

编译原理

Compiler Construction Principles



朱 青

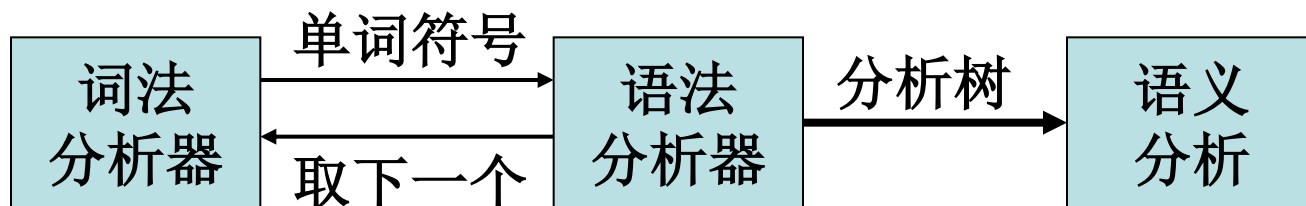
信息学院计算机系，
中国人民大学，

zqruc2012@aliyun.com



第3章: 程序语言的语法描述与分析

语法分析的重要性:



语法分析方法:

- 1) 自顶向下分析法.
- 2) 自底向上分析法.

第3章: 程序语言的语法描述与分析

⌘3.1 上下文无关文法

⌘3.2 自顶向下分析法

⌘3.3 递归下降分析法

⌘3.4 LL(1)分析法

⌘3.5 LL(1)分析器

3.1 上下文无关文法

3.1.1 文法与语言

文法是描述语言的语法结构的形式
规则(即语法规则)

上下文有关文法:

上下文无关文法:

-----在编译中通指此文法.

一.上下文无关文法:

例: 只含*, +和()的算术表达式的文法

$$E \longrightarrow E + E | E * E | (E) | a \quad (3.1)$$

定义: 一个上下文无关文法是一个

四元式 (V_T, V_N, S, P)

其中: 1) V_T 是非空有限集, 每个元素是一个终结符.

2) V_N 是非空有限集, 每个元素是一个非终结符.

3) S是一个非终结符,是开始符号.

4) P是产生式的集合:它的每一个产生式的形式为 $P \rightarrow \alpha$.

其中: P 属于 V_N , α 属于 $(V_NUV_T)^*$.

S在产生式的左部必须至少出现一次.

为简单起见, 文法可以只给出产生式集合, 并指出开始符号.

例如:

$$E \longrightarrow EAE \mid a \qquad V_T = \{a, +, *\}$$

$$A \longrightarrow + \mid * \qquad V_N = \{E, A\}$$

$$E \text{ 为开始符号.} \qquad S = \{E\}$$

$$\text{缩写为: } E \longrightarrow EAE \mid a$$

$$A \longrightarrow + \mid *$$

左部为非终结符号, 右部为后选式.

二.由文法产生语言:

文法(3.1)可以生成句子:(a^*a+a)

$E \Rightarrow (E) \Rightarrow (E^*E) \Rightarrow (a^*E)$

$\Rightarrow (a^*E+E) \Rightarrow (a^*a+E) \Rightarrow (a^*a+a)$

\Rightarrow :表示推导.

- 基本术语定义:

- 1) 称 $\alpha A\beta$ 直接推导出 $\alpha\gamma\beta$, 仅当有 $A \rightarrow \gamma$ 且 $\alpha, \beta \in (V_T \cup V_N)^*$.
表示为 $\alpha A\beta \Rightarrow \alpha\gamma\beta$.

2) 若有 $\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$,则称从 α_1 推导出 α_n .

3) $\alpha_1 \stackrel{*}{\Rightarrow} \alpha_n$ 表示从 α_1 经过0次或若干次可推导出 α_n .

$\alpha_1 \stackrel{+}{\Rightarrow} \alpha_n$ 表示从 α_1 经过1次或若干次可推导出 α_n .

若 $\alpha \stackrel{*}{\Rightarrow} \beta$ 表示 或者 $\alpha = \beta$ 或 $\alpha \stackrel{+}{\Rightarrow} \beta$.

4) 若**G**是一个文法,**S**是开始符号,若有 $S^* \Rightarrow \alpha$,则称 α 为句型.仅含终结符的句型叫句子.

文法**G**所生成的语言为:

$$L(G) = \{\alpha | S \Rightarrow^+ \alpha \text{ 且 } \alpha \text{ 属于 } V_T^*\}$$

5)最左(右)推导

每次直接推导的是最左(右)的非终结符.

6)短语,直接短语:

若**G**是一个文法,**S**是开始符号,
假定 $\alpha\beta\chi$ 是文法**G**的一个句型,如
果有

$$S \xRightarrow{*} \alpha A \chi \quad \text{且} \quad A \xRightarrow{+} \beta$$

则称 β 是一个关于非终结符号**A**的,句型 $\alpha\beta\chi$ 的短语.
如果有

$$S \xRightarrow{*} \alpha A \chi \quad \text{且} \quad A \xRightarrow{=} \beta$$

则称 β 是直接短语.

一个句型的最左直接短语称为
该句型的句柄.

例如:考虑一个文法 **$G=(\{a,b\},\{S\},S,P)$**
其中 **P** :

$$S \longrightarrow aSb$$

$$S \longrightarrow ab$$

$$S \implies aSb \implies aaSbb$$

$$\implies a^3Sb^3 \implies \dots$$

$$\implies a^{n-1}Sb^{n-1}$$

$$\implies a^n b^n (n \geq 1)$$

$$L(G) = \{a^n b^n \mid n \geq 1\}$$

例如：以下述文法为例说明(**a^*a+a**)
是文法的句子,及术语举例.

$$E \longrightarrow T|E+T$$

$$T \longrightarrow F|T^*F \quad (3.2)$$

$$F \longrightarrow (E)|a$$

解: $E \Rightarrow T \Rightarrow F \Rightarrow (E) \Rightarrow (E+T)$

$$\Rightarrow (T+T) \Rightarrow (T^*F+T) \quad (3.3)$$

$$\Rightarrow (F^*F+T) \Rightarrow (a^*F+T)$$

$$\Rightarrow (a^*a+T) \Rightarrow (a^*a+F) \Rightarrow (a^*a+a)$$

所以 $(a*a+a)$ 是文法 **G** 的句子.

而 $T, F, (E), (E+T), (T+T), (T*F+T),$
 $(F*F+T), (a*F+T), (a*a+T), (a*a+F),$
都是文法 **G** 的句型.

最右推导: $E \Rightarrow T \Rightarrow F \Rightarrow (E) \Rightarrow (E+T)$
 $\Rightarrow (E+F) \Rightarrow (E+a)$
 $\Rightarrow (T*F+a) \Rightarrow (T*a+a)$
 $\Rightarrow (F*a+a) \Rightarrow (a*a+a)$

考虑文法(3.2)的一个句型 $a_1 * a_2 + a_3$:

- $a_1, a_2, a_3, a_1 * a_2$ 都是句型 $a_1 * a_2 + a_3$ 的短语,
- a_1, a_2 和 a_3 是直接短语,
- a_1 是最左直接短语,
- $a_2 + a_3$ 不是句型 $a_1 * a_2 + a_3$ 短语,因为
有 $E \Rightarrow a_2 + a_3$ 但不存在从文法的开始符号 E 到
 $a_1 * E$ 的推导.

