

Chapter 4 Top-Down Parsing

- Study Goals:
 - Master
 - Recognition of LL(1) grammar, Construction of Recursive-descent parsing, LL(1) parsing
 - Understand
 - First set, Follow set, LL(1) grammar
 - Know
 - Backtracking Parsing, Error Recovery in Top-Down parsers, Syntax Tree Construction in Top-down parsers

Top-Down Parsing

• Definition

Parsing begins with the start symbol of grammar and tries to find out the derivation of the input string tracing out the steps in a leftmost derivation

The construction of Parse Tree:

Start symbol of grammar is the root, parse tree is constructed from the root to leaves in preorder, the leaves of the parse tree are just input string of tokens

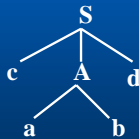
Example G: $S \rightarrow cAd$

$A \rightarrow ab$

$A \rightarrow a$

Top-down parsing of string "cabd"

Derivation: $\underline{S} \Rightarrow cAd \Rightarrow cabd$



• The Key of Top-down Parsing

Another top-down parsing process for input string "cabd" is:



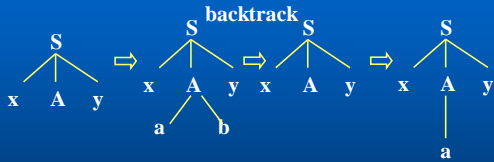
Fail, because a derivation step has selected a wrong production $A \rightarrow a$

- So the problem of top-down parsing is the choice of a production for a derivation step
- If the leftmost nonterminal to be replaced is **B**, there are totally **n** productions of **B**: $B \rightarrow A_1 | A_2 | \dots | A_n$, how can we determine which one to use?

• Categories of Top-down parsing

1. Backtracking Parsing
 - ❖ A nonterminal has more than one productions
 - ❖ and basing on the current input symbol, the parser can't determine which one to choose
 - ❖ it must try different possibilities, backing up an arbitrary amount in the input if one possibility fails.

Example: Grammar : $S \rightarrow xAy$ $A \rightarrow ab|a$
 Top-down parsing of string "xay"



2. Predictive parsing

Parser attempts to predict the next construction in the input string using one or more lookahead tokens

Two Kinds of Predictive Parsing

- Recursive-descent parsing
- LL(1) parsing

4.1 Predictive Parsing

4.2 Recognition of LL(1) Grammar

4.3 Non LL(1) grammar to LL(1) grammar

4.4 Top-Down Parsing by Recursive-Descent

4.5 LL(1) Parsing

4.6 Error Recovery in Top-Down Parsers

4.1 Predictive Parsing

1 The Condition of Predictive Parsing

2 The Definition of Lookahead Sets

3 The Definition of LL(1) Grammar

1 The Condition of Predictive Parsing

- ❖ Parsing of the input string begins with the start symbol of the grammar
- ❖ if it can uniquely determine which production to use next in derivation basing on the current token of the input, parsing is predictive

The Condition of Predictive Parsing

- Predictive parsing require that the grammar must be a **LL(1) grammar**
- What is LL(1) grammar?
 It's definition depends on the definition of lookahead sets--**First set** and **Follow set**

2 The Definition of Lookahead Sets

1) First Sets

Definition

$G=(V_N, V_T, P, S)$ is a grammar, $\beta \in (V_N \cup V_T)^*$
 $FIRST(\beta) = \{ a \in V_T \mid \beta \Rightarrow^* a \dots \}$
 if $\beta \Rightarrow^* \epsilon$ then $\epsilon \in FIRST(\beta)$

Intuitively, the first set of string β is the set of first terminals (including ϵ) that can be derived from β

Example $G[S]$:

$S \rightarrow Ap \quad FIRST(Ap) = \{a, c\}$
 $S \rightarrow Bq \quad FIRST(Bq) = \{b, d\}$
 $A \rightarrow a \quad FIRST(a) = \{a\}$
 $A \rightarrow cA \quad FIRST(cA) = \{c\}$
 $B \rightarrow b \quad FIRST(b) = \{b\}$
 $B \rightarrow dB \quad FIRST(dB) = \{d\}$

If there are more than one productions of nonterminal A : $A \rightarrow \alpha \mid \beta \dots$, but $FIRST(\alpha) \cap FIRST(\beta) \cap \dots = \emptyset$, this grammar can be predictive parsed. Basing on the current input symbol which production to choose is determined

2) Follow Sets

Definition

$G=(V_T, V_N, S, P)$ is a grammar, $A \in V_N$,
 $FOLLOW(A) = \{ a \in V_T \mid S \Rightarrow^* \dots Aa \dots \}$,
 if $S \Rightarrow^* \dots A$, then $\$ \in FOLLOW(A)$
 ($\$$ is used to mark the end of the input)

Intuitively, the follow set of nonterminal A is the set of terminals (include $\$$) following A in all sentential form of the grammar

Example $G_3[S]$:

$S \rightarrow aA \mid d$
 $A \rightarrow bAS \mid \epsilon$
 $\triangleright S \Rightarrow^* aA, \$ \in FOLLOW(A)$
 $S \Rightarrow^* abAS \Rightarrow^* abAaA, a \in FOLLOW(A)$
 $\dots \Rightarrow^* abAd, d \in FOLLOW(A)$
 $FOLLOW(A) = \{ \$, a, d \}$
 $\triangleright S \Rightarrow^* S, \$ \in FOLLOW(S)$
 $S \Rightarrow aA \Rightarrow abAS \Rightarrow abbaSS \Rightarrow abbaSaaS$
 $\dots \Rightarrow abbaSd$
 $FOLLOW(S) = \{ \$, a, d \}$

