

誠朴雄偉
勵學敦行

第四章 语法分析

陈 林





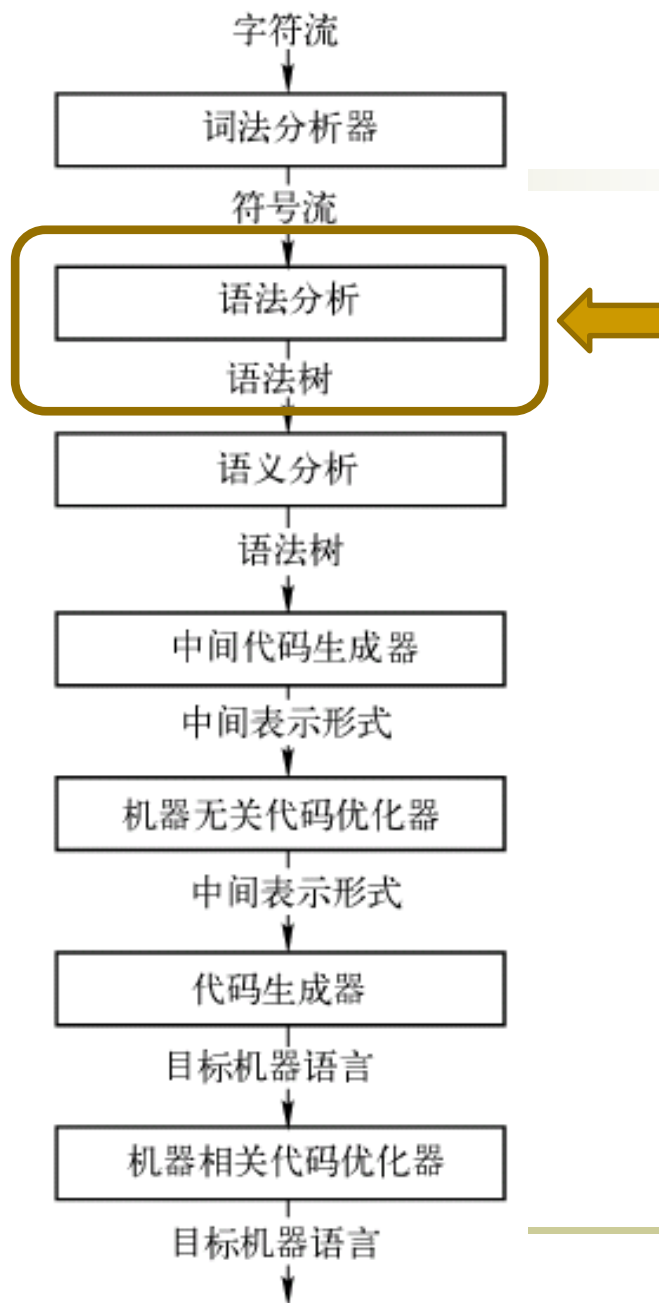
语法分析：提纲



- 引言
- 上下文无关文法
- 语法分析技术
 - 自顶向下
 - 自底向上
- 语法分析器生成工具



符号表



描述

分析



引言



■ 程序设计语言源程序的构成： 语法结构

The screenshot displays a Java source file named `MyClass.java` within an IDE. The code is annotated with labels indicating its syntactic components:

- IPackageDeclaration**: Points to the `package com.mycompany.sample;` statement.
- IImportDeclaration**: Points to the `import java.util.Date;` statement.
- IType**: Points to the `public class MyClass {` declaration.
- IField**: Points to the `static int count = 0;` and `private int localCount;` declarations.
- IMethod**: Points to the `MyClass() {`, `main(String[] args) {`, and `doSomething() {` method declarations.

The Outline window on the right shows the project structure:

- com.mycompany.sample
 - import declarations
 - MyClass
 - count : int
 - localCount : int
 - MyClass()
 - main(String[])
 - doSomething()

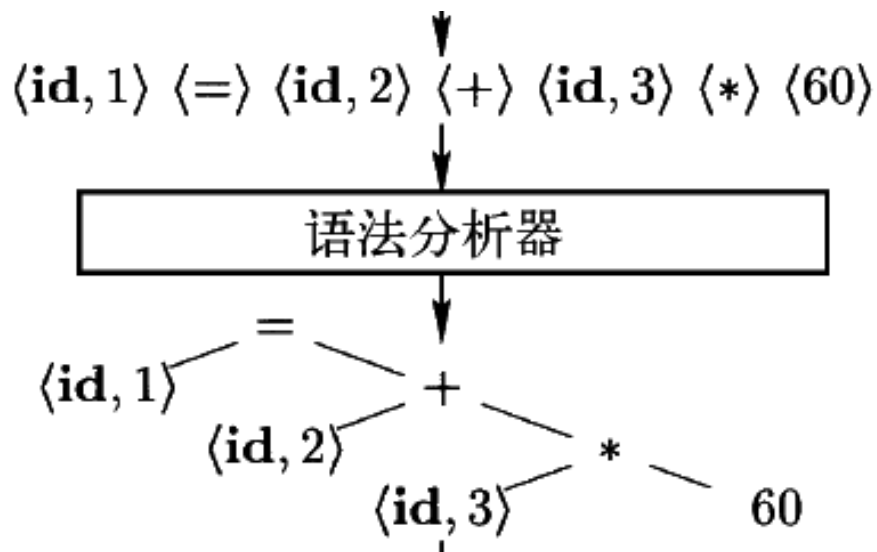


一个实例



○ $\text{position} = \text{initial} + \text{rate} * 60$

<id, 1> <=, > <id, 2> <+, > <id, 3> <*, > <number, 4>



语法分析



引言



- **文法**：一种用于描述程序设计语言语法的表示方法，能够自然地描述程序设计语言构造的**层次化**语法结构
 - 文法给出了一个程序设计语言的精确易懂的**语法规约**
 - 可以基于文法构造语法分析器，帮助确定源程序的语法结构
 - 语法结构有助于把源程序翻译为正确的目标代码
 - 文法的扩展性，有助于迭代的语言演化



语法分析器

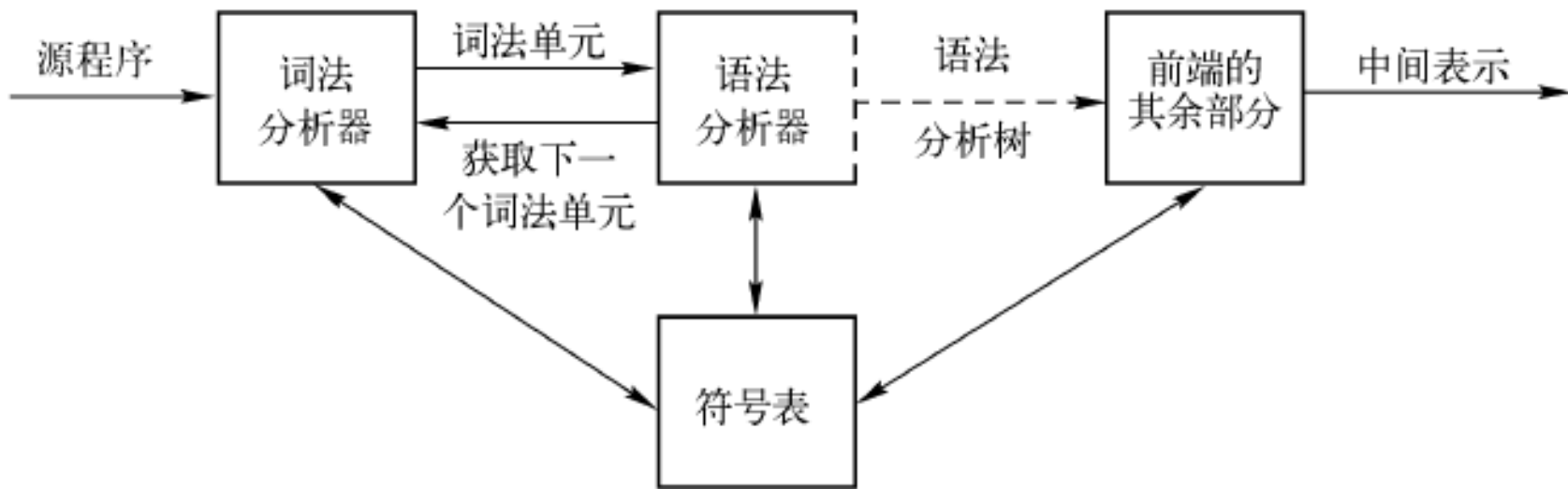


- 输入：词法分析器输出的词法单元序列
- 输出：语法树表示
- 语法分析器功能：
 - 验证输入源程序的合法性，输出良构程序的语法结构
 - 对于病构的程序，能够报告语法错误，进行错误恢复
- 语法分析器的类型：
 - 通用型
 - 自顶向下：通常处理LL文法
 - 自底向上：通常处理LR文法

从左到右扫描字符



语法分析器



- 类型检查，语义分析，翻译生成中间代码等往往和语法分析过程**交错完成**，实践中往往和语法分析放入一个模块，图上用“前端的其余部分”表示上述活动



代表性文法



LR文法类
自底向上分析

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned}$$

LL类文法
(无左递归)
自顶向下分析

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned}$$



上下文无关文法



- 上下文无关文法 (Context Free Grammar, CFG)
- 上下文无关文法是一种能够很好描述程序设计语言的表示方法

$stmt \rightarrow \text{if } (expr) stmt \text{ else } stmt$

$expr \rightarrow \dots\dots$

$stmt \rightarrow \dots\dots$



上下文无关文法的定义



- 一个CFG由以下几个部分构成
 - 终结符号
 - 组成串的基本符号，与“词法单元名字”同义
 - 非终结符号
 - 语法变量，表示特定串的集合
 - 给出了语言的层次结构，这种层次结构是语法分析和翻译的关键
 - 一个开始符号
 - 某个特定的非终结符号，其表示的串集合是这个文法生成的语言
 - 一组产生式
 - 描述将终结符号和非终结符号组合成串的方法
 - 产生式左部（头）是一个非终结符号
 - 符号“ \rightarrow ”
 - 一个由零个或多个终结符号与非终结符号组成的产生式右部（体）



用于描述算术表达式的文法定义



<i>expression</i>	\rightarrow	<i>expression</i> + <i>term</i>
<i>expression</i>	\rightarrow	<i>expression</i> - <i>term</i>
<i>expression</i>	\rightarrow	<i>term</i>
<i>term</i>	\rightarrow	<i>term</i> * <i>factor</i>
<i>term</i>	\rightarrow	<i>term</i> / <i>factor</i>
<i>term</i>	\rightarrow	<i>factor</i>
<i>factor</i>	\rightarrow	(<i>expression</i>)
<i>factor</i>	\rightarrow	id



符号表示约定



- 终结符号: $a\ b\ +\ 3\ \mathbf{id}\dots$
- 非终结符号: $A\ B\ C\dots, S, stmt$
- 文法符号: $X\ Y\dots$
 $E \rightarrow E + T \mid E - T \mid T$
- 终结符号串: $u\ v\ w\dots$
 $T \rightarrow T * F \mid T / F \mid F$
- 文法符号串: $\alpha\ \beta\dots$
 $F \rightarrow (E) \mid \mathbf{id}$
- 不同可选体: $\alpha_1\ \alpha_2\ \alpha_3\dots$
- 开始符号: S



推导



- 产生式又可称为重写规则 (**Rewriting rule**)
- 推导
 - 将待处理的串中的某个非终结符号替换为这个非终结符号的某个产生式的体
 - 从开始符号出发, 不断进行上面的替换, 就可以得到文法的不同句型
- 推导的实例
 - 文法: $E \rightarrow -E \mid E+E \mid E * E \mid (E) \mid id$
 - 从E到 $-(id)$ 的推导序列: $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(id)$



推导



■ 推导的一般性定义

- 如果 $A \rightarrow \gamma$ 是一个产生式, 那么 $\alpha A \beta \Rightarrow \alpha \gamma \beta$

■ 经过零步或者多步推导出: $\xRightarrow{*}$

- 对于任何串 α , $\alpha \xRightarrow{*} \alpha$
- 如果 $\alpha \xRightarrow{*} \beta$ 且 $\beta \Rightarrow \gamma$, 那么 $\alpha \xRightarrow{*} \gamma$

■ 经过一步或者多步推导出: $\xRightarrow{+}$

- $\alpha \xRightarrow{*} \beta$ 且 α 不等于 β 等价于 $\alpha \xRightarrow{+} \beta$

■ 最左(右)推导

- 符号: $\xRightarrow{*}_{lm}$ $\xRightarrow{*}_{rm}$



推导的例子



■ 文法实例

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

从上述文法推导出串 $\mathbf{-(id+id)}$

$$E \Rightarrow - \quad E \Rightarrow - (E) \Rightarrow - (E + E) \Rightarrow - (\mathbf{id} + E) \Rightarrow - (\mathbf{id} + \mathbf{id})$$



句型/句子/语言



- 句型 (sentential form) :
 - 如果 $S \xRightarrow{*} \alpha$, 那么 α 就是文法的一个句型
 - 可能既包含非终结符, 又包含终结符号; 可以是空串
- 句子 (sentence)
 - 文法的句子就是不包含非终结符号的句型
- 语言
 - 文法 G 的语言就是 G 的句子的集合, 记为 $L(G)$
 - ω 在 $L(G)$ 中当且仅当 ω 是 G 的句子, 即 $S \xRightarrow{*} \omega$



推导的问题



- 从推导的角度看，语法分析的任务是：接受一个终结符号串作为输入，找出从文法的开始符号推导出这个串的方法

$$E \rightarrow E+E \mid E * E \mid -E \mid (E) \mid id$$

推导出串 $-(id+id*id)$



推导的问题



- 推导中可能遇到的两个问题
 - 每一步替换哪个非终结符号？
 - 若以这个非终结符号为头的产生式有多个，用哪个产生式的右部替换？

$$E \rightarrow E+E \mid E * E \mid -E \mid (E) \mid id$$

推导出串 $-(id+id*id)$



非终结符号的替换顺序

- 通常使用两种方式进行推导
 - 最左推导：总是选择每个句型的最左非终结符号。记作 \Rightarrow_{lm}
 - 最右推导：总是选择最右边的非终结符号。记作 \Rightarrow_{rm}
- 每个最左推导步骤可以写成 $wA\gamma \Rightarrow_{lm} w\delta\gamma$
应用产生式： $A \rightarrow \delta$
- 如果 α 经过最左推导得到 β ，记作 $\alpha \xRightarrow{*}_{lm} \beta$
- 最左句型：**S** 是文法 **G** 的识别符号，如果 $S \xRightarrow{*}_{lm} \alpha$ 则 α 是 **G** 的最左句型

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

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



语法分析树



- 推导的图形表示形式
 - 根结点的标号是文法的开始符号
 - 每个叶子结点的标号是非终结符号、终结符号或 ϵ
 - 每个内部节点的标号是非终结符号
 - 每个内部结点表示某个产生式的一次应用
 - 内部结点的标号为产生式头，结点的子结点从左到右是产生式的体
- 有时允许树的根不是开始符号（对应于某个短语）
- 树的叶子组成的序列是根的文法符号的句型
- 一棵分析树可对应多个推导序列，但是分析树和最左（右）推导序列之间具有一一对应关系



语法分析树



- 推导的图形表示形式，树上
看不出推导的顺序
- 能够反映串的语法层次结构
- 语法分析树
 - 内部节点：对应于一个非终结符号
 - 子节点：对应于其父节点为头的产生式体
 - 叶子节点：可以是终结符号或非终结符号，从左到右排列可以得到一个句型，称为这棵树的结果

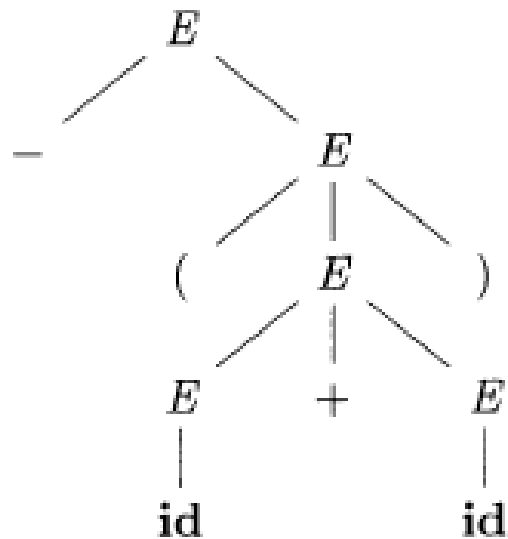


图 4-3 $-(\text{id} + \text{id})$ 的
语法分析树

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

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

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



推导和语法树



- 对于推导中的每个句型 α_i ，我们都可以构造出一个结果为 α_i 的语法树
 - 构造过程是对 i 的一次归纳过程

$$\alpha_{i-1} = X_1 X_2 \cdots X_k$$

- 假设已经构造出 $X_j \rightarrow \beta$

$$\beta = Y_1 Y_2 \cdots Y_m$$

$$\alpha_i = X_1 X_2 \cdots X_{j-1} \beta X_{j+1} \cdots X_k$$

- 在推导的第 i 步，对 α_i 应用规则 $X_j \rightarrow \beta$



从推导序列构造分析树



■ 假设有推导序列:

○ $A \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \cdots \Rightarrow \alpha_n$

■ 算法:

- 初始化: α_1 的分析树是标号为A的单个结点
- 假设已经构造出 $\alpha_{i-1} = X_1 X_2 \cdots X_k$ 的分析树, 且 α_{i-1} 到 α_i 的推导是将 X_j 替换为 β , 那么在当前分析树中找出第 j 个非 ε 结点, 向这个结点增加构成 β 的子结点。(如果 $\beta = \varepsilon$, 则增加一个标号为 ε 的子结点)



推导与语法树的例子



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

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

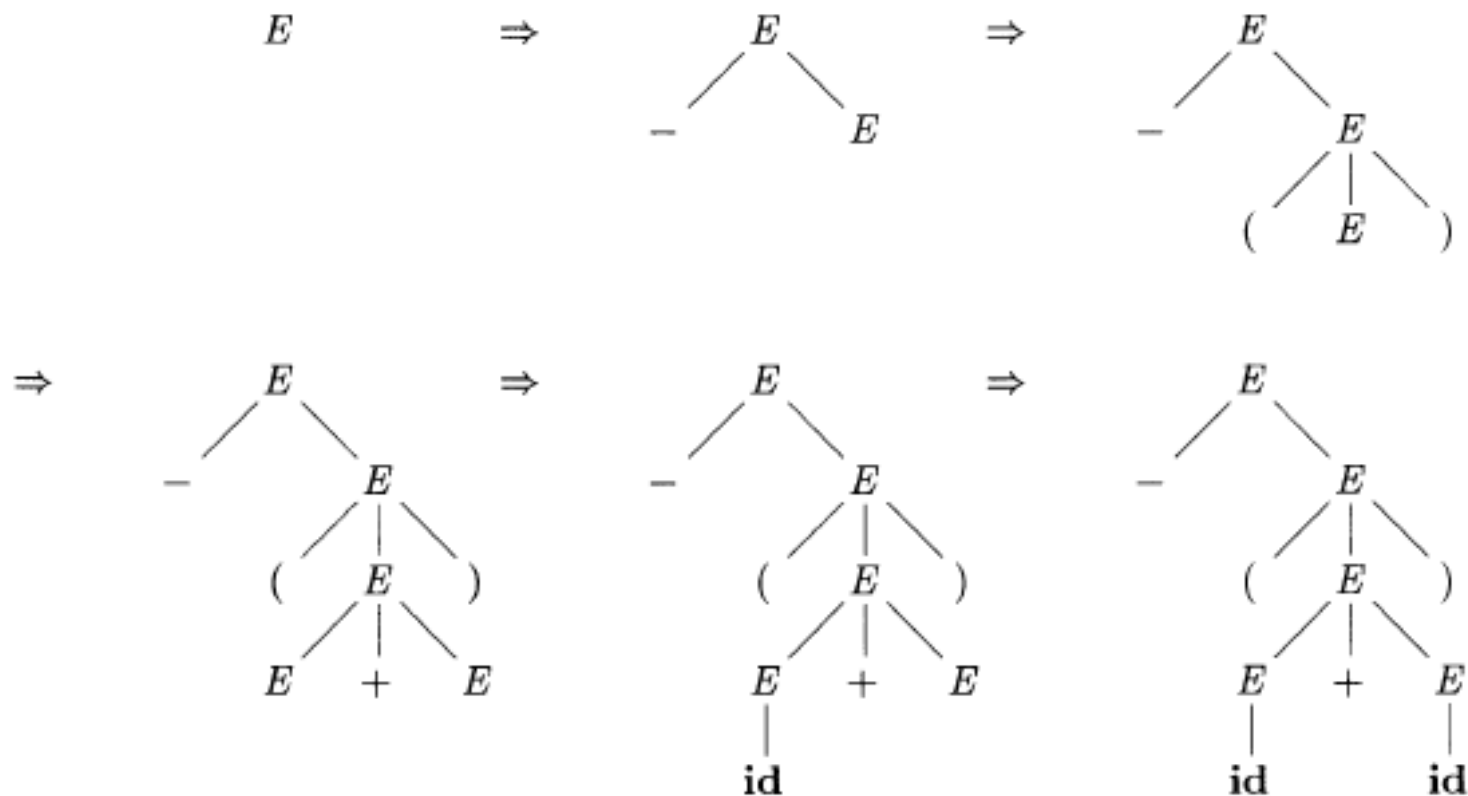


图 4-4 推导(4.8)的语法分析树序列



二义性文法



$E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid id$

- 如果一个文法可以为一个句子生成多棵不同的语法分析树，则该文法为二义性文法

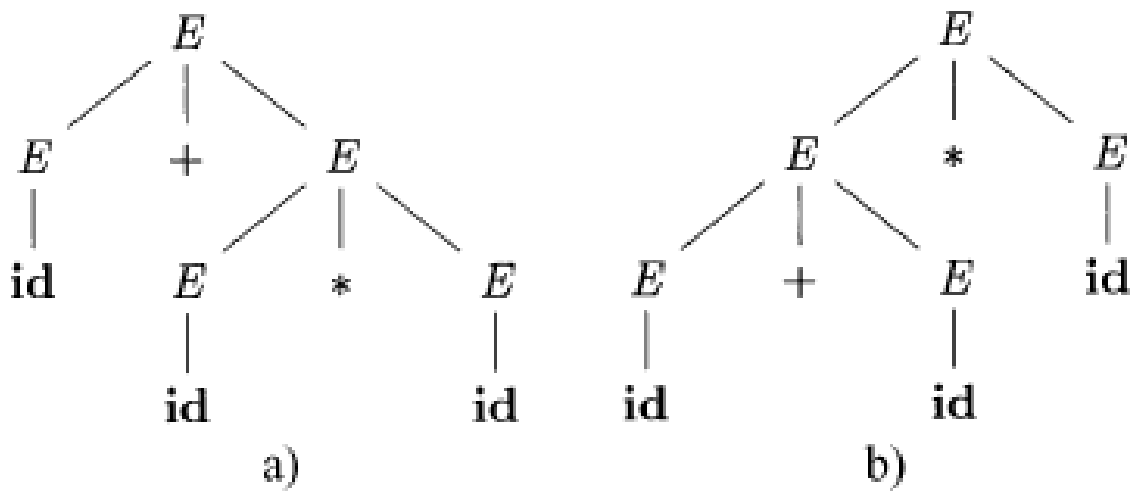


图 4-5 $id + id * id$ 的两棵语法树

- 通常情况下，我们需要无二义性的文法



二义性



- 程序设计语言的文法通常都应该是无二义性的
 - 否则就会导致一个程序有多种“正确”的解释。
 - 比如文法： $E \rightarrow -E \mid E+E \mid E * E \mid (E) \mid \text{id}$ 的句子 $a+b*c$
- 但有些二义性的情况可以方便文法或语法分析器的设计
 - 需要消二义性规则来剔除不要的语法分析树
 - 比如：先乘除后加减