3.1.2 语法树和二义性

上述概念可以用一张图(语法树)表示.
如文法(3.1)的如下推导的语法树可表示:
 $E \Rightarrow (E)$

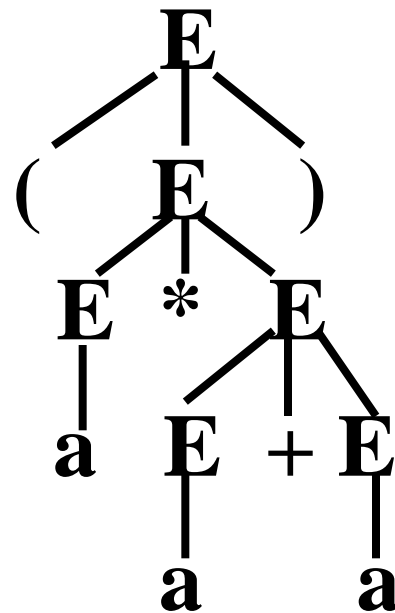
$\Rightarrow (E * E)$

$\Rightarrow (a * E)$

$\Rightarrow (a * E + E)$

$\Rightarrow (a * a + E)$

$\Rightarrow (a * a + a)$



(3.4)

文法(3.1)的如下推导的语法树可表示:

$E \Rightarrow (E)$

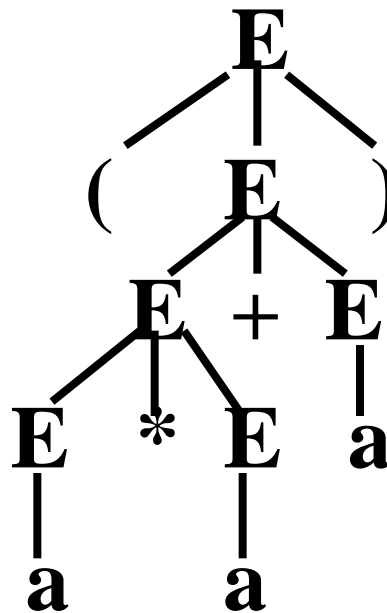
$\Rightarrow (E + E)$

$\Rightarrow (E * E + E)$

$\Rightarrow (a * E + E)$

$\Rightarrow (a * a + E)$

$\Rightarrow (a * a + a)$



(3.5)

若某文法存在一个句子,对应两个不同的语法树,则称该文法是二义的.

已经证明,文法的二义性是不可判定的.可以为某些二义性找到充分条件,使其变成非二义的.

例如: 上述文法规定先乘后加,
同级运算左结合,即:

$$E \longrightarrow T | E + T$$

$$T \longrightarrow F | T * F$$

$$F \longrightarrow (E) | a$$

此文法无二义性.

例如: **PASCAL**语言关于**IF**语句
的文法有下列规则:

(**E**--表达式,**S**--语句,**other**--非**IF**语句)

S \longrightarrow **if E then S**

 | **if E then S else S**

 | **other**

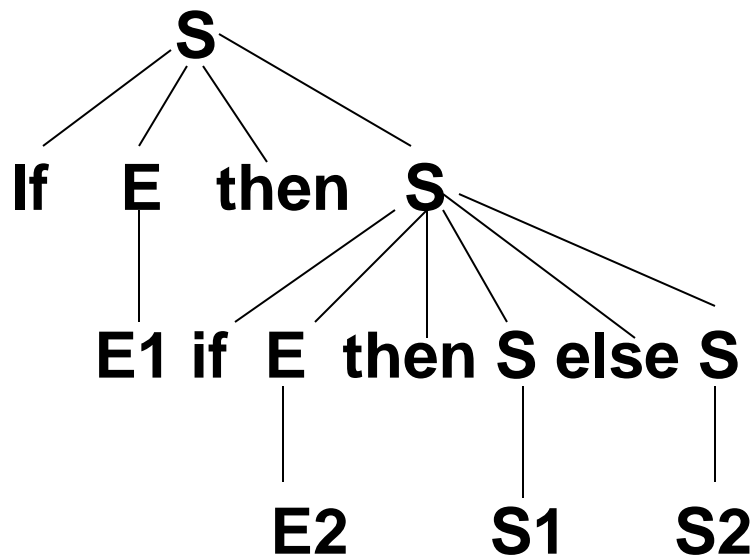
判断此文法是否为二义文法.

解： 对复合条件语句：

if E1 then if E2 then S1 else S2

有两棵语法树(如下所示)

故此文法是二义的。



改写文法为下列的无二义性文法:

(**M_S**为匹配的, **UN_S**为非匹配的)

$$\begin{aligned} S &\longrightarrow M_S \\ &\quad | UN_S \end{aligned}$$
$$\begin{aligned} M_S &\longrightarrow \text{if } E \text{ then } M_S \text{ else } M_S \\ &\quad | \text{other} \end{aligned}$$
$$\begin{aligned} UN_S &\longrightarrow \text{if } E \text{ then } S \\ &\quad | \text{if } E \text{ then } M_S \text{ else } UN_S \end{aligned}$$

3.1.3 形式语言概观

通常将文法分成四类:

0型,1型,2型,3型.

0型>1型>2型>3型. ('>'---强于)

0型语言:递归可枚举语言.

1型语言:上下文有关语言.

2型语言:上下文无关语言.

3型语言:正规语言.

文献已经证明:这四类语言分别
可被四种自动机所识别.

图灵机(**Turing machine---TM**)识
别0型语言(递归可枚举语言).

线性界限自动机(**linear bounded
automata---LBA**)识别上下文有关语
言.

下推自动机 (**push down automata---PDA**) 识别上下文无关语言.

有限自动机 (**finite automata ---FA**) 识别正规语言.

文献已经证明:各类文法所产生的语言恰好与各类自动机所识别的语言相同.

定义: (0型文法)

$G=(V_N, V_T, S, P)$ 是一个0型文法,如果他的每个产生式

$$\alpha \longrightarrow \beta$$

是这样的一种结构: α 属于 $(V_NUV_T)^*$ 且至少含一个非终结符而 β 属于 $(V_NUV_T)^*$.

如果对**0**型文法加上以下的第**i**条限制,就得到**i**型文法.

1. **G**的任何产生式 $\alpha \rightarrow \beta$ 均满足 $|\alpha| < |\beta|$ 或 $|\alpha| = |\beta|$; 仅仅 $S \rightarrow \varepsilon$ 例外, 但**S**不得出现在任何产生式的右部.
2. **G**的任何产生式为 $A \rightarrow \beta$,
A属于 V_N 且 β 属于 $(V_N U V_T)^*$.

**3. G 的任何产生式为 $A \rightarrow \alpha B$ 或 $A \rightarrow \alpha$,
其中, α 属于 V_T^* , A, B 属于 V_N .**

0型文法：称短语文法(无限制文法).

1型文法：称上下文有关文法.

2型文法：称上下文无关文法.

3型文法：称右线性文法,正规文法.

文法与语言举例:

例1: 3型文法G:

含有下列产生式的集合:

$S \longrightarrow 0|C0$

$C \longrightarrow B1$

$B \longrightarrow 0|C0$

$D \longrightarrow 1|C1|D0|D1|B0$

此文法构造的**NFA**能够且只能够识别所产生的正规语言.

