



# 编译原理

## Compiler Principles

### Lecture 4

### Syntax Analysis: Top-Down

---

赵帅

计算机学院  
中山大学

# Before we start...

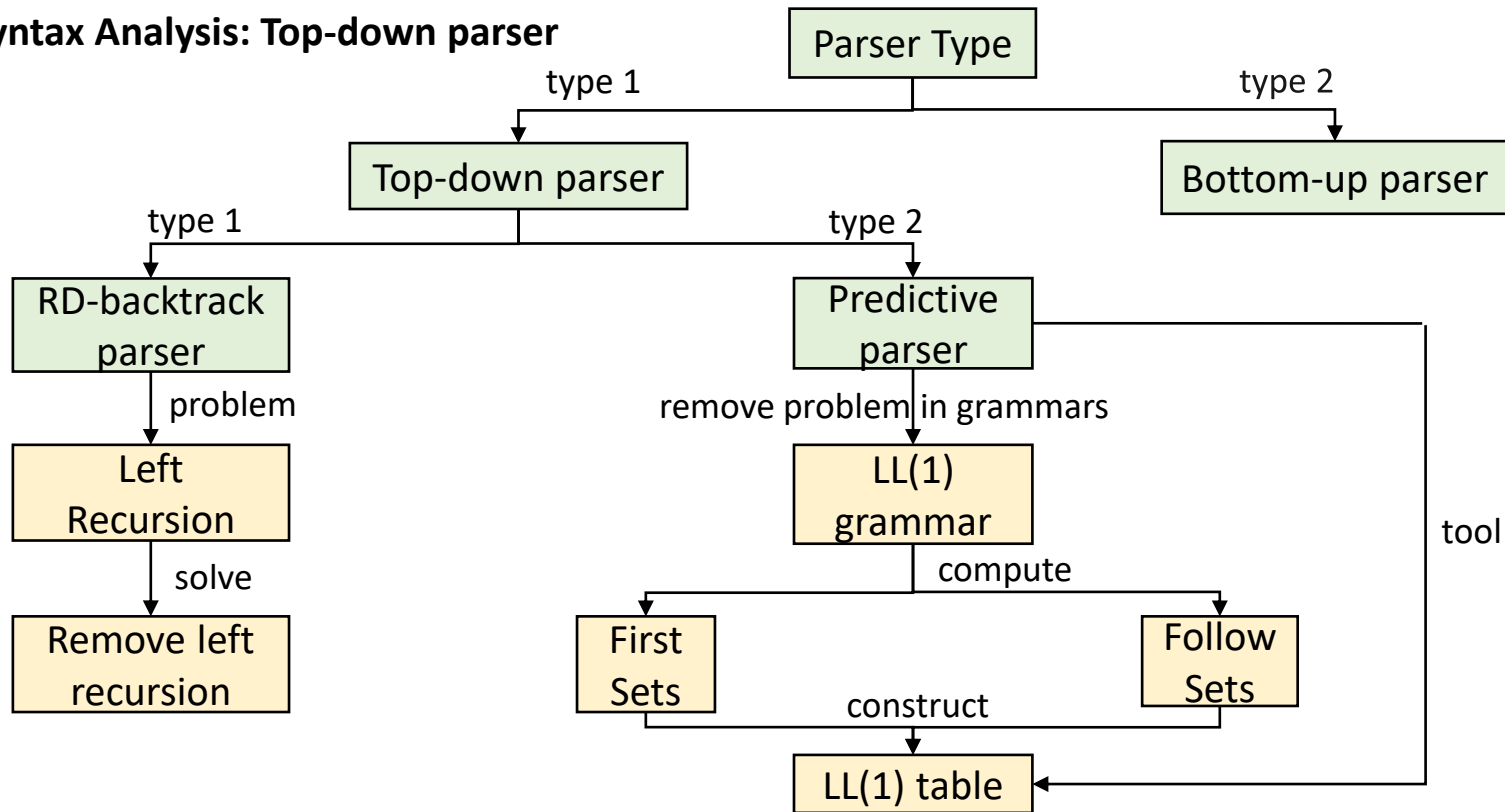


So, please pay attention and keep on track!

# Mind Map[思维导图]



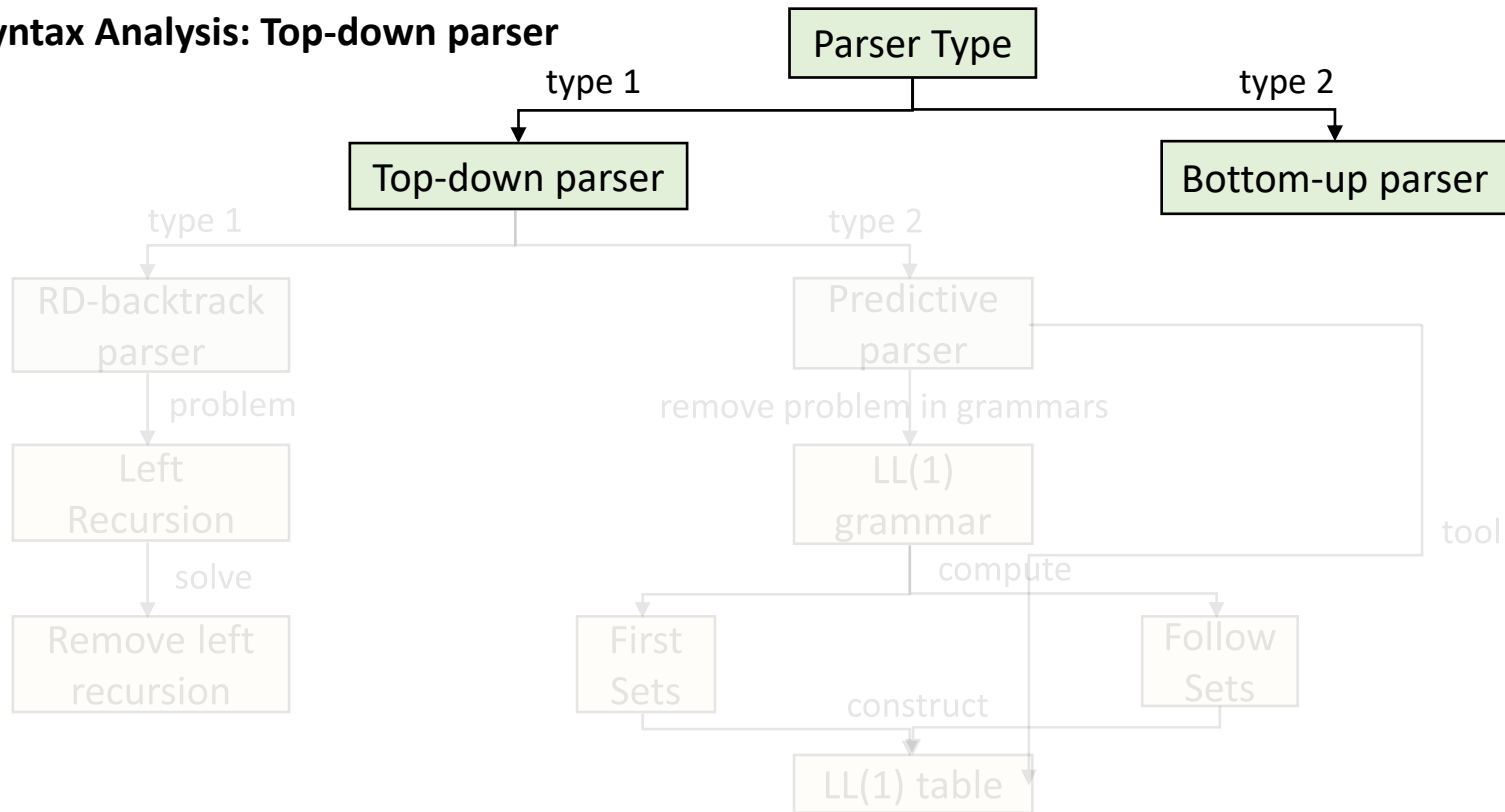
## Syntax Analysis: Top-down parser



# Mind Map[思维导图]



## Syntax Analysis: Top-down parser



- Most compilers use either **Top-Down** or **Bottom-Up** parsers.
- Bottom-up parsing [自底向上分析]
  - ◆ Begin at the **leaves** (the bottom) and working up towards the **root** (the top).
  - ◆ Tries to reduce[规约] **the input string** to **the start symbol**.
  - ◆ Finds reverse order of the rightmost derivation[最右推导的逆过程即最左规约]  
规范推导 规范规约
  - ◆ Parser code structure nothing like grammar.
    - Very difficult to implement manually.
    - Automated tools exist to convert to code (e.g., Yacc, Bison).

- Top-Down parsing [自顶向下分析]
  - ◆ Starting from the **root** (top) and create the leaves (down) of the parse tree in a pre-defined order (depth-first) [深度优先, 先根次序/前序].
  - ◆ Top-down parsing can be viewed as **finding a leftmost derivation** [寻求最左推导] for an input string. *Why not rightmost or arbitrary derivation?*
  - ◆ **Review:** In each step of derivation, the following choices need to be made:
    - Choice of the non-terminal to be replaced. [替换哪个非终结符] **Leftmost!**
    - Choice of the production to be applied for a non-terminal. [使用文法中哪个规则来替换] **Key Problem!**
  - ◆ **Question:** At each step of a top-down parse, What is the key problem?

- Top-Down parsing [自顶向下分析]
  - ◆ Once a production is chosen, we try to match the terminal symbols in the production body with the input string.
  - ◆ Parser code structure closely mimics grammar.
    - Manually implementation is feasible.
    - Automated tools exist to convert to code. (e.g. ANTLR)
- Top-Down **vs.** Bottom-Up [对比]
  - ◆ Top-down: easier to understand and implement manually. (E.g. ANTLR)
  - ◆ Bottom-up: more powerful, can be implemented automatically. (E.g. YACC/Bison)

# Parser Type[语法分析类型]

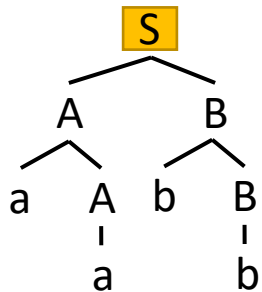


## ◆Example

◆ Grammar  $G(S)$ :  $S \rightarrow AB$  ;  $A \rightarrow aA \mid a$  ;  $B \rightarrow bB \mid b$  ;

◆ Sentence: aabb;

$S \Rightarrow AB$  (1)  
 $\Rightarrow aAB$  (2)  
 $\Rightarrow aaB$  (3)  
 $\Rightarrow aabB$  (4)  
 $\Rightarrow aabb$  (5)



$S \Rightarrow AB$  (5)  $B$  reduce to  $S$ .  
 $\Rightarrow AbB$  (4)  $bB$  reduce to  $B$ .  
 $\Rightarrow Abb$  (3) 2<sup>nd</sup>  $b$  reduce to  $B$ .  
 $\Rightarrow aAbb$  (2)  $aA$  reduce to  $A$ .  
 $\Rightarrow aabb$  (1) 2<sup>nd</sup>  $a$  reduce to  $A$ .

