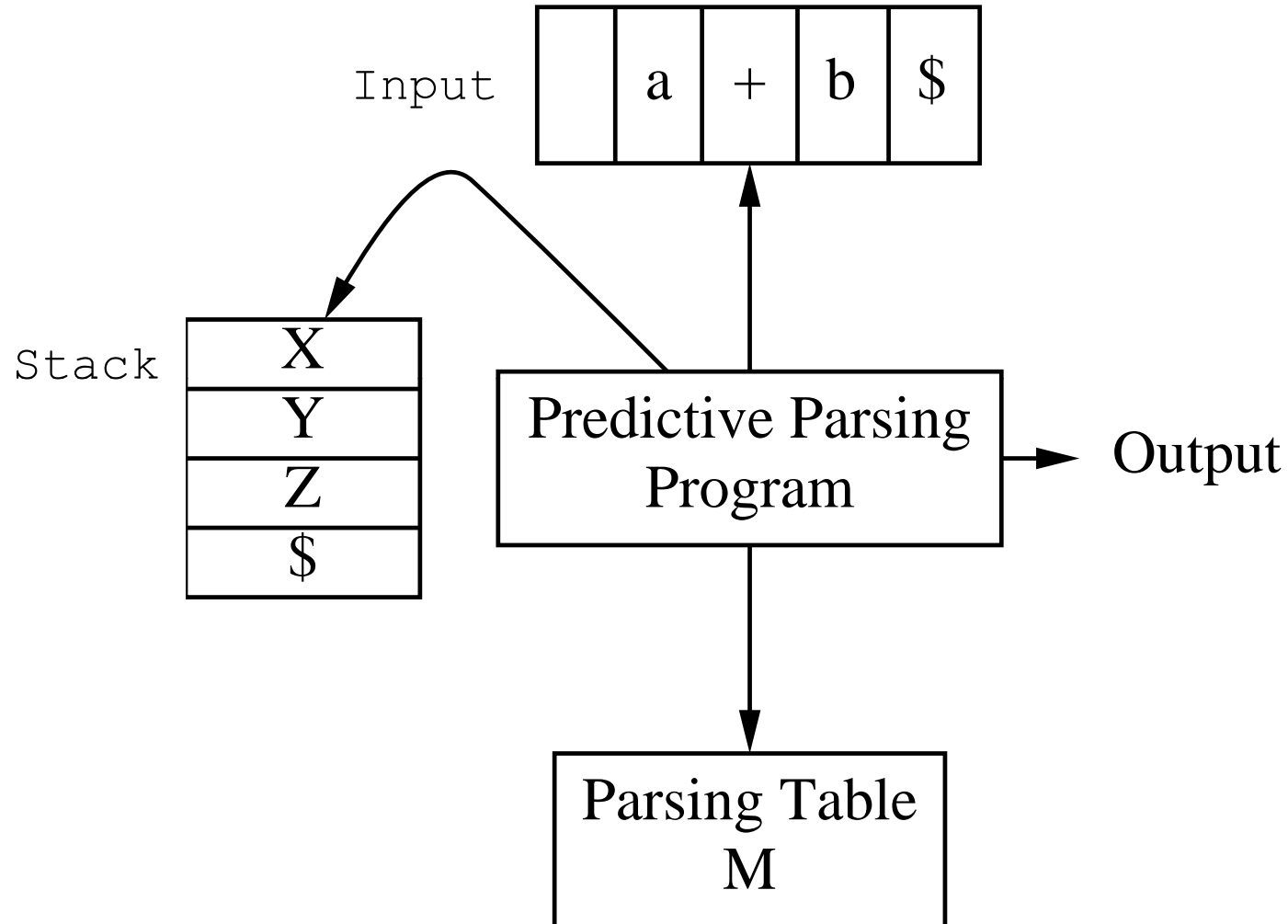          and move to the next token;
        - call B
      e,d:  do nothing
    }
}

**follow set of B**

# Nonrecursive Predictive Parsing (1)

Input | a | + | b | $

Stack
X
Y
Z
$

Predictive Parsing
Program

Output

Parsing Table
M

# Nonrecursive Predictive Parsing (2)

- The symbol at the top of the stack (say X) and the current symbol in the input string (say a) determine the parser action.

- There are four possible parser actions.

1. If X and a are $ ➔ parser halts (successful completion)

2. If X and a are the same terminal symbol (different from $)

   ➔ parser pops X from the stack, and moves the next symbol in the input buffer.