例2: 上下文无关文法:

$$S \longrightarrow aSb|ab$$

此文法产生的语言:

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

例3: 上下文有关文法:

$$S \longrightarrow aSAB$$

$$AA' \longrightarrow AB$$

$$S \longrightarrow abB$$

$$bA \longrightarrow bb$$

$$BA \longrightarrow BA'$$

$$bB \longrightarrow bc$$

$$BA' \longrightarrow AA'$$

$$cB \longrightarrow cc$$

此文法产生的语言:

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

例4: 0型语言

$$L_0 = \{ \alpha c \alpha \mid \alpha \in (a|b)^* \}$$

只有0型文法才能生成。

第3章: 程序语言的语法描述与分析

⌘ 3.1 上下文无关文法

⌘ 3.2 自顶向下分析法

⌘ 3.3 递归下降分析法

⌘ 3.4 LL(1)分析法

⌘ 3.5 LL(1)分析器

3.2 自顶向下分析法

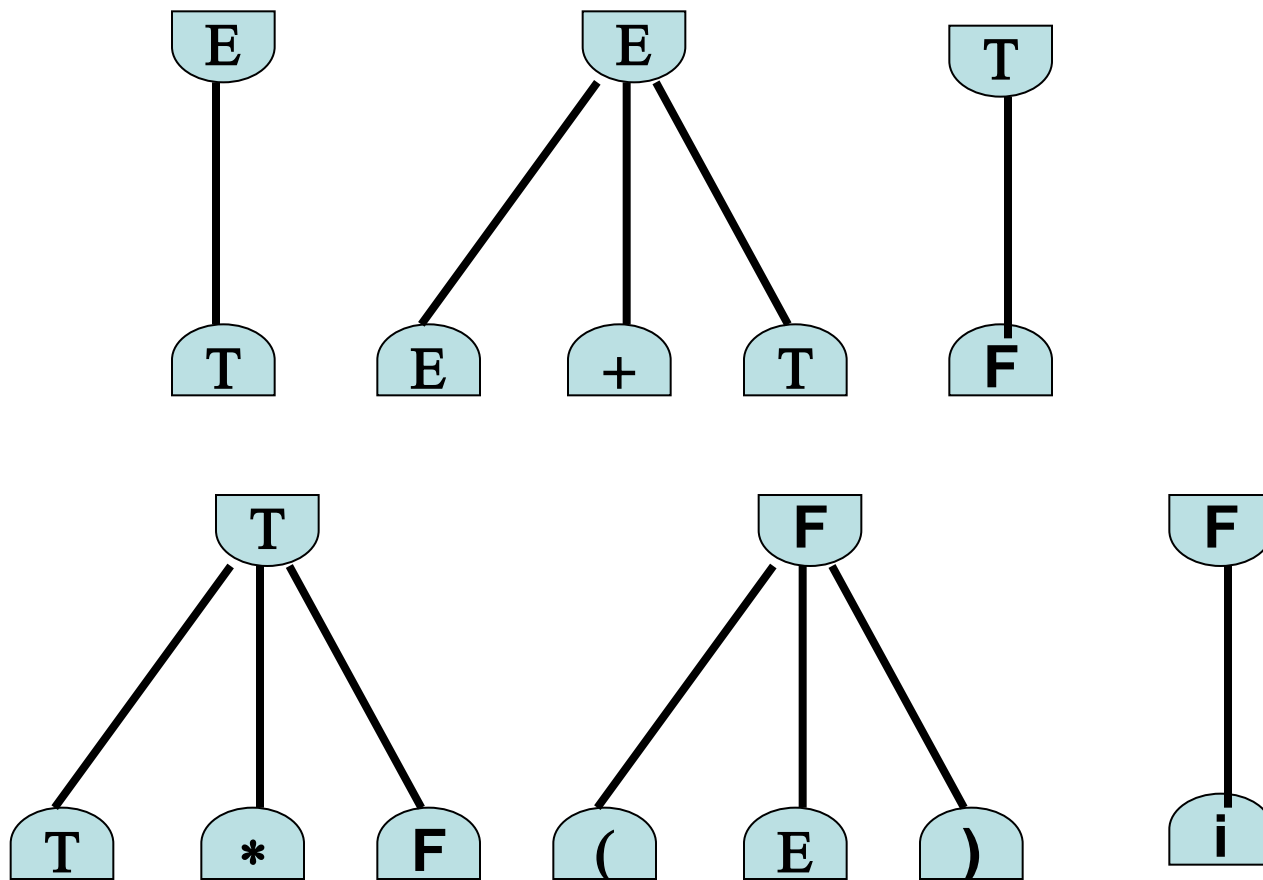
语法分析由两类:自顶向下分析法和自底向上分析法.下面用股子游戏来演示.如:文法

$$E \longrightarrow T|E+T$$

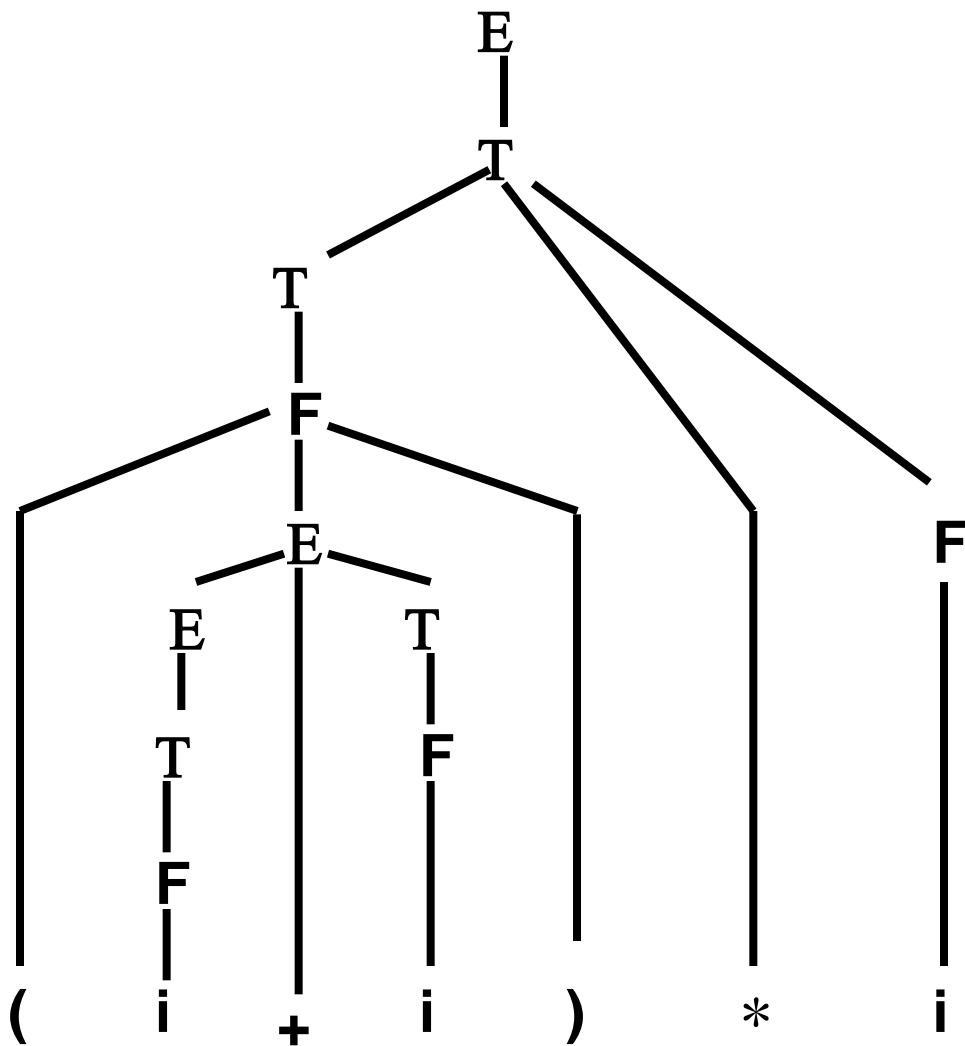
$$T \longrightarrow F|T*F$$

$$F \longrightarrow (E)|i$$

所有的股子为:



若要分析符号串 $(i+i)*i$ 初始态为:



自上而下的分析法:从句子符号到句子,

自下而上的分析法:从句子到句子符号.