Leftmost derivation

Top-Down

Leftmost reduction

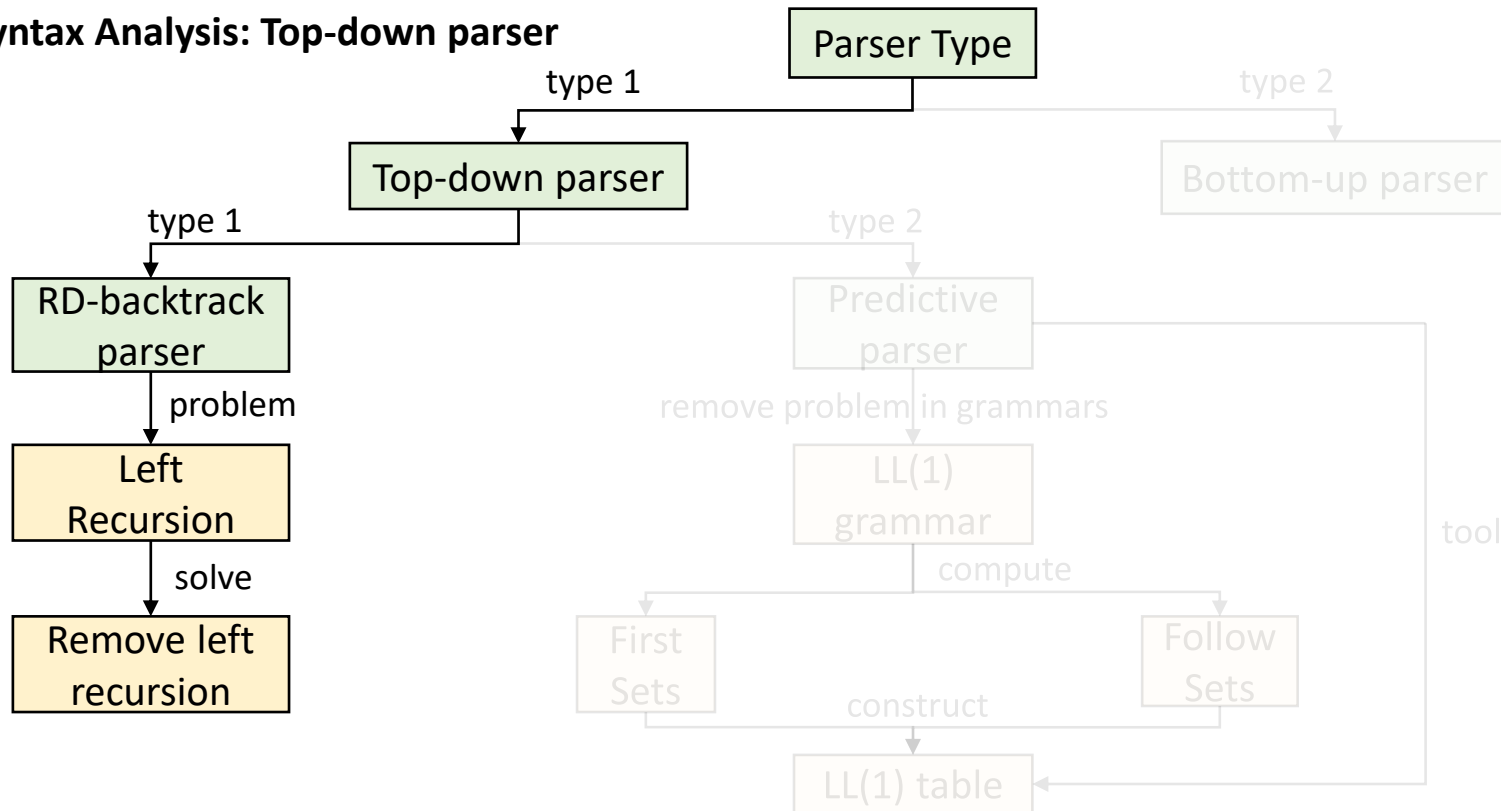
Bottom-Up



# Revisit



## Syntax Analysis: Top-down parser



# Top-down Parsing[自顶向下分析]



- **Recursive-descent parsing**[RDP, 递归下降语法分析]

- ◆ A general form[通用形式] of top-down parsing.
- ◆ A recursive-descent parsing program consists of a set of procedures, one for each non-terminal.
- ◆ Execution begins with the procedure for the start symbol, which halts and announces success if its procedure body scans the entire input string.

```
void A() {  
1)   Choose an A-production,  $A \rightarrow X_1X_2 \cdots X_k$ ;  
2)   for (  $i = 1$  to  $k$  ) {  
3)       if (  $X_i$  is a nonterminal )  
4)           call procedure  $X_i()$ ;  
5)       else if (  $X_i$  equals the current input symbol  $a$  )  
6)           advance the input to the next symbol;  
7)       else /* an error has occurred */;  
    }  
}
```

How to choose A-production?

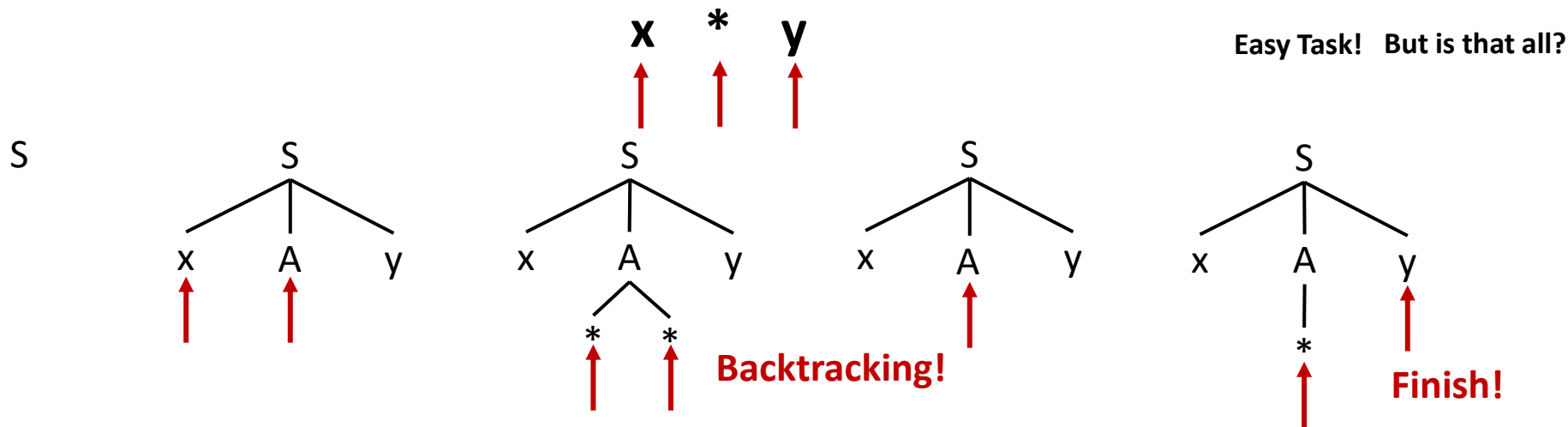
It's not specified, so the pseudocode is nondeterministic[伪代码是不确定的].

# RDP with backtracking[回溯]

- RDP may require *backtracking*.
- **Approach:** for a non-terminal in the derivation, productions are tried in some order until
  - ◆ A production is found that generates a portion of the input, or
  - ◆ No production is found that generates a portion of the input, in which case backtrack to previous non-terminal.
- Terminals of the derivation are compared against input
  - ◆ Match: advance input, continue parsing
  - ◆ Mismatch: backtrack, or fail
- Parsing fails if no derivation generates the entire input.

# RDP with backtracking[回溯]

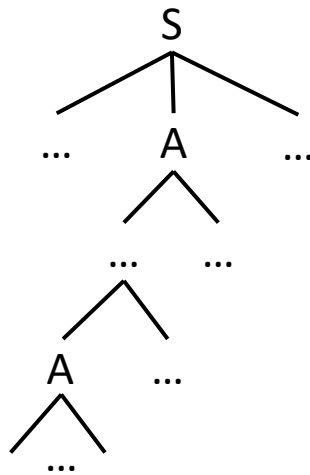
- In the analysis process, when a non-terminal is successfully matched with an alternative, **the match may be temporary**.
- If a mismatch occurs, **backtracking[回溯]** will be performed.
- $G[S]: S \rightarrow xAy; A \rightarrow ** \mid *$ , whether the input string  $x * y$  is its sentence?



# Left Recursion Problem[左递归问题]



- A grammar is **left recursive**[左递归] if it has a non-terminal **A** such that there is a derivation  $A \xRightarrow{+} A\alpha$ .
- Sentence can grow infinitely without consuming input.(Into an infinite loop!)
- **Top-down parsing methods cannot handle left-recursive problems.**[自顶向下语法分析方法不能处理左递归的文法]



# Left Recursion Problem[左递归问题]



- **Immediate** left recursion [直接/立即左递归]
  - ◆ There is a production  $A \rightarrow A\alpha$ .
- **Non-immediate** left recursion [间接/非立即左递归]
  - ◆ Left recursion involving derivation of 2+ step.
    - ◆  $A \rightarrow B\beta; B \rightarrow A\alpha$ .
- A transformation is needed to eliminate left recursion. [需要一个转换方法来消除左递归]
- Rewrite the grammar so that it is right recursive. [改为右递归]

# Remove Left Recursion[消除左递归]



- Immediate left recursion[直接左递归的消除]

- ◆ Grammar:  $A \rightarrow A\alpha \mid \beta$  ( $\alpha \neq \beta$ ,  $\beta$  doesn't start with  $A$ )

- ◆ rewrite the rule of  $A$  as the following form equivalently:

- ◆ Grammar:  $A \rightarrow \beta A'; A' \rightarrow \alpha A' \mid \epsilon$  (right recursion)

$G[A]: A \rightarrow A\alpha \mid \beta$

$A \Rightarrow A\alpha$

$\Rightarrow A\alpha\alpha$

$\Rightarrow A\alpha\alpha\alpha$

.....

$\Rightarrow A\alpha...\alpha$

$\Rightarrow \beta\alpha...\alpha$

Remove Left Recursion



$G[A]: A \rightarrow \beta A'$

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

$A \Rightarrow \beta A'$

$\Rightarrow \beta\alpha A'$

$\Rightarrow \beta\alpha\alpha A'$

.....

$\Rightarrow \beta\alpha...\alpha A'$

$\Rightarrow \beta\alpha...\alpha$

想办法把“ $A\alpha\alpha..$ ”变成“ $\alpha\alpha..A$ ”

$A \rightarrow A\alpha; A \rightarrow \beta$

$\Rightarrow A\alpha \Rightarrow A\alpha\alpha\alpha.. \Rightarrow \beta\alpha\alpha\alpha..$

1. 避免对 $A$ 的左递归

$A \rightarrow \beta x$ , where  $x \rightarrow \alpha\alpha\alpha..$

2. 引入bridging production

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

$\Rightarrow \alpha A' \Rightarrow \alpha\alpha..A' \Rightarrow \alpha\alpha..\alpha$

3. 连接两个productions

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

$\Rightarrow \beta A' \Rightarrow \beta\alpha\alpha..A' \Rightarrow \beta\alpha\alpha..\alpha$

# Remove Left Recursion[消除左递归]



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

◆ First, group the productions as

$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \beta_2 \dots \mid \beta_n$  where no  $\beta_i$  begins with an A.

◆ Then, replace the A-productions by

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

- Exercise-Remove Immediate left recursion.

◆  $G_1[E]: E \rightarrow \overset{\alpha}{\boxed{E+T}} \mid \overset{\beta}{\boxed{T}};$

$T \rightarrow \overset{\alpha}{\boxed{T * F}} \mid \overset{\beta}{\boxed{F}};$

$F \rightarrow (E) \mid i.$

◆  $G_2[E]: E \rightarrow TE'; E' \rightarrow +TE' \mid \epsilon;$

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

$F \rightarrow (E) \mid i.$



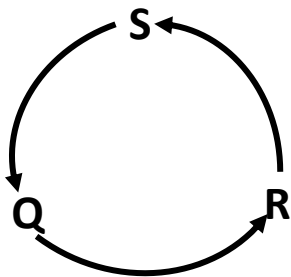
# Remove Left Recursion[消除左递归]



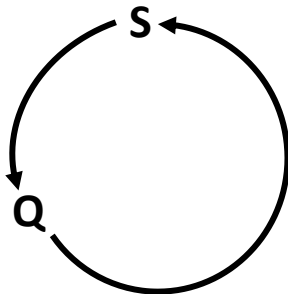
- **Non-Immediate left recursion**[非直接左递归的消除]

- ◆ Grammar:  $S \rightarrow Qc \mid c$ ;  $Q \rightarrow Rb \mid b$ ;  $R \rightarrow Sa \mid a$ .