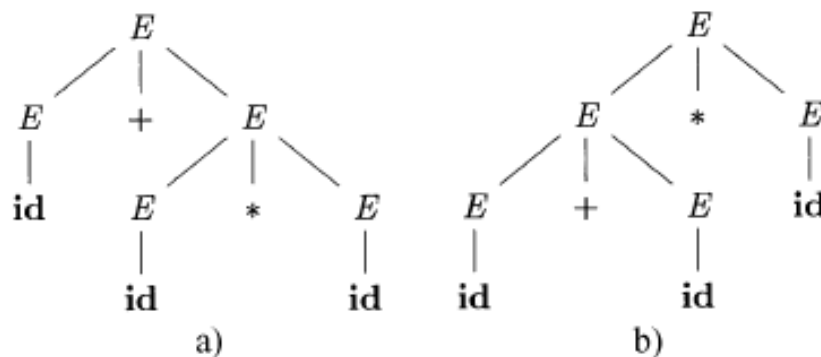


图 4-5 $\text{id} + \text{id} * \text{id}$ 的两棵语法树



文法的二义性

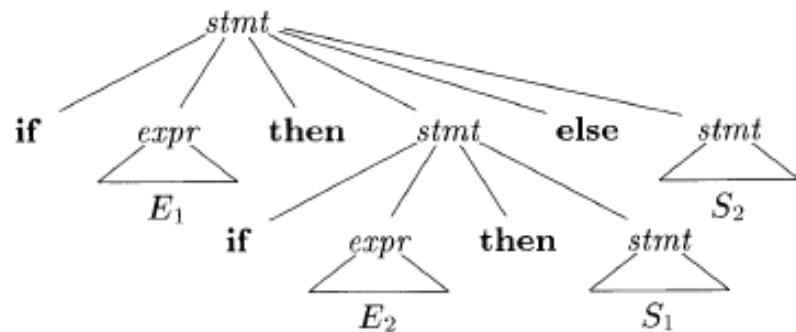
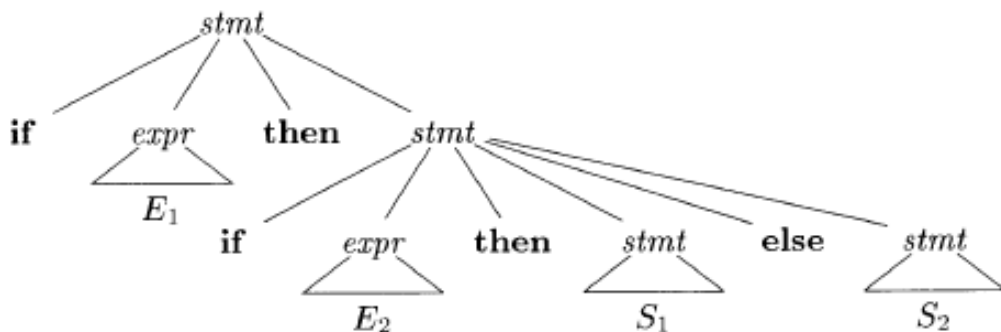


- 例子:

$stmt \rightarrow$

- if** *expr* **then** *stmt*
- if** *expr* **then** *stmt* **else** *stmt*
- other**

- if** E_1 **then** **if** E_2 **then** S_1 **else** S_2 的两棵语法树
- 这个文法是可以被改造成等价的无二义性的文法





文法及其生成的语言



- 语言是由文法的开始符号出发，能够推导得到的所有句子的集合
- 文法G: $S \rightarrow aS|a|b$, $L(G)=\{a^i(a|b), i \geq 0\}$
- 文法G: $S \rightarrow aSb|ab$
 $L(G)=\{a^n b^n, n \geq 1\}$
- 文法G: $S \rightarrow (S)S|\epsilon$
 $L(G)=\{\text{所有具有对称括号对的串}\}$



验证文法生成的语言



- 验证文法 G 生成语言 L 可以帮助我们了解文法可以生成什么样的语言
- 基本步骤:
 - 首先证明 $L(G) \subseteq L$: G 生成的每个串都在 L 中
 - 然后证明 $L \subseteq L(G)$: L 的每个串都能由 G 生成
 - 一般使用数学归纳法
- 证明 $L(G) \subseteq L$: 按照推导序列长度进行数学归纳
- 证明 $L \subseteq L(G)$: 按照符号串的长度来构造推导序列
- 实例: 文法 G : $S \rightarrow (S)S | \varepsilon$, $L(G) = \{\text{所有具有对称括号对的串}\}$



词法分析与语法分析



- 上下文无关文法和正则表达式
- 正则表达式 → 词法
- 上下文无关文法 → 语法



词法分析和语法分析比较



阶段	输入	输出	描述体系
词法分析	源程序符号串	词法单元序列	正则表达式
句法分析	词法单元序列	语法树	上下文无关文法

- 文法比正则表达式描述能力更强
- 正则表达式描述词法单元比较简洁
- 基于正则表达式构造的词法分析器效率更高
- 正则表达式适合描述词法结构，文法适合描述嵌套结构

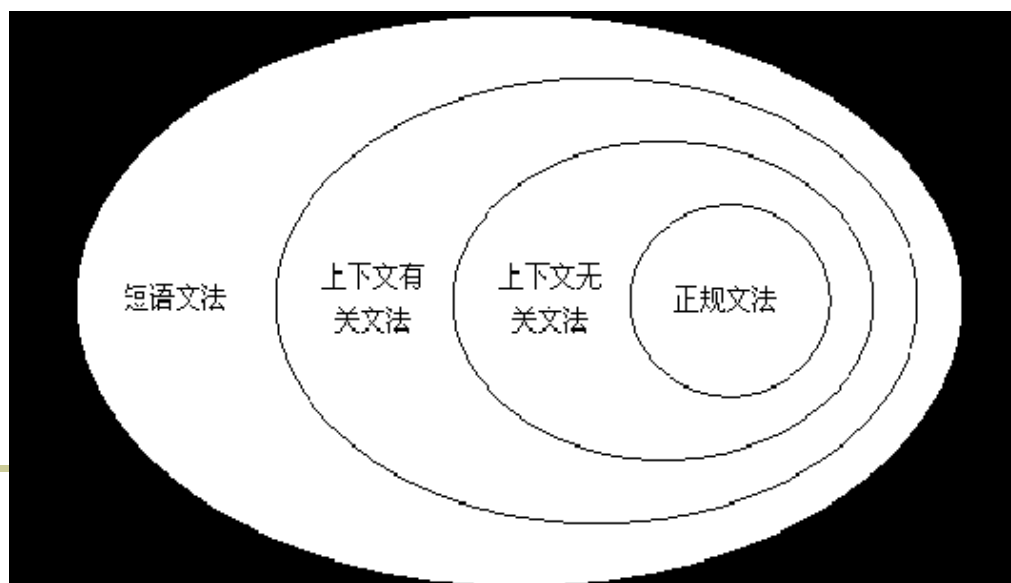


补充



■ 文法的类别, Chomsky文法类

- 0型文法(短语结构文法), $\alpha \rightarrow \beta$
- 1型文法(上下文相关文法), $\alpha A \beta \rightarrow \alpha \gamma \beta$
- 2型文法(上下文无关文法), $A \rightarrow \gamma$
- 3型文法(正则文法), $A \rightarrow t, A \rightarrow tB$





上下文无关文法和正则表达式

- 上下文无关文法的表达能力更强
- 每个正则表达式都可以用一个上下文无关文法来描述，反之不成立
- 每个正则语言都是一个上下文无关语言，反之不成立
- 正则表达式 $(a|b)^*abb$ 等价于文法

$$A_0 \rightarrow a A_0 \mid b A_0 \mid a A_1$$

$$A_1 \rightarrow b A_2$$

$$A_2 \rightarrow b A_3$$

$$A_3 \rightarrow \varepsilon$$



上下文无关文法和正则表达式



- 上下文无关文法比正则表达式的能力更强，即：
 - 一些用文法描述的语言不能用正则文法描述
 - 所有的正则语言都可以使用文法描述
- 证明：
 - 首先 $S \rightarrow aSb \mid ab$ 描述了语言 $\{a^n b^n \mid n > 0\}$ ，但是这个语言无法用DFA识别
 - 假设有DFA识别此语言，且这个文法有K个状态。那么在识别 $a^{k+1} \dots$ 的输入串时，必然两次到达同一个状态。
 - 设自动机在第i个和第j个a时进入同一个状态，那么：因为DFA识别L， $a^i b^j$ 必然到达接受状态，因此 $a^i b^j$ 必然也到达接受状态
 - 直观地讲：有穷自动机不能数（大）数



上下文无关文法和正则表达式



■ 证明（续）

- 其次证明：任何正则语言都可以表示为上下文无关文法的语言
- 任何正则语言都必然有一个NFA。对于任意的NFA构造如下的上下文无关文法
 - 对NFA的每个状态 i ，创建非终结符号 A_i
 - 如果有 i 在输入 a 上到达 j 的转换，增加产生式 $A_i \rightarrow aA_j$
 - 如果 i 在输入 ϵ 上到达 j ，那么增加产生式 $A_i \rightarrow A_j$.
 - 如果 i 是一个接受状态，增加产生式 $A_i \rightarrow \epsilon$
 - 如果 i 是开始状态，令 A_i 为所得文法的开始符号



为NFA构造等价文法



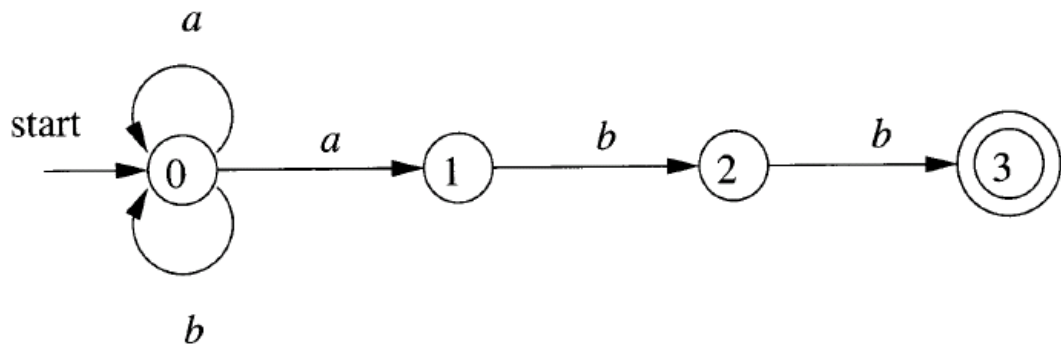
- 1) 对于 NFA 的每个状态 i , 创建一个非终结符号 A_i 。
- 2) 如果状态 i 有一个在输入 a 上到达状态 j 的转换, 则加入产生式 $A_i \rightarrow aA_j$ 。如果状态 i 在输入 ϵ 上到达状态 j , 则加入产生式 $A_i \rightarrow A_j$ 。
- 3) 如果 i 是一个接受状态, 则加入产生式 $A_i \rightarrow \epsilon$ 。
- 4) 如果 i 是自动机的开始状态, 令 A_i 为所得文法的开始符号。

$$A_0 \rightarrow aA_0 \mid aA_1 \mid bA_0$$

$$A_1 \rightarrow bA_2$$

$$A_2 \rightarrow bA_3$$

$$A_3 \rightarrow \epsilon$$



•考虑baabb的推导和接受过程

•NFA接受一个句子的运行过程实际是文法推导出该句子的过程



上下文无关文法和正则表达式



■ 实例：

给出语言 $\{a^n b^j c^i \mid n \geq 0, j \geq 0, 0 \leq i \leq n\}$

- 能用正则文法描述吗？
- 描述它的文法是怎样的？



文法的设计



- 程序设计语言所需要的文法
 - 上下文无关文法足够用来描述语法吗?
 - 标识符必须先声明后使用
 - 语法分析器能够完全按照上下文无关文法来构造吗?
 - 二义性、左递归的文法可能给语法分析器造成很大麻烦
- 可以怎么做?
 - 改造语法分析器，添加语义规则，使它可以做得更多
 - 改造上下文无关文法，消除二义性和左递归



文法的设计

很重要



- 在进行**高效**的语法分析之前，需要对文法做以下处理
 - 消除二义性
 - 文法的二义性：文法可以为一个句子生成多颗不同的树
 - 消除左递归
 - 左递归：文法中一个非终结符号**A**使得对某个串 **α** ，存在一个推导 $A \xRightarrow{+} A\alpha$ ，则称这个文法是左递归的
 - 提取左公因子



消除文法的二义性



$E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid id$

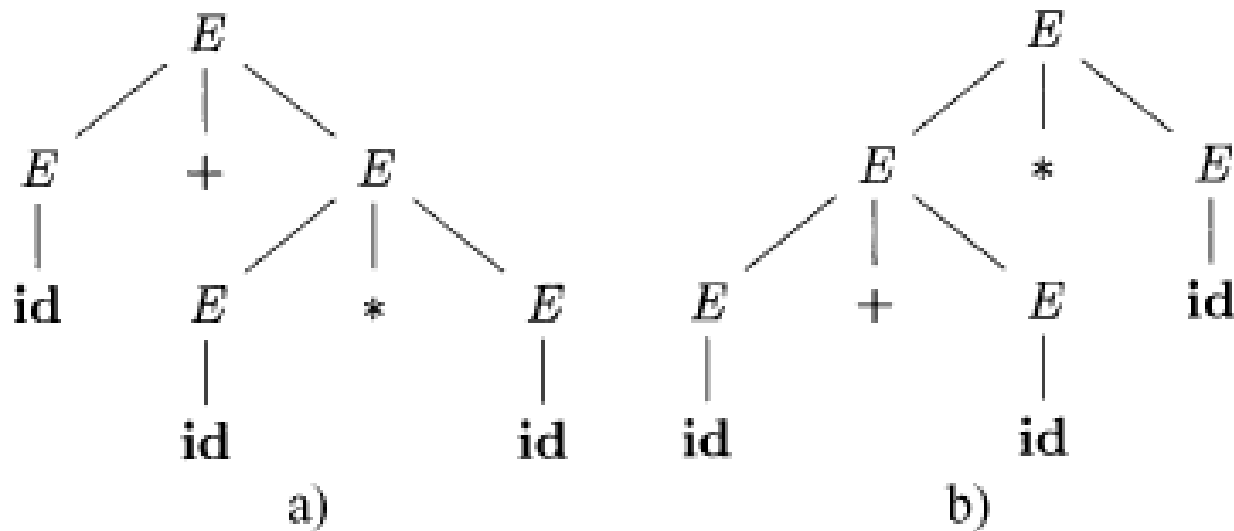


图 4-5 $id + id * id$ 的两棵语法树



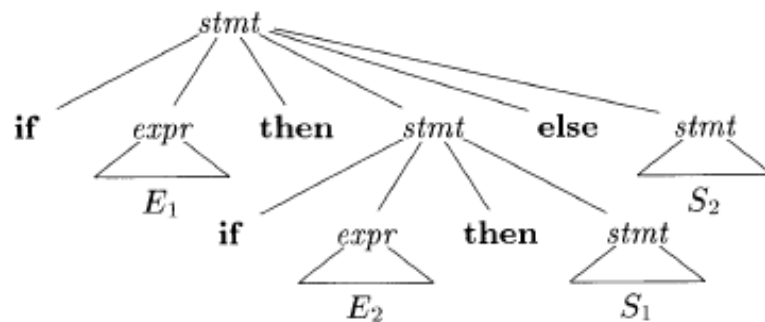
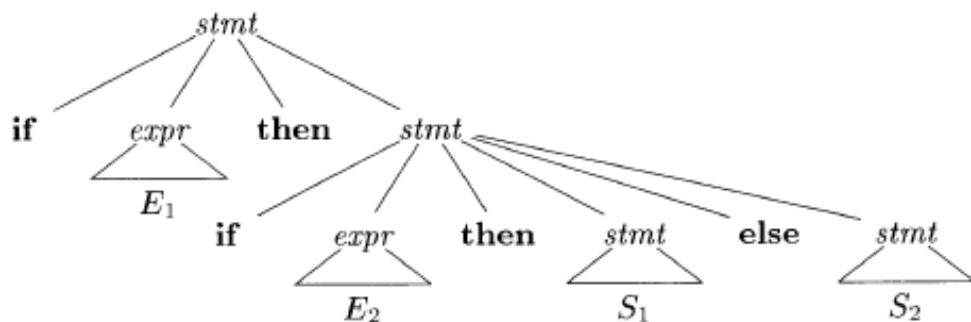
消除文法的二义性



- 一些二义性文法可以被改成等价的无二义性的文法
- 例子:

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

- **if** E_1 **then** **if** E_2 **then** S_1 **else** S_2 的两棵语法树





消除文法的二义性（续）



- 改写文法，基本思想：在一个**then**和一个**else**之间出现的语句必须是“已匹配的”。也就是说**then**和**else**中间的语句不能以一个尚未匹配的**then**结尾
- 解决方案：引入新的非终结符号**matched_stmt**，用来区分是否是成对的**then/else**

<i>stmt</i>	→	<i>matched_stmt</i>
		<i>open_stmt</i>
<i>matched_stmt</i>	→	if <i>expr</i> then <i>matched_stmt</i> else <i>matched_stmt</i>
		other
<i>open_stmt</i>	→	if <i>expr</i> then <i>stmt</i>
		if <i>expr</i> then <i>matched_stmt</i> else <i>open_stmt</i>

if E_1 **then** **if** E_2 **then** S_1 **else** S_2



消除文法的二义性（续）



- 二义性的消除方法没有规律可循
- 通常并不是通过改变文法来消除二义性



消除左递归



■ 左递归的定义

- 文法中一个非终结符号**A**使得对某个串 α , 存在一个推导 $A \Rightarrow^+ A\alpha$ 则称这个文法是左递归的

■ 立即左递归

- 如果存在 $A \rightarrow A\alpha$, 则称为立即左递归

■ 为什么要消除左递归?

- 自顶向下的语法分析技术不能处理左递归的文法

■ 立即左递归的消除实例:

$$\boxed{A \rightarrow A\alpha \mid \beta} \Longrightarrow \boxed{\begin{array}{l} A \rightarrow \beta A' \\ A' \rightarrow \alpha A' \mid \varepsilon \end{array}}$$



立即左递归的消除



- 假设非终结符号A存在立即左递归的情形，假设以A为左部的规则有：

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

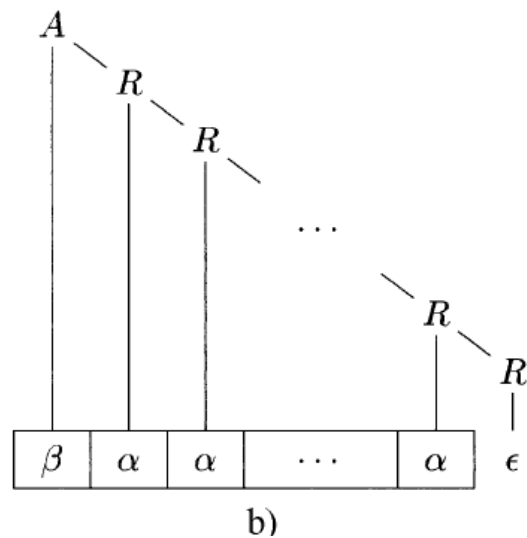
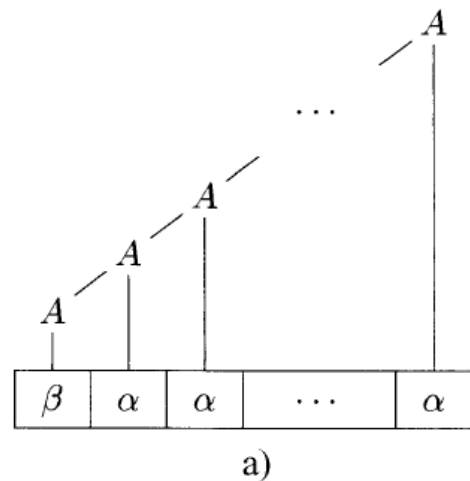
首先对A产生式**分组** (所有 α_i 不等于 ϵ ， β_i 不以A开头)

可以**替换**为

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

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

- 由A生成的串总是以某个 β_i 开头，然后跟上零个或者多个 α_i 的重复





消除立即左递归



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

- 分组

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

- 替换

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

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



立即左递归消除示例



$$\begin{array}{l} E \rightarrow E + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (E) \mid \mathbf{id} \end{array} \quad \Longrightarrow \quad \begin{array}{l} E \rightarrow TE' \\ E' \rightarrow + TE' \mid \epsilon \\ T \rightarrow FT' \\ T' \rightarrow * FT' \mid \epsilon \\ F \rightarrow (E) \mid \mathbf{id} \end{array}$$



消除多步左递归



- 消除立即左递归的方法并不能消除因为两步或多步推导而产生的左递归
- 文法: $S \rightarrow Aa \mid b, \quad A \rightarrow Ac \mid Sd \mid \varepsilon$
- $S \Rightarrow Aa \Rightarrow Sda$
- 如何消除?