分析开始时,句子从左向右读,则最自然,也可以从右向左读,则也同样方便.

注意:不替换已经使用的产生式.如果替换,就是回退.叫回溯.

所谓确定与非确定的分析,取决于是否发生回溯.

第3章: 程序语言的语法描述与分析

⌘3.1 上下文无关文法

⌘3.2 自顶向下分析法

⌘3.3 递归下降分析法

⌘3.4 LL(1)分析法

⌘3.5 LL(1)分析器

3.3 递归下降分析法

是一种易于实现的自上而下的分析法.可以从适当的文法出发写一个语法程序.

考虑下面的文法生成的语言:

PROGRAM \rightarrow begin DECLIST comma STATELIST end

DECLIST \rightarrow d semi DECLIST

DECLIST \rightarrow d

STATELIST \rightarrow s semi STATELIST

STATELIST \rightarrow s

(这个语言的形式结构: **begin d;d,...;s;s,...;s end**).

这个文法稍加改造,可得:

PROGRAM \longrightarrow **begin DECLIST comma STATELIST end**

DECLIST \longrightarrow **dx**

X \longrightarrow **semi DECLIST**

X \longrightarrow ϵ

STATELIST \longrightarrow **sY**

Y \longrightarrow **semi STATELIST**

Y \longrightarrow ϵ

设:有以下程序:

1. **error**是一个出错处理程序.
2. **semi, comma**是终结符.
3. **Lexical**是一个词法分析程序,它提供下一个单词.

4. 语法分析程序执行之前,已经执行过:

symbol:=lexical;

用递归下降分析法写语法分析程序时,只需给每个非终结符写一个过程.

PROC PROGRAM;

BEGIN

IF symbol <> begin THEN error;

symbol:=lexical; DECLIST;

IF symbol <> comma THEN error;

symbol:=lexical; STATELIST;

IF symbol <> end THEN error

END;


```
PROC DECLIST;  
BEGIN  
    IF symbol <> d THEN error;  
    symbol:=lexical; X;  
END;
```

```
PROC X;  
BEGIN  
    IF symbol =semi THEN  
    BEGIN symbol:=lexical; DECLIST END  
    ELSE IF symbol = comma  
        THEN return  
        ELSE error;  
END;
```

```
PROC STATELIST;  
BEGIN  
    IF symbol <> s THEN error;  
    symbol:=lexical; Y;  
END;
```

```
PROC Y;  
BEGIN  
    IF symbol =semi THEN  
        BEGIN symbol:=lexical;  
            STATELIST      END  
    ELSE IF symbol = end  
        THEN return  
        ELSE error;  
END;
```

用这样的语法分析句子,因多次的递归调用,效率很低,可以用重复来代替递归.

DECLIST \longrightarrow **d(semi d)***

STATELIST \longrightarrow **s (semi s)***

递归下降分析程序改写:

PROC PROGRAM 保留.

PROC DECLIST;

BEGIN

IF symbol <> d THEN error;

```
symbol := lexical;  
WHILE symbol =semi DO  
BEGIN  
    symbol := lexical;  
    IF symbol <> d THEN error;  
    symbol := lexical;  
END  
END;
```

```
PROC STATELIST;  
  BEGIN  
    IF symbol <> s THEN error;  
    symbol := lexical;  
    WHILE symbol =semi DO  
      BEGIN  
        symbol := lexical;  
        IF symbol <> s THEN error;  
        symbol := lexical;  
      END  
    END;
```

用重复代替递归可较前的方法提高效率即改善可读性.

文法的另一种表示法:

0 在**LEX**中规定辅助定义式: **$D_i \rightarrow R_i$** 的右部 **R_i** 只允许出现 **Σ** 中的字符和已定义过的简名 **D_1, \dots, D_{i-1}** ,即 **D_i** 不能递归调用.去掉这一条,允许递归调用.这种定义系统就可以表示上下文无关文法.

现在许多程序设计语言的语法是用这种定义系统来描述的.

某些符号的记法:(用重复代替递归)

1)用花括号 $\{a\}$ 表示闭包运算 a^* .

2)用 $\{a\}_0^n$ 表示 a 可任意重复0次至 n 次.

$$\{a\}_0^0 = a^0 = \varepsilon.$$

3)用方括号 $[a]$ 表示 $\{a\}_0^1$,即表示 a 的出现可有可无(等价于 $a|\varepsilon$).

例如:文法

$E \rightarrow T|E+T$ 可表示为: $E \rightarrow T\{+T\}$

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

```
PROCEDURE E;  
BEGIN  
    T;  
    WHILE SYM='+' DO  
        BEGIN ADVANCE; T END;  
END;
```

```
PROCEDURE T;  
BEGIN  
    F;  
    WHILE SYM='*' DO  
        BEGIN ADVANCE; F END;  
END; (ADVANCE 指向下一个输入符号).
```


递归下降分析法:

优点: 1) 从适当的文法出发写一个递归下降分析器很容易.

2) 分析器和文法很相似,故分析器的可靠性很高.

缺点: 1)因不断的递归调用,分析速度慢.

2)分析器本身比较大.

3)会引起分析过程中插入代码生成,而导致编译的不同阶段混在一起.给软件维护带来困难,并过早地依赖于机器.

左递归与公共左因子

1.消除左递归

定义: (左递归)

如果文法有一个非终结符号 **A**,
存在推导 $A \Rightarrow A \alpha$, α 为某个符号串,该文法是左递归的.

自顶向下文法不能处理左递归文法.

简单情况下,如文法的非终结符号**A**的产生式为

$$\mathbf{A} \longrightarrow \mathbf{A} \alpha | \beta$$

其中 β 不以 α 打头,那么可以把**A**的两个规则改写如下的非左递归产生式:

$$\mathbf{A} \longrightarrow \beta \mathbf{A}'$$

$$\mathbf{A}' \longrightarrow \alpha \mathbf{A}' | \varepsilon$$

这里**A**是直接左递归的.

考虑下述表达式文法:

$$E \longrightarrow E+T \mid T$$

$$T \longrightarrow T * F \mid F$$

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

(3.10)

经消除直接左递归后变成:

$$E \longrightarrow TE'$$

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

$$T \longrightarrow FT'$$

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

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

一般而言,假定关于**A**的全部产生式是:

$$\mathbf{A} \longrightarrow \mathbf{A}\alpha_1|\mathbf{A}\alpha_2|\dots|\mathbf{A}\alpha_m|\beta_1|\beta_2|\dots|\beta_n$$

其中 $\beta_i(i=1,2,\dots,n)$ 不以**A**打头;

$\alpha_j(j=1,2,\dots,m)$ 不等于 ε ,

那么可以把上述产生式改写为:

$$\mathbf{A} \longrightarrow \beta_1 \mathbf{A}'|\beta_2 \mathbf{A}'|\dots|\beta_n \mathbf{A}'$$

$$\mathbf{A}' \longrightarrow \alpha_1 \mathbf{A}'|\alpha_2 \mathbf{A}'|\dots|\alpha_m \mathbf{A}'|\varepsilon$$

此法易于消除见于表面的直接左递归, 但不意味着已经消除整个文法的左递归性. 例如文法:

$$S \longrightarrow Aa \mid b$$

$$A \longrightarrow Ac \mid Sd \mid \varepsilon$$

非终结符号S是左递归的, 因为有

$S \Rightarrow Aa \Rightarrow Sda$, 但它不是直接左递归的.