There are two productions of nonterminal A :

$A \rightarrow bAS$ and $A \rightarrow \epsilon$, let the current input symbol is " x "

- \triangleright If $x \in FIRST(bAS) = \{b\}$, then choose $A \rightarrow bAS$ for derivation
- \triangleright If $x \in FOLLOW(A) = \{ \$, a, d \}$, then choose $A \rightarrow \epsilon$ for derivation
- \triangleright Because $FIRST(bAS) \cap FOLLOW(A) = \emptyset$, which production to choose is determined

$S \rightarrow AB \mid bC$

$A \rightarrow a \mid \epsilon$

$B \rightarrow aD \mid \epsilon$

$C \rightarrow AD \mid b$

$D \rightarrow aS \mid c$

aac

3 The Definition of LL(1) Grammar

• LL(1) Grammar

A grammar is LL(1) if the following conditions are satisfied:

1. For each production $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$
 $\text{First}(\alpha_i) \cap \text{First}(\alpha_j) = \emptyset$ for all i and $j, 1 \leq i, j \leq n, i \neq j$
2. For each nonterminal A such that $\text{First}(A)$ contains ϵ , $\text{First}(A) \cap \text{Follow}(A) = \emptyset$

• Example G[S] :

$S \rightarrow aAS$ $\text{First}(aAS) = \{a\}$

$S \rightarrow b$ $\text{First}(b) = \{b\}$

$A \rightarrow bA$ $\text{First}(bA) = \{b\}$

$A \rightarrow \epsilon$ $\text{First}(\epsilon) = \{\epsilon\}$

$\text{First}(A) = \{b, \epsilon\}$ $\text{Follow}(A) = \{a, b\}$

Because $\text{First}(A) \cap \text{Follow}(A) = \{b\} \neq \emptyset$,

G[S] is not a LL(1) grammar, when the leftmost nonterminal to be replaced is A and the current input symbol is " b ", parsing can't determine which productions of A to choose : $A \rightarrow bA$ or $A \rightarrow \epsilon$



4.2 Recognition of LL(1) Grammar

Recognition of LL(1) grammar has the following four steps:

1. Compute the set of nullable nonterminals
2. Compute $\text{FIRST}(\alpha)$ for the right-hand side string α of each production
3. Compute $\text{FOLLOW}(A)$ for each nonterminal A
4. Recognize basing on the definition of LL(1)

1 Compute the set of nullable nonterminals

• Nullable

A nonterminal A is nullable if there exists a derivation $A \Rightarrow^* \epsilon$

• Algorithm

Let S is the set of nullable nonterminals

1. First, $S = \{A_i \mid A_i \rightarrow \epsilon \text{ is a production}\}$
2. For each production $p: A_p \rightarrow X_1 \dots X_n$, if $X_1 \dots X_n \in S$, then $S := S \cup \{A_p\}$
3. Repeat step 2 until there is no change to S

Example G[S]:

$S \rightarrow AB|bC$

$A \rightarrow b|\epsilon$

$B \rightarrow aD|\epsilon$

$C \rightarrow AD|b$

$D \rightarrow aS|c$

| | Set S |
|--------|---------|
| Origin | {A,B} |
| Pass 1 | {A,B,S} |
| Pass 2 | {A,B,S} |

The set of nullable nonterminals is {A,B,S}

2 Compute $\text{FIRST}(\alpha)$ for the right-hand string α of each production

- Algorithm for computing $\text{First}(A)$ for each grammar symbol A ($A \in V_T \cup V_N$)
- Algorithm for computing $\text{First}(\alpha)$ for a string α

Algorithm for computing First(A) for each grammar symbol $A (A \in V_T \cup V_N)$

1. For all $a \in V_T$ do $\text{First}(a) = \{a\}$
2. For all $A \in V_N$, if $A \Rightarrow^* \epsilon$ then $\text{First}(A) = \{\epsilon\}$ else $\text{First}(A) = \{ \}$
3. For each production $A \rightarrow X_1 \dots X_j \dots X_n$, $\text{First}(A) = \text{First}(A) \cup \text{SectionFirst}(X_1 \dots X_j \dots X_n)$
4. Repeat step 3 until there is no change to any First set

SectionFirst($X_1 \dots X_j \dots X_n$)

$$= (\text{First}(X_1) - \{\epsilon\}) \cup (\text{First}(X_2) - \{\epsilon\}) \cup \dots \cup (\text{First}(X_j) - \{\epsilon\}) \cup \text{First}(X_{j+1})$$

X_{j+1} is the first symbol that is not nullable in the right-hand of production

- If X_j is not nullable, then $\text{SectionFirst}(X_1 \dots X_j \dots X_n) = \text{First}(X_j)$

- If $X_1 \dots X_n$ are all nullable, then $\text{SectionFirst}(X_1 \dots X_n) = (\text{First}(X_1) - \{\epsilon\}) \cup (\text{First}(X_2) - \{\epsilon\}) \cup \dots \cup (\text{First}(X_n) - \{\epsilon\}) \cup \{\epsilon\}$