消除算法



- 输入：没有环 $A_i \xRightarrow{+} A_i$ 和 $A \rightarrow \varepsilon$
- 输出：一个等价的无左递归的文法
- 算法原理：
 - 给非终结符号**排序**， A_1, A_2, \dots, A_n
 - 如果只有 $A_i \Rightarrow A_j$ ($i < j$)，则不会有左递归
 - 如果发现 $A_i \Rightarrow A_j$ ($i > j$)，**代入** A_j 的当前产生式，若替换后有 A_i 的直接左递归，再**消除**



通用的左递归消除方法



- 输入：没有环和 ϵ 产生式的文法G
- 输出：等价的无左递归的文法
- 步骤：

将文法的非终结符号任意排序为 A_1, A_2, \dots, A_n

for $i=1$ to n do {

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

 {

 将形如 $A_i \rightarrow A_j \gamma$ 的产生式替换为 $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$,

 其中 $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$ 是以 A_j 为左部的所有产生式;

 }

 消除 A_i 的立即左递归;

}



通用的左递归消除方法（续）



1. 内层循环的每一次迭代保证了 A_i 的产生式的右部首符号都比 A_j 更加靠后
2. 当内层循环结束时， A_i 的产生式的右部首符号不先于 A_i
3. 消除立即左递归就保证了 A_i 的产生式的右部首符号必然比 A_i 后。（不能有环和 ϵ 产生式）
4. 当外层循环结束时，所有的产生式都是如此。使用这些产生式当然不会产生左递归的情形



消除算法（示例）



- $S \rightarrow Aa \mid b, \quad A \rightarrow Ac \mid Sd \mid \varepsilon$
- 排序: $S \ A$
- 替换:
 - $i=1$, 没有处理
 - $i=2$, 替换 $A \rightarrow Sd$ 中的 S , 得到 $A \rightarrow Ac \mid Aad \mid bd \mid \varepsilon$
- 消除立即左递归

$$S \rightarrow A a \mid b$$

$$A \rightarrow b d A' \mid A'$$

$$A' \rightarrow c A' \mid a d A' \mid \epsilon$$

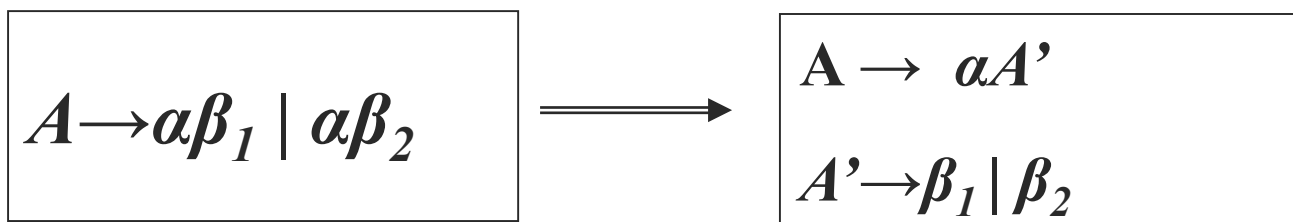


提取左公因子



- 在推导的时候，不知道该如何选择（自顶向下算法会详细描述）

$stmt \rightarrow \mathbf{if} \ expr \ \mathbf{then} \ stmt \ \mathbf{else} \ stmt$
 $\quad | \ \mathbf{if} \ expr \ \mathbf{then} \ stmt$





提取左公因子算法



- 输入：文法 G
- 输出：一个等价的提取了左公因子的文法
- 方法：对于每个非终结符号 A ，找出它的两个或多个可选项之间的最长公共前缀 α ，且 $\alpha \neq \varepsilon$ ，那么将 A 所有的产生式

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

替换为

$$A \rightarrow \alpha A' \mid \gamma$$

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



提取左公因子



$$S \rightarrow i E t S \mid i E t S e S \mid a$$

$$E \rightarrow b$$

对于S而言，最长的前缀是 $i E t S$ ，因此：

$$S \rightarrow i E t S S' \mid a$$

$$S' \rightarrow e S \mid \epsilon$$

$$E \rightarrow b$$



非上下文无关语言的构造



- 抽象语言 $L1 = \{w cw \mid w \text{ 在 } (a|b)^* \text{ 中} \}$
- 这个语言不是上下文无关的语言
- 它抽象地表示了C或者Java中“标识符先声明后使用”的规则
 - 说明了C/Java这些语言不是上下文无关语言，不能使用上下文无关文法描述。
 - 通常用上下文无关文法描述其基本结构，不能用文法描述的特性在语义分析阶段完成。原因是上下文文法具有高效的处理算法。



自顶向下分析技术



- 自顶向下分析可以被看作是为输入串构造语法分析树的问题，也可以看作是一个寻找输入串的**最左推导**的过程
- 问题
 - 在推导的每一步，对非终结符号**A**，应用哪个产生式，以可能产生于输入串相匹配的终结符号串



id+id*id的自顶向下分析



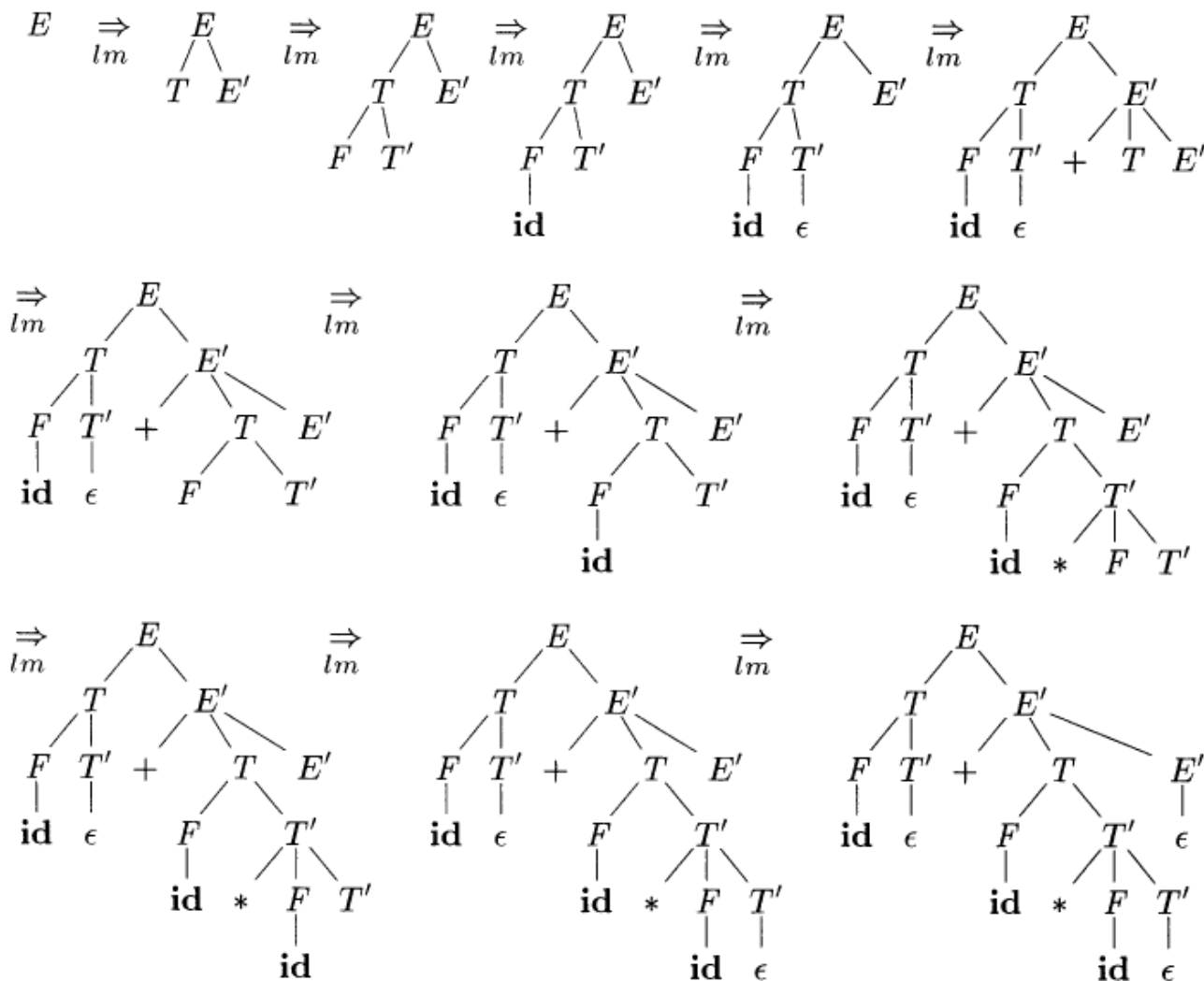
$$E \rightarrow T E'$$

$$E' \rightarrow + T E' \mid \epsilon$$

$$T \rightarrow F T'$$

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

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





自顶向下语法分析



- 问题：对于非终结符号 **A**，选择哪一个产生式
- 一种**通用**的**递归下降**分析框架
 - 由一组过程组成，每个非终结符号对应一个过程
 - 程序的执行从开始符号对应的过程开始
 - 每个过程的功能是：选择一个产生式体，扫描相应的句子。若遇到非终结符号，调用该符号对应的过程

```
void A() {  
1)      选择一个 A 产生式,  $A \rightarrow X_1X_2 \cdots X_k$ ;  
2)      for (  $i = 1$  to  $k$  ) {  
3)          if (  $X_i$  是一个非终结符号 )  
4)              调用过程  $X_i()$ ;  
5)          else if (  $X_i$  等于当前的输入符号  $a$  )  
6)              读入下一个输入符号;  
7)          else /* 发生了一个错误 */;  
          }  
}
```



递归下降分析过程示例



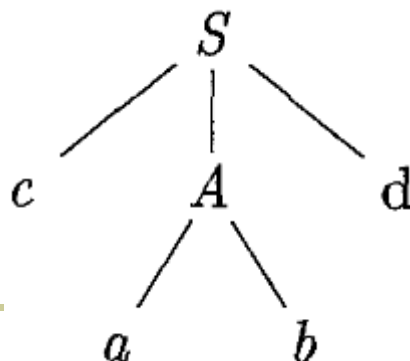
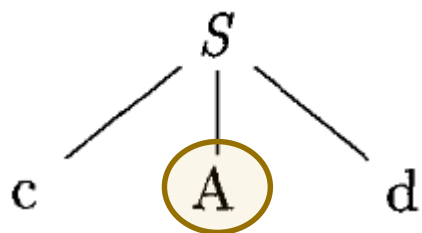
■ 输入串 $w = c a d$

■ 文法:

○ $S \rightarrow cAd$

○ $A \rightarrow ab \mid a$

```
void A() {  
    1) 选择一个  $A$  产生式,  $A \rightarrow X_1X_2 \cdots X_k$ ;  
    2)  for (  $i = 1$  to  $k$  ) {  
    3)      if (  $X_i$  是一个非终结符号 )  
    4)          调用过程  $X_i()$ ;  
    5)      else if (  $X_i$  等于当前的输入符号  $a$  )  
    6)          读入下一个输入符号;  
    7)      else /* 发生了一个错误 */;  
    }  
}
```





递归下降分析过程示例



■ 输入串 $w = c a d$

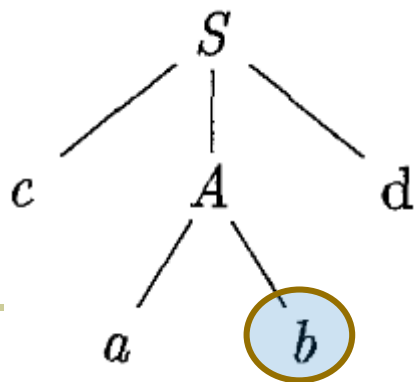
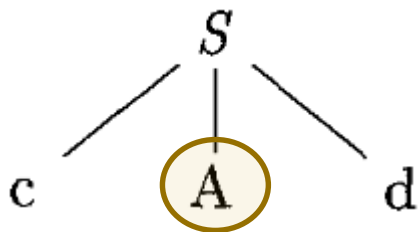
■ 文法:

○ $S \rightarrow cAd$

○ $A \rightarrow ab \mid a$

- 1)
- 2)
- 3)
- 4)
- 5)
- 6)
- 7)

```
void A() {  
    选择一个  $A$  产生式,  $A \rightarrow X_1X_2 \cdots X_k$ ;  
    for (  $i = 1$  to  $k$  ) {  
        if (  $X_i$  是一个非终结符号 )  
            调用过程  $X_i()$ ;  
        else if (  $X_i$  等于当前的输入符号  $a$  )  
            读入下一个输入符号;  
        else /* 发生了一个错误 */;  
    }  
}
```





递归下降分析过程示例



■ 输入串 $w = c a d$

■ 文法:

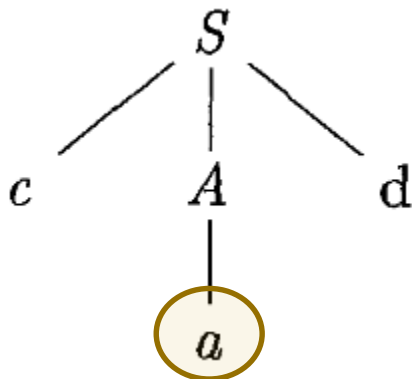
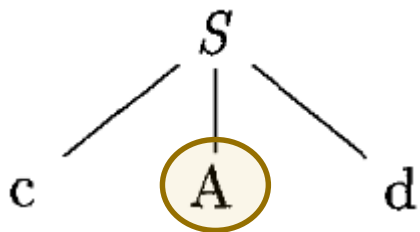
○ $S \rightarrow cAd$

○ $A \rightarrow ab \mid a$



- 1)
- 2)
- 3)
- 4)
- 5)
- 6)
- 7)

```
void A() {  
    选择一个 A 产生式,  $A \rightarrow X_1X_2 \cdots X_k$ ;  
    for (  $i = 1$  to  $k$  ) {  
        if (  $X_i$  是一个非终结符号 )  
            调用过程  $X_i()$ ;  
        else if (  $X_i$  等于当前的输入符号  $a$  )  
            读入下一个输入符号;  
        else /* 发生了一个错误 */;  
    }  
}
```





递归下降分析过程示例



■ 输入串 $w = c a d$

■ 文法:

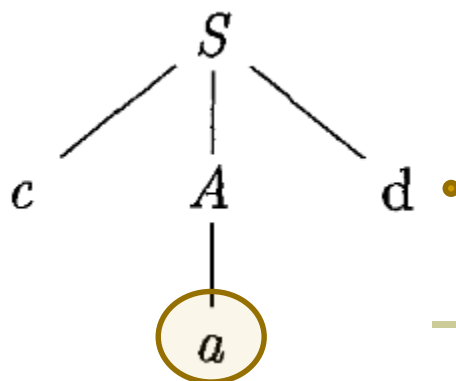
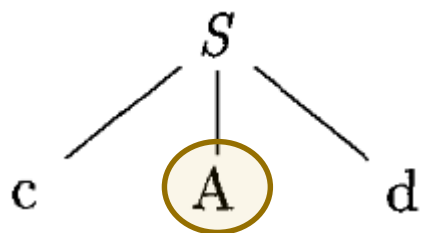
○ $S \rightarrow cAd$

○ $A \rightarrow ab \mid a$



- 1)
- 2)
- 3)
- 4)
- 5)
- 6)
- 7)

```
void A() {  
    选择一个 A 产生式,  $A \rightarrow X_1X_2 \cdots X_k$ ;  
    for (  $i = 1$  to  $k$  ) {  
        if (  $X_i$  是一个非终结符号 )  
            调用过程  $X_i()$ ;  
        else if (  $X_i$  等于当前的输入符号  $a$  )  
            读入下一个输入符号;  
        else /* 发生了一个错误 */;  
    }  
}
```



匹配成功!



递归下降过程示例(续)



- 可能需要回溯，或者说，可能需要重复扫描输入
- “在回到A时，我们必须把输入指针重新设置到位置2，即我们第一次尝试展开A时该指针指向的位置。这意味着A的过程必须将输入指针存放在一个局部变量中。”
- 如果文法中存在左递归.....
- 回溯（穷试 ☺）显然是笨办法



预测分析法简介



- 试图从开始符号推导出输入符号串
- 以开始符号作为初始的当前句型
- 每次为最左边的非终结符号选择适当的产生式
 - 通过查看下一个输入符号来选择这个产生式
 - 有多个可能的产生式时预测分析法无能为力
- 比如文法: $E \rightarrow +EE \mid -EE \mid id$; 输入为 **+ id - id id**
- 当两个产生式具有相同的前缀时无法预测
 - 文法: $stmt \rightarrow \text{if } expr \text{ then } stmt \text{ else } stmt \mid \text{if } expr \text{ then } stmt$
 - 输入: **if a then ...**
 - 新文法: $stmt \rightarrow \text{if } expr \text{ then } stmt \text{ elsePart}$
 - $elsePart \rightarrow \text{else } stmt \mid \varepsilon$
- 需要提取公因子



预测分析技术

消除二义性
消除左递归
提取左公因子



- 一种确定性的、无回溯的分析技术
- 在每一步选择正确的产生式

□ 例1：文法**G(S)**:

$S \rightarrow pA$

$S \rightarrow qB$

$A \rightarrow cAd$

$A \rightarrow a$

□ 输入串 **$W = pccadd$**



预测分析技术（续）



□ 例2：文法**G(S)**为：

$$S \rightarrow Ap$$

$$S \rightarrow Bq$$

$$A \rightarrow a \mid cA$$

$$B \rightarrow b \mid dB$$

□ 输入**W=ccap**



预测分析技术（续）



- 通过在输入中向前看固定多个符号来选择正确的产生式
- 通常情况下，我们只需要向前看一个符号
 - 递归下降分析一般只适合于每个子表达式的第一个终结符能够为产生式选择提供足够信息的那些文法
- 给出两个与文法相关的两个函数
 - FIRST
 - FOLLOW
- 基于上述两个函数，可以根据下一个输入符号来选择应用哪个产生式



FIRST和FOLLOW (1)

很重要



- 在自顶向下的分析技术中，通常使用向前看几个符号来唯一地确定产生式（这里假定只看一个符号）
- 当前句型是 $xA\beta$ ，而输入是 $xa...$ 。那么选择产生式 $A \rightarrow \alpha$ 的必要条件是下列之一：
 - $A \xRightarrow{*} \varepsilon$ ，且 β 以 a 开头；也就是说在某个句型中 a 跟在 A 之后；
 - $A \xRightarrow{*} a...$ ；
 - 如果按照这两个条件选择时能够保证唯一性，那么我们就可以避免回溯。
- 因此，我们定义FIRST和FOLLOW



FIRST和FOLLOW (2)



■ FIRST(α):

- 可以从 α 推导得到的串的首符号的集合;
- 如果 $\alpha \xRightarrow{*} \varepsilon$, 那么 ε 也在 FIRST(α)中;

■ FOLLOW(A):

- 可能在某些句型中紧跟在A右边的终结符号的集合。

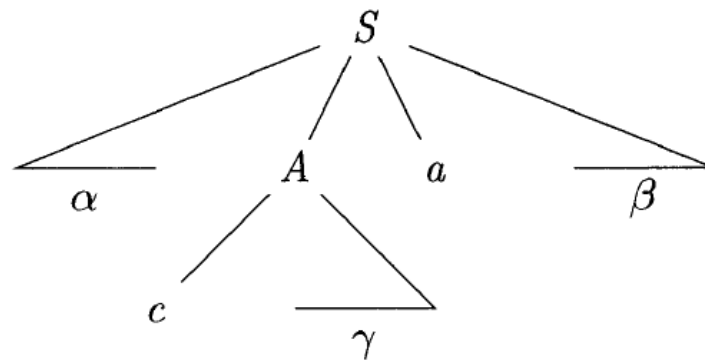


图 4-15 终结符号 c 在 FIRST(A) 中且 a 在 FOLLOW(A) 中



FIRST (α)



- 定义：可从 α 推导得到的串的首符号的集合，其中 α 是任意的文法符号串。如果 $\alpha \in \epsilon$ ，那么 ϵ 也在FIRST(α)中
- FIRST函数的意义
 - 如果两个 A 产生式 $A \rightarrow \alpha \mid \beta$ ，其中First(α)和First(β)是不相交的集合。下一个输入符号是 a ，若 $a \in \text{First}(\alpha)$ ，则选择 $A \rightarrow \alpha$ ，若 $a \in \text{First}(\beta)$ ，则选择 $A \rightarrow \beta$
 - 用于预测产生式的选择



First的计算



- 对于文法符号 X 的 $\text{FIRST}(X)$ ，通过不断应用下列规则，直到没有新的终结符号或者 ϵ 可以加入到任何 FIRST 集合为止
 - 如果 X 是终结符号，那么 $\text{FIRST}(X)=\{X\}$
 - 如果 X 是非终结符号，且有规则 $X \rightarrow a\dots$ ，那么将 a 添加到 $\text{FIRST}(X)$ 中；如果 $X \rightarrow \epsilon$ ，那么 ϵ 也在 $\text{FIRST}(X)$ 中
 - 对于规则 $X \rightarrow Y_1Y_2\dots Y_n$ ，把 $\text{FIRST}(Y_1)$ 中的非 ϵ 符号添加到 $\text{FIRST}(X)$ 中。如果 ϵ 在 $\text{FIRST}(Y_1)$ 中，把 $\text{FIRST}(Y_2)$ 中的非 ϵ 符号添加到 $\text{FIRST}(X)$ 中...；如果 ϵ 在 $\text{FIRST}(Y_n)$ 中，把 ϵ 添加到 $\text{FIRST}(X)$ 中



First的计算（续）



- 对于文法符号串 $X_1X_2...X_n$ 的First集合
 - 向 $\text{First}(X_1X_2...X_n)$ 加入 $\text{First}(X_1)$ 中所有的非 ϵ 符号。
 - 如果 ϵ 在 $\text{First}(X_1)$ 中，再加入 $\text{First}(X_2)$ 中的所有非 ϵ 符号。如果 ϵ 在 $\text{First}(X_1)$ 和 $\text{First}(X_2)$ 中，再加入 $\text{First}(X_3)$ 中的所有非 ϵ 符号。依次类推。
 - 如果对所有的 i （1到 n ）， ϵ 都在 $\text{First}(X_i)$ 中，则 ϵ 加入 $\text{First}(X_1X_2...X_n)$ 中。



First的计算示例



$$E \rightarrow T E'$$

$$E' \rightarrow + T E' \mid \epsilon$$

$$T \rightarrow F T'$$

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

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

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

- $\text{First}(E') = \{ +, \epsilon \}$

- $\text{First}(T') = \{ *, \epsilon \}$

- $\text{First}(TE') = \text{First}(T) = \{ (, \text{id} \}$

- $\text{First}(+TE') = \{ + \}$

- $\dots\dots\dots$



一个实例



■ 文法G[S], 输入bcd

- $S \rightarrow AB \mid CD$
- $A \rightarrow aD \mid \varepsilon$
- $C \rightarrow cD$
- $B \rightarrow bC$
- $D \rightarrow d$
- $\text{First}(S) = \{a, b, c\}$
- $\text{First}(A) = \{a, \varepsilon\}$
- $\text{First}(B) = \{b\}$
- $\text{First}(C) = \{c\}$
- $\text{First}(D) = \{d\}$



FOLLOW函数

文法G[S], 输入bcd

$S \rightarrow AB|CD$

$A \rightarrow aD| \epsilon$

$C \rightarrow cD$

$B \rightarrow bC$

$D \rightarrow d$



- 对于非终结符号A, FOLLOW(A)定义为可能在某些句型中紧跟在A右边的终结符号的集合
 - 例如 $S \xRightarrow{*} \alpha A a \beta$, 终结符号 $a \in \text{Follow}(A)$
- 如果A是某些句型的最右符号, 那么 $\$ \in \text{Follow}(A)$ 。\$是特殊的输入串“结束标记”
- FOLLOW函数的意义:
 - 如果 $A \rightarrow \alpha$, 当 $\alpha \rightarrow \epsilon$ 或 $\alpha \xRightarrow{*} \epsilon$ 时, FOLLOW(A)可以帮助我们做出选择恰当的产生式
 - 例如: 如果 $A \rightarrow \alpha$, b属于FOLLOW(A), 如果 $\alpha \xRightarrow{*} \epsilon$, 则若当前输入符号是b, 可以选择 $A \rightarrow \alpha$, 因为A最终到达了 ϵ , 而且后面跟着b



FOLLOW计算



- 计算各个非终结符号 A 的 $FOLLOW(A)$ 集合，不断应用下列规则，直到没有新的终结符号可以被加入到任意 $FOLLOW$ 集合中
 - 将 $\$$ 放入 $FOLLOW(S)$ ， S 是开始符号，而 $\$$ 是输入串的结束标记
 - 如果存在产生式 $A \rightarrow \alpha B \beta$ ，那么 $First(\beta)$ 中除 ϵ 之外的所有符号都在 $Follow(B)$ 中
 - 如果存在一个产生式 $A \rightarrow \alpha B$ ，或存在产生式 $A \rightarrow \alpha B \beta$ 且 $First(\beta)$ 包含 ϵ ，那么 $Follow(A)$ 中的所有符号都在 $Follow(B)$ 中