3. If X is a non-terminal

   ➔ parser looks at the parsing table entry M[X,a]. If M[X,a] holds a production rule $X \rightarrow Y_1Y_2...Y_k$, it pops X from the stack and pushes $Y_k,Y_{k-1},...,Y_1$ into the stack. The parser also outputs the production rule $X \rightarrow Y_1Y_2...Y_k$ to represent a step of the derivation.

4. none of the above ➔ error

   – all empty entries in the parsing table are errors.

   – If X is a terminal symbol different from a, this is also an error case.

# Predictive Parsing Table

| Nonter-minal | Input Symbol | | | | | |
|---|---|---|---|---|---|---|
| | id | + | * | ( | ) | $ |
| E | E→TE′ | | | E→TE′ | | |
| E′ | | E′→+TE′ | | | E′→ε | E′→ε |
| T | T→FT′ | | | T→FT′ | | |
| T′ | | T′→ε | T′→*FT′ | | T′→ε | T′→ε |
| F | F→id | | | F→(E) | | |

# Using a Predictive Parsing Table

| Stack | Input | Output |
|-------|-------|--------|
| $E | id+id*id$ | |
| $E'T | id+id*id$ | E→TE' |
| $E'T'F | id+id*id$ | T→FT' |
| $E'T'id | id+id*id$ | F→id |
| $E'T' | +id*id$ | |
| $E' | +id*id$ | T'→ε |
| $E'T+ | +id*id$ | E'→+TE' |
| $E'T | id*id$ | |
| $E'T'F | id*id$ | T→FT' |

| Stack | Input | Output |
|-------|-------|--------|
| ... | ... | ... |
| $E'T'id | id*id$ | F→id |
| $E'T' | *id$ | |
| $E'T'F* | *id$ | T'→*FT' |
| $E'T'F | id$ | |
| $E'T'id | id$ | F→id |
| $E'T' | $ | |
| $E' | $ | T'→ ε |
| $ | $ | E'→ ε |

# FIRST

- $\text{FIRST}(\alpha)$ is the set of all terminals that begin any string derived from $\alpha$

- Computing `FIRST`:

  - If `X` is a terminal, $\text{FIRST(X)} = \{X\}$

  - If $X \rightarrow \varepsilon$ is a production, add $\varepsilon$ to $\text{FIRST(X)}$

  - If `X` is a nonterminal and $X \rightarrow Y_1 Y_2 ... Y_n$ is a production:

    - For all terminals a, add a to FIRST(X) if a is a member of any FIRST($Y_i$) and ε is a member of FIRST($Y_1$), FIRST($Y_2$), … FIRST($Y_{i-1}$)

    - If ε is a member of FIRST($Y_1$), FIRST($Y_2$), … FIRST($Y_n$), add ε to FIRST(X)

# FOLLOW

- FOLLOW(A), for any nonterminal A, is the set of terminals a that can appear immediately to the right if A in some sentential form

- More formally, a is in FOLLOW(A) if and only if there exists a derivation of the form S $^*$=>αAaβ

- $ is in FOLLOW(A) if and only if there exists a derivation of the form S $^*$=> αA

## Computing FOLLOW

- Place `$` in `FOLLOW(S)`

- If there is a production `A` → αBβ, then everything in `FIRST`(β) (except for ε) is in `FOLLOW(B)`

- If there is a production `A` → αB, or a production `A` → αBβ where `FIRST`(β) contains ε,then everything in `FOLLOW(A)` is also in `FOLLOW(B)`

# FIRST and FOLLOW Example

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

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

# Creating a Predictive Parsing Table

- For each production `A` → α :

  - For each terminal `a` in `FIRST(`α`)` add `A` → α to `M[A, a]`

  - If ε is in `FIRST(`α`)` add `A` → α to `M[A, b]` for every terminal `b` in `FOLLOW(A)`

  - If ε is in `FIRST(`α`)` and `$` is in `FOLLOW(A)` add `A` → α to `M[A, $]`

- Mark each undefined entry of M as an error entry (use some recovery strategy)

# Example

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

FIRST($E$) = FIRST($T$) = FIRST($F$) = $\{(, \mathbf{id}\}$.