算法3.1:消除左递归:

输入: 无环路无 ε -产生式的文法**G**.

输出: 不带有左递归的一个等价文法.

1. 把文法**G**的所有非终结符号按一种顺序排列成 **A1,A2, ... ,An**.

2. for $i:= 1$ to n do

for $j:=1$ to $i-1$ do

begin

把每一个形如 $A_i \rightarrow A_j \gamma$ 的规则改写成

$$A_i \longrightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$$

其中 $A_j \longrightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ 是关于当前 A_j 的所有产生式;

消除关于 A_i -产生式的直接左递归性.

end

3. 化简由 (2) 所得到的文法, 即去除那些从开始符号出发永远无法到达的非终结符号的产生式.

例如:文法如下,消除左递归.

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd \mid \varepsilon$$

解: 改写文法为 ε -无关文法.

$S \rightarrow Aa$	$S \rightarrow Aa \mid \varepsilon a$ (即 $S \rightarrow Aa \mid a$)
$S \rightarrow b$	$S \rightarrow b$
$A \rightarrow Ac$	$A \rightarrow Ac \mid \varepsilon c$
$A \rightarrow Sd$	$A \rightarrow Sd$
$A \rightarrow \varepsilon$	无

现在考虑对改写后的文法:

$$S \longrightarrow Aa|a|b$$
$$A \longrightarrow Ac|c|Sd$$

根据算法4.1

1)非终结符安排为 S, A

2) S 不存在左递归, 故 $i=1$ 不做工作.

3) $i=2$ 时, 将 $S \longrightarrow Aa|a|b$ 代入到 A 的有关候选式后, 得到

$$A \longrightarrow Ac \mid c \mid Aad \mid ad \mid bd$$

再消除A-产生式中的直接左递归:

$$S \longrightarrow Aa|a|b$$
$$A \longrightarrow cA'|adA'|bdA'$$
$$A' \longrightarrow cA'|adA'|\varepsilon$$

例如:消除文法(3.17)的左递归

$$S \longrightarrow Qc|c$$
$$Q \longrightarrow Rb|b$$
$$R \longrightarrow Sa|a$$

(P83)

2. 提取公共左因子:

提取公共左因子是对文法改写的一种方法,为适应预测分析法.

例如:考虑以下两个产生式:

stmt \rightarrow if expr then stmt else stmt
| if expr then stmt

当输入 **if** 时,选择哪个产生式展开 **stmt**. 故提取公共左因子.

一般地,如果有产生式 $A \longrightarrow \alpha\beta_1 | \alpha\beta_2$
是两个产生式,那么可以提取左因子.
产生式变为:

$$A \longrightarrow \alpha A'$$

$$A' \longrightarrow \beta_1 | \beta_2$$

同样,若有 $A \longrightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_n | \gamma$,
其中 γ 是不以 α 开头的候选式,那么这
些产生式可用如下的式子代替:

$$A \longrightarrow \alpha A' | \gamma$$

$$A' \longrightarrow \beta_1 | \beta_2 | \dots | \beta_n$$

A'是一个新的非终结符号.

例如: 下面的文法是“未定的**else**”
文法的简写形式.

$$S \longrightarrow iEtS | iEtSeS | a$$

$$E \longrightarrow b \quad (3.13)$$

其中: **i,t,e**分别表示**if then else**,

提取公共左因子后,文法为:

$$\begin{aligned} S &\longrightarrow iEtSS'|a \\ S' &\longrightarrow eS|\varepsilon \\ E &\longrightarrow b \end{aligned} \quad (3.14)$$

文法(3.13) 和 (3.14)都是二义的,需进一步处理.

例如: 文法:

$$\begin{aligned} S &\longrightarrow aSb \\ S &\longrightarrow aSc \\ S &\longrightarrow \varepsilon \end{aligned}$$

提取公共左因子后文法形式

S \longrightarrow **aSX**

X \longrightarrow **b**

X \longrightarrow **c**

S \longrightarrow ϵ

第3章: 程序语言的语法描述与分析

⌘3.1 上下文无关文法

⌘3.2 自顶向下分析法

⌘3.3 递归下降分析法

⌘3.4 LL(1)分析法

⌘3.5 LL(1)分析器

3.4 LL(1)分析法

LL(1)文法是上下文无关文法中的一种.用它可以写确定的自上而下的语法分析器.LL(1)是S-文法的一般化.一个S-文法满足下面的二个条件:

- 1)每个产生式右部的第一个符号为终结符.
- 2)若每个非终结符有若干个候选式,

则这些候选式的第一个符号彼此不同.

(这可以帮助我们做确定的分析.)

例如:看下面的文法是否为**S**-文法.

$$S \longrightarrow pX$$
$$S \longrightarrow qY$$
$$X \longrightarrow aXb$$
$$X \longrightarrow x$$
$$Y \longrightarrow aYd \quad \text{是S-文法.}$$
$$Y \longrightarrow y$$

例如: 下面的文法不是**S-文法**:

$$S \longrightarrow R$$
$$S \longrightarrow T$$
$$R \longrightarrow pX$$
$$T \longrightarrow q$$
$$X \longrightarrow x$$
$$Y \longrightarrow aYd \quad \text{不是**S-文法**.}$$
$$Y \longrightarrow y$$

虽然他们产生的是同一种语言.但不是同一个文法.

在某些情况下,一个非**S**-文法可以改造成**S**-文法,而不影响所生成的语言.

下面对串**paaaxbbb**进行分析,看是否为上述文法的句子:

$S \Rightarrow pX \Rightarrow paXb \Rightarrow paaXbb$

$\Rightarrow paaaXbbb \Rightarrow paaaxbbb.$

可见,用**S**-文法可以很容易地写确定的自上而下的语法分析器.

LL(1)文法虽然是**S**-文法的一般化,但从**LL(1)**文法写确定的自上而下的语法分析器也很容易.

LL(1) :**L**--从左向右读符号串.

L--左推导.

1--在分析时,只要超前搜索一个符号,来决定要选着哪个候选式.

在LL(1)文法中,某些产生式右部的第一个符号虽然不是终结符,但从右部的第一个符号(非终结符)可推出一组以终结符开头的串来,这些终结符组成一个集合.

例如: 对文法

$$S \longrightarrow Ry$$
$$S \longrightarrow Tz$$
$$R \longrightarrow pacb$$
$$T \longrightarrow qayd$$

当从**S**出发推导时,

1)若在被分析的第一个符号是**p**,则应使用 $S \rightarrow Ry$.

2)若在被分析的第一个符号是**q**,则应使用 $S \rightarrow Tz$.

定义: (非终结符的首符号集合)

a 属于 $\text{first}(A)$ $\iff A \xRightarrow{*} a\alpha$

其中: $\text{first}(A)$ 是 A 的首符号集合, A 为非终结符, α 为字串.