Follow的计算示例



□ 文法

$$E \rightarrow T E'$$

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

$$T \rightarrow F T'$$

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

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

$$\square \text{ FOLLOW}(E)=\{ \} , \$\}$$

$$\leftarrow \text{FIRST}()) \cup \{ \$ \}$$

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

$$\leftarrow \text{FOLLOW}(E)$$

$$\square \text{ FOLLOW}(T)=\{ + ,) , \$ \}$$

$$\leftarrow \text{FIRST}(E') \cup \text{FOLLOW}(E')$$

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

$$\leftarrow \text{FOLLOW}(T)$$

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

$$\leftarrow \text{FIRST}(T') \cup \text{FOLLOW}(T)$$

- 将\$放入FOLLOW(S), S是开始符号, 而\$是输入串的结束标记
- 如果存在产生式 $A \rightarrow \alpha B \beta$, 那么First(β)中除 ε 之外的所有符号都在Follow(B)中
- 如果存在一个产生式 $A \rightarrow \alpha B$, 或存在产生式 $A \rightarrow \alpha B \beta$ 且First(β)包含 ε , 那么Follow(A)中的所有符号都在Follow(B)中



课堂练习



- 有如下文法:

- $E \rightarrow id\ X$
- $X \rightarrow \epsilon \mid (A) \mid [E]$
- $A \rightarrow EY$
- $Y \rightarrow \epsilon \mid ;A$

- 给出各个非终结符号的First集和Follow集合。



思考： 为什么



- 消除二义性

- 消除左递归

改造文法

- 提取左公因子

- 求First和Follow集合



LL(1) 文法 (1)



- LL(1)文法，第一个L表示自左向右扫描，第二个L指产生最左推导，而1则表示向前看一个输入符号
- LL(1)文法可以利用不回溯的确定性的预测分析技术，因为只需要检查当前输入符号就可以为一个非终结符号选择正确的产生式



LL(1) 文法 (2)



- 定义:对于文法的任意两个不同的产生式 $A \rightarrow \alpha | \beta$
 - 不存在终结符号 a 使得 α 和 β 都可以推导出以 a 开头的串
 - α 和 β 最多只有一个可以推导出空串
 - 如果 β 可以推导出空串, 那么 α 不能推导出以 FOLLOW 中任何终结符号开头的串;
- 等价于:
 - $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \Phi$; (条件一、二)
 - 如果 $\epsilon \in \text{FIRST}(\beta)$, 那么 $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \Phi$; 反之亦然。(条件三)



LL(1) 文法 (3)



- 思考：对于一个LL(1)文法，我们期望从A推导出以终结符号a开头的符号串
 - 我们如何选择产生式？
 - 有多个选择吗？



LL(1) 文法 (4)



- 对于LL(1)文法，可以在自顶向下分析过程中，根据当前输入符号来确定使用的产生式
- 例如：
 - 产生式: $stmt \rightarrow \mathbf{if (exp) stmt \ else stmt} \mid \mathbf{while (exp) stmt} \mid \mathbf{a}$
 - 输入: $\mathbf{if (exp) while (exp) a \ else a}$
 - 首先选择产生式 $stmt \rightarrow \mathbf{if (exp) stmt \ else stmt}$ 得到句型 $\mathbf{if (exp) stmt \ else stmt}$
 - 然后发现要把句型中的第一个stmt展开，选择产生式 $stmt \rightarrow \mathbf{while (exp) stmt}$ ，得到句型 $\mathbf{if (exp) while (exp) stmt \ else stmt}$ 。
 - 然后再展开第一个stmt，得到 $\mathbf{if (exp) while (exp) a \ else stmt}$
 - ...



LL(1) 文法的预测分析技术



- LL(1)文法
- 将**first**和**follow**集合中的信息放入一个预测分析表 $M[A, a]$ ，该预测表告诉我们当非终结符号为 A ，当前输入符号为 a 时，要选择哪条产生式
- 预测分析表的构造
 - 若当前输入符号 a 在 $\text{First}(\alpha)$ 中，选择产生式 $A \rightarrow \alpha$
 - 若 $\alpha \xRightarrow{*} \varepsilon$ ，如果 a 在 $\text{Follow}(A)$ 中，选择产生式 $A \rightarrow \alpha$
 - 如果当前符号 a 是 $\$$ ，且 $\$$ 在 $\text{Follow}(A)$ 中，则选择产生式 $A \rightarrow \alpha$



预测分析表构造



- 输入：文法 G
- 输出：预测分析表 M
- 方法：对于文法 G 的每个产生式 $A \rightarrow \alpha$ ，进行如下处理
 - 对于 $\text{First}(\alpha)$ 中的每个终结符号 a ，将 $A \rightarrow \alpha$ 加入到 $M[A, a]$
 - 如果 ϵ 在 $\text{First}(\alpha)$ 中，那么对于 $\text{Follow}(A)$ 中的每个终结符号 b ，将 $A \rightarrow \alpha$ 加入到 $M[A, b]$ 中
 - 如果 ϵ 在 $\text{First}(\alpha)$ 中，且 $\$$ 在 $\text{Follow}(A)$ 中，将 $A \rightarrow \alpha$ 加入到 $M[A, \$]$ 中
- 完成上述操作后，若 $M[A, a]$ 中没有产生式，填为error



预测分析表的例子



■ 文法:

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

这个例子恰巧使得每个产生式的右部的第一个符号的FIRST集合就等于产生式右部的FIRST集合。但是在一般情况下不总是这样的。

■ FIRST集:

- $F: \{ (, \text{id} \}; \quad E, T: \{ (, \text{id} \}; \quad E': \{ +, \epsilon \}; \quad T': \{ *, \epsilon \}$

■ FOLLOW集:

- $E: \{ \$,) \}; \quad E': \{ \$,) \}; \quad T, T': \{ +,), \$ \}; \quad F: \{ +, *,), \$ \}$

非终结符号	输入符号					
	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	$F \rightarrow \text{id}$			$F \rightarrow (E)$		



预测分析方法



- 对于任何文法**G**，都可以构造预测分析表
- 对于**LL(1)**文法，预测表中每个条目都唯一地指定了一个产生式，或者标明**error**



二义性文法的预测分析表



- 对于某些文法，表中可能会有有一些多重定义的条目，比如左递归或二义性文法。

$S \rightarrow iEtSS' \mid a$ -FIRST(eS)= $\{e\}$ ，使得 $S' \rightarrow eS$ 在 $M[S',e]$ 中；
 $S' \rightarrow eS \mid \epsilon$ -FOLLOW(S')= $\{\$,e\}$ 使得 $S' \rightarrow \epsilon$ 也在 $M[S',e]$ 中
 $E \rightarrow b$

非终结符号	输入符号					
	a	b	e	i	t	$\$$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				



文法和自顶向下分析



- 设计文法(描述能力, 适用于自顶向下分析)
 - 消除二义性
 - 消除左递归
 - 提取左公因子
- 自顶向下分析
 - 递归下降分析 (通用的可能需要回溯)
 - 预测分析
- LL(1)文法及其预测分析



非递归的预测分析（1）



- 在自顶向下分析的过程中，我们总是
 - 匹配掉句型中左边的所有终结符号
 - 对于最左边的非终结符号，我们选择适当的产生式展开
 - 匹配成功的终结符号不会再被考虑，因此我们只需要记住句型的余下部分，以及尚未匹配的输入终结符号串
 - 由于展开的动作总是发生在余下部分的左端，我们可以用栈来存放这些符号



非递归的预测分析（2）



- 分析处理过程
 - 初始化时，栈中仅包含开始符号
 - 如果栈顶元素是终结符号，那么进行匹配
 - 如果栈顶元素是非终结符号
 - 使用预测分析表来**选择**适当的产生式
 - 在栈顶用产生式右部**替换**产生式左部
- 对所有文法的预测分析都可以共用同样的驱动程序



非递归的预测分析技术

- 文法符号栈
- 输入缓冲区
- 控制程序

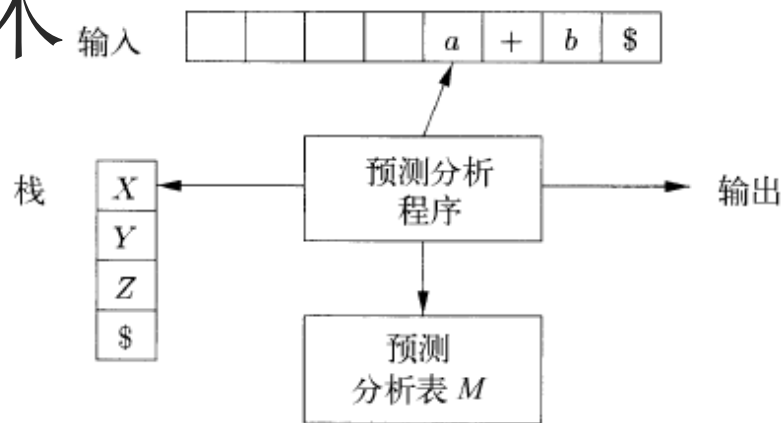


图 4-19 一个分析表驱动的分析器的模型

- 在任何时刻，栈顶符号 X 与当前输入符号 a 决定了控制程序所应执行的分析动作。四种可能的动作
 - 如果 $X=a=\$$ ，则分析成功结束
 - 如果 $X=a\neq \$$ ，则从栈中退去 X ，并把输入指针推进到指向下一个输入符号
 - 如果 X 是一个非终结符号，且分析表 A 的元素 $A[X][a]=“X \rightarrow X_1X_2\dots X_k”$ ，则把栈顶的 X 替换为 $X_kX_{k-1}\dots X_1$ （反向下推入栈，使得 X_1 在栈顶）
 - 除上述情况外的其它情况，调用出错程序



分析表驱动预测分析器



- 性质1: 如果栈中的符号序列为 α , 而 w 是已经被读入的部分输入, w' 是尚未处理的输入; 那么
 - S推导出 $w\alpha$
 - 我们试图从 α 推导出余下的输入终结符号串 w'
- 预测分析程序使用 $M[X,a]$ 来扩展 X , 将此产生式的右部按倒序压入栈中
- 请注意: 这样的操作使得分析过程中性质1得到保持。

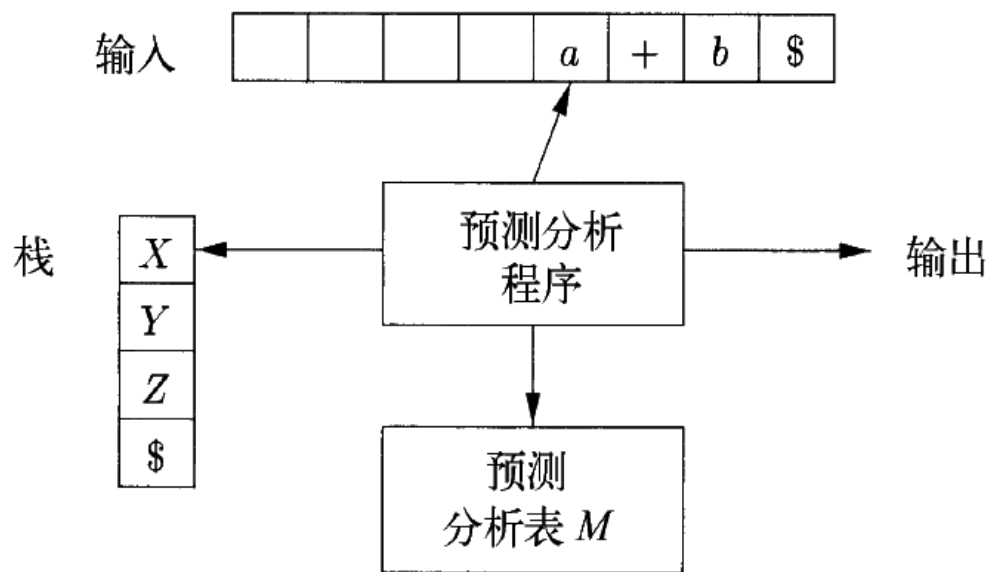


图 4-19 一个分析表驱动预测分析器的模型



表驱动的非递归的预测语法分析



- 输入：一个串 w ，文法 G 的预测分析表 M 。
- 输出：如果 w 在 $L(G)$ 中，输出 w 的一个最左推导；否则报错。
- 方法：初始化输入缓冲区为 $w\$$ ，栈顶是 G 的开始符号 S ，下面是 $\$$ 。

```
设置  $ip$  使它指向  $w$  的第一个符号，其中  $ip$  是输入指针；
令  $X$  = 栈顶符号；
while (  $X \neq \$$  ) { /* 栈非空 */
    if (  $X$  等于  $ip$  所指向的符号  $a$  ) { 弹出栈的弹出操作，将  $ip$  向前移动一个位置；
    else if (  $X$  是一个终结符号 )  $error()$ ；
    else if (  $M[X, a]$  是一个报错条目 )  $error()$ ；
    else if (  $M[X, a] = X \rightarrow Y_1 Y_2 \cdots Y_k$  ) {
        输出产生式  $X \rightarrow Y_1 Y_2 \cdots Y_k$ ；
        弹出栈顶符号；
        将  $Y_k, Y_{k-1}, \dots, Y_1$  压入栈中，其中  $Y_1$  位于栈顶。
    }
    令  $X$  = 栈顶符号；
}
```




分析表驱动预测分析的例子



$$E \rightarrow T E'$$

$$E' \rightarrow + T E' \mid \epsilon$$

$$T \rightarrow F T'$$

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

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

非终结符号	输入符号					
	id	+	*	()	\$
E	$E \rightarrow T E'$			$E \rightarrow T E'$		
E'		$E' \rightarrow + T E'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow F T'$			$T \rightarrow F T'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow * F T'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

已匹配	栈	输入	动作
	$E\$$	$\text{id} + \text{id} * \text{id}\$$	
	$T E'\$$	$\text{id} + \text{id} * \text{id}\$$	输出 $E \rightarrow T E'$
	$F T' E'\$$	$\text{id} + \text{id} * \text{id}\$$	输出 $T \rightarrow F T'$
	$\text{id } T' E'\$$	$\text{id} + \text{id} * \text{id}\$$	输出 $F \rightarrow \text{id}$
id	$T' E'\$$	$+ \text{id} * \text{id}\$$	匹配 id
id	$E'\$$	$+ \text{id} * \text{id}\$$	输出 $T' \rightarrow \epsilon$
id	$+ T E'\$$	$+ \text{id} * \text{id}\$$	输出 $E' \rightarrow + T E'$
$\text{id} +$	$T E'\$$	$\text{id} * \text{id}\$$	匹配 $+$

已经匹配部分加上栈中符号必然是一个**最左**句型

- 文法4.28, 输入: $\text{id} + \text{id} * \text{id}$;



自底向上语法分析



- 概述
- 移进-规约技术
- LR语法分析技术
 - 简单LR技术
 - LR技术



自底向上的语法分析



- 为一个输入串构造语法分析树的过程;
- 从叶子（输入串中的终结符号，将位于分析树的底端）开始，向上到达根结点
- 重要的自底向上语法分析的通用框架
 - 移入-归约
- LR: 最大的可以构造出移入-归约语法分析器的语法类



自底向上语法分析概述 - 归约



■ Bottom-Up Parsing

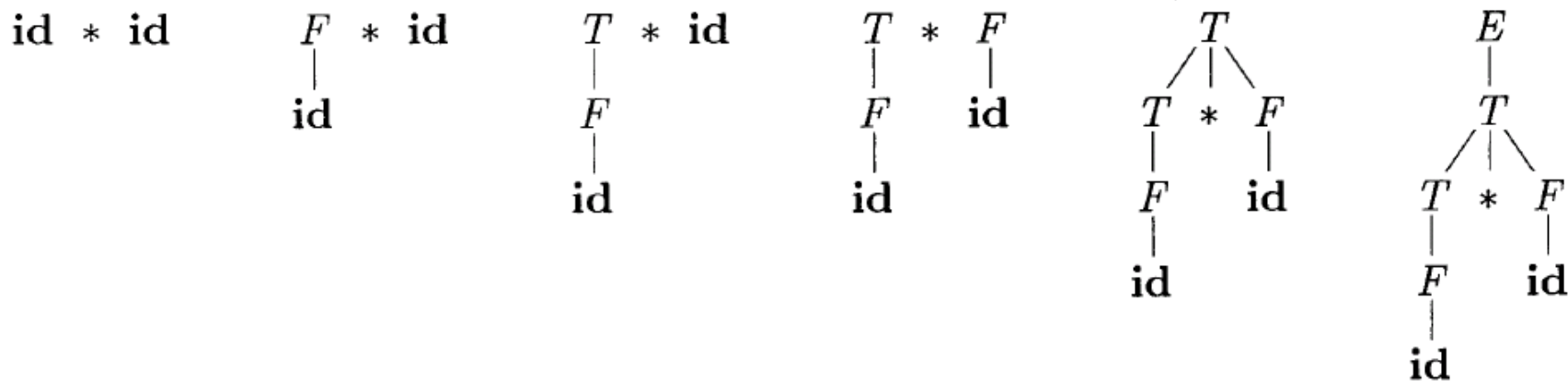


图 4-25 $\text{id} * \text{id}$ 的自底向上分析过程



自底向上语法分析概述 - 归约



- **Bottom-up parsing** – 将一个串 w 归约为文法符号的过程
- 在每一步的归约中，一个与某产生式体相匹配的特定子串被替换为该产生式头部的非终结符号，一次归约实质上是一个推导的反向操作



归约的例子



- $id * id$ 的归约过程
 - $id * id$, $F * id$, $T * id$, $T * F$, T , E
- 对于句型 $T * id$, 有两个子串和某产生式右部匹配
 - T 是 $E \rightarrow T$ 的右部
 - id 是 $F \rightarrow id$ 的右部
 - 为什么选择将 id 归约为 F , 而不是将 T 归约为 E ?
 - T 归约为 E 之后, $E * id$ 不再是句型

$F \rightarrow id$

$T \rightarrow F$

$F \rightarrow id$

$E \rightarrow T * F$

$id * id$

$F * id$
|
 id

$T * id$
|
 F
|
 id

$T * F$
| |
 F id
|
 id

T
/ | \
 T $*$ F
| |
 F id
|
 id

E
|
 T
/ | \
 T $*$ F
| |
 F id
|
 id



自底向上分析- 规约



- 自底向上分析的过程也是规约的过程
- 规约的问题:
 - 选择哪部分进行归约?
 - 应用哪个产生式进行归约?



句型/句子/语言 (回顾)



- 句型 (sentential form) :
 - 如果 $S \xRightarrow{*} \alpha$, 那么 α 就是文法的句型
 - 可能既包含非终结符号, 又包含终结符号; 可以是空串
- 句子 (sentence)
 - 文法的句子就是不包含非终结符号的句型
- 语言
 - 文法 G 的语言就是 G 的句子的集合, 记为 $L(G)$
 - w 在 $L(G)$ 中当且仅当 w 是 G 的句子, 即 $S \xRightarrow{*} w$



句柄

重要



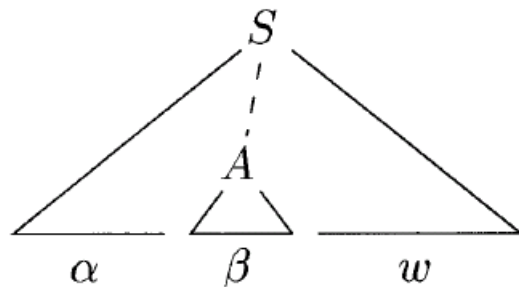
- 最右句型 γ 的一个句柄
 - 满足下述条件的产生式 $A \rightarrow \beta$ 及串 β 在 γ 中出现的位置
 - 条件：将这个位置上的 β 替换为 A 之后得到的串是 γ 的某个最右推导序列中出现在位于 γ 之前的最右句型
- 通俗地说
 - 句柄是最右推导的反向过程中被规约的那些部分
- 句柄的作用
 - 对句柄的规约，代表了相应的最右推导中的一个反向步骤



自底向上句法分析概述: 句柄剪枝



- Bottom-Up Parsing – 归约过程 – 反向最右推导过程 – 句柄剪枝过程
- 句柄: 和某个产生式体相匹配的子串, 对它的归约代表了相应的最右推导的一个反向步骤
- 句柄的形式定义: $S \xRightarrow{*}_{rm} \alpha Aw \xRightarrow{rm} \alpha \beta w$, 那么紧跟 α 的 β 可以一步归约到 A 。 w 是所有的终结符号串。



- 注意:
 - 归约前后都是最右句型
 - 和某个产生式匹配的最左子串不一定是句柄
 - 无二义性的文法, 其每个最右句型都有且只有一个句柄



自底向上语法分析概述: 句柄剪枝 (续)

- 通过“句柄剪枝”可以得到一个反向的最右推导。从句子 w 开始, 令 $w = \gamma_n$, γ_n 是未知最右推导的第 n 个最右句型

$$S = \gamma_0 \underset{rm}{\Rightarrow} \gamma_1 \underset{rm}{\Rightarrow} \gamma_2 \underset{rm}{\Rightarrow} \cdots \underset{rm}{\Rightarrow} \gamma_{n-1} \underset{rm}{\Rightarrow} \gamma_n = w$$

- 以相反的顺序重构这个推导（实际上是重构归约序列），我们在 γ_n 中寻找句柄 β_n ，并将替换为相应产生式 $A \rightarrow \beta_n$ 的头部，得到前一个最右句型 γ_{n-1} ；重复这个过程，直到 S
- 和推导类似，如果能重现这个序列，则可以完成语法分析任务（其中可能存在什么问题）
 - 如何找到句柄



句柄的例子



■ 输入: $id * id$

最右句型	句柄	归约用的产生式
$id_1 * id_2$	id_1	$F \rightarrow id$
$F * id_2$	F	$T \rightarrow F$
$T * id_2$	id_2	$F \rightarrow id$
$T * F$	$T * F$	$T \rightarrow T * F$

$id * id$

$\begin{array}{c} F \\ | \\ id \end{array} * id$

$\begin{array}{c} T \\ | \\ F \\ | \\ id \end{array} * id$

$\begin{array}{c} T \\ | \\ F \\ | \\ id \end{array} * \begin{array}{c} F \\ | \\ id \end{array}$

$\begin{array}{c} T \\ / \quad | \quad \backslash \\ T \quad * \quad F \\ | \quad \quad | \\ F \quad \quad id \\ | \\ id \end{array}$

$\begin{array}{c} E \\ | \\ T \\ / \quad | \quad \backslash \\ T \quad * \quad F \\ | \quad \quad | \\ F \quad \quad id \\ | \\ id \end{array}$



句柄的例子



- 已知文法

$E \rightarrow E + T \mid E - T \mid T$

$T \rightarrow T * F \mid T / F \mid F$

$F \rightarrow (E) \mid id$

- 句型 $T / (F - id)$ 的句柄是什么？

- 句柄是**最右推导**的**反向**过程中**被规约**的那些部分



移进归约语法分析框架



- 一种自底向上的分析形式
- 使用一个栈保存文法符号，一个输入缓冲区存放将要进行分析的剩余符号
- 初始栈：\$ 初始输入 $w\$$
- 对输入串的一次从左到右的扫描过程中，语法分析器将零个或多个输入符号移到栈的顶端，直到它可以**对栈顶的一个文法符号串**进行**归约**为止。它将归约为某个产生式的头。不断重复这个循环，直到它检测到一个语法错误，或者栈中包含了开始符号且输入缓冲区为空。
- 分析成功时： 栈 $\$ S$ ， 输入 $\$$