G[S]: $S \rightarrow AB|bC$ The set of nullable nonterminals is $\{A, B, S\}$
 $A \rightarrow b|\epsilon$
 $B \rightarrow aD|\epsilon$
 $C \rightarrow AD|b$
 $D \rightarrow aS|c$



| | Origin | Pass 1 | Pass 2 | Pass 3 |
|---|------------|--------------|----------------|----------------|
| S | ϵ | ϵ b | ϵ b a | ϵ b a |
| A | ϵ | ϵ b | ϵ b | ϵ b |
| B | ϵ | ϵ a | ϵ a | ϵ a |
| C | | b | b a c | b a c |
| D | | a c | a c | a c |
| a | a | a | a | a |
| b | b | b | b | b |

Algorithm for computing First(α) for a string $\alpha = X_1 X_2 \dots X_n$

1. If X_1 is not nullable, then $\text{FIRST}(\alpha) = \text{FIRST}(X_1)$
2. if X_j ($1 \leq j < n$) is nullable then $\text{FIRST}(\alpha) = (\text{FIRST}(X_1) - \{\epsilon\}) \cup \dots \cup (\text{FIRST}(X_j) - \{\epsilon\}) \cup \text{FIRST}(X_{j+1})$
3. if X_i ($1 \leq i \leq n$) are all nullable, then $\text{FIRST}(\alpha) = (\text{FIRST}(X_1) - \{\epsilon\}) \cup \dots \cup (\text{FIRST}(X_n) - \{\epsilon\}) \cup \{\epsilon\}$



Example

G[S] $S \rightarrow AB|bC$ First set for nonterminals
 $A \rightarrow b|\epsilon$ $\text{First}(S) = \{a, b, \epsilon\}$ $\text{First}(A) = \{b, \epsilon\}$
 $B \rightarrow aD|\epsilon$ $\text{First}(B) = \{a, \epsilon\}$ $\text{First}(C) = \{a, b, c\}$
 $C \rightarrow AD|b$ $\text{First}(D) = \{a, c\}$
 $D \rightarrow aS|c$

First sets for strings in the right-hand side of productions

$S \rightarrow AB$ $\text{FIRST}(AB) = (\text{FIRST}(A) - \{\epsilon\}) \cup (\text{FIRST}(B) - \{\epsilon\}) \cup \{\epsilon\} = \{a, b, \epsilon\}$

$S \rightarrow bC$ $\text{FIRST}(bC) = \{b\}$

$A \rightarrow \epsilon$ $\text{FIRST}(\epsilon) = \{\epsilon\}$

$A \rightarrow b$ $\text{FIRST}(b) = \{b\}$

$C \rightarrow AD$ $\text{FIRST}(AD) = (\text{FIRST}(A) - \{\epsilon\}) \cup \text{FIRST}(D) = \{b, a, c\}$

$D \rightarrow aS$ $\text{FIRST}(aS) = \{a\}$

3 Compute FOLLOW(A) for nonterminal A

1. S is the start symbol, $\text{Follow}(S) = \{\$ \}$; for all $A \in V_N$, and $A \neq S$, $\text{Follow}(A) = \{ \}$;

2. For each production $B \rightarrow \alpha A \gamma$, for each A that is a nonterminal do

$\text{Follow}(A) = \text{Follow}(A) \cup (\text{First}(\gamma) - \{\epsilon\})$

if $\epsilon \in \text{First}(\gamma)$ then add $\text{Follow}(B)$ to $\text{Follow}(A)$

if $b \in \text{FOLLOW}(B)$, then $S \Rightarrow^* \dots Bb \dots$,
 because $B \rightarrow \alpha A \gamma$, and $\gamma \Rightarrow^* \dots$,
 so $S \Rightarrow^* \dots Bb \dots \Rightarrow^* \dots \alpha A \gamma b \dots \Rightarrow^* \dots \alpha A b \dots$,
 that is $S \Rightarrow^* \dots \alpha A b \dots$, $b \in \text{FOLLOW}(A)$

because $S \Rightarrow^* S$,
 $\$ \in \text{FOLLOW}(S)$

G[S]:
 [1] S → AB
 [2] S → bC
 [3] A → b
 [4] A → ε
 [5] B → aD
 [6] B → ε
 [7] C → AD
 [8] C → b
 [9] D → aS
 [10] D → c

First sets of nonterminals are:

First(S) = {a, b, ε} First(A) = {b, ε}

First(B) = {a, ε} First(C) = {a, b, c}

First(D) = {a, c}

| | Origin | Pass 1 | Pass 2 |
|---|--------|--------|--------|
| S | \$ | \$ | \$ |
| A | | a \$ c | a \$ c |
| B | | \$ | \$ |
| C | | \$ | \$ |
| D | | \$ | \$ |

• Compare First set with Follow set

- ε is an element of First set but never of Follow set
- First set is defined for nonterminals and strings of terminals and nonterminals, while Follow set is defined only for nonterminals
- The definition of Follow set works “on the right” of production, while the definition of the First set works “on the left”

4 Recognize basing on the definition of LL(1)

Condition of LL(1) :