例如:文法:

$$P \longrightarrow Ac$$

$$P \longrightarrow Bd$$

$$A \longrightarrow a$$

$$A \longrightarrow Aa$$

$$B \longrightarrow bB$$

$$\text{first}(P) = \{a, b\}$$

定义: (一个串的首符号集合)

$$a \text{ 属于 } \text{first}(\alpha) \iff \alpha \overset{*}{\Longrightarrow} a\beta$$

判断LL(1)文法的必要条件是:

一个非终结符,若有多余一个的候选式,
则这些后选式的首符号集合,彼此不相交.

注意: 当某些串中能推出 ϵ 时要特别留神.

例如:

$P \longrightarrow AB$

$P \longrightarrow BG$

不是LL(1)文法.

$A \longrightarrow Aa$

因为

$A \longrightarrow \epsilon$

$S(AB)=\{a,c,b\}$

$B \longrightarrow c$

$S(BG)=\{c,b\}$

$B \longrightarrow bB$

再考虑文法:

$T \rightarrow AB$

$A \rightarrow PQ$

$A \rightarrow BC$

$P \rightarrow pP$

$P \rightarrow \varepsilon$

$Q \rightarrow qQ$

$Q \rightarrow \varepsilon$

$B \rightarrow bB$

$B \rightarrow e$

$C \rightarrow d$

$C \rightarrow f$

$S(PQ) = \{p, q\}$

$S(BC) = \{e, b\}$

看上去彼此不相交,
但当分析的串的首符
号为: b, e 时推导就不唯一.

$T \Rightarrow AB \Rightarrow BCB \Rightarrow \dots$

$T \Rightarrow AB \Rightarrow PQB \Rightarrow B \dots$

都可以成立.这就是不确定性.
是不确定的分析.

定义: (引导符号集)

一个非终结符**P**的一个给定候选 **α** 的引导符号集为:

$$\mathbf{DS(P,\alpha)=\{a|a\in first(\alpha) \text{ or } (\alpha\Rightarrow\varepsilon, a \in F(P))\}}$$

其中: **F(P)**为文法所有产生式中能跟在**P**后面的符号集合。如上例中

$$\mathbf{DS(A, PQ)=\{p, q, b, e\}}$$

$$\mathbf{DS(A, BC)=\{b, e\}}$$

定义LL(1)文法:

若一个文法,它的每一个非终结符,有多余一个的候选式,则这些候选式的引导符号集,彼此不相交.这样的文法就是LL(1)文法.

LL(1)文法能进行确定的自上而下的文法分析.

下面给出形式化描述及判定LL(1)的算法:

判定LL(1)的算法:

步骤:

- 1)求文法符号串的首符号集**first**.
- 2)求文法符号串的后随符号集**follow**.
- 3)求文法产生式的引导符号集(选用集)**select**.
- 4)判定一个文法是否为LL(1)的.

首先,形式化定义首符号集**first** 和后随符号集**follow**.

定义1:(首符号集**first**)

假定 α 是文法**G**的任意符号串,定义

$$\text{first}(\alpha) = \{a \mid \alpha \xRightarrow{*} a \dots, a \in V_T\}$$

特别是, 若 $\alpha \xRightarrow{*} \epsilon$ 则规定 $\epsilon \in \text{first}(\alpha)$.

定义2:(后随符号集**follow**)

假定**S**是文法**G**的开始符号,对于**G**的任何非终结符**A**,定义

$$\text{follow}(A) = \{a \mid S \xRightarrow{*} \dots Aa \dots, a \in V_T\}$$

特别是,若 $S \xRightarrow{*} \dots A$,则规定 $\# \in \text{follow}(A)$.

一.求文法符号串的首符号集**first** 的算法:

应用下列规则,直到再没有任何终结符号或 ε 能加到该首符号集为止.

(1)如果 X 是终结符,则**first**(X)= $\{X\}$.

(2)如果 $X \in V_N$,且有产生式 $X \rightarrow a...$,则把 a 加入到**first**(X)中;若 $X \rightarrow \varepsilon$ 也是一个产生式,则 $\varepsilon \in \text{first}(X)$.

(3)若 $X \rightarrow Y \dots$ 是一个产生式且 X 是非终结符,则把 $\text{first}(Y)$ 中的所有非 ϵ -元素都加到 $\text{first}(X)$ 中;

若 $X \rightarrow Y_1 Y_2 \dots Y_i$ 是一个产生式,
 Y_1, \dots, Y_{i-1} 都是非终结符,而且对任何 j ,

$1 \leq j \leq i-1$, $\text{first}(Y_j)$ 都含有 ϵ
(即 $Y_1 Y_2 \dots Y_{i-1} \Rightarrow \epsilon$) , 则把 $\text{first}(Y_j)$ 中的所有非 ϵ -元素都加到 $\text{first}(X)$ 中;

特别是,若所有的 $\text{first}(Y_j)$ 均含有 ϵ ,
 $j=1,2,\dots,k$, 则把 ϵ 加到 $\text{first}(X)$ 中.

例如： 给定文法,求其各非终结符相关的串的首符号集.

$$\begin{aligned} E &\longrightarrow TE' \\ E' &\longrightarrow +TE' \mid \varepsilon \\ T &\longrightarrow FT' \\ T' &\longrightarrow *FT' \mid \varepsilon \\ F &\longrightarrow (E) \mid id \end{aligned} \quad (3.16)$$

解：按上述规则：

$$\text{first}(F) = \{ (, id \} \quad \text{first}(T') = \{ *, \varepsilon \}$$

$$\text{first}(E') = \{ +, \varepsilon \}$$

$$\text{first}(E) = \text{first}(T) = \text{first}(F) = \{ (, id \}$$

二.求文法符号串的后随符号集**follow** 的算法:

应用下列规则,直到每个后随符号集**follow**不再增大为止.

1) 对于文法的开始符号**S**,置**#**于**follow(S)**中;

2) 若 **$A \rightarrow \alpha B \beta$** 是一个产生式,则把 **$\text{first}(\beta) \setminus \{\epsilon\}$** 加至**follow(B)**中;

3) 若 **$A \rightarrow \alpha B$** 是一个产生式,或 **$A \rightarrow \alpha B \beta$** 是一个产生式而 **$\beta \Rightarrow \epsilon$** (即 **$\epsilon \in \text{first}(\beta)$**),则把**follow(A)**加至**follow(B)**中.

例如: 对文法(3.16),求每个非终结符的后随符号集.

解: (1)求 $\text{follow}(E)$: E 为开始符号,所以 $\#$ 属于 $\text{follow}(E)$; 由 $F \longrightarrow (E)$, 得 $) \in \text{follow}(E)$, 故 $\text{follow}(E) = \{), \# \}$.