- ◆ Although there is no immediate left recursion,  **$S$ ,  $Q$  and  $R$  are all left recursion.**

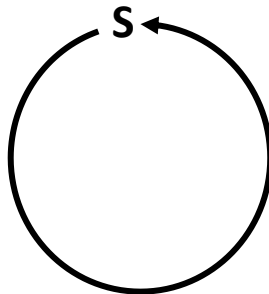


$S \rightarrow Qc \mid c$   
 $Q \rightarrow Rb \mid b$   
 $R \rightarrow Sa \mid a$



$S \rightarrow Qc \mid c$   
 $Q \rightarrow Sab \mid ab \mid b$   
 $R \rightarrow Sa \mid a$

**Left Recursion!**



$S \rightarrow Sabc \mid abc \mid bc \mid c$   
 $Q \rightarrow Sab \mid ab \mid b$   
 $R \rightarrow Sa \mid a$

# Remove Left Recursion[消除左递归]



- The following algorithm **systematically eliminates left recursion** from a grammar[直接/间接]. It is guaranteed to work if:
  - ◆ the grammar has no cycles (derivations of the form  $A \xRightarrow{+} A$ )
  - ◆ the grammar has no  $\epsilon$ -productions (productions of the form  $A \rightarrow \epsilon$ ).
  - ◆ These two can be eliminated systematically from a grammar.
- **Algorithm:** Eliminating left recursion.
  - ◆ **INPUT:** Grammar  $G$  with no cycles or  $\epsilon$ -productions.
  - ◆ **OUTPUT:** An equivalent grammar with no left recursion.
  - ◆ **METHOD:** Apply the following 3 steps to  $G$ . Note that the resulting non-left-recursive grammar may have  $\epsilon$ -productions.

# Remove Left Recursion[消除左递归]



◆ **Step 1** : Arrange all non-terminals of grammar G **in some order**  $A_1, A_2, \dots, A_n$ ;

◆ **Step 2** : Execute in order obtained in Step 1:

FOR  $i:=1$  TO  $n$  DO

FOR  $j:=1$  TO  $i-1$  DO

Replace each production of the form  $A_i \rightarrow A_j \gamma$  by the productions

$A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \delta_k \gamma$ , where  $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \delta_k$  are all current  $A_j$ -productions;  
[间接左递归问题 转换为 直接左递归问题]

**eliminate the immediate left recursion among the  $A_i$ -productions;**

[解决直接左递归问题]

END

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

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

END

◆ **Step 3** : Simplify the grammar obtained from Step 2 --- remove the production rules of non-terminal that can never be reached from the start symbol.

[去掉无效productions]

# Remove Left Recursion[消除左递归]



- **Example:** Consider Grammar  $G(S)$ : (**ILR = immediate left recursion**)

$$R \rightarrow Sa \mid a$$

$$Q \rightarrow Rb \mid b$$

$$S \rightarrow Qc \mid c$$

- **Step 1** : non-terminal order: **R,Q,S**; ( $A_1, A_2, A_3$ ) [顺序决定了哪个non-terminal可能会被消除 (如果符合算法要求), 最后的non-terminal会留下]  
[ $A_1=R$ ]
- **Step 2** : ( $i=1$ ) For **R**, there is no **ILR**;

( $i=2, j=1$ ) Replace production  **$Q \rightarrow Rb$**  by the productions **R** [ $A_i=Q, A_j=R$ ]  
 **$\rightarrow Sa \mid a$** , which generates  **$Q \rightarrow Sab \mid ab \mid b$** ; **not ILR**

( $i=3, j=1$ ) No operations. [ $A_i=S, A_j=R$ ]

( $i=3, j=2$ ) Replace production  **$S \rightarrow Qc$**  by the productions  **$Q \rightarrow Sab \mid ab \mid b$**  [ $A_i=S, A_j=Q$ ]  
 **$S \rightarrow Sabc \mid abc \mid bc \mid c$** ; **contain ILR**;

(continue...)

# Remove Left Recursion[消除左递归]



(continue...)

( $\alpha$ )

( $\beta$ )

( $i=3, j=2$ )  $S \rightarrow \text{Sabc} \mid \text{abc} \mid \text{bc} \mid \text{c}$  contains ILR.

**Eliminate ILR about S:**

$S \rightarrow \text{abcS}' \mid \text{bcS}' \mid \text{cS}'$

$S' \rightarrow \text{abcS}' \mid \epsilon$

$Q \rightarrow \text{Sab} \mid \text{ab} \mid \text{b}$

$R \rightarrow \text{Sa} \mid \text{a}$

- **Step 3** : Simplify the grammar and get **Final Grammar**:

$S \rightarrow \text{abcS}' \mid \text{bcS}' \mid \text{cS}' \quad S' \rightarrow \text{abcS}' \mid \epsilon$

$R \rightarrow \text{Sa} \mid \text{a}$

$Q \rightarrow \text{Rb} \mid \text{b}$

$S \rightarrow \text{Qc} \mid \text{c}$

(Q&R's production is included by S)

# Remove Left Recursion[消除左递归]



- **Question:** What will happen if the order in step 1 is different
- Again, consider  $G(S): S \rightarrow Qc \mid c \quad Q \rightarrow Rb \mid b \quad R \rightarrow Sa \mid a$
- **Exercise:** If the order in step 1 is **S,Q,R**, the final grammar without ILR is?  
(was R, Q, S in the previous example)

$S \rightarrow Qc \mid c$

$Q \rightarrow Rb \mid b$

$R \rightarrow bcaR' \mid caR' \mid aR'$

$R' \rightarrow bcaR' \mid \epsilon$

**S,Q,R**

equivalent



$S \rightarrow abcS' \mid bcS' \mid cS'$

$S' \rightarrow abcS' \mid \epsilon$

**R,Q,S**

- Different order in Step1 may cause the final grammar different in form[可能在形式上不同], but it is not difficult to prove that they are equivalent.

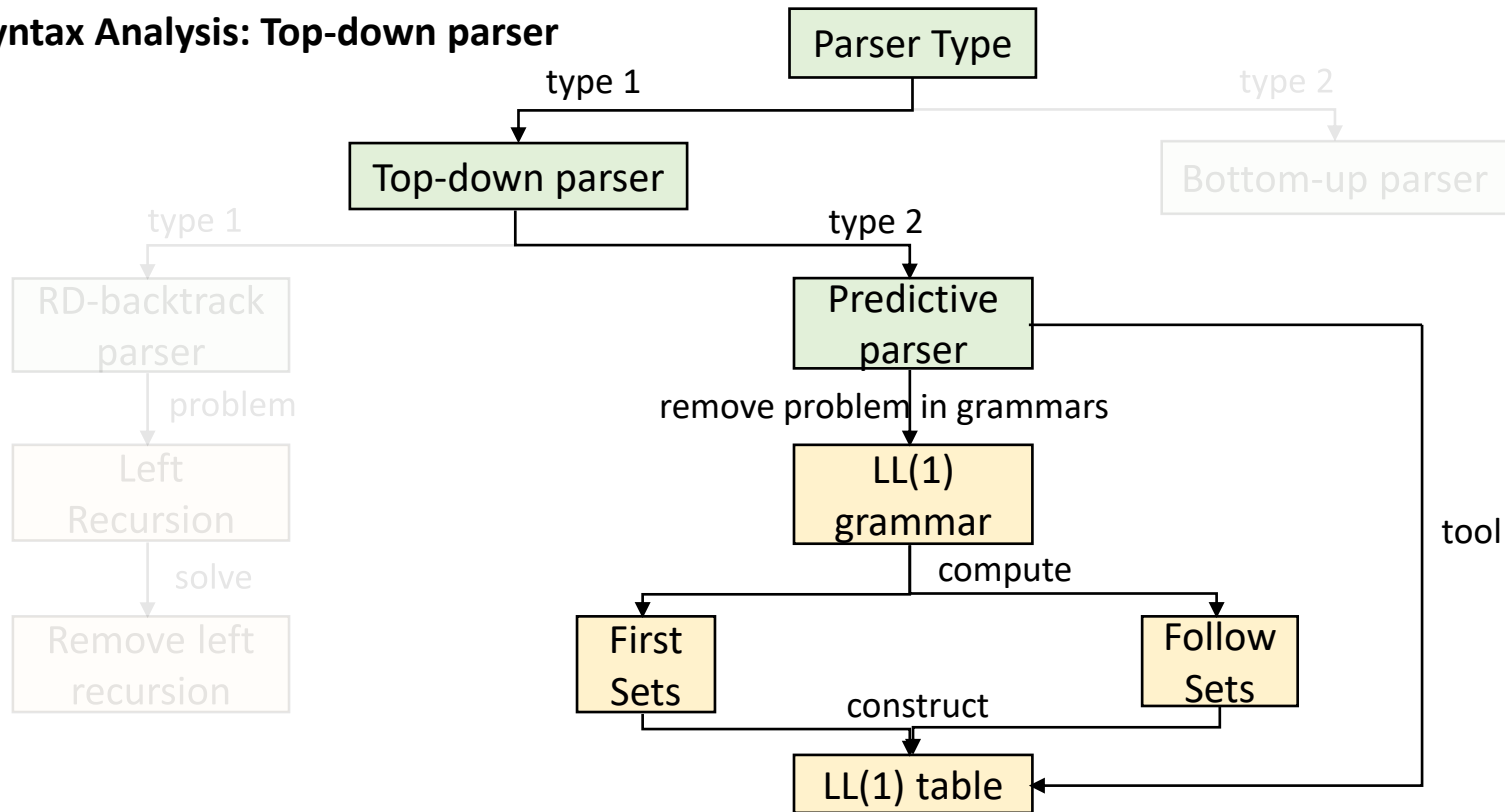
# Recursive-descent parsing[递归下降语法分析]

- **Recursive-descent parsing** [RDP, 递归下降语法分析]
  - ◆ A **general form** of top-down parsing.
  - ◆ RDP is a simple and general parsing strategy.
    - ▢ **Left-recursion must be eliminated first.**(Can be eliminated automatically using some algorithm)
  - ◆ However, it is not popular because of **backtracking**.
    - ▢ Backtracking requires re-parsing the same string.
    - ▢ Tried to solve problem in any possible ways[穷尽一切可能的试探法] which is inefficient.(can take exponential time)
    - ▢ Also removing already added nodes in parse tree is troublesome.