1. For each production $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$,
 $\text{First}(\alpha_i) \cap \text{First}(\alpha_j) = \emptyset$, for all i and j, $1 \leq i, j \leq n, i \neq j$
2. For every nonterminal A such that
 $\text{First}(A)$ contains ε,
 $\text{First}(A) \cap \text{Follow}(A) = \emptyset$

G[S]:

S → AB | bC

A → b | ε

B → aD | ε

C → AD | b

D → aS | c

| | nullable | First | Follow |
|---|----------|-----------|------------|
| S | yes | {a, b, ε} | { \$ } |
| A | yes | {b, ε} | {a, c, \$} |
| B | yes | {a, ε} | { \$ } |
| C | no | {a, b, c} | { \$ } |
| D | no | {a, c} | { \$ } |

➤ S → AB | bC : First(AB) = {a, b, ε}, First(bC) = {b}

First(AB) ∩ First(bC) = {b} ≠ ∅

➤ C → AD | b : First(AD) = {b, a, c}, First(b) = {b}

First(AD) ∩ First(b) = {b} ≠ ∅

➤ This grammar is not LL(1)



4.3 Non LL(1) grammar to LL(1) grammar

- Non LL(1) grammar
 - If a grammar has either **left factor** or **left recursion**, or both, then it must be non LL(1) grammar
 - Whereas, a grammar that doesn't have **left factor** and **left recursion** is not always a LL(1) grammar

1. Left Factor

- Left factor is two or more grammar rule choices sharing a common prefix string, as in the rule
 $A \rightarrow \alpha\beta | \alpha\gamma$
- Because $\text{First}(\alpha\beta) \cap \text{First}(\alpha\gamma) \neq \emptyset$, so it's not LL(1)

2. Left Recursion

A grammar is left recursive if its productions have the following forms:

a) $A \rightarrow A\beta$

b) $A \rightarrow B\beta \quad B \rightarrow A\alpha$

- ❖ a) is called **immediate left recursion**, where the left recursion occurs only within the production of a single nonterminal
- ❖ b) is called **indirect left recursion**, where $A \Rightarrow B\beta \Rightarrow A\alpha\beta$, that is $A \Rightarrow^* A\dots$

Take **immediate left recursion** for example,

if there are productions: $A \rightarrow A\alpha \mid A \rightarrow \beta$

where α and β are arbitrary strings

because $\text{First}(A\alpha) \supseteq \text{First}(\beta)$

so it is not a LL(1) grammar

Techniques for rewriting non LL(1) grammar into LL(1)

- Left recursion removal
- Left factoring

Notice: There is no guarantee that the application of these techniques will turn a grammar into LL(1) grammar

1 Left Recursion Removal

Immediate left recursion removal

➤ Simple case

$A \rightarrow A\alpha \mid \beta$ rewrite this rule into

$A \rightarrow \beta A', A' \rightarrow \alpha A' \mid \epsilon$

➤ General case

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

rewrite this rule into:

$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$

$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$

2 Left Factoring

➤ Simple case

Rewrite the rule $A \rightarrow \alpha\beta \mid \alpha r$ as

$A \rightarrow \alpha(\beta \mid r)$, let A' represents $\beta \mid r$, we get

$A \rightarrow \alpha A' \quad A' \rightarrow \beta \mid r$

α must be the longest string shared by the right-hand sides.

➤ General case

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

rewrite the rule as:

$A \rightarrow \alpha A' \quad A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

➤ Example

G1: $S \rightarrow aSb \mid aS \mid \epsilon$

Left factoring:

$S \rightarrow aS(b \mid \epsilon) \mid \epsilon$

The same as :

$S \rightarrow aSS' \mid \epsilon$

$S' \rightarrow b \mid \epsilon$



4.4 Top-Down Parsing by Recursive-Descent

- Main Idea of Recursive-Descent
 - Define a procedure for each nonterminal **A** that will recognize **A**
 - The right-hand side of the grammar rule for **A** specifies the structure of the code for this procedure
 - **Terminals** correspond to matches of the input
 - **Nonterminals** correspond to calls to other procedures
 - **Choices** correspond to alternatives (case or if statement) within the code

We will talk about:

- General method for constructing Recursive-Descent parser
- A simpler method based on **EBNF**

• General method for constructing Recursive-Descent parser

1. Determine whether the grammar is LL(1)

Example

Grammar G

$E \rightarrow E+T \mid T$

$T \rightarrow T^*F \mid F$

$F \rightarrow (E) \mid i$

1) After left recursion removal we get grammar G':

$E \rightarrow TE'$

$E' \rightarrow +TE' \mid \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' \mid \epsilon$

$F \rightarrow (E) \mid i$

2) First and Follow Set

| | nullable | First | Follow |
|----|----------|-------------------|-----------------|
| E | no | { (, i } | {), \$ } |
| E' | yes | { +, ϵ } | {), \$ } |
| T | no | { (, i } | { +,), \$ } |
| T' | yes | { *, ϵ } | { +,), \$ } |
| F | no | { (, i } | { *, +,), \$ } |

$E \rightarrow TE'$

$E' \rightarrow +TE' \mid \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' \mid \epsilon$

$F \rightarrow (E) \mid i$

G' is LL(1) grammar

2 Construct Recursive-Descent Parser for G'

- When a grammar is LL(1), we can construct recursive-descent parser for it.
- The parser consists of a main procedure and a group of recursive procedures, each corresponds to a nonterminal of the grammar

