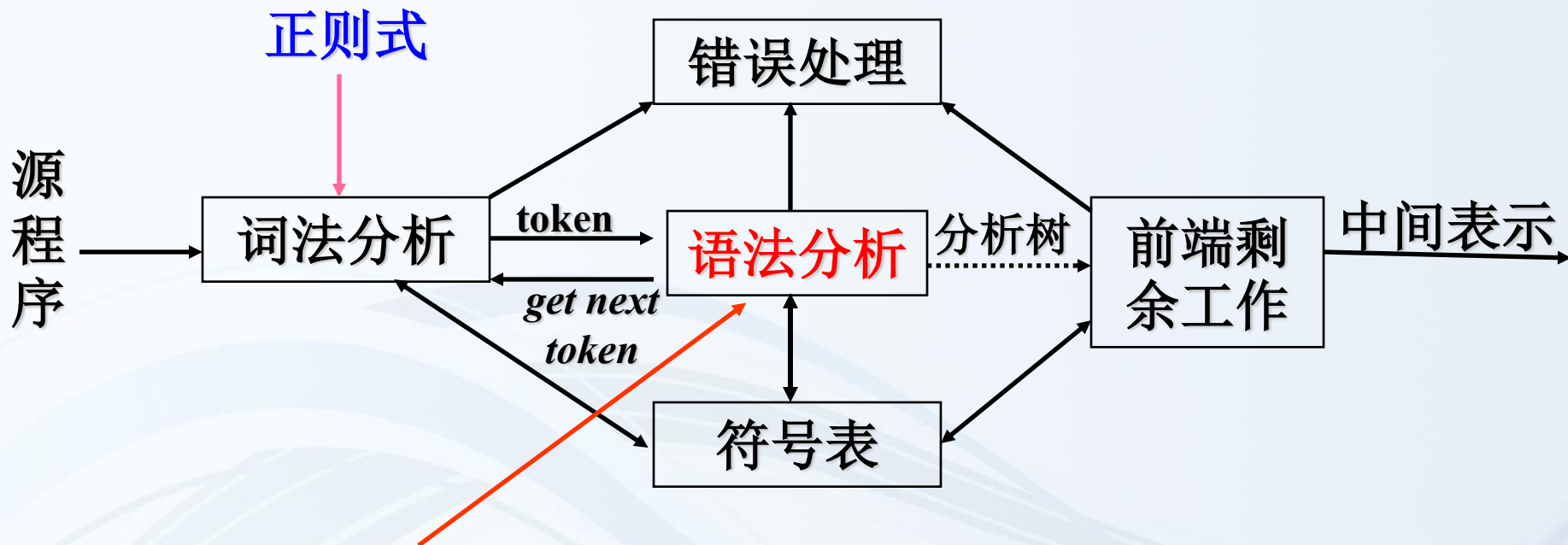


语法分析

4.1 引论



基于文法检查记号流的构成
生成分析树
发现、报告语法错误，并恢复

语法分析器的位置

语法分析和文法

❖ 为什么通过文法形式化地描述语言？

- 精确易懂的表示方法
- 编译器构建工具可以基于文法来自动构建编译器
- 使语言的演化变得容易

❖ 为什么不继续使用正则式描述语言？

正则式的描述能力不足（正则式能描述正则语言，但不足以描述上下文无关语言）

证明： 不存在识别 $L = \{x^n y x^n \mid n = 0, 1, 2, \dots\}$ 的 DFA

反证：假设有一个状态数为 K 的 DFA 能够识别 L

则对于 $K+1$ 个不同的输入 $x, x^2, x^3, \dots, x^K, x^{K+1}$

必然存在 a, b ($1 \leq a < b \leq K+1$)