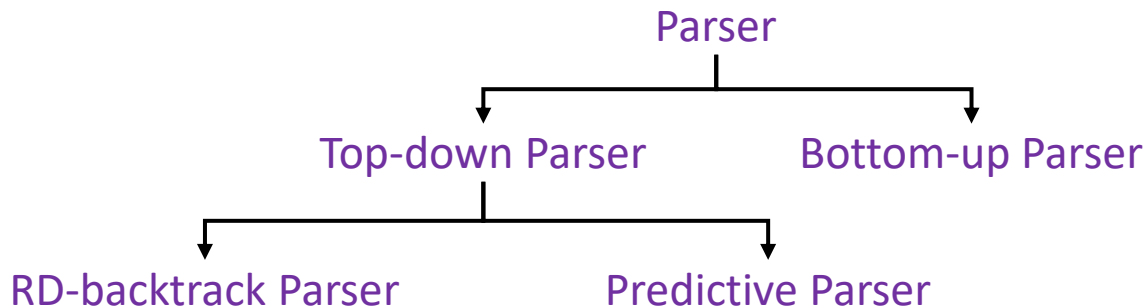
# Mind Map[思维导图]



## Syntax Analysis: Top-down parser







## • Predictive Parsing[预测分析法]

- ◆ A special case of recursive-descent parsing **without backtracking**[无回溯].
- ◆ Predictive parsing **chooses the correct production by looking ahead** at the input a fixed number of symbols, typically we may look only at one.(that is, the next input symbol) [通过向前看一步输入符号，预测正确的产生式]
- ◆ Restrictions on the grammar to **avoid backtracking**. [LL(k)]

- **A parser with no backtracking** [无回溯]: select the correct alternative through given next input terminal(s)[下一个输入符号/终结符]
- A predictive parser chooses the production to apply solely on the basis of [选取产生式的依据]
  - ◆ Next input symbol(s).
  - ◆ Current non-terminal being processed.

$G(S): S \rightarrow aBD \mid bBB; \quad B \rightarrow c \mid bce; \quad D \rightarrow d$   
parsing input “**abced**” requires no backtracking

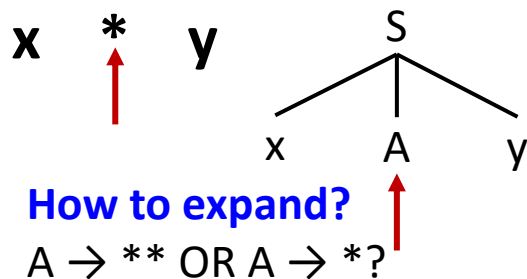
**What if**  $S \rightarrow aBD \mid aBB$      **Common Prefix**[共同前缀]

Given input terminal(s)  $a$ , cannot choose between two rules.

# Common Prefix[共同前綴]



- $G[S]: S \rightarrow xAy; A \rightarrow ** \mid *$ , If the current matching symbol is  $*$ , the next step is to expand  $A$ , and  $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ . How to choose  $\alpha_i$ ?



- ◆ (1) If there is a unique  $\alpha_i$  with  $*$  as the head, replace using this unique  $\alpha_i$ .
- ◆ (2) If there are multiple  $\alpha_i$  with  $*$  as the head, the substitution is not unique[替换是不唯一的], backtracking may happen.

- **Left factoring**[提取左公因子]: Rewrite the productions to defer the decision until enough of the input has been seen that we can make the right choice. [推后决定直至可选择]

# Left factoring[提取公共左因子]



- For each non-terminal  $A$ , find the longest prefix  $\alpha$  common to two or more of its alternatives. If  $\alpha \neq \epsilon$ , **replace** all of the  $A$ -productions  $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma$ , where  $\gamma$  represents all alternatives that do not begin with  $\alpha$ , by  $A \rightarrow \alpha A' \mid \gamma; A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$ .
- Repeatedly apply this transformation until no two alternatives for a non-terminal have a common prefix.
- **Example:**  $G[S]: S \rightarrow abc \mid abd \mid ae;$ 
  - ◆ **Step1:**  $S \rightarrow aA'; A' \rightarrow bc \mid bd \mid e;$
  - ◆ **Step2:**  $S \rightarrow aA'; A' \rightarrow bB' \mid e; B' \rightarrow c \mid d;$

- Patterns in grammars that prevent predictive parsing [并非总是能预测分析]:

- ◆ **Left Recursion Problem** [左递归问题]

- Lookahead symbol changes **only** when a terminal is matched.
- Q: How to solve? A: Remove Left Recursion.

- ◆ **Common Prefix**[共同前缀] —Cause **Backtracking Problem**.

- Q: How to solve? A: Left factoring.

- ◆ **Question:** After left factoring, Can we completely avoid the backtracking and do predictive parsing? **What if**  $S \rightarrow EBD \mid FBB$ ;  $E \rightarrow a\dots$ ;  $F \rightarrow a\dots$

$G(S): S \rightarrow aBD \mid bBB; B \rightarrow c \mid bce; D \rightarrow d$  **FIRST && FOLLOW!**  
parsing input “**abced**” requires no backtracking

- During top-down parsing, **FIRST** and **FOLLOW** allow us to choose which production to apply, based on the next input symbol.
- **FIRST( $\alpha$ )** [终结首符集] ( $\alpha$  is any string of grammar symbols)
  - ◆ Define FIRST( $\alpha$ ) to be the set of terminals that are the first symbol in strings [串的终结首符的集合] derived from  $\alpha$ .  **$\text{FIRST}(\alpha)^* = \{a \mid \alpha \Rightarrow a..., a \in V_T\}$**
  - ◆ AND if  $\alpha \xRightarrow{*} \varepsilon$ , then  $\varepsilon$  is also in FIRST( $\alpha$ ).
- Consider two(or more) A-productions  $A \rightarrow \alpha \mid \beta$ , **where FIRST( $\alpha$ ) and FIRST( $\beta$ ) are disjoint sets** [不相交的集合]. We can then choose between these A-productions by looking at the next input symbol  $a$ , **since  $a$  can be in at most one of FIRST( $\alpha$ ) and FIRST( $\beta$ ), not both.**

- **Question:** And then, if for an input symbol  $a$  and the non-terminal  $A$  ( $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ ), **but for any  $\alpha_i$ ,  $a \notin \text{FIRST}(\alpha_i)$** , in this case, how to choose  $\alpha_i$ , or draw a conclusion of grammatical errors?
- **FOLLOW(A):** For nonterminal  $A$ , the set of terminals  $a$  that can appear immediately to the right of  $A$  in some sentential form; that is,  **$\text{FOLLOW}(A) = \{a \mid S \xRightarrow{*} \dots Aa\dots, a \in V_T\}$** . [在某些句型中紧跟在A右边的终结符号的集合]
  - ◆ If  $A$  is the rightmost symbol in some sentential form, then  $\$$  is in  $\text{FOLLOW}(A)$ ; recall that  $\$$  is a special “endmarker”[结束标记] symbol that is assumed not to be a symbol of any grammar.

[First与Follow的区别]

First( $\alpha$ ):  $\alpha$ 展开的串的第一个终结符

Follow(A): A后面紧跟的终结符 (不包括A展开的串)

- **Question:** for an input symbol  $a$  and the non-terminal  $A$  ( $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ ),  $a \notin \text{FIRST}(\alpha_i)$ , in this case, how to choose  $\alpha_i$ , or draw a conclusion of grammatical errors?

- **Answer:** if  $\epsilon \in \text{FIRST}(\alpha_i)$ , then when  $a \in \text{FOLLOW}(A)$ , select  $A \rightarrow \alpha_i$ , otherwise, report grammatical error.   
 [如果A的所有候选式都推不到a, 看看有没有候选式能推到 $\epsilon$  (即 $\epsilon \in \text{FIRST}(\alpha_i)$ ), 并且A后直接跟着a (即 $a \in \text{FOLLOW}(A)$ ) ]

- **Example:**  $G(S): S \rightarrow aA \mid d; \quad A \rightarrow bAS \mid \epsilon. \quad \text{String} : abd.$

$\text{FIRST}(aA) = \{a\}$        $\text{FOLLOW}(A) = \{\$, a, d\}$

$\text{FIRST}(d) = \{d\}$

$\text{FIRST}(bAS) = \{b\}$

$\text{FIRST}(A) = \{b, \epsilon\}$

$\text{FIRST}(S) = \{a, d\}$

**HOW TO compute  
FIRST & FOLLOW?**

$S$

$\Rightarrow aA$

$\Rightarrow abAS$

$\Rightarrow abS$

$\Rightarrow abd$

$abd$

$abd$

$abd$  (First cannot tell, check Follow)

$abd$



# Compute FIRST [计算FIRST集合的方法]



- To compute **FIRST(X)** for all grammar symbol **X**, apply the following rules **until no more terminals or  $\epsilon$  can be added to any FIRST set**:

- ◆ **Rule1**: If  $X \in V_T$ , then  $\text{FIRST}(X) = \{X\}$ . [X是终结符]
- ◆ **Rule2**: If  $X \in V_N$  and  $X \rightarrow \epsilon$  exists, then add  $\epsilon$  to  $\text{FIRST}(X)$ . [非终结符,空式]
- ◆ **Rule3**: If  $X \in V_N$  and  $X \rightarrow Y_1 Y_2 Y_3 \dots Y_k$ , then [非终结符,非空式]
  - If for some  $Y_i$  and a terminal  $a$ : ①  $\epsilon \in \text{all of } \text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$ , i.e.,  $Y_1 \dots Y_{i-1} \xRightarrow{*} \epsilon$ ; ②  $a \in \text{FIRST}(Y_i) \setminus \{\epsilon\}$ . Then,  $a \in \text{FIRST}(X)$ .
  - E.g.,
    - Everything in  $\text{FIRST}(Y_1) \setminus \{\epsilon\}$  is surely in  $\text{FIRST}(X)$ .
    - If  $Y_1$  doesn't derive  $\epsilon$ , then we add nothing more.
    - But if  $Y_1 \xRightarrow{*} \epsilon$ , then we add  $\text{FIRST}(Y_2) \setminus \{\epsilon\}$ , and so on...
  - Add  $\epsilon$  to  $\text{FIRST}(X)$ , if  $\epsilon$  is in  $\text{FIRST}(Y_i)$  for all  $i=1,2,\dots,k$ .

# Compute FIRST [计算FIRST集合的方法]



- Next, we can compute FIRST for any string  $\alpha = X_1X_2...X_n$  [符号串]
  - ◆ Add all non- $\epsilon$  symbols of  $FIRST(X_1)$  to  $FIRST(\alpha)$ .
  - ◆ Add non- $\epsilon$  symbols of  $FIRST(X_i)$ ,  $2 \leq i \leq n$ , to  $FIRST(\alpha)$ , if  $FIRST(X_1), \dots, FIRST(X_{i-1})$  all contain  $\epsilon$ . [前i-1个都包含空串]
  - ◆ Add  $\epsilon$  to  $FIRST(\alpha)$ , if  $FIRST(X_1), \dots, FIRST(X_n)$  all contain  $\epsilon$ .

[First(X)集合的计算本质上就是看X能推导出什么, First(X)包括所有X能推导出的句子的首个字符 (必须是终结符或空串)]

# Compute FIRST [计算FIRST集合的方法]



- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$ , compute the **FIRST** set of each non-terminal.

◆  $FIRST(E): \{$

◆  $FIRST(T): \{$

◆  $FIRST(E'): \{ +, \varepsilon$

◆  $FIRST(T'): \{ *, \varepsilon$

◆  $FIRST(F): \{ (, id$

**Apply rules for the first time:**

$E \rightarrow TE'$   $FIRST(T)$  add to  $FIRST(E)$ ,  $FIRST(T)$  doesn't contain  $\varepsilon$

$E' \rightarrow +TE'$   $FIRST(+)$  add to  $FIRST(E')$ ,  $FIRST(+)$  doesn't contain  $\varepsilon$

$E' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(E')$

$T \rightarrow FT'$   $FIRST(F)$  add to  $FIRST(T)$ ,  $FIRST(F)$  doesn't contain  $\varepsilon$

$T' \rightarrow *FT'$   $FIRST(*)$  add to  $FIRST(T')$ ,  $FIRST(*)$  doesn't contain  $\varepsilon$

$T' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(T')$

$F \rightarrow (E)$   $FIRST('(')$  add to  $FIRST(F)$ ,  $FIRST('(')$  doesn't contain  $\varepsilon$

$F \rightarrow id$   $FIRST(id)$  add to  $FIRST(F)$

# Compute FIRST [计算FIRST集合的方法]



- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$ , compute the **FIRST** set of each non-terminal.