Construction of procedure for a nonterminal

- Sub procedures used :
 - **match** is a procedure that matches the current next token with its parameter, advances the input if it succeeds, and declares error if it does not
 - **error** is a procedure that prints an error message and exit
- Variable used:
 - **TOKEN** is a variable that keeps the current next token in the input

- 1) If productions of nonterminal U are $U \rightarrow x_1 \mid x_2 \mid \dots \mid x_n$, and $x_1, \dots, x_n \neq \epsilon$, then the code for procedure U is as follow:
- ```

if TOKEN in First(x_1) then p_x1
else if TOKEN in First(x_2) then p_x2
 else ...

 if TOKEN in First(x_n) then p_xn
 else ERROR

```

- 2) If a production of  $U$  is  $U \rightarrow \epsilon$ , then rewrite code  
 if TOKEN in First( $x_n$ ) then p\_xn else ERROR  
 into  
 if TOKEN in First( $x_n$ ) then p\_xn  
 else if TOKEN not in Follow( $U$ ) then ERROR
- 3) The code for p\_x where  $x = y_1 y_2 \dots y_n$  is:  
 begin p\_y1; p\_y2; ...; p\_yn end  
 if  $y_i \in V_N$  then p\_yi is the call of procedure yi;  
 otherwise, if  $y_i \in V_T$  then p\_yi is match(yi)

#### Example: Recursive-descent parser for $G'$

$E \rightarrow TE'$        $E' \rightarrow +TE' \mid \epsilon$   
 $T \rightarrow FT'$        $T' \rightarrow *FT' \mid \epsilon$   
 $F \rightarrow (E) \mid i$

- (1) program MAIN;      /\* main \*/  
 begin  
 GETNEXT (TOKEN);  
 E ;      /\* call E \*/  
 if TOKEN  $\neq$  '\$' then ERROR  
 end.

- (2) procedure E;      /\*  $E \rightarrow TE'$  \*/  
 begin  
 T;      /\* call T \*/  
 E'      /\* call E' \*/  
 end;
- (3) procedure T;      /\*  $T \rightarrow FT'$  \*/  
 begin  
 F;      /\* call F \*/  
 T'      /\* call T' \*/  
 end;

- (4) procedure E';      /\*  $E' \rightarrow +TE' \mid \epsilon$  \*/  
 begin  
 if TOKEN = '+' then      /\*  $E' \rightarrow +TE'$  \*/  
 begin  
 match('+');  
 T;      /\* call T \*/  
 E'      /\* call E' \*/  
 end  
 else      /\*  $E' \rightarrow \epsilon$  \*/  
 if TOKEN  $\neq$  ')' and TOKEN  $\neq$  '\$' then ERROR  
 end;

- (5) procedure T';      /\*  $T' \rightarrow *FT' \mid \epsilon$  \*/  
 begin  
 if TOKEN = '\*' then      /\*  $T' \rightarrow *FT'$  \*/  
 begin  
 match('\*');  
 F;      /\* call F \*/  
 T'      /\* call T' \*/  
 end  
 else      /\*  $T' \rightarrow \epsilon$  \*/  
 if TOKEN  $\neq$  '+' and TOKEN  $\neq$  ')' and TOKEN  $\neq$  '\$'  
 then ERROR  
 end;

```

(6) procedure F; /* F → (E) | i */
begin
 if TOKEN = '(' then /* F → (E) */
 begin
 match('(');
 E; /* call E */
 match(')')
 end
 else /* F → i */
 if TOKEN = 'i' then match('i')
 else error;
end;

```

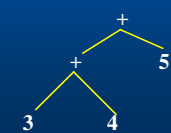
#### Note:

- Left factoring and left recursion removal change the grammar, these changes cause the complication of the parser
- Left factoring and left recursion removal can obscure the semantics of the language structure (for example, they obscure the associativity in arithmetic expressions)

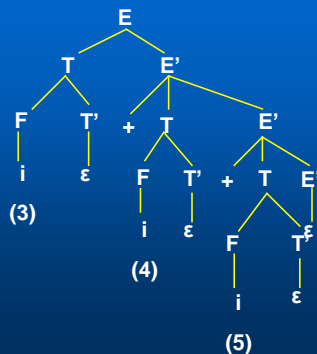
#### Example

$E \rightarrow TE'$   
 $E' \rightarrow +TE' \mid \epsilon$   
 $T \rightarrow FT'$   
 $T' \rightarrow *FT' \mid \epsilon$   
 $F \rightarrow (E) \mid i$

The syntax tree for "3+4+5" is



The parse tree for "3+4+5" is



#### • A simpler method based on EBNF

- Notations of EBNF (extended BNF)

- ❖  $\{ \}$  is used to express repetition

Example:

$A \rightarrow A\alpha \mid \beta$  is written as

$A \rightarrow \beta\{\alpha\}$

- ❖  $[ ]$  is used to surround options

Example:

if-stmt  $\rightarrow$  if (exp) stmt

          | if (exp) stmt else stmt

is written as

if-stmt  $\rightarrow$  if (exp) stmt [else stmt]