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

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

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

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

FOLLOW($F$) = $\{+, *, ), \$\}$

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

# Constructing LL(1) Parsing Table -- Example

E → TE'        FIRST(TE')={(,id}        ➔ E → TE' into M[E,(] and M[E,id]

E' → +TE'      FIRST(+TE' )={+}          ➔ E' → +TE' into M[E',+]

E' → ε         FIRST(ε)={ε}             ➔ none
               but since ε in FIRST(ε)
               and FOLLOW(E')={$,)}      ➔ E' → ε   into M[E',$]  and M[E',)]

T → FT'        FIRST(FT')={(,id}         ➔ T → FT' into M[T,(] and M[T,id]

T' → *FT'      FIRST(*FT' )={*}          ➔ T' → *FT' into M[T',*]

T' → ε         FIRST(ε)={ε}             ➔ none
               but since ε in FIRST(ε)
               and FOLLOW(T')={$,),+}    ➔ T' → ε  into M[T',$], M[T',)] and M[T',+]

F → (E)        FIRST((E) )={(}           ➔ F → (E) into M[F,(]

F → id         FIRST(id)={id}            ➔ F → id  into M[F,id]

# LL(1) Grammars

- A grammar whose parsing table has no multiply-defined entries is said to be LL(1) grammar.

one input symbol used as a look-head symbol do determine parser action

$\downarrow$

LL(1) — left most derivation

$\uparrow$

input scanned from left to right

- The parsing table of a grammar may contain more than one production rule. In this case, we say that it is not a LL(1) grammar.

# A Grammar which is not LL(1)

S → i C t S E | a          FOLLOW(S) = { $,e }

E → e S | ε                FOLLOW(E) = { $,e }

C → b                      FOLLOW(C) = { t }

FIRST(iCtSE) = {i}

FIRST(a) = {a}

FIRST(eS) = {e}

FIRST(ε) = {ε}

FIRST(b) = {b}

|   | **a** | **b** | **e** | **i** | **t** | **$** |
|---|---|---|---|---|---|---|
| **S** | S → a |  |  | S → iCtSE |  |  |
| **E** |  |  | E → e S<br>E → ε |  |  | E → ε |
| **C** |  | C → b |  |  |  |  |

two production rules for M[E,e]

# A Grammar which is not LL(1) (cont.)

- What do we have to do it if the resulting parsing table contains multiply defined entries?

  - If we didn't eliminate left recursion, eliminate the left recursion in the grammar.

  - If the grammar is not left factored, we have to left factor the grammar.

  - If its (new grammar's) parsing table still contains multiply defined entries, that grammar is ambiguous or it is inherently not a LL(1) grammar.

- A left recursive grammar cannot be a LL(1) grammar.

  - $A \rightarrow A\alpha \mid \beta$

    ➔ any terminal that appears in FIRST($\beta$) also appears FIRST($A\alpha$) because $A\alpha \Rightarrow \beta\alpha$.

    ➔ If $\beta$ is $\varepsilon$, any terminal that appears in FIRST($\alpha$) also appears in FIRST($A\alpha$) and FOLLOW(A).

- A grammar is not left factored, it cannot be a LL(1) grammar

  - $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$

    ➔ any terminal that appears in FIRST($\alpha\beta_1$) also appears in FIRST($\alpha\beta_2$).

- An ambiguous grammar cannot be a LL(1) grammar.

# Properties of LL(1) Grammars

- A grammar G is LL(1) if and only if the following conditions hold for two distinctive production rules $A \to \alpha$ and $A \to \beta$

  1. Both $\alpha$ and $\beta$ cannot derive strings starting with same terminals.

  2. At most one of $\alpha$ and $\beta$ can derive to $\varepsilon$.

  3. If $\beta$ can derive to $\varepsilon$, then $\alpha$ cannot derive to any string starting with a terminal in FOLLOW(A).

- A Grammar to be LL(1), following conditions must satisfied:

  For every pair of productions A → α Ι β

  {

  $\qquad$ FIRST(α) ∩ FIRST(β) = Φ

  $\qquad$ and if FIRST(β) contains ε then

  $\qquad$ FIRST(α) ∩ FOLLOW(A) = Φ

  }

# Example

- Test the Following Grammar is LL(1) or not ?

    S $\rightarrow$ 1AB | ε

    A $\rightarrow$ 1AC | 0C

    B $\rightarrow$ 0S

    C $\rightarrow$ 1

    For Production S $\rightarrow$ 1AB | ε

    FIRST(1AB) ∩ FIRST(ε) = {1} ∩ {ε} = Φ and

    FIRST(1AB) ∩ FOLLOW(S) = {1} ∩ {$} = Φ

    Similarly   A $\rightarrow$ 1AC | 0C

    FIRST(1AC) ∩ FIRST(0C) = {1} ∩ {0} = Φ

    Hence The Grammar is LL(1)