- ◆  $FIRST(E): \{$
- ◆  $FIRST(T): \{ (, id$
- ◆  $FIRST(E'): \{ +, \varepsilon$
- ◆  $FIRST(T'): \{ *, \varepsilon$
- ◆  $FIRST(F): \{ (, id$

It is necessary to determine whether the FIRST set has changed after each rule application. **first time YES!**

**Apply rules for the second time:**

$E \rightarrow TE'$   $FIRST(T)$  add to  $FIRST(E)$ ,  $FIRST(T)$  doesn't contain  $\varepsilon$

$E' \rightarrow +TE'$   $FIRST(+)$  add to  $FIRST(E')$ ,  $FIRST(+)$  doesn't contain  $\varepsilon$

$E' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(E')$

$T \rightarrow FT'$   $FIRST(F)$  add to  $FIRST(T)$ ,  $FIRST(F)$  doesn't contain  $\varepsilon$

$T' \rightarrow *FT'$   $FIRST(*)$  add to  $FIRST(T')$ ,  $FIRST(*)$  doesn't contain  $\varepsilon$

$T' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(T')$

$F \rightarrow (E)$   $FIRST('(')$  add to  $FIRST(F)$ ,  $FIRST('(')$  doesn't contain  $\varepsilon$

$F \rightarrow id$   $FIRST(id)$  add to  $FIRST(F)$

# Compute FIRST [计算FIRST集合的方法]



- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$ , compute the FIRST set of each non-terminal.

◆  $FIRST(E): \{ (, id$

◆  $FIRST(T): \{ (, id$

◆  $FIRST(E'): \{ +, \varepsilon$

◆  $FIRST(T'): \{ *, \varepsilon$

◆  $FIRST(F): \{ (, id$

It is necessary to determine whether the FIRST set has changed after each rule application. **second time YES!**

**Apply rules for the third time:**

$E \rightarrow TE'$   $FIRST(T)$  add to  $FIRST(E)$ ,  $FIRST(T)$  doesn't contain  $\varepsilon$

$E' \rightarrow +TE'$   $FIRST(+)$  add to  $FIRST(E')$ ,  $FIRST(+)$  doesn't contain  $\varepsilon$

$E' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(E')$

$T \rightarrow FT'$   $FIRST(F)$  add to  $FIRST(T)$ ,  $FIRST(F)$  doesn't contain  $\varepsilon$

$T' \rightarrow *FT'$   $FIRST(*)$  add to  $FIRST(T')$ ,  $FIRST(*)$  doesn't contain  $\varepsilon$

$T' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(T')$

$F \rightarrow (E)$   $FIRST('(')$  add to  $FIRST(F)$ ,  $FIRST('(')$  doesn't contain  $\varepsilon$

$F \rightarrow id$   $FIRST(id)$  add to  $FIRST(F)$

# Compute FIRST [计算FIRST集合的方法]



- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$ , compute the FIRST set of each non-terminal.

◆  $FIRST(E): \{ (, id \}$

◆  $FIRST(T): \{ (, id \}$

◆  $FIRST(E'): \{ +, \varepsilon \}$

◆  $FIRST(T'): \{ *, \varepsilon \}$

◆  $FIRST(F): \{ (, id \}$

It is necessary to determine whether the FIRST set has changed after each rule application. **third time YES!**

**Apply rules for the 4th time:**

$E \rightarrow TE'$   $FIRST(T)$  add to  $FIRST(E)$ ,  $FIRST(T)$  doesn't contain  $\varepsilon$

$E' \rightarrow +TE'$   $FIRST(+)$  add to  $FIRST(E')$ ,  $FIRST(+)$  doesn't contain  $\varepsilon$

$E' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(E')$

$T \rightarrow FT'$   $FIRST(F)$  add to  $FIRST(T)$ ,  $FIRST(F)$  doesn't contain  $\varepsilon$

$T' \rightarrow *FT'$   $FIRST(*)$  add to  $FIRST(T')$ ,  $FIRST(*)$  doesn't contain  $\varepsilon$

$T' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(T')$

$F \rightarrow (E)$   $FIRST('(')$  add to  $FIRST(F)$ ,  $FIRST('(')$  doesn't contain  $\varepsilon$

$F \rightarrow id$   $FIRST(id)$  add to  $FIRST(F)$

It is necessary to determine whether the FIRST set has changed after each rule application. **4th time NO!**

# Compute FOLLOW [计算FOLLOW集合的方法]

- To compute FOLLOW(A) for all non-terminals A, **apply the following rules until nothing can be added to any FOLLOW set.**
  - ◆ **Rule1:** Place \$ in FOLLOW(S), where S is the start symbol.
  - ◆ **Rule2:** If there is a production  $A \rightarrow \alpha B \beta$ , then everything in FIRST( $\beta$ ) except  $\epsilon$  is in FOLLOW(B).
  - ◆ **Rule3:** If there is a production  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha B \beta$ , where FIRST( $\beta$ ) contains  $\epsilon$ , then everything in FOLLOW(A) is in FOLLOW(B).
- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \epsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \epsilon ; F \rightarrow (E) \mid id$ , compute the FOLLOW set of each non-terminal.
  - FIRST(E):{ (, id }; FIRST(T):{ (, id }; FIRST(E'):{ +,  $\epsilon$  }; FIRST(T'):{ \*,  $\epsilon$  }; FIRST(F):{ (, id }.

[FOLLOW(A)计算的本质是在看这个A后 (不包含A本身的推导结果) 的符号 (必须是终结符或结束符号)]

# Compute FOLLOW [计算FOLLOW集合的方法]

- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$  , compute the **FOLLOW** set of each non-terminal.

◆  $FIRST(E): \{ (, id \}$ ;  $FIRST(T): \{ (, id \}$ ;  $FIRST(E'): \{ +, \varepsilon \}$ ;  $FIRST(T'): \{ *, \varepsilon \}$ ;  $FIRST(F): \{ (, id \}$ .

◆  $FOLLOW(E): \{ \$$

◆  $FOLLOW(T): \{ +, \$$

◆  $FOLLOW(E'): \{ \$$

◆  $FOLLOW(T'): \{$

◆  $FOLLOW(F): \{$

## Apply rules for the first time:

- ◆ Place  $\$$  in  $FOLLOW(E)$ , since  $E$  is the start symbol.

- ◆  $E \rightarrow TE'$

- $FIRST(E')$  except  $\varepsilon$  is in  $FOLLOW(T)$ .
- Everything in  $FOLLOW(E)$  is in  $FOLLOW(E')$ .
- Since  $FIRST(E')$  contains  $\varepsilon$ , then everything in  $FOLLOW(E)$  is in  $FOLLOW(T)$ .

- ◆  $E' \rightarrow +TE' \mid \varepsilon$

- $FIRST(E')$  except  $\varepsilon$  is in  $FOLLOW(T)$ .
- Everything in  $FOLLOW(E')$  is in  $FOLLOW(E')$ .
- Since  $FIRST(E')$  contains  $\varepsilon$ , then everything in  $FOLLOW(E')$  is in  $FOLLOW(T)$ .



# Compute FOLLOW [计算FOLLOW集合的方法]

- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$  , compute the **FOLLOW** set of each non-terminal.

◆  $FIRST(E): \{ (, id \}$ ;  $FIRST(T): \{ (, id \}$ ;  $FIRST(E'): \{ +, \varepsilon \}$ ;  $FIRST(T'): \{ *, \varepsilon \}$ ;  $FIRST(F): \{ (, id \}$ .

◆  $FOLLOW(E): \{ \$$

**Apply rules for the first time:**

◆  $FOLLOW(T): \{ +, \$$

◆  $T \rightarrow FT'$

□  $FIRST(T')$  except  $\varepsilon$  is in  $FOLLOW(F)$ .

□ Everything in  $FOLLOW(T)$  is in  $FOLLOW(T')$ .

□ Since  $FIRST(T')$  contains  $\varepsilon$ , then everything in  $FOLLOW(T)$  is in  $FOLLOW(F)$ .

◆  $FOLLOW(E'): \{ \$$

◆  $FOLLOW(T'): \{ +, \$$

◆  $FOLLOW(F): \{ *, +, \$$

◆  $T' \rightarrow *FT' \mid \varepsilon$

□  $FIRST(T')$  except  $\varepsilon$  is in  $FOLLOW(F)$ .

□ Everything in  $FOLLOW(T')$  is in  $FOLLOW(T')$ .

□ Since  $FIRST(T')$  contains  $\varepsilon$ , then everything in  $FOLLOW(T')$  is in  $FOLLOW(F)$ .

# Compute FOLLOW [计算FOLLOW集合的方法]

- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$  , **compute the FOLLOW set of each non-terminal.**

◆  $FIRST(E): \{ (, id \}$ ;  $FIRST(T): \{ (, id \}$ ;  $FIRST(E'): \{ +, \varepsilon \}$ ;  $FIRST(T'): \{ *, \varepsilon \}$ ;  $FIRST(F): \{ (, id \}$ .

◆  $FOLLOW(E): \{ \$, )$

**Apply rules for the first time:**

◆  $FOLLOW(T): \{ +, \$, )$

◆  $F \rightarrow (E) \mid id$

□  $FIRST('')$  except  $\varepsilon$  is in  $FOLLOW(E)$ .

◆  $FOLLOW(E'): \{ \$, )$

Since FOLLOW sets has changed at the first time, **Apply rules for**

◆  $FOLLOW(T'): \{ +, \$$

**the second time:**

◆  $FOLLOW(F): \{ *, +, \$$

◆  $E \rightarrow TE'$

□  $FIRST(E')$  except  $\varepsilon$  is in  $FOLLOW(T)$ .

□ Everything in  $FOLLOW(E)$  is in  $FOLLOW(E')$ .

□ Since  $FIRST(E')$  contains  $\varepsilon$ , then everything in  $FOLLOW(E)$  is in  $FOLLOW(T)$ .

# Compute FOLLOW [计算FOLLOW集合的方法]

- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$  , **compute the FOLLOW set of each non-terminal.**

- ◆  $FIRST(E): \{ (, id \}$ ;  $FIRST(T): \{ (, id \}$ ;  $FIRST(E'): \{ +, \varepsilon \}$ ;  $FIRST(T'): \{ *, \varepsilon \}$ ;  $FIRST(F): \{ (, id \}$ .
- ◆  $FOLLOW(E): \{ \$, )$       ◆  $E' \rightarrow +TE' \mid \varepsilon$ 
  - $FIRST(E')$  except  $\varepsilon$  is in  $FOLLOW(T)$ .
  - Everything in  $FOLLOW(E')$  is in  $FOLLOW(E')$ .
  - Since  $FIRST(E')$  contains  $\varepsilon$ , then everything in  $FOLLOW(E')$  is in  $FOLLOW(T)$ .
- ◆  $FOLLOW(T): \{ +, \$, )$
- ◆  $FOLLOW(E'): \{ \$, )$
- ◆  $FOLLOW(T'): \{ +, \$, )$
- ◆  $FOLLOW(F): \{ *, +, \$, )$       ◆  $T \rightarrow FT'$ 
  - $FIRST(T')$  except  $\varepsilon$  is in  $FOLLOW(F)$ .
  - Everything in  $FOLLOW(T)$  is in  $FOLLOW(T')$ .
  - Since  $FIRST(T')$  contains  $\varepsilon$ , then everything in  $FOLLOW(T)$  is in  $FOLLOW(F)$ .