#### Example G

$E \rightarrow E+T \mid T$

$T \rightarrow T*F \mid F$

$F \rightarrow (E) \mid i$

rewrite G into EBNF as G':

$E \rightarrow T \{+ T\}$

$T \rightarrow F \{* F\}$

$F \rightarrow (E) \mid i$

### ➤ Construct Recursive-Descent Parsing for G'

- **[ ]** is translated into a test in the code
- **{ }** is translated into the code for a while loop

**if-stmt → if (exp) stmt [else stmt]**

can be translated into the procedure:

```

Procedure ifStmt;
begin
 match('if');
 match('(');
 Exp;
 match(')');
 Stmt;
 if token='else' then
 match('else');
 Stmt;
 end if;
end ifStmt;

```

**E → T {+ T}**

```

(1) procedure E;
begin
 T;
 while token='+' do
 match('+');
 T;
 end while;
end

```

**T → F { \* F }**

```

(2) procedure T;
begin
 F;
 while token='*' do
 match('*');
 F;
 end while;
end

```

### ➤ Action for Constructing a Syntax Tree

#### Sub functions used

- **MakeOpNode** that receives an operator token as a parameter and returns a newly constructed syntax tree node
- **leftChild(t):=p** or **rightChild(t):=p** indicates the assignment of a syntax tree **p** as a left or right child of a syntax tree **t**

```

function E:syntaxTree; /* E → T {+ T} */
var temp,newtemp:syntaxTree;
begin
 temp:=T;
 while token='+' do
 newtemp:=makeOpNode(token);
 match('+');
 leftChild(newtemp):=temp;
 rightChild(newtemp):=T;
 temp:=newtemp;
 end while;
 return temp;
end

```

```

function ifStmt: syntaxTree;
/* if-stmt → if (exp) stmt [else stmt] */
var temp :syntaxTree;
begin
 match('if');
 match('(');
 temp:=makeStmtNode('if');
 testChild(temp):=Exp;
 match(')');
 thenChild(temp):=Stmt;
 if token='else' then
 match('else');
 elseChild(temp):=Stmt;
 else
 elseChild(temp):=nil;
 end if;
end
end

```

**Note:**

- ❖ The method of turning grammar rules in EBNF into code is quite powerful
- ❖ But it may be difficult to convert a grammar originally written in BNF into EBNF form

**Character of recursive-descent parsing****Merit:**

- ❖ Versatile, powerful
- ❖ Flexible by allowing the programmer to adjust the scheduling of the actions
- ❖ Suitable for hand-generated parsers

**Pitfall:**

- ❖ Care must be taken in scheduling the actions within the code
- ❖ Recursive calls result in slowness and much more space

**Quiz**

What are the first and follow sets of S?

$$S \rightarrow A ( S ) B \mid \epsilon$$

$$A \rightarrow S \mid S B \mid x \mid \epsilon$$

$$B \rightarrow S B \mid y$$

- ☐ A First: {x, '('}, Follow: {\$, y, x}
- ☒ B First: {x, y, '(',  $\epsilon$ }, Follow: {\$, y, x, '(', ')'
- ☐ C First: {x, y, '(',  $\epsilon$ }, Follow: {y, x, '(', ')'
- ☐ D First: {x, y, '(',  $\epsilon$ }, Follow: {\$, '(', y}
- ☐ E First: {x,  $\epsilon$ }, Follow: {\$, y, x, '(', ')'
- ☐ F First: {x, y, '(', ')'}, Follow: {\$, y, x, '(', ')'

**Questions: Predictive Parsing****Left-factor the grammar** $S \rightarrow aSb|aS|\epsilon$  $S \rightarrow aS(b|\epsilon)|\epsilon$ 

$$S \rightarrow aSS'|\epsilon$$

$$S' \rightarrow b|\epsilon$$
**In-video Quiz: Left Recursion**

Choose the grammar that correctly

eliminates left recursion from the given grammar:  $E \rightarrow E + T \mid T$   
 $T \rightarrow id \mid (E)$ 

- ☐ A  $E \rightarrow E + id \mid E + (E)$   
 $\mid id \mid (E)$
- ☒ B  $E \rightarrow TE'$   
 $E' \rightarrow +TE' \mid \epsilon$   
 $T \rightarrow id \mid (E)$
- ☐ C  $E \rightarrow E' + T \mid T$   
 $E' \rightarrow id \mid (E)$   
 $T \rightarrow id \mid (E)$
- ☐ D  $E \rightarrow id + E \mid E + T \mid T$   
 $T \rightarrow id \mid (E)$

**4.5 LL(1) Parsing****The meaning of LL(1)**

- The first “L” refers to the fact that it processes the input from left to right
- The second “L” refers to the fact it traces out a leftmost derivation for the input string
- The number “1” means that it uses only one symbol of input to predict the direction of the parse

### 4.5.1 The Basic Method of LL(1) Parsing

- LL(1) parsing uses an explicit stack rather than recursive calls to perform a parse
  - Stack stores the symbols waiting to be matched in parsing
  - Parsing begins by pushing the start symbol into the stack
  - The stack and the input will be empty at the end of a successful parsing

- A schematic for a successful LL(1) parsing is:

**\$StartSymbol**    **InputString\$**

...                      ...