(2) 求 $\text{follow}(E')$: 由产生式 $E \longrightarrow TE'$, 得 $\text{follow}(E) \in \text{follow}(E')$, 故 $\text{follow}(E') = \text{follow}(E) = \{), \# \}$

(3) 求 $\text{follow}(T)$: 由产生式 $E \longrightarrow TE'$ 和 $E' \longrightarrow \varepsilon$, 得

$\text{follow}(T) = \text{first}(E') \cup \text{follow}(E)$,
即 $\text{follow}(T) = \{ +,), \# \}$;

(4) 求 $\text{follow}(T')$: 由产生式 $T \rightarrow FT'$, 得
 $\text{follow}(T) \in \text{follow}(T')$, 故
 $\text{follow}(T') = \text{follow}(T) = \{), +, \# \}$

(5) 求 $\text{follow}(F)$: 由产生式 $T \rightarrow FT'$ 和
 $T' \rightarrow \varepsilon$, 得
 $\text{follow}(F) = \text{first}(T') \setminus \varepsilon \cup \text{follow}(T)$,
即 $\text{follow}(F) = \{ *, +,), \# \}$;

三. 产生式 $A \rightarrow \alpha$ 的选用集定义如下:

(1) 若 $\alpha \neq \varepsilon$, 且不存在推导 $\alpha \xrightarrow{+} \varepsilon$, 则产生式 $A \rightarrow \alpha$ 的选用集
 $\text{select}(A \rightarrow \alpha) = \text{first}(\alpha)$

(2) 若 $\alpha \neq \varepsilon$, 且存在推导 $\alpha \xRightarrow{+} \varepsilon$, 则产生式 $A \rightarrow \alpha$ 的选用集

$$\text{select}(A \rightarrow \alpha) = \text{first}(\alpha) \cup \text{follow}(A)$$

(3) 若 $\alpha = \varepsilon$, 即产生式 $A \rightarrow \varepsilon$ 的选用集

$$\text{select}(A \rightarrow \varepsilon) = \text{follow}(A)$$

例如: 求文法(3.16)各产生式的选用集.

$$\text{select}(E \rightarrow TE') = \text{first}(TE') = \{ (, \text{id} \}$$

$$\text{select}(E' \rightarrow +TE') = \text{first}(+TE') = \{ + \}$$

$$\text{select}(E' \rightarrow \varepsilon) = \text{follow}(E') = \{ \#,) \}$$

$\text{select}(F \longrightarrow (E)) = \text{first}((E)) = \{ (\}$

$\text{select}(F \longrightarrow \text{id}) = \{\text{id}\}$

$\text{select}(T \longrightarrow FT') = \text{first}(FT') = \{ (, \text{id} \}$

$\text{select}(T' \longrightarrow *FT') = \text{first}(*FT') = \{ * \}$

$\text{select}(T' \longrightarrow \varepsilon) = \text{follow}(T') = \{ \#,), + \}$

四.LL(1)文法的判定.

设有文法**G**,如果其中任意产生式

$A \longrightarrow \alpha \mid \beta$, 存在

$\text{select}(A \rightarrow \alpha) \cap \text{select}(A \rightarrow \beta) = \Phi$

则称**G**为**LL(1)**文法.

LL(1)文法是一种无二义性,非左递归的上下文无关文法,并且关于任意非终结符的不同产生式,其选用集不相交.

例如: 文法(3.16)

$\text{select}(E \rightarrow +TE') \cap \text{select}(E' \rightarrow \varepsilon) = \Phi$

$\text{select}(T' \rightarrow *FT') \cap \text{select}(T' \rightarrow \varepsilon) = \Phi$

$\text{select}(F \rightarrow (E)) \cap \text{select}(F \rightarrow \text{id}) = \Phi$

所以,文法(3.16)是LL(1)文法.

例如: 验证下列文法是否为LL(1)文法.

(1) $S \longrightarrow \text{if } E \text{ then } S S'$

(2) $S \longrightarrow \text{others}$

(3) $S' \longrightarrow \text{else } S$

(4) $S' \longrightarrow \varepsilon$

(5) $E \longrightarrow \text{expression}$

解: 求非终结符的首符号集和后随符号集

$\text{first}(S) = \{\text{if}, \text{other}\}$ $\text{follow}(S) = \{\#, \text{else}\}$

$\text{first}(S') = \{\text{else}, \varepsilon\}$ $\text{follow}(S') = \{\#, \text{else}\}$

$\text{first}(E) = \{\text{expression}\}$ $\text{follow}(E) = \{\text{then}\}$

求每个产生式的选用集

select(1)={if} select(2)={others}

select(3)={else} select(4)={ # , else }

select(5)={expression}

S的产生式的两个选用集

select(1) n select(2) = Φ

S'的产生式的两个选用集

select(3) n select(4) = {else}

由于文法同一终结符的两个产生式的选用集相交,所以,该文法不是LL(1)的.

LL(1)语言

已经可以判断一个文法是否是LL(1)的,但有两个问题需解决.

(1)是否每个语言都有一个LL(1)文法?

答案是否定的.

例如: 语言L(G)

$L(G) = \{aa^t \mid a \in \{0,1\}^* \text{ 中} \}$. (a^t 表示 a 的转置). 由文法G生成.

$S \longrightarrow 0S0$

$S \longrightarrow 1S1$ 不是LL(1)的.

$S \longrightarrow \varepsilon$

(2)能否找一个算法判断一个语言是否是LL(1)的?

答案是否定的.

上述问题的本质,在于它说明不是所有的程序设计语言都对应一个LL(1)文法;但是,程序设计语言的绝大部分上下文无关性质可被LL(1)文法描述.

问题的关键: 对于给定的非LL(1)文法, 如何找出一个等价的LL(1)文法. 从上述结论知: 实现此事无通用的算法. 但在大多情况下,

文法的转换并不困难.有经验的编译工作者可做手工转换.

文法转换涉及的问题:

1).消除左递归:

消除直接左递归

消除循环左递归

2).提取公共左因子:

但上述方法没有普遍性.否则,于判断一个语言是否为LL(1)的问题相矛盾.