# Compute FOLLOW [计算FOLLOW集合的方法]

- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$  , **compute the FOLLOW set of each non-terminal.**

- ◆  $FIRST(E): \{ (, id \}; FIRST(T): \{ (, id \}; FIRST(E'): \{ +, \varepsilon \}; FIRST(T'): \{ *, \varepsilon \}; FIRST(F): \{ (, id \}.$
- ◆  $FOLLOW(E): \{ \$, ) \}$  ◆  $T' \rightarrow *FT' \mid \varepsilon$ 
  - $FIRST(T')$  except  $\varepsilon$  is in  $FOLLOW(F)$ .
  - Everything in  $FOLLOW(T')$  is in  $FOLLOW(T')$ .
  - Since  $FIRST(T')$  contains  $\varepsilon$ , then everything in  $FOLLOW(T')$  is in  $FOLLOW(F)$ .
- ◆  $FOLLOW(T): \{ +, \$, ) \}$
- ◆  $FOLLOW(E'): \{ \$, ) \}$
- ◆  $FOLLOW(T'): \{ +, \$, ) \}$
- ◆  $FOLLOW(F): \{ *, +, \$, ) \}$  ◆  $F \rightarrow (E) \mid id$ 
  - $FIRST('')$  except  $\varepsilon$  is in  $FOLLOW(E)$ .

Since FOLLOW sets has changed at the second time, **Apply rules for the third time:**

- ◆ **EXERCISE! –Ans: no FOLLOW set will be changed.**

## • Top-down Parsing

- **Recursive-descent parsing** [RDP, 递归下降语法分析]
  - (immediate and indirect) left recursion -> **remove left recursion**
  - **Backtracking** -> brutal force, not preferred due to complications and high overheads
- **Predicative parsing – LL(k)** [预测分析]
  - Predict the right production based on the input symbol (*i.e.*, **lookahead**), **no backtracking!**
  - Preliminaries: (1) **remove left recursion**, (2) **left factoring**
  - To enable prediction: **FIRST** and **FOLLOW**
    - $\text{FIRST}(\alpha) = \{ a \mid \alpha \Rightarrow a\dots, a \in V_T \}$ ,  $\epsilon$  can be included
    - $\text{FOLLOW}(A) = \{ a \mid S \Rightarrow \dots Aa\dots, a \in V_T \}$ ,  $\$$  can be included
  - Prediction based on the input symbol **a** when expanding **A**:
    - For **A** with  $A \rightarrow \alpha \mid \beta$ , where  $\text{FIRST}(\alpha)$  and  $\text{FIRST}(\beta)$  are disjoint, *i.e.*, **chose the one with a inside its FIRST set.**
    - If all FIRST sets do not contain **a**, but  $\epsilon \in \text{FIRST}(\alpha)$  while  $a \in \text{FOLLOW}(A)$ , **chose FIRST( $\alpha$ )**

# Compute FIRST (Revisit)



- To compute **FIRST(X)** for all grammar symbol **X**, apply the following rules **until no more terminals or  $\epsilon$  can be added to any FIRST set**:

- ◆ **Rule1**: If  $X \in V_T$ , then  $\text{FIRST}(X) = \{X\}$ . [X是终结符]
- ◆ **Rule2**: If  $X \in V_N$  and  $X \rightarrow \epsilon$  exists, then add  $\epsilon$  to  $\text{FIRST}(X)$ . [非终结符,空式]
- ◆ **Rule3**: If  $X \in V_N$  and  $X \rightarrow Y_1 Y_2 Y_3 \dots Y_k$ , then [非终结符,非空式]
  - If for some  $Y_i$  and a terminal  $a$ : ①  $\epsilon \in \text{all of } \text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$ , i.e.,  $Y_1 \dots Y_{i-1} \xRightarrow{*} \epsilon$ ; ②  $a \in \text{FIRST}(Y_i) \setminus \{\epsilon\}$ . Then,  $a \in \text{FIRST}(X)$ .
  - E.g.,
    - Everything in  $\text{FIRST}(Y_1) \setminus \{\epsilon\}$  is surely in  $\text{FIRST}(X)$ .
    - If  $Y_1$  doesn't derive  $\epsilon$ , then we add nothing more.
    - But if  $Y_1 \xRightarrow{*} \epsilon$ , then we add  $\text{FIRST}(Y_2) \setminus \{\epsilon\}$ , and so on...
  - Add  $\epsilon$  to  $\text{FIRST}(X)$ , if  $\epsilon$  is in  $\text{FIRST}(Y_i)$  **for all**  $i=1,2,\dots,k$ .

# Compute FIRST (Revisit)



- Next, we can compute FIRST for any string  $\alpha = X_1X_2...X_n$  [符号串]
  - ◆ Add all non- $\epsilon$  symbols of  $FIRST(X_1)$  to  $FIRST(\alpha)$ .
  - ◆ Add non- $\epsilon$  symbols of  $FIRST(X_i)$ ,  $2 \leq i \leq n$ , to  $FIRST(\alpha)$ , if  $FIRST(X_1), \dots, FIRST(X_{i-1})$  all contain  $\epsilon$ . [前i-1个都包含空串]
  - ◆ Add  $\epsilon$  to  $FIRST(\alpha)$ , if  $FIRST(X_1), \dots, FIRST(X_n)$  all contain  $\epsilon$ .

[First( $\alpha$ )集合的计算本质上就是看 $\alpha$ 能推导出什么, First( $\alpha$ )包括所有 $\alpha$ 能推导出的句子的首个字符 (必须是终结符或空串)]

[给定 $\alpha \rightarrow X_1X_2...X_n$ , First( $\alpha$ )的计算基于单个First( $X_i$ )的计算 (方法见上页)]

# Compute FIRST (Revisit)



- **Example:**  $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$ , compute the FIRST set of each non-terminal.

◆  $FIRST(E): \{ (, id \}$

◆  $FIRST(T): \{ (, id \}$

◆  $FIRST(E'): \{ +, \varepsilon \}$

◆  $FIRST(T'): \{ *, \varepsilon \}$

◆  $FIRST(F): \{ (, id \}$

It is necessary to determine whether the FIRST set has changed after each rule application.

$E \rightarrow TE'$   $FIRST(T)$  add to  $FIRST(E)$ ,  $FIRST(T)$  doesn't contain  $\varepsilon$

$E' \rightarrow +TE'$   $FIRST(+)$  add to  $FIRST(E')$ ,  $FIRST(+)$  doesn't contain  $\varepsilon$

$E' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(E')$

$T \rightarrow FT'$   $FIRST(F)$  add to  $FIRST(T)$ ,  $FIRST(F)$  doesn't contain  $\varepsilon$

$T' \rightarrow *FT'$   $FIRST(*)$  add to  $FIRST(T')$ ,  $FIRST(*)$  doesn't contain  $\varepsilon$

$T' \rightarrow \varepsilon$   $\varepsilon$  add to  $FIRST(T')$

$F \rightarrow (E)$   $FIRST('(')$  add to  $FIRST(F)$ ,  $FIRST('(')$  doesn't contain  $\varepsilon$

$F \rightarrow id$   $FIRST(id)$  add to  $FIRST(F)$

$\alpha = X_1X_2...X_n$

- (1) all non- $\varepsilon$  symbols of  $FIRST(X_1)$  to  $FIRST(\alpha)$
- (2) If  $FIRST(X_1), \dots, FIRST(X_{i-1})$  all contain  $\varepsilon$ , add non- $\varepsilon$  symbols of  $FIRST(X_i)$ ,  $2 \leq i \leq n$ , to  $FIRST(\alpha)$
- (3) If  $FIRST(X_1), \dots, FIRST(X_n)$  all contain  $\varepsilon$ , add  $\varepsilon$  to  $FIRST(\alpha)$



# Compute FOLLOW (Revisit)

- To compute FOLLOW(A) for all non-terminals A, **apply the following rules until nothing can be added to any FOLLOW set.**
  - ◆ **Rule1:** Place \$ in FOLLOW(S), where S is the start symbol.
  - ◆ **Rule2:** If there is a production  $A \rightarrow \alpha B \beta$ , then everything in FIRST( $\beta$ ) except  $\epsilon$  is in FOLLOW(B).
  - ◆ **Rule3:** If there is a production  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha B \beta$ , where FIRST( $\beta$ ) contains  $\epsilon$ , then everything in FOLLOW(A) is in FOLLOW(B).

[FOLLOW(A)计算的本质是在看这个A后 (不包含A本身的推导结果) 的符号 (必须是终结符或结束符号)]

[给定  $A \rightarrow \alpha B \beta$ , FOLLOW(B)计算基于FIRST( $\beta$ )的计算]

# Compute FOLLOW (Revisit)

- $G[E]: E \rightarrow TE' ; E' \rightarrow +TE' \mid \varepsilon ; T \rightarrow FT' ; T' \rightarrow *FT' \mid \varepsilon ; F \rightarrow (E) \mid id$ 
    - ◆  $FIRST(E):{ (, id } ; FIRST(T):{ (, id } ; FIRST(E'):{ +, \varepsilon } ; FIRST(T'):{ *, \varepsilon } ; FIRST(F):{ (, id } .$
    - ◆  $FOLLOW(E):{ \$, ) }$
    - ◆  $FOLLOW(T):{ +, \$, ) }$
    - ◆  $FOLLOW(E'):{ \$, ) }$
    - ◆  $FOLLOW(T'):{ +, \$, ) }$
    - ◆  $FOLLOW(F):{ *, +, \$, ) }$
  - ◆ Place  $\$$  in  $FOLLOW(E)$ , since  $E$  is the start symbol.
  - ◆  $E \rightarrow TE'$ 
    - $FIRST(E')$  except  $\varepsilon$  is in  $FOLLOW(T)$ .
    - Everything in  $FOLLOW(E)$  is in  $FOLLOW(E')$ .
    - Since  $FIRST(E')$  contains  $\varepsilon$ , then everything in  $FOLLOW(E)$  is in  $FOLLOW(T)$ .
  - ◆  $E' \rightarrow +TE' \mid \varepsilon$ 
    - $FIRST(E')$  except  $\varepsilon$  is in  $FOLLOW(T)$ .
    - Everything in  $FOLLOW(E')$  is in  $FOLLOW(E')$ .
    - Since  $FIRST(E')$  contains  $\varepsilon$ , then everything in  $FOLLOW(E')$  is in  $FOLLOW(T)$ .
- $A \rightarrow \alpha B \beta$
- (1) If  $A$  is the start symbol,  $\$ \in FOLLOW(A)$
  - (2) Everything in  $FIRST(\beta)$  except  $\varepsilon$  is in  $FOLLOW(B)$ .
  - (3) If  $\beta$  does not exist or  $FIRST(\beta)$  contains  $\varepsilon$ , then everything in  $FOLLOW(A)$  is in  $FOLLOW(B)$ .