DFA 在处理 x^a 和 x^b 后转向同一状态 S

因为 DFA 接受 $x^a y x^a$ ，所以在状态 S 继续输入 $y x^a$ 将到达终态，

因此该 DFA 接受 $x^b y x^a$ 。

由于 a, b 不等，所以此 DFA 识别的语言不是 L 。

所以原假设不成立，即不存在识别 L 的 DFA。

4.2 上下文无关文法（复习）

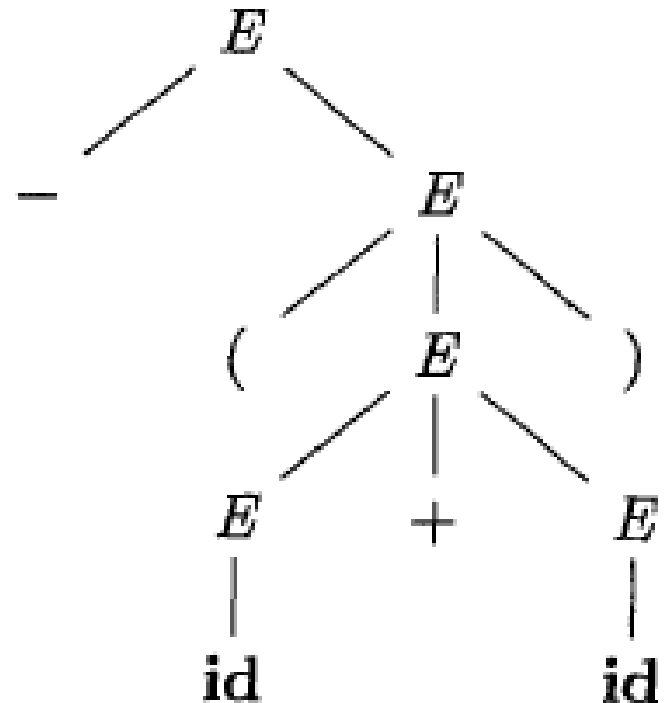
❖ 文法

$E \rightarrow E + E \mid E * E \mid (E) \mid -E \mid \text{id}$

❖ 推导

$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E)$
 $\Rightarrow -(\text{id}+E) \Rightarrow -(\text{id}+\text{id})$

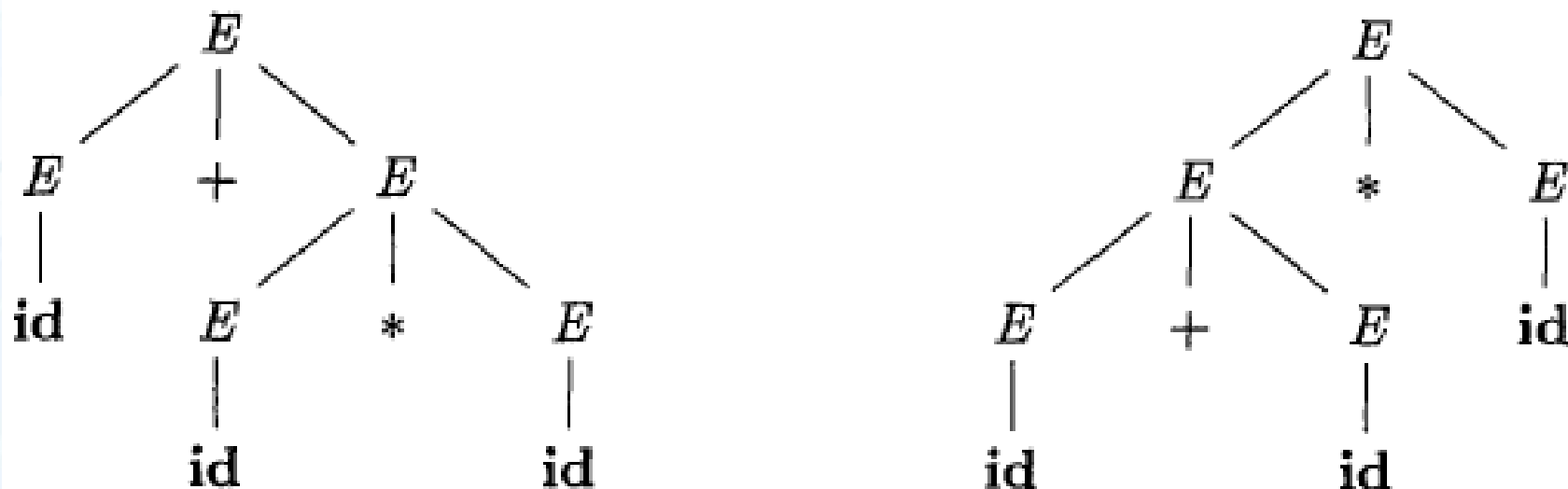
❖ 语法分析树



文法的二义性（复习）

❖ 一个文法是二义的，如果：

- 能为一个句子生成多棵语法分析树（推荐）
- 能为一个句子生成多个最左推导
- 能为一个句子生成多个最右推导



4.3 设计文法

二义性的去除

if E_1 then if E_2 then S_1 else S_2

if E_1 then

if E_2 then

S_1

else

S_2

vs

if E_1 then

if E_2 then

S_1

else

S_2

stmt* → if *expr* then *stmt

| if *expr* then *stmt* else *stmt*

| other (any other statement)