自顶向下的非递归预测分析中栈是如何工作的？



移入-归约分析技术



- 使用一个栈来保存归约/扫描移入的文法符号
- 栈中符号（从底向上）和待扫描的符号组成了一个最右句型
- 开始时刻：栈中只包含 $\$$ ，而输入为 $w\$$
- 成功结束时刻：栈中 $\$S$ ，而输入 $\$$
- 在分析过程中，不断地移入符号，并在识别到句型时进行归约



主要分析动作



- 移入：将下一个输入符号移动到栈顶
- 归约：将句柄归约为相应的非终结符号
 - 句柄总是在栈顶
 - 具体操作时弹出句柄，压入被归约到的非终结符号
- 接受：宣布分析过程成功完成
- 报错：发现语法错误，调用错误恢复子程序



归约分析过程的例子



栈	输入	动作
\$	$\text{id}_1 * \text{id}_2 \$$	移入
$\$ \text{id}_1$	$* \text{id}_2 \$$	按照 $F \rightarrow \text{id}$ 归约
$\$ F$	$* \text{id}_2 \$$	按照 $T \rightarrow F$ 归约
$\$ T$	$* \text{id}_2 \$$	移入
$\$ T *$	$\text{id}_2 \$$	移入
$\$ T * \text{id}_2$	$\$$	按照 $F \rightarrow \text{id}$ 归约
$\$ T * F$	$\$$	按照 $T \rightarrow T * F$ 归约
$\$ T$	$\$$	按照 $E \rightarrow T$ 归约
$\$ E$	$\$$	accept

图 4-28 一个移入 - 归约语法分析器
在处理输入 $\text{id}_1 * \text{id}_2$ 时经历的格局



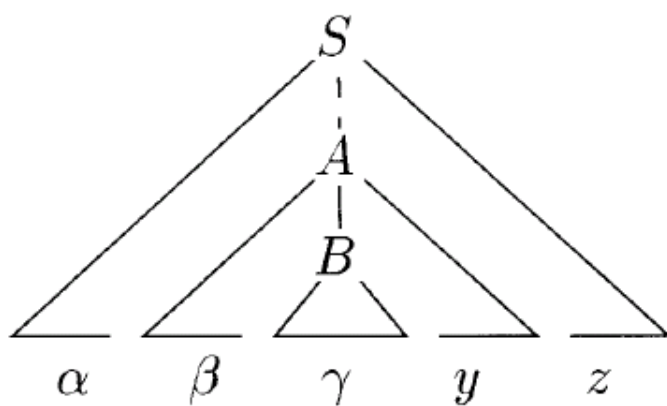
移进归约语法分析框架（续）



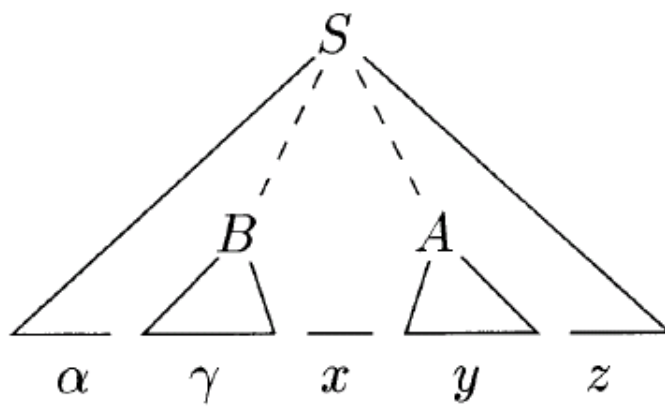
- 可行性分析：句柄总是出现在栈的顶端，绝不会出现在栈的中间

$$1) S \xRightarrow{*}_{rm} \alpha Az \xRightarrow{rm} \alpha \beta B y z \xRightarrow{rm} \alpha \beta \gamma y z$$

$$2) S \xRightarrow{*}_{rm} \alpha B x A z \xRightarrow{rm} \alpha B x y z \xRightarrow{rm} \alpha \gamma x y z$$



情况 (1)

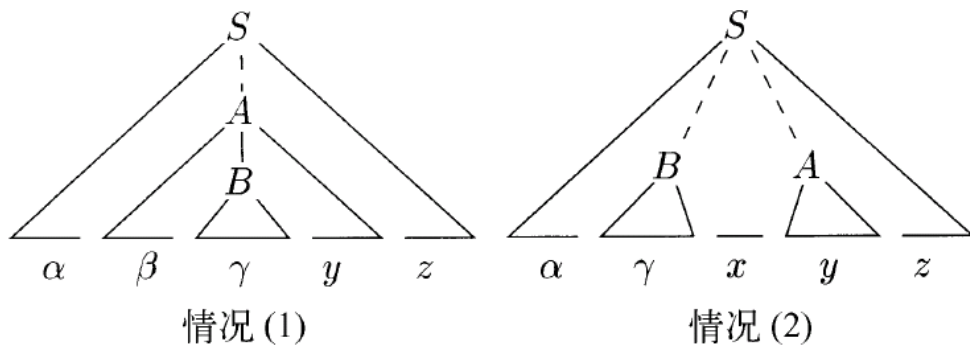


情况 (2)

图 4-29 一个最右推导中两个连续步骤的两种情况



移进归约语法分析框架（续）



$$1) S \xRightarrow{rm}^* \alpha Az \Rightarrow \alpha \beta B y z \Rightarrow \alpha \beta \gamma y z$$

$$2) S \xRightarrow{rm}^* \alpha B x A z \Rightarrow \alpha B x y z \Rightarrow \alpha \gamma x y z$$

图 4-29 一个最右推导中两个连续步骤的两种情况

STACK	INPUT	STACK	INPUT
$\$ \alpha \beta \gamma$	$y z \$$	$\$ \alpha \gamma$	$x y z \$$
$\$ \alpha \beta B$	$y z \$$	$\$ \alpha B x y$	$z \$$
$\$ \alpha \beta B y$	$z \$$		

- 语法分析器进行一次归约以后，都必须接着移入零个或多个符号才能在栈顶找到下一个句柄
- 不需要去栈中间去寻找句柄



移进/归约冲突



- 移入-归约技术并不能处理所有上下文无关文法
- 某些上下文无关文法（比如二义性文法）
 - 移入/归约冲突：栈中的内容和接下来的 k 个输入符号，都不能确定进行移入还是归约操作(不能确定是否是句柄)
 - 归约/归约冲突：存在多个可能的归约到不同非终结符号的归约(不能确定句柄归约到那个非终结符号)



移进/归约冲突



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

栈

... **if** $expr$ **then** $stmt$

输入

else ... \$

规约为**stmt**?
还是移入**else**?



归约/归约冲突



- | | | | |
|-----|-----------------------|---------------|--|
| (1) | <i>stmt</i> | \rightarrow | id (<i>parameter_list</i>) |
| (2) | <i>stmt</i> | \rightarrow | <i>expr</i> := <i>expr</i> |
| (3) | <i>parameter_list</i> | \rightarrow | <i>parameter_list</i> , <i>parameter</i> |
| (4) | <i>parameter_list</i> | \rightarrow | <i>parameter</i> |
| (5) | <i>parameter</i> | \rightarrow | id |
| (6) | <i>expr</i> | \rightarrow | id (<i>expr_list</i>) |
| (7) | <i>expr</i> | \rightarrow | id |
| (8) | <i>expr_list</i> | \rightarrow | <i>expr_list</i> , <i>expr</i> |
| (9) | <i>expr_list</i> | \rightarrow | <i>expr</i> |

?

图 4-30 有关过程调用和数组引用的产生式

栈

... **id** (**id**

输入

, **id**) ...



解决冲突问题



- 让词法分析区分id和procid
- 用[]表示数组，用()表示函数
-
- 但是并不能完全解决冲突问题
- 有一类文法可以解决句柄查找及归约唯一性的问题



自底向上分析-回顾



- 从输入串出发构造一个反向最右推导序列（归约）
- 每一步是对**句柄**进行归约 -- 寻找句柄
- 一种移进-归约的分析框架
 - 一个栈存放文法符号，一个输入缓冲区
 - 移进 or 归约
 - 总是能够在栈的顶端找到句柄
 - 通过证明是可行的
- 有时候会有冲突，是移进还是归约？ 怎么归约？



LR语法分析技术



- LR(k) 的语法分析概念
 - L表示最左扫描，R表示反向构造出最右推导
 - k表示最多向前看k个符号
- 当k的数量增大时，相应的语法分析器的规模急剧增大
 - $K=2$ 时，程序设计语言的语法分析器的规模通常非常庞大
 - 当 $k=0$ 、 1 时已经可以解决很多语法分析问题，因此具有实践意义
 - 因此，我们只考虑 $k \leq 1$ 的情况



LR语法分析器的优点



- 表格驱动
 - 虽然手工构造表格工作量大，但表格可以自动生成
- 对于几乎所有的程序设计语言，只要写出上下文无关文法，就能够构造出识别该构造的**LR**语法分析器
- 最通用的无回溯移入-归约分析技术，且和其它技术一样高效
- 可以尽早检测到错误
- 能分析的文法集合是**LL(k)**文法的超集



LR分析相关概念



- 移入-归约语法分析器如何知道何时该移入、何时该归约呢？
 - LR语法分析器试图用一些**状态**来表明我们在移进归约语法分析过程中所处的位置，从而做出移入-归约决定
- 项的意义
 - 指明在语法分析过程中的给定点上，我们已经看到了一个产生式的哪些部分。或者说，如果我们想用这个产生式进行归约，还需要看到哪些文法符号
- 项的集合（项集）对应于一个状态



LR(0) 项



- 文法的一个产生式加上在其产生式体中某处的一个点
 - $A \rightarrow .XYZ$, $A \rightarrow X.YZ$, $A \rightarrow XY.Z$, $A \rightarrow XYZ.$
 - 注意: $A \rightarrow \epsilon$ 只对应一个项 $A \rightarrow .$
- 直观含义
 - 项 $A \rightarrow \alpha.\beta$ 表示已经扫描/归约到了 α , 并期望接下来的输入中经过扫描/归约得到 β , 然后把 $\alpha\beta$ 归约到 A
 - 如果 β 为空, 表示我们可以把 α 归约为 A
- 项也可以用一对整数表示
 - (i,j) 表示第 i 条规则, 点位于右部第 j 个位置



规范LR(0)项集族



- 规范LR(0)项集族提供构建LR(0)自动机的基础
 - LR(0)自动机可用于做出语法分析决定
 - LR(0)自动机中每个状态代表了LR(0)项集族中的一个项集



以项为基础构造自动机



- 构造以项为状态的自动机
 - 开始状态 $S' \rightarrow .S$
 - 转换
 - $A \rightarrow \alpha.B\beta$ 到 $B \rightarrow .\gamma$ 有一个 ϵ 转换
 - 从 $A \rightarrow \alpha.X\beta$ 到 $A \rightarrow \alpha X.\beta$ 有一个 X 转换
 - 接受状态: $A \rightarrow \alpha.$, 即点在最后的项



规范LR(0)项集族的构造



- 三个概念
 - 增广文法
 - 项集闭包: **CLOSURE**
 - **GOTO**



规范LR(0)项集族的构造



■ 增广文法

- **G**的增广文法**G'**是在**G**中增加新开始符号**S'**，并加入产生式**S' → S**而得到的
- **G'**和**G**接受相同的语言，且按照**S' → S**进行归约实际上就表示已经将输入符号串归约成为开始符号
- 引入的目的是告诉语法分析器何时宣布接受输入符号串，即用**S' → S**进行归约时，表明分析结束。



规范LR(0)项集族的构造



- 构造过程中用到的子函数
 - **CLOSURE(I)**: I的项集闭包
 - 对应于DFA化算法的 ϵ -CLOSURE
 - **GOTO(I,X)**: I的X后继
 - 对应于DFA化算法的MOVE(I,X)



CLOSURE (I) 的构造算法



- 如果I是文法G的一个项集，那么CLOSURE(I)就是根据下列规则从I构造得到的项集
 - 将I中的各个项加入到CLOSURE(I)中
 - 如果 $A \rightarrow \alpha.B\beta$ 在CLOSURE(I)中，那么对B的任意产生式 $B \rightarrow \gamma$ ，将 $B \rightarrow \gamma$ 加到CLOSURE(I)中
 - 不断重复第二步，直到收敛
- 第二步的意义
 - 项 $A \rightarrow \alpha.B\beta$ 表示期望在接下来的输入中归约到B
 - 显然，要归约到B，首先要扫描归约到B的某个产生式的右部
 - 因此对每个产生式 $B \rightarrow \gamma$ ，加入 $B \rightarrow \gamma$
 - 表示它期望能够扫描归约到 γ



项集闭包构造的例子



■ 增广文法:

- $E' \rightarrow E$ $E \rightarrow E+T \mid T$ $T \rightarrow T * F \mid F$ $F \rightarrow (E) \mid id$

■ 项集 $I = \{[E' \rightarrow \cdot E]\}$ 的闭包

- $[E' \rightarrow \cdot E]$ 在闭包中
- $[E \rightarrow \cdot E+T]$, $[E \rightarrow \cdot T]$ 在闭包中
- $[T \rightarrow \cdot T * F]$, $[T \rightarrow \cdot F]$ 在闭包中
- $[F \rightarrow \cdot (E)]$, $[F \rightarrow \cdot id]$ 在闭包中

```
SetOfItems CLOSURE( $I$ ) {  
     $J = I$ ;  
    repeat  
        for ( $J$  中的每个项  $A \rightarrow \alpha \cdot B \beta$ )  
            for ( $G$  的每个产生式  $B \rightarrow \gamma$ )  
                if (项  $B \rightarrow \cdot \gamma$  不在  $J$  中)  
                    将  $B \rightarrow \cdot \gamma$  加入  $J$  中;  
    until 在某一轮中没有新的项被加入到  $J$  中;  
    return  $J$ ;  
}
```



LR(0) 项集中的内核项和非内核项



- 内核项：初始项 $S' \rightarrow .S$ 、以及所有点不在最左边的项
- 非内核项：除了 $S' \rightarrow .S$ 之外、点在最左边的项
- 由于表可能很庞大，实现算法时可以考虑只保存相应的非终结符号；甚至可以只保存内核项，而在要使用非内核项时调用 **CLOSURE** 函数重新计算



GOTO函数



- I 是一个项集， X 是一个文法符号， $GOTO(I, X)$ 定义为 I 中所有形如的项 $[A \rightarrow \alpha \cdot X \beta]$ 所对应的项 $[A \rightarrow \alpha X \cdot \beta]$ 的集合的闭包
 - 根据项的历史-期望的含义， $GOTO(I, X)$ 表示读取输入中的 X 或者归约到一个 X 之后的情况
 - $GOTO(I, X)$ 定义了LR(0)自动机中状态 I 在 X 之上的转换
- 例如：
 - $I = \{[E' \rightarrow E \cdot], [E \rightarrow E \cdot + T]\}$
 - $GOTO(I, +)$ 计算如下
 - I 中只有一个项的点后面跟着 $+$ ，对应的项为 $[E \rightarrow E + \cdot T]$
 - $CLOSURE(\{[E \rightarrow E + \cdot T]\}) = \{[E \rightarrow E + \cdot T], [T \rightarrow \cdot T * F], [T \rightarrow \cdot F], [F \rightarrow \cdot (E)], [F \rightarrow \cdot id]\}$



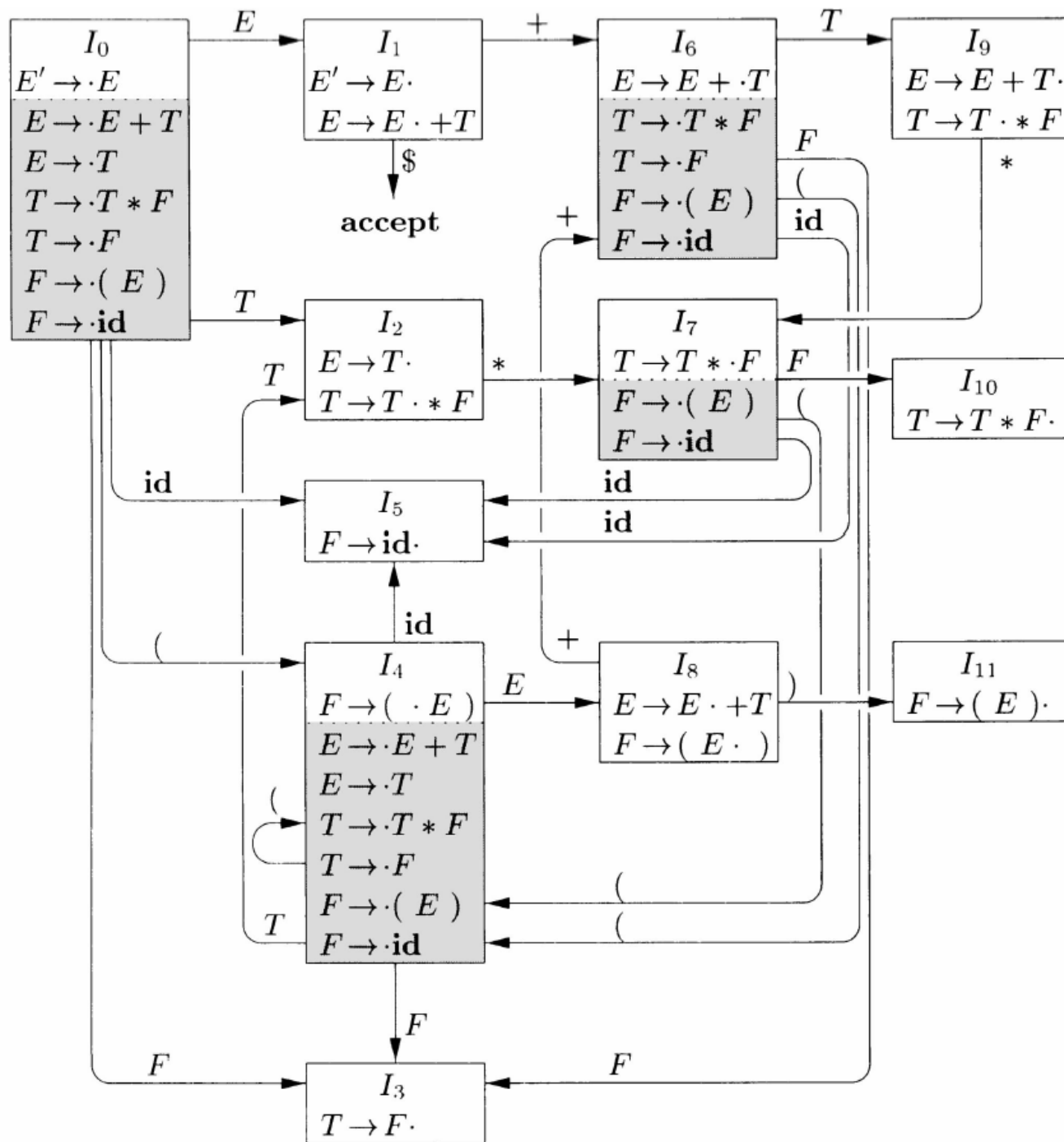
求LR(0)项集规范族的算法



■ 为文法G构造LR(0)项集规范族C

- 步骤一：构造增广文法 $G': S' \rightarrow S \dots\dots$
- 步骤二：构造 $I_0 = \text{CLOSURE}(S' \rightarrow \cdot S)$ ， I_0 是C的初始项集
- 步骤三：对C中的每个项集 I_i 及每个文法符号 X_j ，将 $\text{GOTO}(I_i, X_j)$ 加入到C中
- 重复步骤三，直到没有新的项集可以加入

```
void items( $G'$ ) {  
     $C = \text{CLOSURE}(\{[S' \rightarrow \cdot S]\})$ ;  
    repeat  
        for ( $C$ 中的每个项集 $I$ )  
            for (每个文法符号  $X$ )  
                if (  $\text{GOTO}(I, X)$  非空且不在  $C$ 中)  
                    将  $\text{GOTO}(I, X)$  加入  $C$  中;  
    until 在某一轮中没有新的项集被加入到  $C$  中;  
}
```



- (1) $E \rightarrow E + T$
- (2) $E \rightarrow T$
- (3) $T \rightarrow T * F$
- (4) $T \rightarrow F$
- (5) $F \rightarrow (E)$
- (6) $F \rightarrow id$

笑和GOTO函



LR(0)自动机的构造



■ 构造方法

- 规范LR(0)项集族中的项集可以作为LR(0)自动机的状态
- $GOTO(I, X) = J$, 则从I到J有一个标号为X的转换
- 初始状态为 $CLOSURE(\{S' \rightarrow .S\})$ 对应的项集
- 接受状态: 包含形如 $A \rightarrow \alpha.$ 的项集对应的状态



LR(0)自动机的作用



- 根据文法构造出**LR(0)**自动机，通过自动机的运行进行语法分析
- 假设文法符号串 γ 使**LR(0)**自动机从开始状态**0**运行到某个状态**j**，**LR(0)**自动机按照如下方式决定移入或归约
 - 如果下一个输入符号为**a**，且**j**上有**a**的转换，就移入**a**
 - 否则就归约，状态**j**中的项会告诉我们使用那个产生式进行归约



LR(0) 自动机的作用 (1)



- 假设文法符号串 γ 使 LR(0) 自动机从开始状态运行到状态（项集） j
 - 如果 j 中有一个形如 $A \rightarrow \alpha.$ 的项，那么
 - 在 γ 之后添加一些终结符号可以得到一个最右句型
 - α 是 γ 的后缀，且 $A \rightarrow \alpha$ 是这个句型的句柄
 - 表示可能找到了当前最右句型的句柄
 - 如果 j 中存在一个项 $B \rightarrow \alpha.X\beta$ ，那么
 - 在 γ 之后添加 $X\beta$ ，然后再添加一个终结符号串可得到一个最右句型
 - 在这个句型中 $B \rightarrow \alpha X\beta$ 是句柄
 - 此时表示还没有找到句柄，需要移入



LR(0) 自动机的作用 (2)