# LL(1) Grammar[LL(1)文法]



- Predictive parsers, that is, recursive-descent parsers needing no backtracking, can be constructed for a class of grammars called **LL(1)**.
- A grammar **G** is LL(1) **if and only if** whenever **any two distinct productions of G** [G的任意两个不同的产生式],  $A \rightarrow \alpha \mid \beta$ , the following conditions hold:

So,  $S \rightarrow EBD \mid FBB$ ;  $E \rightarrow a\dots$ ;  $F \rightarrow a\dots$  is not LL(1)

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**. [条件1和2不满足则直接产生二义性, 即使用 $\alpha$ 还是 $\beta$ ?]
3. If  $\beta \xRightarrow{*} \epsilon$ , then  $\alpha$  does not derive any string beginning with a terminal in **FOLLOW(A)**. Likewise, if  $\alpha \xRightarrow{*} \epsilon$ , then  $\beta$  does not derive any string beginning with a terminal in **FOLLOW(A)**. [条件3不满足也会产生两种可能的推导, 如推导 $\beta$ 为 $\epsilon$ 还是推导 $\alpha$ ]

- The first two conditions are equivalent to the statement that  $\text{FIRST}(\alpha)$  and  $\text{FIRST}(\beta)$  are disjoint sets [不相交的集合].
- The third condition is equivalent to stating that if  $\epsilon$  is in  $\text{FIRST}(\beta)$ , then  $\text{FIRST}(\alpha)$  and  $\text{FOLLOW}(A)$  are disjoint sets, and likewise if  $\epsilon$  is in  $\text{FIRST}(\alpha)$ . [确保无等价的推导]  

反例: input symbol为**b**, 当前句型aa**A**b, 产生式**A**-> **b** |  $\epsilon$
- LL(1) Grammar does not contain left recursion. [LL(1)文法不含左递归]
- LL(1) Grammar is not ambiguous. [LL(1)文法不是二义的]

# LL(1)/LL(k) Grammar [LL(1)/LL(k) 文法]



- **LL (1) grammar.**
  - ◆ **L:** The first "L" in LL(1) stands for scanning the input from left to right.
  - ◆ **L:** The second "L" for producing a leftmost derivation.
  - ◆ **1:** The "1" for using one input symbol of lookahead at each step to make parsing action decisions.
- **LL (k) grammar.**
  - ◆ **k:** using **k** input symbols of lookahead at each step to make parsing action decisions.
- Many languages are LL(k), in fact many are LL(1).
- **Is LL(0) useful at all?**
  - ◆ Grammar where rules can be predicted with no lookahead.
  - ◆  $\Rightarrow$  That means, there can only be one rule per non-terminal.
  - ◆  $\Rightarrow$  That means, this language can have only one string.

# LL(1) Parser Implementation[实现]



- Recursive LL(1) parser for  $G[S]: S \rightarrow A \mid B;$   
 $A \rightarrow a; B \rightarrow b.$

- ◆ Maintaining a stack implicitly. [隐式维护栈]
- ◆ What is the stack for? Derivation results

- Is there a way to express above code more concisely? [简洁]

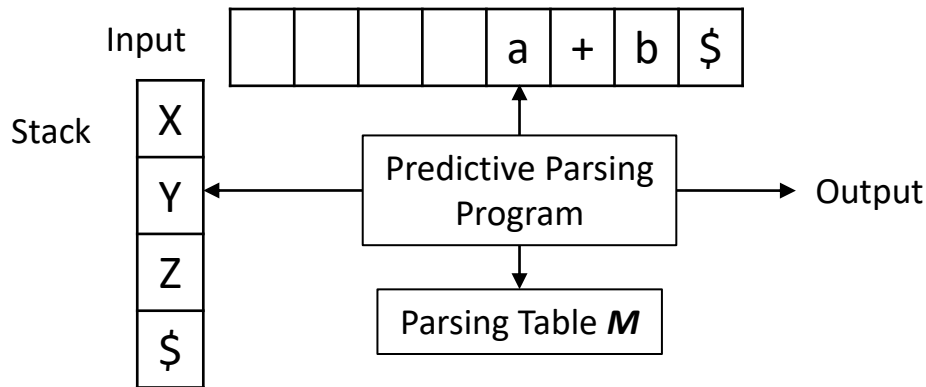
- ◆ Non-recursive LL(1) parsers
  - ▣ Use a predictive parsing table. [预测分析表]
  - ▣ maintaining a stack explicitly. [显式维护栈]
  - ▣ Table-driven parser. [表驱动]

```
void S(){
    token = Next(); // lookahead
    if(token == a) // 'A' starts with 'a'
        A(); // call procedure A()
    else if (token == b) // 'B' starts
        with 'b'
        B(); // call procedure B()
    else
        return; // error, reject.
}
```

# Non-recursive LL(1) Parser [非递归]



- **Input buffer** [輸入串]: contains the string to be parsed, followed by the “end marker” \$.
- **Stack** [推导的中间结果]: holds a sequence of grammar symbols and the symbol \$ to mark the bottom of the stack. It may contain:
  - ◆ **Terminals** that have yet to be matched against the input symbol.
  - ◆ **Non-terminals** that have yet to be expanded.



- **Parsing Table  $M[A,a]$**  [语法分析表]: an entry containing production “ $A \rightarrow \dots$ ” or error.
- **Predictive Parsing Program** [语法分析程序]: Execute the action according to **<stack top, current input symbol>**

# LL(1) Parse Table [预测分析表]



**G[E]:**

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

	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	$F \rightarrow id$			$F \rightarrow (E)$		

Input symbol,  
lists all possible  
terminals and \$

all non-terminals  
in the grammar

One action for each **<non-terminal, next input>**, it  
“predicts” the correct action based on **one lookahead**

- **Reject** on reaching error state
- **Accept** on end of input & empty stack

FIRST(E):{ (, id }

FIRST(T):{ (, id }

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

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

FIRST(F):{ (, id }

FOLLOW(E):{ \$, ) }

FOLLOW(T):{ +, \$, ) }

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

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

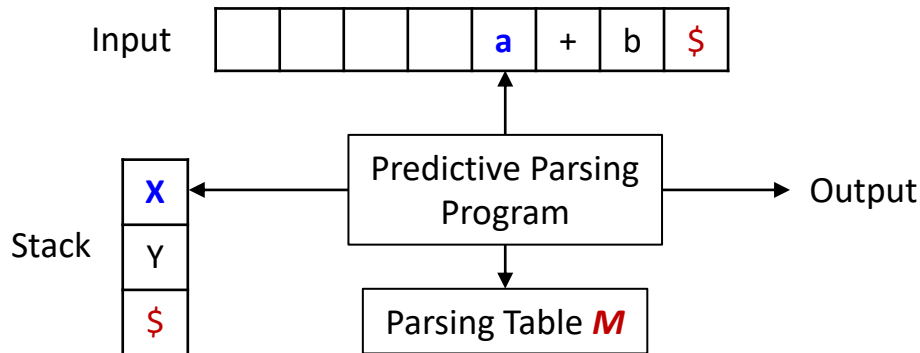
FOLLOW(F):{ \*, +, \$, ) }



# LL(1) Parsing Algorithm [非递归算法]



- Initial state [初始态]
  - ◆ **Input**: A string  $w$  and a parsing table  $M$  for grammar  $G$ . **Input Buffer**:  $w\$$ .
  - ◆ **Stack**: start symbol followed by '\$' at bottom.
  - ◆ Assume  $X$ : symbol at the top of the stack,  $a$ : current input symbol.
- General idea [总体思路] : repeat one of two actions
  - ◆ **Expand non-terminal symbol** at top of stack by applying a production
  - ◆ **Match terminal symbol** at top of stack with input token

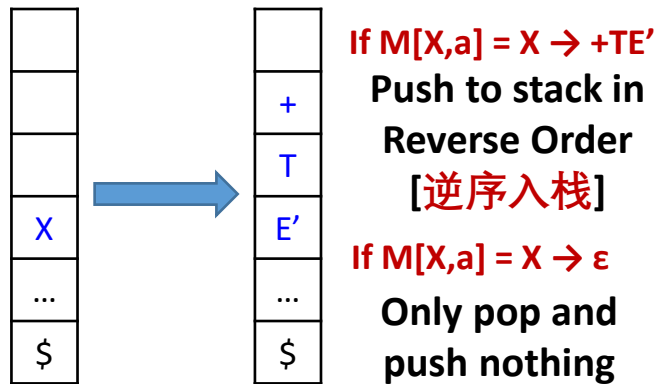
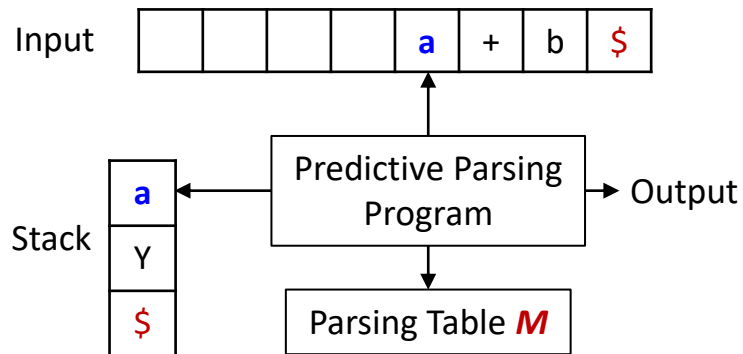


# LL(1) Parsing Algorithm [非递归算法]



- Algorithm Step-by-Step based on  $\langle X, a \rangle$ :

- ◆  $X$ : symbol at the top of the stack
- ◆  $a$ : current input token
- ◆ If  $X \in V_T$ , [栈顶符号为终结符] and
  - $X == a == \$$ , declare **SUCCESS**, stop parsing.
  - $X == a != \$$ , **pop  $X$  from stack** and move the current input symbol forward one.
  - $X != a$ , declare **ERROR**, input is **rejected**, stop parsing.
- ◆ If  $X \in V_N$ , [栈顶符号为非终结符] and
  - $M[X, a]$  has a production about  $X$ , **pop  $X$**  and **push right side of production to stack**.
  - $M[X, a] == \text{empty}$ , declare **ERROR**, input is **rejected**, stop parsing.



# Use the Parse Table [使用表进行预测分析]



- When using parse table for predictive parsers, if the grammar does not conform to the specification of LL (1) grammar:
  - the grammar should be **rewritten into LL (1) grammar** by **removing left recursion** and **backtracking**, then the parse table should be constructed.
- LL(1) parser example:** consider  $G[E]$  to recognize “**id + id \* id**”

$G[E]$ :

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid \text{id}$

	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	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

E
\$

Predictive Parsing Program

	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)		

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id** + **id** \* **id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

T
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

**G[E]:**

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id** + **id** \* **id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

F
T'
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id** + **id** \* **id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

id
T'
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

**G[E]:**

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

T'
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$



# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

+
T
E'
\$

Predictive Parsing Program

	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	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

**G[E]:**

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid \text{id}$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id** + **id** \* **id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

T
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

**G[E]:**

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id** + **id** \* **id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

F
T'
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id** + **id** \* **id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

id
T'
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

**G[E]:**

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

T'
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

*
F
T'
E'
\$

Predictive Parsing Program

**G[E]:**

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid \text{id}$

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

# Use the Parse Table [使用表进行预测分析]



- recognize “**id** + **id** \* **id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

F
T'
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$



# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

id
T'
E'
\$

Predictive Parsing Program

	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	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

**G[E]:**

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid \text{id}$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

T'
E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

E'
\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



- recognize “**id + id \* id**”

Input Buffer

		id	+	id	*	id	\$
--	--	----	---	----	---	----	----

Stack

\$

Predictive Parsing Program

	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	$F \rightarrow id$			$F \rightarrow (E)$		

ACCEPT!!

G[E]:

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

# Use the Parse Table [使用表进行预测分析]



Matched	Stack	Input	Action
	E\$	id + id * id\$	$E \rightarrow TE'$
	TE'\$	id + id * id\$	$T \rightarrow FT'$
	FT'E'\$	id + id * id\$	$F \rightarrow id$
	idT'E'\$	id + id * id\$	match id
id	T'E'\$	+ id * id\$	$T' \rightarrow \epsilon$
id	E'\$	+ id * id\$	$E' \rightarrow +TE'$
id	+TE'\$	+ id * id\$	match +
id +	TE'\$	id * id\$	$T \rightarrow FT'$
id +	FT' E'\$	id * id\$	$F \rightarrow id$
id +	idT' E'\$	id * id\$	match id
id + id	T' E'\$	* id\$	$T' \rightarrow *FT'$
id + id	*FT'E'\$	* id\$	match *

Matched	Stack	Input	Action
id + id *	FT'E'\$	id\$	$F \rightarrow id$
id + id *	idT'E'\$	id\$	match id
id + id * id	T'E'\$	\$	$T' \rightarrow \epsilon$
id + id * id	E'\$	\$	$E' \rightarrow \epsilon$
id + id * id	\$	\$	ACCEPT

The parser mimics a leftmost derivation.