二义性的去除（续）

原来的文法:

$stmt \rightarrow$ **if** *expr* **then** *stmt*
| **if** *expr* **then** *stmt* **else** *stmt*
| **other** (any other statement)

修改后的文法:

$stmt \rightarrow$ *matched_stmt* | *unmatched_stmt*
 $matched_stmt \rightarrow$ **if** *expr* **then** *matched_stmt* **else** *matched_stmt* | **other**
 $unmatched_stmt \rightarrow$ **if** *expr* **then** *stmt*
| **if** *expr* **then** *matched_stmt* **else** *unmatched_stmt*

递归

❖ 直接递归 (规则递归) : 产生式体的左部和右部具有相同非终结符。

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

❖ 间接递归 (文法递归) : 文法中的某个非终结符, 能建立一个推导过程, 得到的符号串中又出现该非终结符本身。

$$A \stackrel{+}{\Rightarrow} \alpha A \beta$$

通过递归, 我们能有限规则去定义无穷集合的语言

直接左递归(*direct left recursion*)

❖ 产生式体的最左边的符号和产生式头部相同。

❖ 左递归的消除:

将 $A \rightarrow A\alpha \mid \beta$

改成: $A \rightarrow \beta R \quad R \rightarrow \alpha R \mid \varepsilon$

例: 消除文法 $\text{expr} \rightarrow \text{expr} + \text{term} \mid \text{term}$ 的左递归

解: 引入非终结符 R , 将文法改写为

$\text{expr} \rightarrow \text{term } R$

$R \rightarrow + \text{term } R \mid \varepsilon$

直接左递归的消除

for (each non-terminal A) {

(1) 将A为左部的规则组织成:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_m \mid \delta_1 \mid \delta_2 \mid \dots \mid \delta_n$$

(2) 引入非终结符A', 将以上规则改写成:

$$A \rightarrow \delta_1 A' \mid \delta_2 A' \mid \dots \mid \delta_n A'$$

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

}

间接左递归(*indirect left recursion*)

❖ 定义：文法G存在推导 $A \xRightarrow{+} A\alpha$

由多步推导带来的左递归

$$\left. \begin{array}{l} S \rightarrow Aa \mid b \\ A \rightarrow Sd \mid \varepsilon \end{array} \right\} S \Rightarrow Aa \Rightarrow Sda$$

例：以下文法加入哪一条规则可导致左递归（多选）

$S \rightarrow A$ $A \rightarrow B|C$ $B \rightarrow (C)$

$C \rightarrow B+C|D$ $D \rightarrow 1|0$

(A) $C \rightarrow C+B$

(B) $D \rightarrow B$

(C) $C \rightarrow 1C$

(D) $B \rightarrow C$

(E) $A \rightarrow D$

(F) $D \rightarrow A$

间接左递归的消除

设定非终结符顺序，假设依次记为 A_1, A_2, \dots, A_n .

for (each i from 1 to n) {

(1) 消除与 A_i 相关的间接左递归:

for (each j from 1 to $i-1$) {

向 $A_i \rightarrow A_j \gamma$ 中代入以 A_j 为头的规则，得 $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$,

其中 $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ 为所有以 A_j 为头的规则

}

(2) 消除与 A_i 相关的直接左递归

}

输入文法有环 ($A \overset{\pm}{\Rightarrow} A$)
或 ϵ 规则 ($A \rightarrow \epsilon$) 时,
无法保证消除左递归。