■ LR(0)自动机的使用

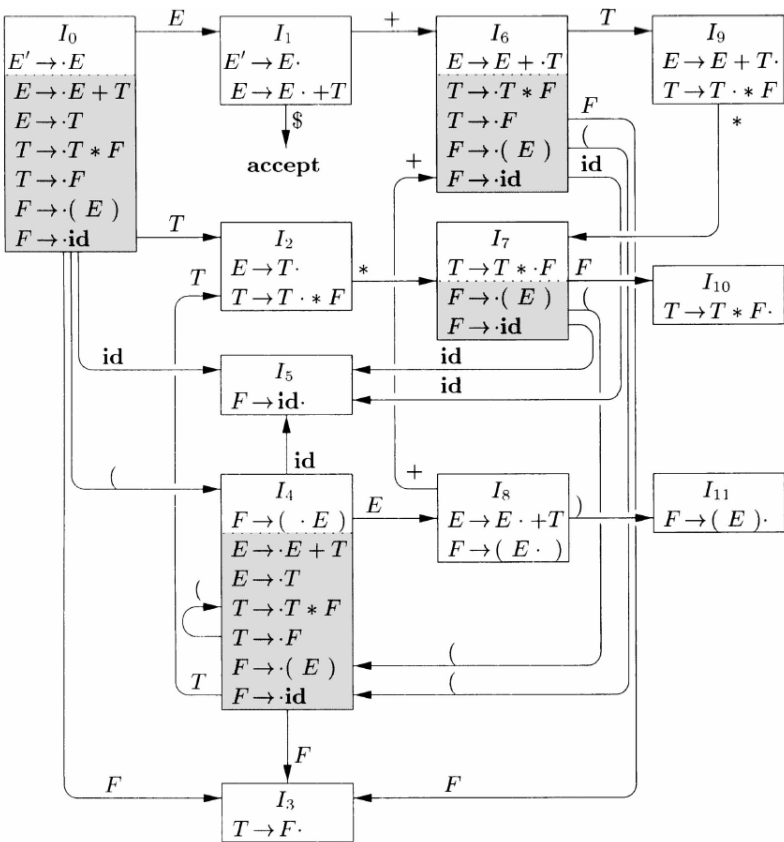
- 移入-归约时，LR(0)自动机被用于识别句型
- 已经归约/移入得到的文法符号序列对应于LR(0)自动机的一条路径
- 如果路径到达接受状态，表明栈上端的某个符号串可能是句柄

■ 不需要每次用归约/移入得到的串来运行LR(0)自动机

- 路径被放到栈中；且文法符号可以省略，因为由LR(0)状态可以确定相应文法符号
- 在移入后，根据原来的栈顶状态即可知道新的状态
- 在归约时，根据归约产生式的右部长长度弹出相应状态，仍然可以根据此时的栈顶状态计算得到新状态



LR(0) 自动机作用示例：分析id*id

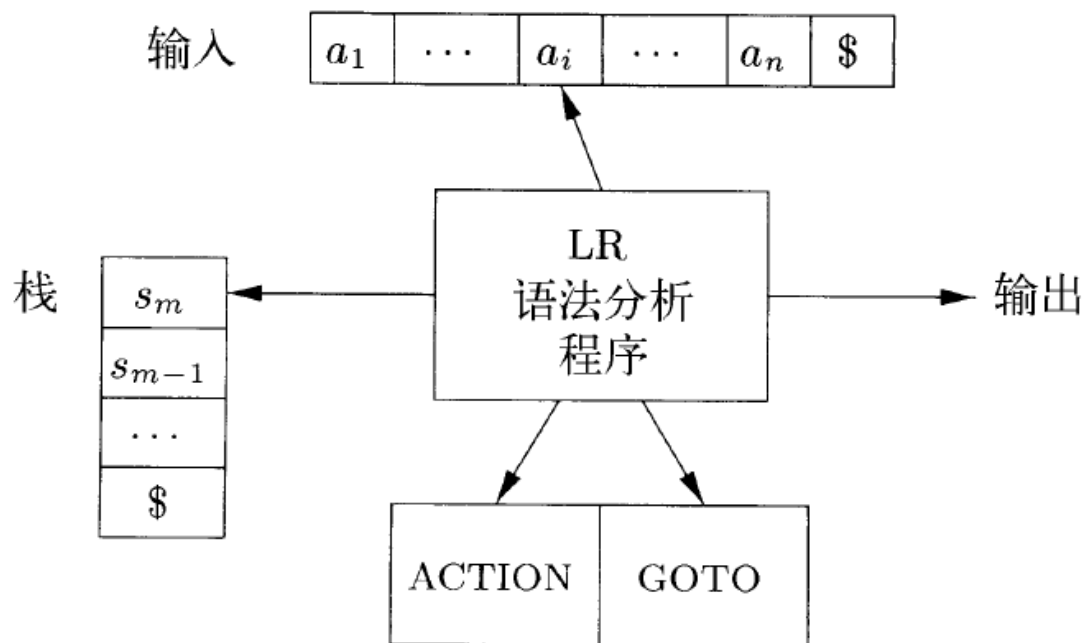


行号	栈	符号	输入	动作
(1)	0	\$	id * id \$	移入到 5
(2)	0 5	\$ id	* id \$	按照 $F \rightarrow \text{id}$ 归约
(3)	0 3	\$ F	* id \$	按照 $T \rightarrow F$ 归约
(4)	0 2	\$ T	* id \$	移入到 7
(5)	0 2 7	\$ T *	id \$	移入到 5
(6)	0 2 7 5	\$ T * id	\$	按照 $F \rightarrow \text{id}$ 归约
(7)	0 2 7 10	\$ T * F	\$	按照 $T \rightarrow T * F$ 归约
(8)	0 2	\$ T	\$	按照 $E \rightarrow T$ 归约
(9)	0 1	\$ E	\$	接受

- 根据LR(0)自动机状态和符号的对应关系，可以得到符号栏中的符号串
- 路径被放到栈中；且文法符号可以省略，因为由LR(0)状态可以确定相应文法符号



LR语法分析器的结构



- 所有的分析器都使用相同的驱动程序
- 分析表随文法以及LR分析技术不同而不同
- 栈中存放的是状态序列；可以由状态序列求出符号序列
- 分析程序根据栈顶状态、当前输入，通过分析表确定语法分析动作



LR语法分析表的结构



- 两个部分：动作ACTION，转换GOTO
- ACTION有两个参数：状态i、终结符号a
 - 移入j：j是一个状态。把j压入栈
 - 归约 $A \rightarrow \beta$ ：把栈顶的 β 归约为A
 - 接受：接受输入、完成分析
 - 报错：在输入中发现语法错误
- 状态集上的GOTO函数
 - 如果 $GOTO[I_i, A] = I_j$ ，那么 $GOTO[i, A] = j$



LR语法分析器的格局



- LR语法分析器的格局包含了栈中内容和余下输入
 $(s_0 s_1 \dots s_m, a_i a_{i+1} \dots a_n \$)$
 - 第一个分量是栈中的内容（右侧是栈顶）
 - 第二个分量是余下输入
- LR语法分析器的每一个状态都对应一个文法符号（ s_0 除外）
 - 如果进入状态 s 的边的标号为 X ，那么 s 就对应于 X
- 令 X_i 为 s_i 对应的符号，那么
 - $X_1 X_2 \dots X_m a_i a_{i+1} \dots a_n$ 对应于一个最右句型



LR语法分析器的行为



- 对于格局($s_0s_1\dots s_m, a_ia_{i+1}\dots a_n\$$), LR语法分析器查询条目 $\text{ACTION}[s_m, a_i]$ 确定相应动作
 - 移入s: 执行移入动作, 将s (对应a) 移入栈中, 新格局($s_0s_1\dots s_ms, a_{i+1}\dots a_n\$$)
 - 归约 $A \rightarrow \beta$, 将栈顶的 β 归约为A, 压入状态s ($s_0s_1\dots s_{m-r}s, a_ia_{i+1}\dots a_n\$$), 其中r是 β 的长度, $s = \text{GOTO}[s_{m-r}, A]$
 - 接受: 语法分析过程完成
 - 报错: 发现语法错误, 调用错误恢复例程



LR语法分析算法



- 输入：文法 G 的LR语法分析表，输入串 w
- 输出：如果 w 在 $L(G)$ 中，输出最左归约步骤（最右推导的反向过程），否则输出错误指示
- 算法

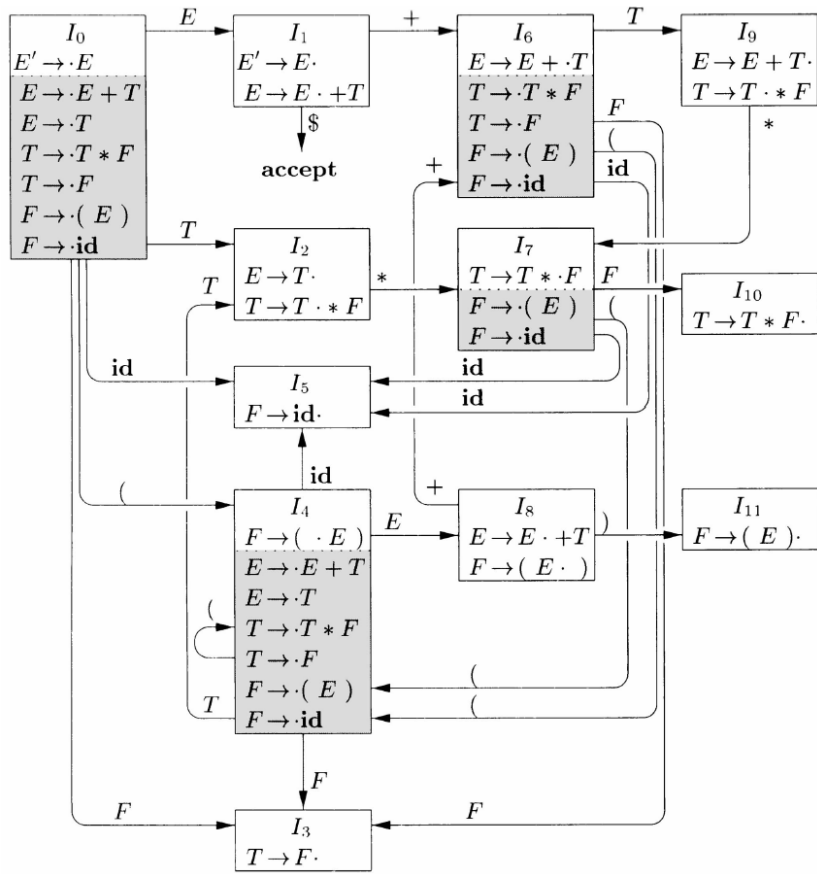
```
令  $a$  为  $w\$$  的第一个符号;  
while(1) { /* 永远重复 */  
    令  $s$  是栈顶的状态;  
    if ( ACTION[ $s, a$ ] = 移入  $t$  ) {  
        将  $t$  压入栈中;  
        令  $a$  为下一个输入符号;  
    } else if ( ACTION[ $s, a$ ] = 归约  $A \rightarrow \beta$  ) {  
        从栈中弹出  $|\beta|$  个符号;  
        令  $t$  为当前的栈顶状态;  
        将 GOTO[ $t, A$ ] 压入栈中;  
        输出产生式  $A \rightarrow \beta$ ;  
    } else if ( ACTION[ $s, a$ ] = 接受 ) break; /* 语法分析完成 */  
    else 调用错误恢复例程;  
}
```



LR分析表的例子



- 文法:
- (1) $E \rightarrow E+T$
 - (2) $E \rightarrow T$
 - (3) $T \rightarrow T*F$
 - (4) $T \rightarrow F$
 - (5) $F \rightarrow (E)$
 - (6) $F \rightarrow id$



状态	ACTION						GOTO		
	id	+	*	()	\$	E	T	F
0	s5			s4			1	2	3
1		s6				acc			
2		r2	s7		r2	r2			
3		r4	r4		r4	r4			
4	s5			s4			8	2	3
5		r6	r6		r6	r6			
6	s5			s4				9	3
7	s5			s4				10	
8		s6			s11				
9		r1	s7		r1	r1			
10		r3	r3		r3	r3			
11		r5	r5		r5	r5			



LR分析

■ 输入: $id * id + id$

状态	ACTION						GOTO		
	id	+	*	()	\$	E	T	F
0	s5				s4		1	2	3
1		s6				acc			
2		r2	s7			r2			
3		r4	r4			r4			
4	s5				s4		8	2	3
5		r6	r6			r6			
6	s5				s4			9	3
7	s5				s4				10
8		s6				s11			
9		r1	s7			r1			
10		r3	r3			r3			
11		r5	r5			r5			

	栈	符号	输入	动作
(1)	0		$id * id + id \$$	移入
(2)	0 5	id	$* id + id \$$	根据 $F \rightarrow id$ 归约
(3)	0 3	F	$* id + id \$$	根据 $T \rightarrow F$ 归约
(4)	0 2	T	$* id + id \$$	移入
(5)	0 2 7	$T *$	$id + id \$$	移入
(6)	0 2 7 5	$T * id$	$+ id \$$	根据 $F \rightarrow id$ 归约
(7)	0 2 7 10	$T * F$	$+ id \$$	根据 $T \rightarrow T * F$ 归约
(8)	0 2	T	$+ id \$$	根据 $E \rightarrow T$ 归约
(9)	0 1	E	$+ id \$$	移入
(10)	0 1 6	$E +$	$id \$$	移入
(11)	0 1 6 5	$E + id$	$\$$	根据 $F \rightarrow id$ 归约
(12)	0 1 6 3	$E + F$	$\$$	根据 $T \rightarrow F$ 归约
(13)	0 1 6 9	$E + T$	$\$$	根据 $E \rightarrow E + T$ 归约
(14)	0 1	E	$\$$	接受



SLR分析表构造

很重要



- 以LR(0)项和LR(0)自动机为基础，基于文法 G 的规范项集族 C 和GOTO函数就可以构造出SLR的语法分析表



SLR分析表构造算法



- 输入：一个增广文法 G'
- 输出： G' 的SLR语法分析表函数ACTION和GOTO
- 方法：
 - 1) 构造 G' LR(0)项集规范族 $C=\{I_0, I_1, \dots, I_n\}$
 - 2) 对于状态 i 中的项：
 - $[A \rightarrow \alpha \cdot a \beta]$ 且 $\text{GOTO}[I_i, a]=I_j$, 那么 $\text{ACTION}[i, a]=\text{移入}j$
 - $[A \rightarrow \alpha \cdot]$, 那么对于 $\text{FOLLOW}(A)$ 中的所有 a , $\text{ACTION}[i, a]=\text{归约}A \rightarrow \alpha$
 - $[S' \rightarrow S \cdot]$, 那么 $\text{ACTION}[i, \$]=\text{接受}$
 - 3) 状态 i 对于非终结符号 A 的转换规则：

如果 $\text{GOTO}[I_i, A]=I_j$, 那么 $\text{GOTO}[i, A]=j$
 - 4) 规则2)和3)没有定义的所有条目设置为“报错”



SLR(1)



- 根据上一算法构造的语法分析表，若表中各位置没有多个条目，则称为文法**G**的**SLR(1)**分析表。使用该分析表的分析器，称为**G**的**SLR(1)**语法分析器。**G**称为**SLR(1)**文法
- 若表中某位置存在多个条目（多个可选的动作，冲突），则该文法是非**SLR(1)**文法

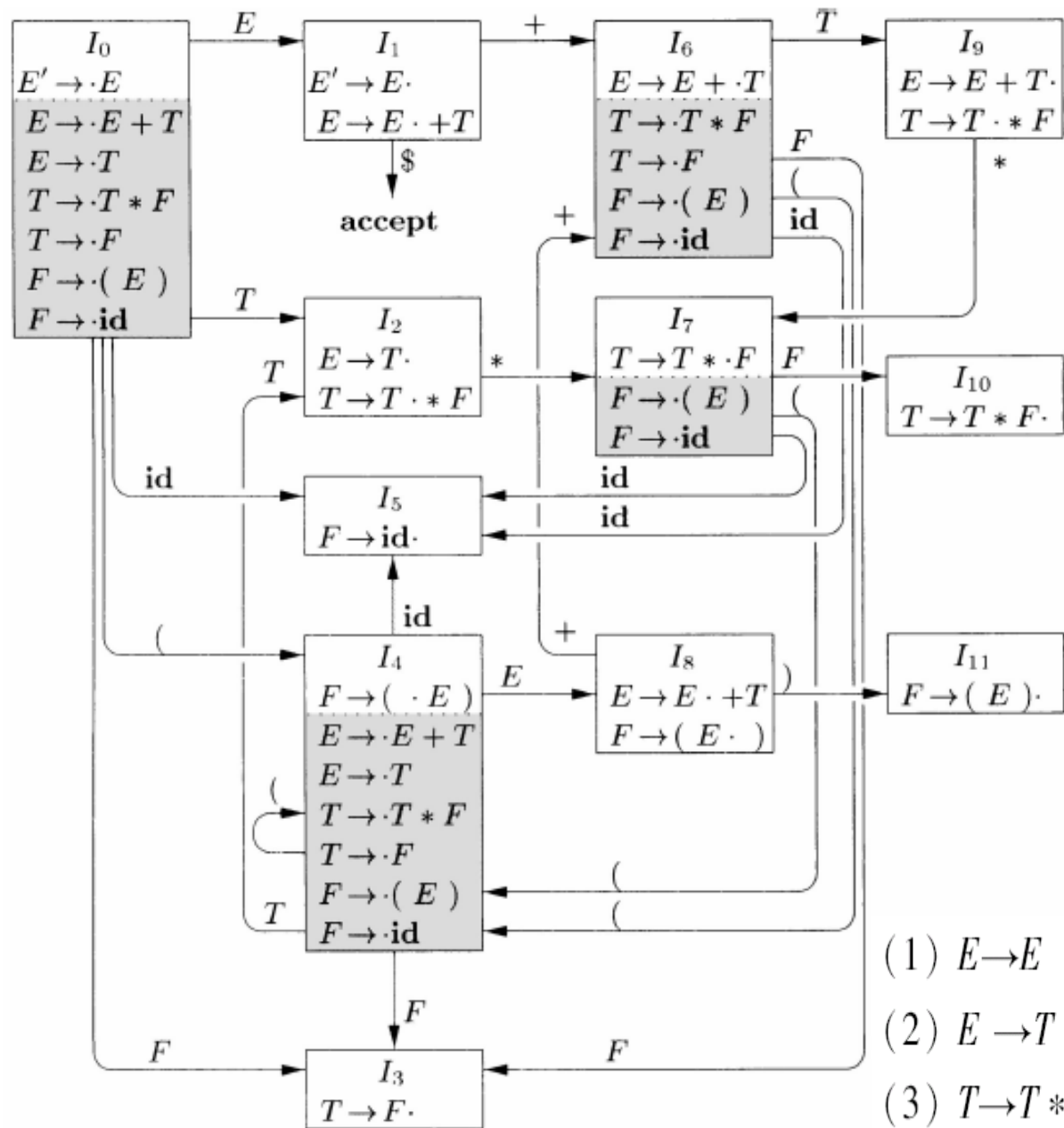


图 4-31 表达式文法(4.1)的 LR(0) 自动机



SLR的原理：可行前缀（1）



- LR(0)自动机刻画了可能出现在语法分析栈中的文法符号串 (能够识别可行前缀)
 - 直接把项当作状态，可以构造得到一个NFA
 - 然后确定化得到DFA就是LR(0)自动机
- 可行前缀
 - 可以出现在移入-归约语法分析器栈中的最右句型的前缀
 - 定义：某个最右句型的前缀，且没有越过该句型的句柄的右端

$$E \xRightarrow[rm]{*} F * \mathbf{id} \Rightarrow_{rm} (E) * \mathbf{id}$$

可能的前缀：
(, (E, (E), (E)*



SLR的原理：可行前缀（2）



■ 有效项

- 如果存在 $S' \xRightarrow{*}_{rm} \alpha Aw \Rightarrow \alpha \beta_1 \beta_2 w$ ，那么我们说项 $A \rightarrow \beta_1 \cdot \beta_2$ 对 $\alpha \beta_1$ 有效

■ 有效项的意义

- 帮助我们决定是进行规约还是移入
- 如果我们知道 $A \rightarrow \beta_1 \cdot \beta_2$ 对 $\alpha \beta_1$ 有效
 1. 当 β_2 不等于空，表示句柄尚未出现在栈中，应该移入或者等待归约
 2. 如果 β_2 等于空，表示句柄出现在栈中，应该归约



SLR的原理：可行前缀（3）



- 如果某个时刻存在两个有效项要求执行不同的动作，那么就应该设法解决冲突
 - 冲突实际上表示可行前缀可能是两个最右句型的前缀，第一个包含了句柄，而另一个尚未包含句柄
 - SLR解决冲突的思想：假如要按照 $A \rightarrow \beta$ 进行归约，那么得到的新句型中A后面跟随下一个输入符号。因此只有当下一个输入在 $FOLLOW(A)$ 中时才可以归约



SLR的原理：可行前缀（4）



- 如果在文法 G 的 $LR(0)$ 自动机中，从初始状态出发，沿着标号为 γ 的路径到达一个状态，那么这个状态对应的项集就是 γ 的有效项集
- 回顾确定分析动作的方法，就可以知道我们实际上是按照有效项来确定的
 - 为了避免冲突，归约时要求下一个输入符号在 $FOLLOW(A)$ 中



可行前缀/有效项的例子



可行前缀 $E+T^*$

对应的LR(0)项

$T \rightarrow T^*.F$

$F \rightarrow \cdot(E)$

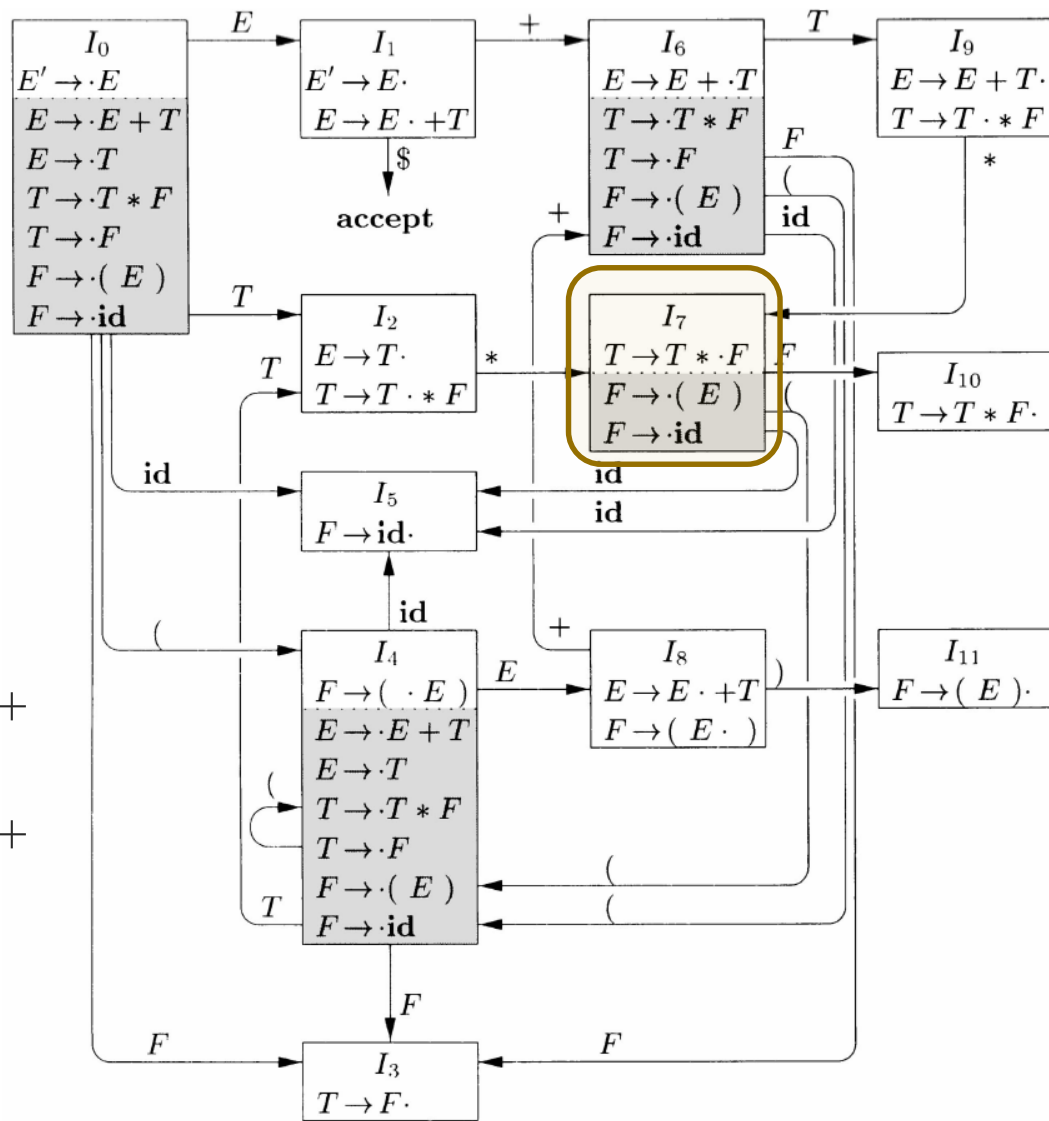
$F \rightarrow \cdot id$

对应的最右推导

$E' \rightarrow E \rightarrow E+T \rightarrow E+T^*F$

$E' \rightarrow E \rightarrow E+T \rightarrow E+T^*F \rightarrow E+T^*(E)$

$E' \rightarrow E \rightarrow E+T \rightarrow E+T^*F \rightarrow E+T^*id$





示例二



$I_0:$ $S' \rightarrow \cdot S$
 $S \rightarrow \cdot L = R$
 $S \rightarrow \cdot R$
 $L \rightarrow \cdot * R$
 $L \rightarrow \cdot \text{id}$
 $R \rightarrow \cdot L$