- Top-down parsing mimics the grammar

## How to construct LL(1) Parse Table?

# Construct LL(1) Parse Table [构建预测分析表]

- Use **FIRST** and **FOLLOW** sets for a predictive parsing table  $M[A, a]$ , and the algorithm is based on the following idea:
  - ◆ The production  $A \rightarrow \alpha$  is chosen if the next input symbol  $a$  is in  $\text{FIRST}(\alpha)$ .
  - ◆ The only complication occurs when  $\alpha = \epsilon$  or, more generally,  $\alpha \xRightarrow{*} \epsilon$ :
    - we should again choose  $A \rightarrow \alpha$ , if (1)  $a$  is in  $\text{FOLLOW}(A)$  or (2) the  $\$$  on the input has been reached and  $\$$  is in  $\text{FOLLOW}(A)$ .

# Construct LL(1) Parse Table [构建预测分析表]

• Algorithm: For each production  $A \rightarrow \alpha$  of the grammar:

◆ For each terminal  $a \in \text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$ .

$G[E]$ :  $E \rightarrow TE'$     $E' \rightarrow +TE' \mid \epsilon$     $T \rightarrow FT'$     $T' \rightarrow *FT' \mid \epsilon$     $F \rightarrow (E) \mid id$

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

Production	FIRST( $\alpha$ )
$E \rightarrow TE'$	(, id
$E' \rightarrow +TE'$	+
$T \rightarrow FT'$	(, id
$T' \rightarrow *FT'$	*
$F \rightarrow (E)$	(
$F \rightarrow id$	id
$E' \rightarrow \epsilon$	FOLLOW
$T' \rightarrow \epsilon$	FOLLOW

# Construct LL(1) Parse Table [构建预测分析表]

• Algorithm: For each production  $A \rightarrow \alpha$  of the grammar:

- ◆ If  $\epsilon \in \text{FIRST}(\alpha)$ , then for each terminal  $b$  in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$ . If  $\epsilon \in \text{FIRST}(\alpha)$  and  $\$ \in \text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \$]$  as well.

$G[E]: E \rightarrow TE' \quad E' \rightarrow +TE' \mid \epsilon \quad T \rightarrow FT' \quad T' \rightarrow *FT' \mid \epsilon \quad F \rightarrow (E) \mid \text{id}$

	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	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

Production	$\text{FIRST}(\alpha)$
$E' \rightarrow \epsilon$	FOLLOW
$T' \rightarrow \epsilon$	FOLLOW
Symbol	$\text{FOLLOW}(A)$
E'	\$, )
T'	+, \$, )



# Construct LL(1) Parse Table [构建预测分析表]

**G[E]:**

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid \text{id}$

Symbol	FIRST	FOLLOW
E	(, i	\$, )
E'	+, $\varepsilon$	\$, )
T	(, i	+, \$, )
T'	*, $\varepsilon$	+, \$, )
F	(, i	*, +, \$, )

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

Production	FIRST( $\alpha$ )
$E \rightarrow TE'$	(, i
$E' \rightarrow +TE'$	+
$T \rightarrow FT'$	(, i
$T' \rightarrow *FT'$	*
$F \rightarrow (E)$	(
$F \rightarrow \text{id}$	id
$E' \rightarrow \varepsilon$	FOLLOW
$T' \rightarrow \varepsilon$	FOLLOW

# Determine If Grammar is LL(1) [判断LL(1)文法]

- Observation [直观依据]

- ◆ If a grammar is LL(1), each of its LL(1) table entry **contains at most one rule**.
- ◆ Otherwise, it is not LL(1).

	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	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

- Two methods to determine if a grammar is LL(1) or not
  - ◆ Construct LL(1) table, and check if there is a multi-rule entry. [看语法分析表]
  - ◆ Check each rule as if the table is getting constructed, grammar G is LL(1) if and only if for any two distinct productions  $A \rightarrow \alpha | \beta$ : [看文法]
    - ◆  $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \phi$
    - ◆ If  $\epsilon \in \text{FIRST}(\beta)$ ,  $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \phi$  (Mentioned before)

# Non-LL(1) Grammar [非LL(1)文法]

- Assume that a grammar is not LL(1). How to solve?
  - ◆ Case1- the language may still be LL(1)
    - Try to **rewrite grammar** to LL(1) grammar [remove left-recursion & left-factoring]
    - Try to **remove ambiguity** in grammar.
  - ◆ Case2- If Case-1 fails, language may not be LL(1)
    - It's impossible to resolve conflict at the grammar level.
    - Programmer chooses which rule to use for conflicting entry (if choosing that rule is always semantically correct)
    - Otherwise, use a more powerful parser (e.g. **LL(k)**, **LR(1)**)

# LL(1) Time and Space Complexity[复杂度]

- **Linear** time and space relative to the length of input.
- **Time:** each input symbol is consumed within a constant number of steps.
  - ◆ If symbol at top of stack is a **terminal**: Matched immediately in **one step**.
  - ◆ If symbol at top of stack is a **non-terminal**:
    - ▢ Matched in at most **N steps**, where **N = number of rules**.
    - ▢ Since LL(1) is not left-recursive, we do not apply same rule twice without consuming input.

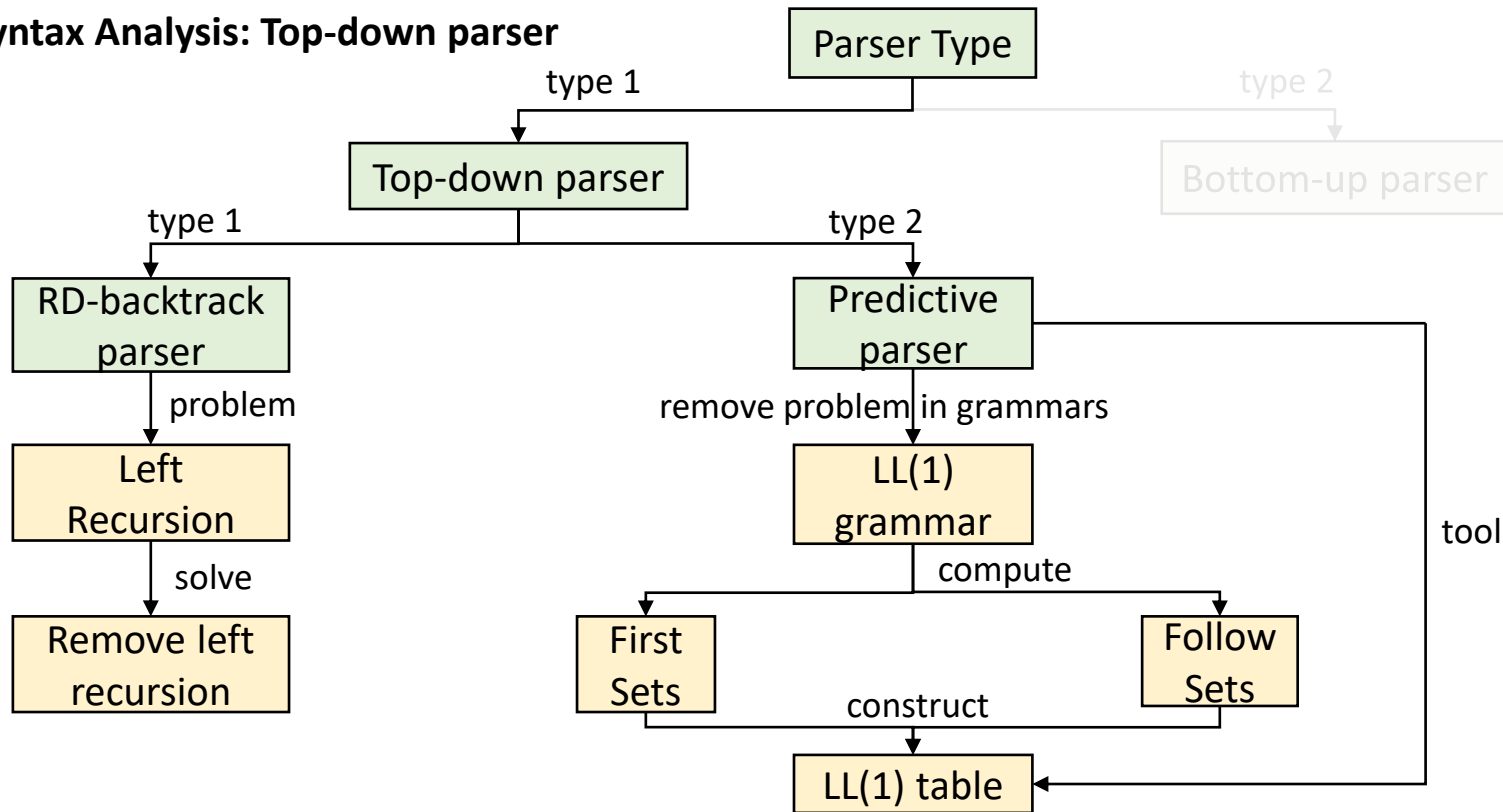
# LL(1) Time and Space Complexity[复杂度]

- **Space:** smaller than input (after removing  $X \rightarrow \epsilon$ )
  - ◆ Right side of production is always longer or equal to left side of production
    - Derivation string **expands monotonically**.
    - Derivation string is always **shorter than final input string**.
  - ◆ Stack is a subset of derivation string (unmatched portion)
- LL(k)'s size of parse table =  $O(|N| * |T|^k)$ 
  - ◆ N = number of non-terminals, T = number of terminals
  - ◆ prevent LL(2) ... LL(k) from wide usage

# Summary



## Syntax Analysis: Top-down parser



# Summary



- Top-down Parsing; RDP with backtracking, predictive parsing.
- (Immediate / Non-immediate) Left Recursion and how to resolve
- Left-factoring.
- FIRST, FOLLOW.
- LL(1)/LL(k) Grammar.
- Recursive / Non-recursive LL(1) parser implementation.
- The use of LL(1) Parse Table.
- The construction of LL(1) Parse Table.

# Summary [自顶向下分析]

---

- RDP [递归下降分析]

- left recursion [左递归]
  - Remove Left Recursion [消除直接/间接左递归]
- Backtracking [回溯]
  - left factoring [提取左公因子]

- Predictive Parsing [预测分析]

- FIRST & FOLLOW [终结首符集 & 后继终结符号集]
- Definition of LL(1)/LL(k) [LL(1)/LL(k)文法的定义]
- LL(1) Parse Table [LL(1)的分析表]
  - The use of the Parse Table [分析表的使用]
  - The Construction of the Parse Table [分析表的构建]
- Table-driven LL(1) Parser Implementation [分析表驱动的LL(1)语法分析实现]
- Complexity of LL(1) [LL(1)的时间与空间复杂度]



# Further Reading



- Dragon Book

- ◆ Comprehensive Reading:

- Section 4.1.2 and 4.4.1 for the introduction to top-down parsing.
    - Section 4.4.2 for function FIRST and FOLLOW.
    - Section 4.1.4, 4.4.3–4.4.4, and 4.4.5 for LL(1) parsing and error recovery.

- ◆ Skip Reading:

- Section 4.1.3–4.1.4 for error recovery in top-down parsing.
    - Section 4.2.7 for differences between CFGs and regular expressions.