消除间接左递归（例）

消除文法 $G[S]$ 的左递归: $S \rightarrow Aa|b$ $A \rightarrow Ac|Sd|\epsilon$

解: 设定非终结符顺序为 S, A

S 不存在左递归。

消除 A 存在的间接左递归, 变为: $A \rightarrow Ac | Aad | bd | \epsilon$

消除 A 存在的直接左递归: 引入非终结符 A' , 变为:

$A \rightarrow bdA' | A'$ $A' \rightarrow cA' | adA' | \epsilon$

所得文法为:

$S \rightarrow Aa | b$ $A \rightarrow bdA' | A'$ $A' \rightarrow cA' | adA' | \epsilon$

消除间接左递归（例2）

消除文法G[S]的左递归: $S \rightarrow Qc|c$ $Q \rightarrow Rb|b$ $R \rightarrow Sa|a$

解: 设定非终结符顺序为S,Q,R

S和Q不存在左递归。

R存在间接左递归 $R \rightarrow Sa$, 将S代入得: $R \rightarrow Qca|ca|a$

R存在间接左递归 $R \rightarrow Qca$, 将Q代入得: $R \rightarrow Rbca|bca|ca|a$

R存在直接左递归 $R \rightarrow Rbca$, 引入R', 消除左递归得:

$R \rightarrow bcaR' | caR' | aR'$ $R' \rightarrow bcaR' | \epsilon$

所得文法为:

$S \rightarrow Qc|c$ $Q \rightarrow Rb|b$ $R \rightarrow bcaR' | caR' | aR'$

$R' \rightarrow bcaR' | \epsilon$

提取左公共因子

某个非终结符的多个产生式具有相同前缀时，称存在**左公共因子**

$stmt \rightarrow \underline{\text{if expr then stmt else stmt}}$

$|\ \underline{\text{if expr then stmt}}$

$A \rightarrow \underline{\alpha}\beta_1 | \underline{\alpha}\beta_2$

左公共因子

转化为:

$A \rightarrow \alpha A'$

$A' \rightarrow \beta_1 | \beta_2$

Choose the alternative that correctly left factors “if” statements in the given grammar

```
EXPR → if BOOL then { EXPR }
      | if BOOL then { EXPR } else { EXPR }
      | ...
BOOL → true | false
```

A EXPR → if true then { EXPR }
| if false then { EXPR }
| if true then { EXPR } else { EXPR }
| if false then { EXPR } else { EXPR }
| ...

C EXPR → EXPR' | EXPR' else { EXPR }
EXPR' → if BOOL then { EXPR }
| ...
BOOL → true | false

B EXPR → if BOOL EXPR'
| ...
EXPR' → then { EXPR }
| then { EXPR } else { EXPR }
BOOL → true | false

D EXPR → if BOOL then { EXPR } EXPR'
| ...
EXPR' → else { EXPR } | ε
BOOL → true | false

提取左公共因子（算法）

对每个非终结符A

找到最长公共前缀 α （并非针对所有产生式）

若 $\alpha \neq \varepsilon$

将 $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \gamma$ (γ 代表所有不以 α 开头的产生式组)

替换为: $A \rightarrow \alpha A' \mid \gamma$ $A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

反复进行以上替换，直至 $\alpha = \varepsilon$

提取左公共因子（例）

改造以下文法G，以利于无回溯的自顶向下分析：

$$S \rightarrow abcd \mid abce \mid abf$$

解：S为左部的产生式体的最长公共前缀为abc, 引入A，将文法变为： $S \rightarrow abcA \mid abf$ $A \rightarrow d \mid e$

S为左部的产生式体的最长公共前缀为ab, 引入B，得：

$$S \rightarrow abB \quad B \rightarrow cA \mid f$$

所得文法为：

$$S \rightarrow abB \quad A \rightarrow d \mid e \quad B \rightarrow cA \mid f$$

改造以下文法，使其**可能**用于无回溯的自顶向下分析。

$$S \rightarrow SS+ \mid SS^* \mid a$$

解：（1）提取左公共因子SS得：

$$S \rightarrow SSA \mid a \quad A \rightarrow + \mid *$$

（2）设定文法符号顺序S, A。

仅S存在直接左递归，消除左递归得： $S \rightarrow aR$ $R \rightarrow SAR \mid \epsilon$

最终文法为：

$$S \rightarrow aR \quad R \rightarrow SAR \mid \epsilon \quad A \rightarrow + \mid *$$