**\$**                      **\$**                      **accept**

'\$' is used to mark the bottom of the stack and the end of the input string

- Example  
 $S \rightarrow (S)S \mid \epsilon$

The actions of LL(1) parser for string "( )"

|   | Parsing stack | Input | Action of parser         |
|---|---------------|-------|--------------------------|
| 1 | \$ S          | ( )\$ | $S \rightarrow (S)S$     |
| 2 | \$ S ) S (    | ( )\$ | match                    |
| 3 | \$ S ) S      | )\$   | $S \rightarrow \epsilon$ |
| 4 | \$ S )        | )\$   | match                    |
| 5 | \$ S          | \$    | $S \rightarrow \epsilon$ |
| 6 | \$            | \$    | accept                   |

### 4.5.2 The LL(1) Parsing Table

- Function  
Expressing production choices for a nonterminal to use at the appropriate parsing step
- Example

| $M[N,T]$ | (                    | )                        | \$                       |
|----------|----------------------|--------------------------|--------------------------|
| S        | $S \rightarrow (S)S$ | $S \rightarrow \epsilon$ | $S \rightarrow \epsilon$ |

### • LL(1) Parsing table

- A two-dimensional array indexed by nonterminals and terminals
- $M[N,t]$  is a production choice of **N** to use at the appropriate parsing step, based on the current input token **t** (terminal or '\$')
- If  $M[N,t]$  remain empty after construction, it represent potential errors that may occur during a parse

### • The Construction of Parsing Tables

Repeat the following two steps for each nonterminal **A** and production choice  $A \rightarrow \alpha$

1. For each token '**a**' in  $\text{First}(\alpha)$ , add  $A \rightarrow \alpha$  to the entry  $M[A,a]$
2. if  $\epsilon$  is in  $\text{First}(\alpha)$ , for each element '**a**' of  $\text{Follow}(A)$  (token or \$), add  $A \rightarrow \alpha$  to  $M[A,a]$

### Example :Grammar of arithmetic expression

$E \rightarrow E + T \mid T$   
 $T \rightarrow T * F \mid F$   
 $F \rightarrow (E) \mid i$

(1) After left recursive removal gets  $G'$

$E \rightarrow TE'$   
 $E' \rightarrow +TE' \mid \epsilon$   
 $T \rightarrow FT'$   
 $T' \rightarrow *FT' \mid \epsilon$   
 $F \rightarrow (E) \mid i$

### (2) First and Follow set for nonterminals

|    | nullable | First            | Follow      |
|----|----------|------------------|-------------|
| E  | no       | {(, i}           | {), \$}     |
| E' | yes      | {+, $\epsilon$ } | {), \$}     |
| T  | no       | {(, i}           | {+), \$}    |
| T' | yes      | {*, $\epsilon$ } | {+), \$}    |
| F  | no       | {(, i}           | {*, +), \$} |

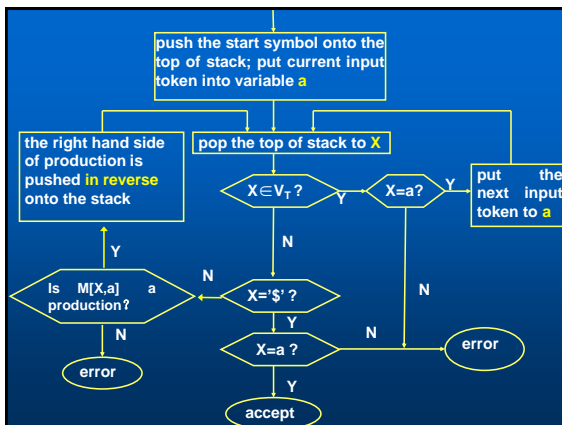
### (3) Construct parsing table

|    | i                   | +                         | *                     | (                   | )                         | \$                        |
|----|---------------------|---------------------------|-----------------------|---------------------|---------------------------|---------------------------|
| 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  | $F \rightarrow i$   |                           |                       | $F \rightarrow (E)$ |                           |                           |

### • LL(1) grammar

A grammar is a LL(1) grammar if the associated LL(1) parsing table has at most one production in each table entry.

## 4.5.3 The LL(1) parsing algorithm



### Parsing process of string "i+i\*i\$"

|    | i                   | +                         | *                     | (                   | )                         | \$                        |
|----|---------------------|---------------------------|-----------------------|---------------------|---------------------------|---------------------------|
| 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  | $F \rightarrow i$   |                           |                       | $F \rightarrow (E)$ |                           |                           |

| Step | Stack   | Input   | Action                    |
|------|---------|---------|---------------------------|
| 1    | \$E     | i+i*i\$ | $E \rightarrow TE'$       |
| 2    | \$E'T   | i+i*i\$ | $T \rightarrow FT'$       |
| 3    | \$E'T'F | i+i*i\$ | $F \rightarrow i$         |
| 4    | \$E'T'i | i+i*i\$ | match i                   |
| 5    | \$E'T'  | +i*i\$  | $T' \rightarrow \epsilon$ |
| 6    | \$E'    | +i*i\$  | $E' \rightarrow +TE'$     |
| 7    | \$E'T+  | +i*i\$  | match +                   |

|    | i                   | +                         | *                     | (                   | )                         | \$                        |
|----|---------------------|---------------------------|-----------------------|---------------------|---------------------------|---------------------------|
| 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  | $F \rightarrow i$   |                           |                       | $F \rightarrow (E)$ |                           |                           |