移进/规约
冲突

$I_1:$ $S' \rightarrow S \cdot$

$I_2:$ $S \rightarrow L \cdot = R$
 $R \rightarrow L \cdot$

$I_3:$ $S \rightarrow R \cdot$

$I_4:$ $L \rightarrow * \cdot R$
 $R \rightarrow \cdot L$
 $L \rightarrow \cdot * R$
 $L \rightarrow \cdot \text{id}$

$I_5:$ $L \rightarrow \text{id} \cdot$

$I_6:$ $S \rightarrow L = \cdot R$
 $R \rightarrow \cdot L$
 $L \rightarrow \cdot * R$
 $L \rightarrow \cdot \text{id}$

$I_7:$ $L \rightarrow * R \cdot$

$I_8:$ $R \rightarrow L \cdot$

$I_9:$ $S \rightarrow L = R \cdot$

$S \rightarrow L = R \mid R$

$L \rightarrow * R \mid \text{id}$

$R \rightarrow L$

状态2，输入符号为“=”时

- 移进
- Follow(R)=? , 规约



SLR语法分析器的弱点 (1)



■ SLR技术解决冲突的方法

- 项集中包含 $[A \rightarrow \alpha.]$ 时, 按照 $A \rightarrow \alpha$ 进行归约的条件是下一个输入符号 a 可以在某个句型中跟在 A 之后
 - 如果对于 a 还有其它的移入/归约操作, 则出现冲突

$$S \rightarrow L = R \mid R$$

$$L \rightarrow * R \mid \mathbf{id}$$

$$R \rightarrow L$$

$$I_2: \quad S \rightarrow L \cdot = R \\ R \rightarrow L \cdot$$

- 假设此时栈中的符号串为 $\beta\alpha$
 - 如果 $\beta A a$ 不可能是某个最右句型的前缀, 那么即使 a 在某个句型中跟在 A 之后, 仍然不应该按照 $A \rightarrow \alpha$ 归约
 - 进行归约的条件更加严格可以降低冲突的可能性



SLR语法分析器的弱点（2）



- $[A \rightarrow \alpha.]$ 出现在项集中的条件：
 - 首先 $[A \rightarrow . \alpha]$ 出现在某个项集中，然后逐步移入/归约到 α 中的符号，点不断后移，到达末端
 - 而 $[A \rightarrow . \alpha]$ 出现的条件是 $B \rightarrow \beta.A\gamma$ 出现在项中
 - 期望首先按照 $A \rightarrow \alpha$ 归约，然后将 $B \rightarrow \beta.A\gamma$ 中的点后移到A之后
 - 显然，在按照 $A \rightarrow \alpha$ 归约时要求下一个输入符号是 γ 的第一个符号
 - 但是从LR(0)项集中不能确定这个信息



习题



- 给出如下文法的语法分析表，判断该文法是否**SLR**文法。

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



更强大的LR语法分析器



- 规范LR方法。充分利用向前看符号。这种方法是用了一个更大的项集，称为LR(1)项集
- 向前看LR，又称LALR



LR(1)项



- 在SLR中，如果项集 I_i 中包含项 $[A \rightarrow \alpha \cdot]$ ，且当前符号 a 在 $FOLLOW(A)$ 中，那么就可以按照 $A \rightarrow \alpha$ 进行归约
- 有时，在任何最右句型中 a 都不可能跟在 βA 之后，这时输入为 a 不能按照 $A \rightarrow \alpha$ 进行归约
- $FOLLOW$ 集提供的可归约条件过松

$$S \rightarrow L = R \mid R$$

$$L \rightarrow * R \mid \mathbf{id}$$

$$R \rightarrow L$$

$$I_2: \quad S \rightarrow L \cdot = R \\ R \rightarrow L \cdot$$

- 每个状态能否明确指出哪些输入符号可以跟在句柄 α 后面，从而使句柄 α 可能被归约为 A



规范LR(1)项



- 在状态中包含更多的信息
 - 让LR分析器的每个状态精确指明哪些输入符号可以跟在句柄 α 后面，使 α 可能被规约为A
- 让项包含第二个分量
 - 形式如 $[A \rightarrow \alpha \cdot \beta, a]$ ，其中 $A \rightarrow \alpha \beta$ 是一条产生式， a 是一个终结符号或者\$符
 - a 是向前看符号，长度为1



规范LR(1)项



- 项 $[A \rightarrow \alpha \cdot \beta, a]$
 - 如果 $\beta \neq \epsilon$ 时， a 没有任何作用
 - a 的意义是若栈顶状态中存在LR(1)项 $[A \rightarrow \alpha \cdot, a]$ ，在输入符号为 a 时才会进行归约
 - a 的集合总是**FOLLOW(A)**的子集

- **a 指出能够进行归约的准确判断**

向前看符号
与确定归约
有关



规范LR(1)项



- 从可行前缀和有效项的角度看
- LR(1)项 $[A \rightarrow \alpha \cdot \beta, a]$ 对于一个可行前缀 $\delta\alpha$ 有效的条件是存在一个推导 $S \xRightarrow{*} \delta A w$
 $\delta\alpha\beta w$, 其中 $a \in \text{First}(\vec{w})$ 。当 \vec{w} 为 ϵ , $a = \$$

$$S \rightarrow L = R \mid R$$

$$L \rightarrow * R \mid \mathbf{id}$$

$$R \rightarrow L$$

$$I_2: \quad \begin{array}{l} S \rightarrow L \cdot = R \\ R \rightarrow L \cdot \end{array}$$

LR(0): 如果存在 $S' \xRightarrow{*}_{rm} \alpha A w \Rightarrow_{rm} \alpha \beta_1 \beta_2 w$, 那么我们说项 $A \rightarrow \beta_1 \cdot \beta_2$ 对 $\alpha\beta_1$ 有效



LR(1)项集的构造



- 过程类似于LR(0)项集构造
- 多了一个向前看符号分量的计算
- CLOSURE
- GOTO



LR(1)项集闭包CLOSURE



- $\text{CLOSURE}(I) = I \cup \{[B \rightarrow \cdot \eta, \mathbf{b}] \mid [A \rightarrow \alpha \cdot B \beta, \mathbf{a}] \in \text{CLOSURE}(I), b \in \text{First}(\beta a), \text{ 且 } B \rightarrow \eta \text{ 为语法规则}\}$
为什么**b**必须在 $\text{First}(\beta a)$ 中?
 - 对于可行前缀 γ 的有效项 $[A \rightarrow \alpha \cdot B \beta, a]$ ，存在最右推导 $S \xRightarrow{*}_{rm} \delta A a x \Rightarrow_{rm} \delta \alpha B \beta a x$ ，其中 $\gamma = \delta \alpha$
 - 假设 $\beta a x$ 推导出终结符号串 $\mathbf{b}\mathbf{y}$ ，那么对于某个形如 $B \rightarrow \eta$ 的产生式，有推导 $S \xRightarrow{*}_{rm} \gamma B b y \Rightarrow_{rm} \gamma \eta b y$ ，因此， $\{[B \rightarrow \cdot \eta, b]\}$ 是 γ 的有效项
 - **b**是 $\text{First}(\beta a)$ 中的终结符号
 - **b**可能是从 β 推导出的第一个终结符号
 - 也可能 β 推导出了 ϵ ，**b**就是**a**



LR(1)项集闭包CLOSURE



```
SetOfItems CLOSURE( $I$ ) {  
    repeat  
        for (  $I$  中的每个项  $[A \rightarrow \alpha \cdot B \beta, a]$  )  
            for (  $G'$  中的每个产生式  $B \rightarrow \gamma$  )  
                for ( FIRST( $\beta a$ ) 中的每个终结符号  $b$  )  
                    将  $[B \rightarrow \cdot \gamma, b]$  加入到集合  $I$  中;  
    until 不能向  $I$  中加入更多的项;  
    return  $I$ ;  
}
```




LR(1)项集GOTO函数



```
SetOfItems GOTO( $I, X$ ) {  
    将  $J$  初始化为空集;  
    for ( $I$  中的每个项  $[A \rightarrow \alpha \cdot X \beta, a]$ )  
        将项  $[A \rightarrow \alpha X \cdot \beta, a]$  加入到集合  $J$  中;  
    return CLOSURE( $J$ );  
}
```



LR(1)项集规范族构造算法

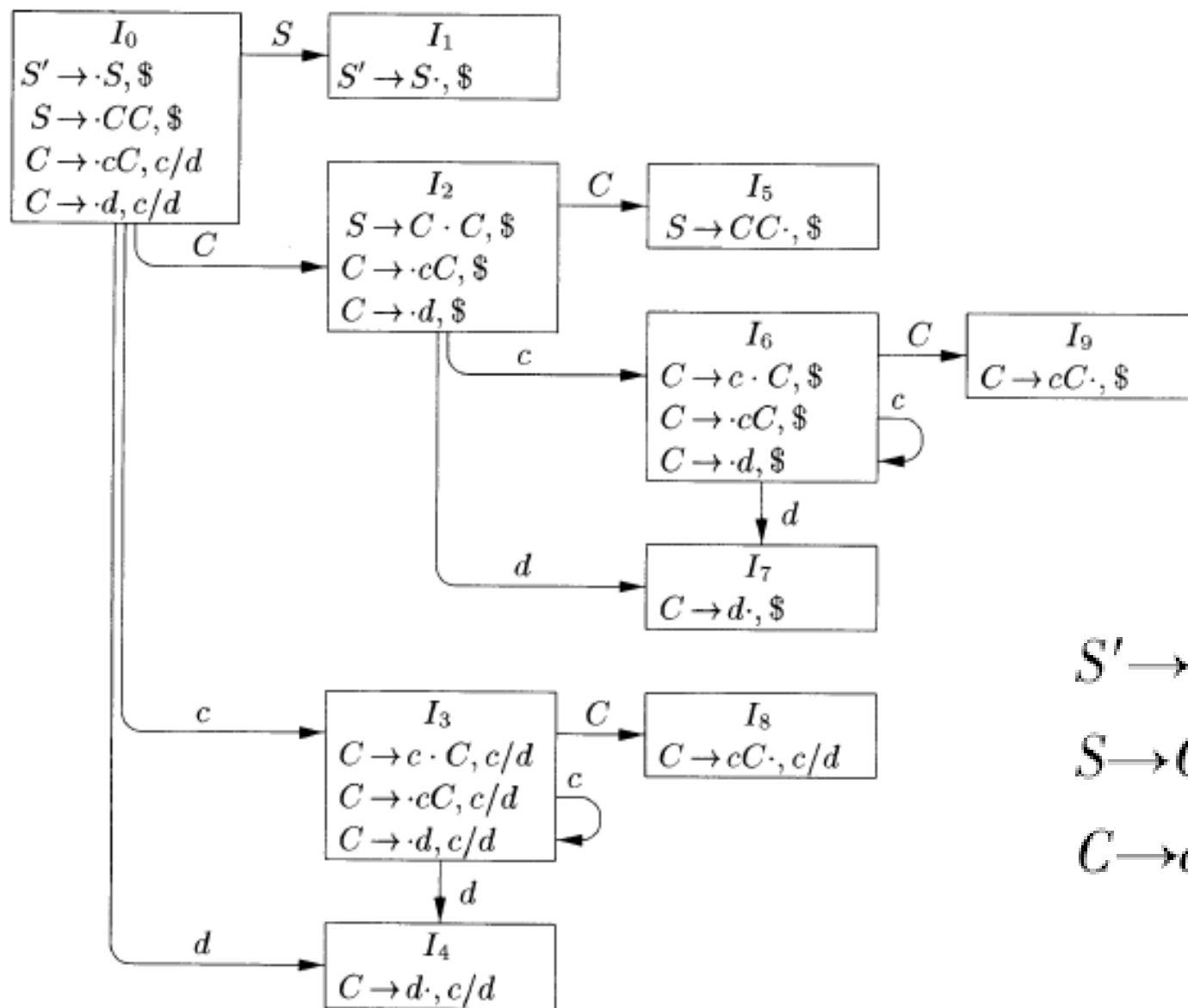


```
SetOfItems CLOSURE( $I$ ) {  
    repeat  
        for (  $I$  中的每个项  $[A \rightarrow \alpha \cdot B\beta, a]$  )  
            for (  $G'$  中的每个产生式  $B \rightarrow \gamma$  )  
                for ( FIRST( $\beta a$ ) 中的每个终结符号  $b$  )  
                    将  $[B \rightarrow \cdot \gamma, b]$  加入到集合  $I$  中;  
    until 不能向  $I$  中加入更多的项;  
    return  $I$ ;  
}  
  
SetOfItems GOTO( $I, X$ ) {  
    将  $J$  初始化为空集;  
    for (  $I$  中的每个项  $[A \rightarrow \alpha \cdot X\beta, a]$  )  
        将项  $[A \rightarrow \alpha X \cdot \beta, a]$  加入到集合  $J$  中;  
    return CLOSURE( $J$ );  
}  
  
void items( $G'$ ) {  
    将  $C$  初始化为 CLOSURE( $\{[S' \rightarrow \cdot S, \$]\}$ );  
    repeat  
        for (  $C$  中的每个项集  $I$  )  
            for ( 每个文法符号  $X$  )  
                if ( GOTO( $I, X$ ) 非空且不在  $C$  中 )  
                    将 GOTO( $I, X$ ) 加入  $C$  中;  
    until 不再有新的项集加入到  $C$  中;  
}
```

图 4-40 为文法 G' 构造 LR(1)项集族的算法



LR(1)规范项集族构造示例



$S' \rightarrow S$

$S \rightarrow C C$

$C \rightarrow c C \mid d$



规范LR(1)语法分析表



- 输入：一个增广文法 G'
- 输出： G' 的规范LR语法分析表函数 **ACTION** 和 **GOTO**
- 方法：
 - 1) 构造 G' LR(1) 项集规范族 $C = \{I_0, I_1, \dots, I_n\}$
 - 2) 对于状态 i 中的项：
 - $[A \rightarrow \alpha \cdot a \beta, b]$ 且 $\text{GOTO}[I_i, a] = I_j$, 那么 $\text{ACTION}[i, a] = \text{移入 } j$
 - $[A \rightarrow \alpha \cdot, a]$, 且 $A \neq S'$, $\text{ACTION}[i, a] = \text{归约 } A \rightarrow \alpha$
 - $[S' \rightarrow S \cdot, \$]$, 那么 $\text{ACTION}[i, \$] = \text{接受}$
 - 3) 状态 i 对于非终结符号 A 的转换规则：如果 $\text{GOTO}[I_i, A] = I_j$, 那么 $\text{GOTO}[i, A] = j$
 - 4) 规则2)和3)没有定义的所有条目设置为“报错”



语法分析表示例

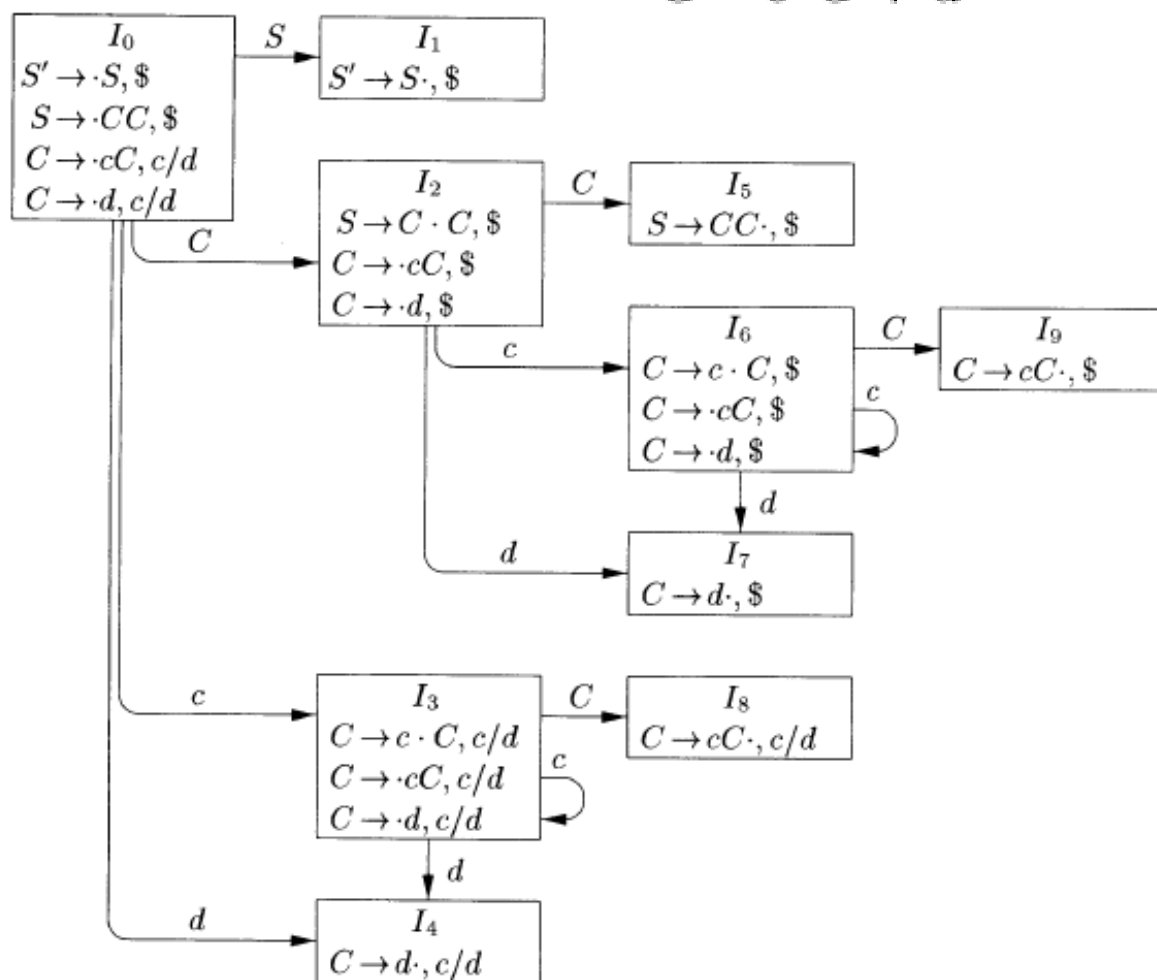


$$S' \rightarrow S$$

$$S \rightarrow C C$$

$$C \rightarrow c C \mid d$$

状态	ACTION			GOTO	
	<i>c</i>	<i>d</i>	\$	<i>S</i>	<i>C</i>
0					
1					
2					
3					
4					
5					
6					
7					
8					
9					





LR(1)文法



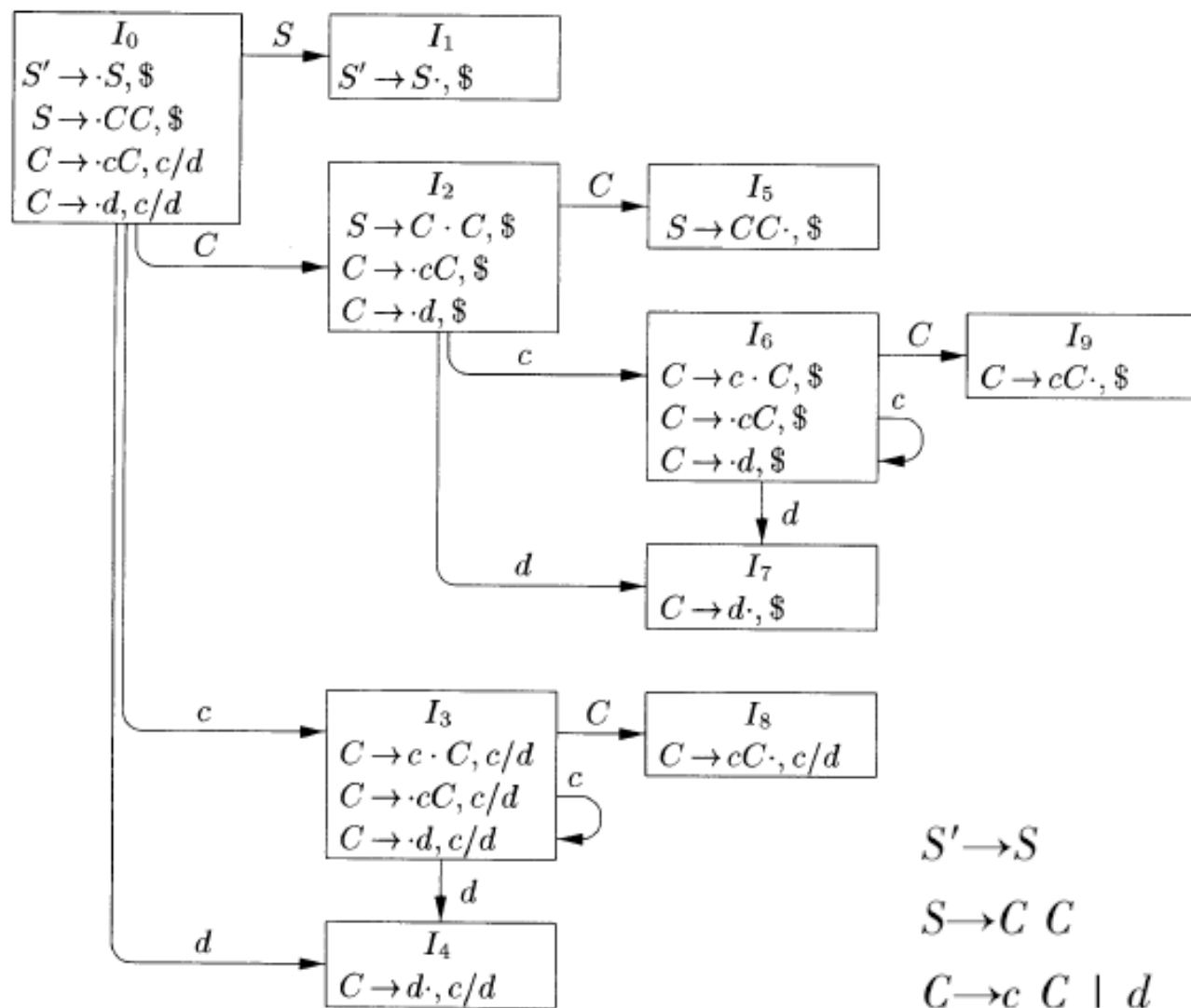
- 通过上述算法构造的LR(1)语法分析表中若不存在冲突条目，则给定的文法称为LR(1)文法