第3章: 程序语言的语法描述与分析

⌘ 3.1 上下文无关文法

⌘ 3.2 自顶向下分析法

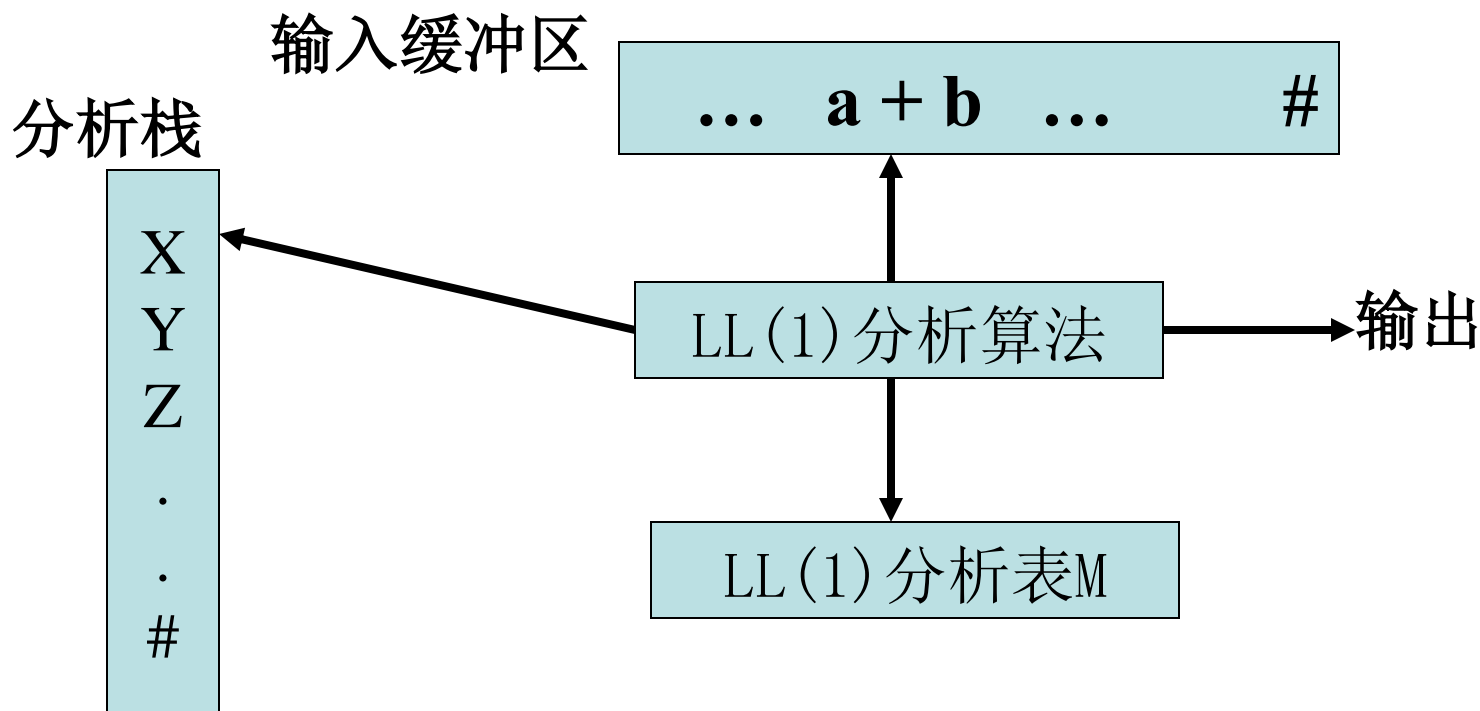
⌘ 3.3 递归下降分析法

⌘ 3.4 LL(1)分析法

⌘ 3.5 LL(1)分析器

3.5 LL(1)分析器

LL(1)分析器的模型.



在LL(1)分析模型中.

输入缓冲区:存放要分析的符号串,
在串的末尾加符号#表示输入结束.

分析栈:存放当前句型中尚待分析的文法符号,其中包括终结符号和非终结符号.栈底用#标记结束.初始化时,把#和文法的开始符号放在栈顶.

分析表M:是一个矩阵.其每行对应文法的一个非终结符号A,其每一列对应文法的一个终结符号a,结束标记符#作为一个终结符号

矩阵 $M[A,a]$ 表示当前栈顶为 A 输入符号为 a 时,应该选用的产生式,或者为空.构造分析表时,按照文法的选用集,若有 $a \in \text{select}(A \rightarrow \alpha)$,则产生式 $A \rightarrow \alpha$ 填入 $M[A, a]$;最后,把某些空白项作为处理错误的入口。

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

LL(1)分析算法控制分析器的执行。
初态，分析栈顶为文法的开始符号，输入缓冲指向输入串的第一个符号。对当前的栈顶符号**X**，输入符号**a**，算法执行的操作如下：

(1)如果 $X=a=\#$ ，算法结束 并且报告分析成功,接受输入串.

(2)如果 $X=a\neq\#$ ，则从栈中弹出X,并使缓冲区指针前进到下一个输入符号.

**(3)如果 X 是一个非终结符,则查分析表;
若 $M[A,a]$ 为 X 的产生式,用该产生式右部的
逆替换栈顶符号,否则出错.**

算法: LL(1)分析算法

输入:串 w ,文法 G 的分析表 M .

**输出:如果 $w \in L(G)$,则产生 w 的最左推导,否则
输出错误信息.**

**方法: 初态时,分析栈为 $\#S$,栈顶 S 是文法的开
始符号;缓冲区为 $w\#$.分析器按下列操作进行
语法分析:**

- ```

(1) push # ,S;指针ip指向串w#的第一个符号;
(2) repeat
(3) 令X为栈顶符号,a为ip所指的符号;
(4) if X 为终结符或# then
(5) if X=a=# then 接受w
(6) else if X=a<># then 弹出X,并使ip前进
(7) else error
(8) else /* X为非终结符 */
(9) if M[X,a]=X Y1Y2...Yk then
(10) begin
(11) 弹出X;
(12) 将Y1Y2...Yk压入栈;
(13) ip指向Y1;
(14) end
(15) until 接受w或error;

```

```
(11) push Yk,...,Y2,Y1 /* Y1在栈顶 */
(12) end
(13) else error
(14) until X=#; /* 栈为空 */
```

例如： 按照前面的文法所示的**LL(1)**分析表，  
对输入串**id+id\*id**进行语法分析。

解： 分析器分析过程中,栈中的符号,剩余的  
输入串和选用的产生式,如下表所列：

| 分析栈(顶)  | 输入串       | 产生式                          |
|---------|-----------|------------------------------|
| #E      | id+id*id# | $E \rightarrow TE'$          |
| #E'T    | id+id*id# | $T \rightarrow FT'$          |
| #E'T'F  | id+id*id# | $F \rightarrow id$           |
| #E'T'id | id+id*id# |                              |
| #E'T'   | +id*id#   | $T' \rightarrow \varepsilon$ |
| #E'     | +id*id#   | $E' \rightarrow +TE'$        |
| #E'T+   | +id*id#   |                              |
| #E'T    | id*id#    | $T \rightarrow FT'$          |
| #E'T'F  | id*id#    | $F \rightarrow id$           |
| #E'T'id | id*id#    |                              |

| 分析栈(顶)  | 输入串   | 产生式                          |
|---------|-------|------------------------------|
| #E'T'   | *id # | $T' \rightarrow *FT'$        |
| #E'T'F* | *id # |                              |
| #E'T'F  | id #  | $F \rightarrow id$           |
| #E'T'id | id #  |                              |
| #E'T'   | #     | $T' \rightarrow \varepsilon$ |
| #E'     | #     | $E' \rightarrow \varepsilon$ |
| #       | #     |                              |

## LL(1)分析表的优点:

- 1)是确定的,永远不需要回溯.
- 2)分析表的体积比其他方法小.
- 3)LL(1)分析表使用广泛.

可以把LL(1) 文法推广到LL(K)文法,即:超前搜索K个符号,来决定采用哪个候选式.一般地,很少使用.



# LL(1)分析中的错误处理

---

- 发现错误
  - 栈顶的终结符与当前输入符不匹配
  - 非终结符A于栈顶，面临的输入符为a，但分析表M的 $M[A, a]$ 为空
- 应急恢复策略
  - 跳过输入串中的一些符号直至遇到“同步符号”为止。

# LL(1)分析中的错误处理

- 同步符号集的选择
  - 把FOLLOW(A)中的所有符号作为A的同步符号。跳过输入串中的一些符号直至遇到这些“同步符号”，把A从栈中弹出，可使分析继续
  - 把FIRST(A)中的符号加到A的同步符号集，当FIRST(A)中的符号在输入中出现时，可根据A恢复分析应急恢复策略

# LL(1) 分析中的错误处理

文法G[ E]:

(1)  $E \rightarrow TE'$

(2)  $E' \rightarrow +TE'$

(3)  $E' \rightarrow \varepsilon$

(4)  $T \rightarrow FT'$

(5)  $T' \rightarrow *FT'$

(6)  $T' \rightarrow \varepsilon$

(7)  $F \rightarrow (E)$

(8)  $F \rightarrow i$

|    | i   | +   | *   | (   | )   | #   |
|----|-----|-----|-----|-----|-----|-----|
| E  | (1) |     |     | (1) | syn | syn |
| E' |     | (2) |     |     | (3) | (3) |
| T  | (4) | syn |     | (4) | syn | syn |
| T' |     | (6) | (5) |     | (6) | (6) |
| F  | (8) | syn | syn | (7) | syn | syn |