|    |            |         |                           |
|----|------------|---------|---------------------------|
| 8  | $\$E'T$    | $i'i\$$ | $T \rightarrow FT'$       |
| 9  | $\$E'T'F$  | $i'i\$$ | $F \rightarrow i$         |
| 10 | $\$E'T'i$  | $i'i\$$ | match i                   |
| 11 | $\$E'T'$   | $*i\$$  | $T' \rightarrow *FT'$     |
| 12 | $\$E'T'F*$ | $*i\$$  | match *                   |
| 13 | $\$E'T'F$  | $i\$$   | $F \rightarrow i$         |
| 14 | $\$E'T'i$  | $i\$$   | match i                   |
| 15 | $\$E'T'$   | $\$$    | $T' \rightarrow \epsilon$ |
| 16 | $\$E'$     | $\$$    | $E' \rightarrow \epsilon$ |
| 17 | $\$$       | $\$$    | accept                    |

#### 4.6 Error Recovery in Top-down Parsers

- A parser must determine whether a program is syntactically correct or not
- The response of a parser to syntax errors is often a critical factor in the usefulness of a compiler

#### • Normal syntactical errors

##### ➤ Start and following symbols error

The start and following symbol of program, expressions and statements for example, the TINY language

|           | start symbol                        | following symbol |
|-----------|-------------------------------------|------------------|
| statement | if, repeat, identifier, read, write | ; else end until |
| exp       | (, number, identifier               | ) then ;         |

##### ➤ Identifiers and constants error

Such as “const”, “var”, “procedure” is not followed by identifiers

##### ➤ Parenthesis matching error

Such as, begin—end, if—then don't match

##### ➤ Symbols error

Such as symbol in assignment is not ‘:=’

#### • Error handling

- Give a meaningful error message
- Pick a likely place to resume the parse.  
A parser should always try to parse as much of the code as possible, in order to find as many real errors as possible during a single translation

#### • Categories of Error Handling

- ✦ **Error Recovery**: after an error has occurred, the parser picks a likely place to resume the parsing
- ✦ **Error Repair**: the parser attempts to infer a correct program from the incorrect one given

#### 4.6.1 Error Recovery in Recursive-Descent Parsers

##### • Panic Mode

In complex situations, the error handler will consume a possibly large number of tokens in an attempt to find a place to resume parsing

##### • The Basic Mechanism of Panic Mode

- Provide each recursive procedure with an extra parameter consisting of a set of **synchronizing tokens**
- As parsing proceeds, tokens that may function as synchronizing tokens are added to this set as each call occurs
- If an error is encountered, the parser **scans ahead**, throwing away tokens until one of the synchronizing set of tokens is seen in the input

Follow sets are important candidates for synchronizing tokens

| Example                                                            | Assume <b>SynchSetS</b> is the synchronizing set of <b>S</b> |
|--------------------------------------------------------------------|--------------------------------------------------------------|
| $S \rightarrow id := E_1$                                          |                                                              |
| $S \rightarrow \text{if } E_2 \text{ then } S_1 \text{ else } S_2$ | $\text{SynchSet}E_1 = \text{SynchSet}S$                      |
| $S \rightarrow \text{while } E_3 \text{ do } S_3$                  | $\text{SynchSet}E_2 = \{\text{then}\}$                       |
| $S \rightarrow \text{repeat } S_4 \text{ until } E_4$              | $\text{SynchSet}S_1 = \{\text{else}\}$                       |
| $S \rightarrow \text{begin SL end}$                                | $\text{SynchSet}S_2 = \text{SynchSet}S$                      |
|                                                                    | $\text{SynchSet}E_3 = \{\text{do}\}$                         |
|                                                                    | $\text{SynchSet}S_3 = \text{SynchSet}S$                      |
|                                                                    | $\text{SynchSet}S_4 = \{\text{until}\}$                      |
|                                                                    | $\text{SynchSet}E_4 = \text{SynchSet}S$                      |
|                                                                    | $\text{SynchSetSL} = \{\text{end}\}$                         |

- ❖ First sets may also be used to prevent the error handler from skipping important tokens that begin major new constructs
- ❖ First sets are also important, in that they allow a recursive descent parser to detect errors early in the parse

##### • Realization of Panic Mode

- 1) At the beginning of each procedure, check whether the current input token is in **FirstSet**. If not, scan ahead until find the first token in  $\text{FirstSet} \cup \text{SynchSet}$ , parsing is resumed from this token
- 2) Before returning from a procedure, check whether the next token is in **SynchSet**. If not, scan ahead until find the first token in  $\text{FirstSet} \cup \text{SynchSet}$ , parsing is resumed from this token

##### Example

```

E -> T {+ T}
T -> F {* F}
F -> (E) | i

```

```

procedure checkinput(firstset, followset);
begin
 if not (token in firstset) then
 error;
 getToken;
 while not (token in (firstset \cup followset))
 do getToken;
 end if;
 end;
end;

```



```
E-> T {+ T}
procedure E(synchset);
begin
 checkinput({ (, i }, synchset);
 if not (token in synchset) then
 T(synchset);
 while token='+' do
 match('+');
 T(synchset);
 end while;
 checkinput(synchset, { (, i });
end if;
end
```