LR(1)语法分析器状态再观察



在(3,6)
(4,7) (8,9)
中, 项的第一个分量是一样的, 实质上它们来自于同样的LR(0)状态





LALR分析



- LR(1)语法分析器的状态过多
- 但LR(1)的分析能力强于SLR(1)
- LALR分析
 - 状态和SLR一样多
 - 能力又比SLR稍微强一些



项集合并



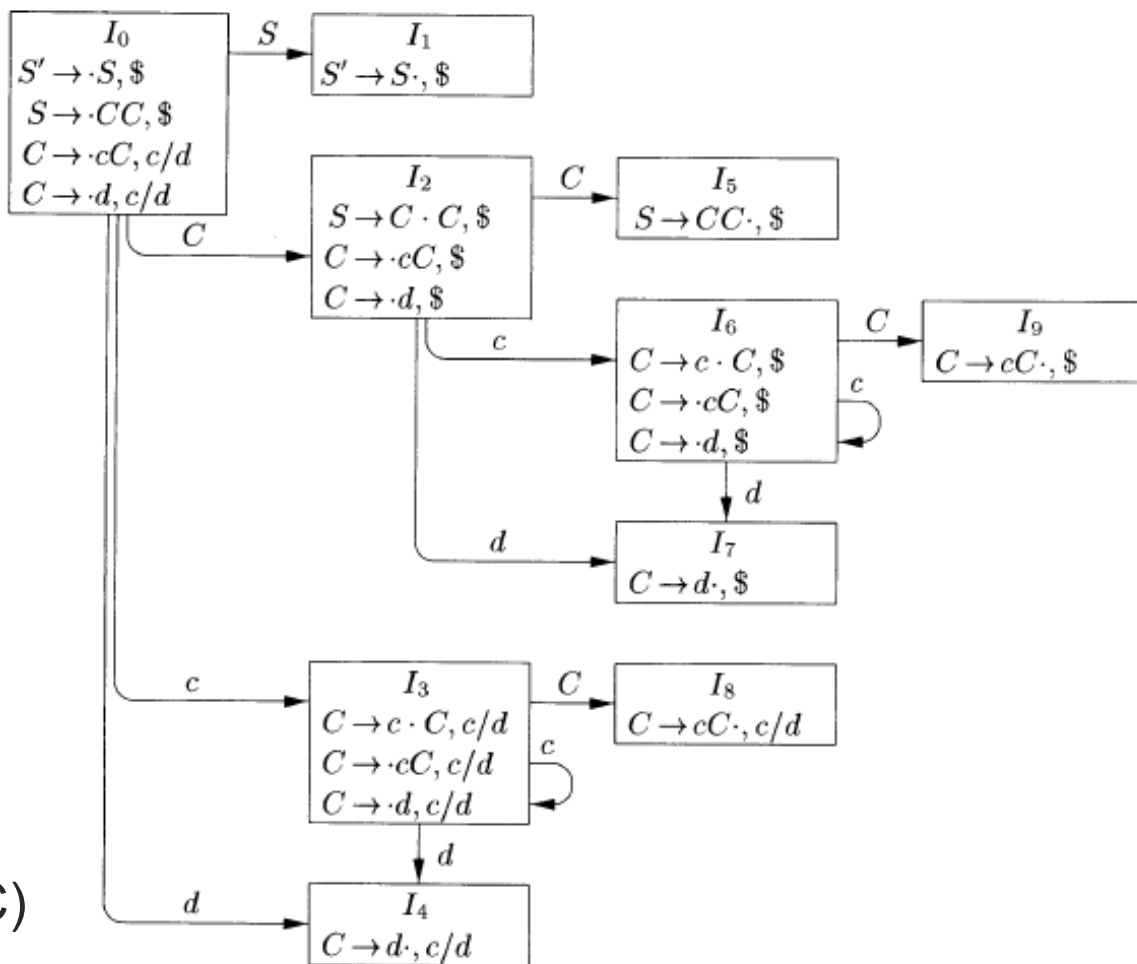
- 将LR(1)项集中具有相同核心的项集合并为一个项集
- 项集核心是指LR(1)项集中第一个分量的集合
- 被合并项集的GOTO目标显然也可以被合并



项集合并



- $I_3 = \{ [C \rightarrow c \cdot C, c/d], [C \rightarrow \cdot cC, c/d], [C \rightarrow \cdot d, c/d] \}$
- $I_6 = \{ [C \rightarrow c \cdot C, \$], [C \rightarrow \cdot cC, \$], [C \rightarrow \cdot d, \$] \}$
- 合并后:
 - $I_{36} = \{ [C \rightarrow c \cdot C, c/d/\$], [C \rightarrow \cdot cC, c/d/\$], [C \rightarrow \cdot d, c/d/\$] \}$
- GOTO(I_3, C)和GOTO(I_6, C)显然也可以合并.....





合并可能产生的冲突

很重要



- 合并引起的冲突是指：本来的LR(1)项集没有冲突，而合并具有相同核心的项集后有冲突
- 不可能引入归约-移入型冲突
 - 假定合并后有移入_归约冲突，就是说：有项 $[A \rightarrow \alpha \cdot, a]$ 和项 $[B \rightarrow \beta \cdot a \gamma, b]$ 。显然，原来的项集中都有 $[B \rightarrow \beta \cdot a \gamma, c]$ 。而 $[A \rightarrow \alpha \cdot, a]$ 必然也在某个原来的项集中。这样，合并前的LR(1)项集已经存在移入-归约冲突
- 但是可能引起归约_归约冲突



归约-归约冲突

$$S' \rightarrow S$$

$$S \rightarrow a A d \mid b B d \mid a B e \mid b A e$$

$$A \rightarrow c$$

$$B \rightarrow c$$

- 语言{acd, ace, bcd, bce}
- 可行前缀ac的有效项集 $\{[A \rightarrow c \cdot, d], [B \rightarrow c \cdot, e]\}$
- 可行前缀bc的有效项集 $\{[A \rightarrow c \cdot, e], [B \rightarrow c \cdot, d]\}$
- 合并之后的项集为 $\{[A \rightarrow c \cdot, d/e], [B \rightarrow c \cdot, d/e]\}$
 - 当输入为d或e时，用哪个归约？
- 从引起新的冲突可以看出：LALR的分析能力比规范LR(1)弱一些



LALR语法分析表构造方法



- 朴素的方法
- 高效的方法



LALR语法分析表构造方法 – 朴素的方法



- 基本思想：
 - 先构造出LR(1)项集，若没有出现冲突，就将核心的项集进行合并。然后根据合并后得到的项集规范组构造语法分析表
- 输入：一个增广文法 G'
- 输出：文法 G' 的LALR语法分析表
- 方法：
 - 得到的分析表成为LALR语法分析表。构造得到LR(1)项集族 $C = \{I_0, I_1, \dots, I_n\}$
 - 对于LR(1)项集中的每个核心，找出所有具有这个核心的项集，并把这些项集替换为它们的并集
 - 令 $C' = \{J_0, J_1, \dots, J_n\}$ 是得到的LR(1)项集族。按照LR(1)分析表的构造方法得到ACTION表。（注意检查若存在冲突，则这个文法不是LALR的）
 - GOTO表的构造：若 J 是一个或者多个LR(1)项集的并集，即 $J = I_1 \cup I_2 \cup \dots \cup I_K$ ，令 K 是所有和 $GOTO(I_1, X)$ 具有相同核心的项集的并集，那么 $GOTO(J, X) = K$



LALR分析表构造示例



$$S' \rightarrow S$$

$$S \rightarrow C C$$

$$C \rightarrow c C \mid d$$

$$I_{36} : C \rightarrow c \cdot C, c/d/ \$$$

$$C \rightarrow \cdot cC, c/d/ \$$$

$$C \rightarrow \cdot d, c/d/ \$$$

$$I_{47} : C \rightarrow d \cdot, c/d/ \$$$

$$I_{89} : C \rightarrow cC \cdot, c/d/ \$$$

状态	ACTION			GOTO	
	<i>c</i>	<i>d</i>	<i>\$</i>	<i>S</i>	<i>C</i>
0	s36	s47		1	2
1			acc		
2	s36	s47			5
36	s36	s47			89
47	r3	r3	r3		
5			r1		
89	r2	r2	r2		



LALR分析器和LR分析器



- 若构造出的LALR分析表中没有冲突，则可以用LALR分析表进行语法分析
- 如果是正确的输入串
 - LALR和LR分析器执行完全相同的移入和归约动作序列，只是栈中的状态名字有所不同
- 如果是错误的输入串
 - 则LALR分析器可能会比LR分析器晚一点报错
 - LALR分析器不会在LR分析器报错后再移入任何符号，但是可能会继续执行一些归约



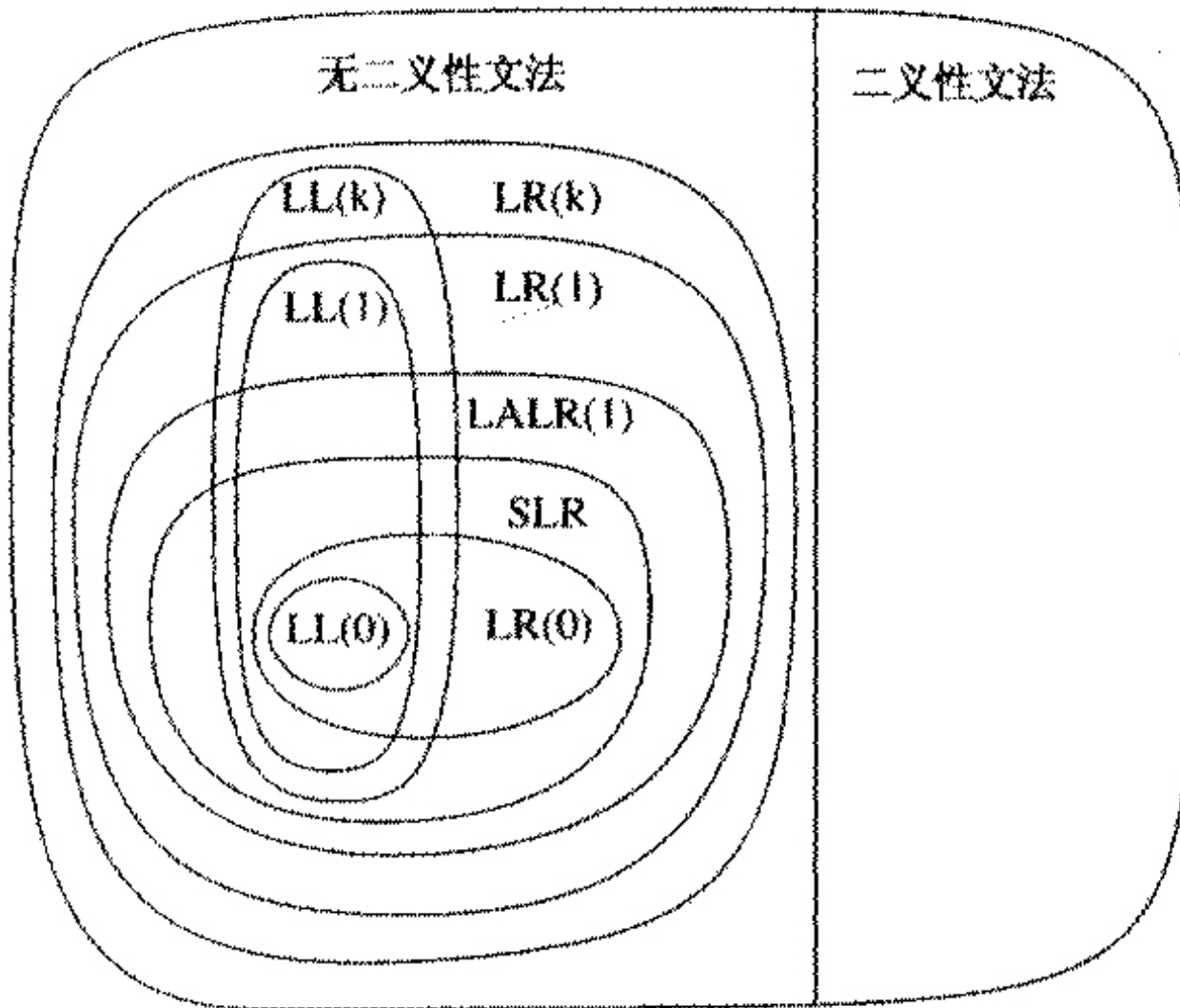
LALR技术本质



- 对LR(1)项集规范族中所谓的同心项集进行合并，从而使得分析表既保持了LR(1)项中向前看符号的信息，又使状态数减少到与SLR分析表的一样多



文法比较





习题



■ 已知文法

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T F \mid F$$

$$F \rightarrow F * \mid (E) \mid a \mid b$$

- ① 给出LR(1)项集 $\{ [F \rightarrow (\cdot E), +] \}$ 的闭包;
- ② 给出①中项集闭包相对于文法符号E的后继LR(1)项集闭包。



二义性文法的自底向上分析



- 二义性文法都不是**LR**的
- 二义性文法却有其存在的必要
- 对于某些二义性文法
 - 可以通过消除二义性规则来保证每个句子只有一棵语法分析树
 - 且可以在**LR**分析器中实现这个规则



优先级/结合性消除冲突



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

$$E \rightarrow E + T \mid T$$

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

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

■ 二义性文法的优点

- 容易修改算符的优先级和结合性
- 简洁：较少的非终结文法符号
- 高效：不需要处理 $E \rightarrow T$ 这样的归约



二义性表达式文法的LR(0)项集



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

- I_7, I_8 中有冲突，在输入+或*时，不能确定是归约还是移入，且不可能通过向前看符号解决

$I_0:$ $E' \rightarrow \cdot E$
 $E \rightarrow \cdot E + E$
 $E \rightarrow \cdot E * E$
 $E \rightarrow \cdot (E)$
 $E \rightarrow \cdot \mathbf{id}$

$I_5:$ $E \rightarrow E * \cdot E$
 $E \rightarrow \cdot E + E$
 $E \rightarrow \cdot E * E$
 $E \rightarrow \cdot (E)$
 $E \rightarrow \cdot \mathbf{id}$

$I_1:$ $E' \rightarrow E \cdot$
 $E \rightarrow E \cdot + E$
 $E \rightarrow E \cdot * E$

$I_6:$ $E \rightarrow (E \cdot)$
 $E \rightarrow E \cdot + E$
 $E \rightarrow E \cdot * E$

$I_2:$ $E \rightarrow (\cdot E)$
 $E \rightarrow \cdot E + E$
 $E \rightarrow \cdot E * E$
 $E \rightarrow \cdot (E)$
 $E \rightarrow \cdot \mathbf{id}$

$I_7:$ $E \rightarrow E + E \cdot$
 $E \rightarrow E \cdot + E$
 $E \rightarrow E \cdot * E$

$I_3:$ $E \rightarrow \mathbf{id} \cdot$

$I_8:$ $E \rightarrow E * E \cdot$
 $E \rightarrow E \cdot + E$
 $E \rightarrow E \cdot * E$

$I_4:$ $E \rightarrow E + \cdot E$
 $E \rightarrow \cdot E + E$
 $E \rightarrow \cdot E * E$
 $E \rightarrow \cdot (E)$
 $E \rightarrow \cdot \mathbf{id}$

$I_9:$ $E \rightarrow (E) \cdot$



基于优先级解决冲突



前缀	栈	输入
$E + E$	0 1 4 7	* id \$

- 如果*的优先级大于+, 且+是左结合的
 - 下一个符号为+时, 我们应该将 $E+E$ 归约为 E
 - 下一个符号为*时, 我们应该移入*, 期待移入下一个符号



解决冲突之后的SLR(1)分析表



■ 对于状态7,
输入

○ +时归约

○ *时移入

■ 对于状态8

$I_7:$ $E \rightarrow E + E \cdot$
 $E \rightarrow E \cdot + E$
 $E \rightarrow E \cdot * E$

$I_8:$ $E \rightarrow E * E \cdot$
 $E \rightarrow E \cdot + E$
 $E \rightarrow E \cdot * E$

状态	ACTION						GOTO
	id	+	*	()	\$	E
0	s3			s2			1
1		s4	s5			acc	
2	s3			s2			6
3		r4	r4		r4	r4	
4	s3			s2			7
5	s3			s2			8
6		s4	s5		s9		
7		r1	s5		r1	r1	
8		r2	r2		r2	r2	
9		r3	r3		r3	r3	



悬空else的二义性



- 栈中内容if *expr* then *stmt*, 是输入else, 还是归约?

- 答案是移入

$S' \rightarrow S$

$S \rightarrow i S e S \mid i S \mid a$

$\text{Follow}(S) = \{e, \$\}$

$I_0:$ $S' \rightarrow \cdot S$
 $S \rightarrow \cdot i S e S$
 $S \rightarrow \cdot i S$
 $S \rightarrow \cdot a$

$I_1:$ $S' \rightarrow S \cdot$

$I_2:$ $S \rightarrow i \cdot S e S$
 $S \rightarrow i \cdot S$
 $S \rightarrow \cdot i S e S$
 $S \rightarrow \cdot i S$
 $S \rightarrow \cdot a$

$I_3:$ $S \rightarrow a \cdot$

$I_4:$ $S \rightarrow i S \cdot e S$
 $S \rightarrow i S \cdot$

$I_5:$ $S \rightarrow i S e \cdot S$
 $S \rightarrow \cdot i S e S$
 $S \rightarrow \cdot i S$
 $S \rightarrow \cdot a$

$I_6:$ $S \rightarrow i S e S \cdot$



语法错误的处理



- 错误难以避免
- 编译器需要具有处理错误的能力
- 程序中可能存在不同层次的错误
 - 词法错误
 - 语法错误
 - 语义错误
 - 逻辑错误
- 语法错误相同容易发现，语义和逻辑错误较难精确的检测到
- 语法分析器中错误处理程序的设计目标
 - 清晰准确地**报告**出现错误，并指出错误的**位置**
 - 能从当前错误中**恢复**，以继续检测后面的错误
 - 尽可能减少处理正确程序的开销



预测分析中的错误恢复



■ 错误恢复

- 当预测分析器报错时，表示输入的串不是句子
- 对于使用者而言，希望预测分析器能够进行恢复处理后继续语法分析过程，以便在一次分析中找到更多的语法错误
- 但是有可能恢复得并不成功，之后找到的语法错误有可能是假的
- 进行错误恢复时可用的信息：栈里面的符号，待分析的符号

■ 两类错误恢复方法

- 恐慌模式
- 短语层次的恢复



恐慌模式



■ 基本思想

- 语法分析器忽略输入中的一些符号，直到出现由设计者选定的某个同步词法单元
- 解释
 - 语法分析过程总是试图在输入的前缀中找到和某个非终结符号对应的串；出现错误表明现在已经不可能找到这个非终结符号（程序结构）的串
 - 如果编程错误仅仅局限于这个程序结构，那么我们可以考虑跳过这个程序结构，假装已经找到了这个结构；然后继续进行语法分析处理
- 同步词法单元就是这个程序结构结束的标志



同步词法单元的确定



- 非终结符号A的同步集合的启发式规则
 - FOLLOW(A)的所有符号在非终结符A的同步集合中
 - 这些符号的出现可能表示之前的那些符号是错误的A的串
 - 将高层次的非终结符号对应串的开始符号加入到较低层次的非终结符号的同步集合
 - 比如语句的开始符号, **if**, **while**等可以作为表达式的同步集合
 - FIRST(A)中的符号加入到非终结符号A的同步集合
 - 碰到这些符号时, 可能意味着前面的符号是多余的符号
 - 如果A可以推导出空串, 那么把此产生式当作默认值
 - 在栈顶的终结符号出现匹配错误时, 可以直接弹出相应的符号, 并且发出消息称已经插入了这个终结符号
- 哪些符号需要确定同步集合需要根据特定的应用而定



恐慌模式的例子（1）



- 对文法4.28对表达式进行语法分析时，可以直接使用FIRST、FOLLOW作为同步集合
 - 我们为E、T和F设定同步集合
 - 空条目表示忽略输入符号；synch表示忽略到同步集合，然后弹出非终结符号
 - 终结符号不匹配时，弹出栈中终结符号

非终结符号	输入符号					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	synch	synch
E'		$E \rightarrow +TE'$			$E \rightarrow \epsilon$	$E \rightarrow \epsilon$
T	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \text{id}$	synch	synch	$F \rightarrow (E)$	synch	synch



恐慌模式 的例子 (2)

■ 错误输入:
 $+id * + id$

栈	输入	说明
$E \$$) $id * + id \$$	错误, 略过)
$E \$$	$id * + id \$$	id 在 $FIRST(E)$ 中
$TE' \$$	$id * + id \$$	
$FT'E' \$$	$id * + id \$$	
$id T'E' \$$	$id * + id \$$	
$T'E' \$$	$* + id \$$	
$* FT'E' \$$	$* + id \$$	
$FT'E' \$$	$+ id \$$	错误, $M[F, +] = synch$
$T'E' \$$	$+ id \$$	F 已经被弹出栈
$E' \$$	$+ id \$$	
$+ TE' \$$	$+ id \$$	
$TE' \$$	$id \$$	
$FT'E' \$$	$id \$$	
$id T'E' \$$	$id \$$	
$T'E' \$$	$\$$	
$E' \$$	$\$$	
$\$$	$\$$	

非终结符号	输入符号					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	synch	synch
E'		$E \rightarrow +TE'$			$E \rightarrow \epsilon$	$E \rightarrow \epsilon$
T	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$	synch	synch	$F \rightarrow (E)$	synch	synch



短语层次的恢复



- 在预测语法分析表的空白条目中插入错误处理例程的函数指针
- 例程可以改变、插入或删除输入中的符号；发出适当的错误消息



LR语法分析中的错误恢复（1）



- LR语法分析器查询GOTO表时不会发现错误
- 但是可能在查询ACTION表时发现报错条目
 - 假设栈中的符号串为 α ，当前输入符号为 a ；报错表示不可能存在终结符号串 x 使得 αax 是一个最右句型
- 恐慌模式的错误恢复策略
 - 从栈顶向下扫描，找到状态 s ， s 有一个对应于非终结符号 A 的GOTO目标；（ s 之上的状态被丢弃）
 - 在输入中读入并丢弃一些符号，直到找到一个可以合法跟在 A 之后的符号 a （不丢弃 a ）
 - 将GOTO(s, A)压栈；继续进行正常的语法分。
 - **基本思想：**假定当前正在试图归约到 A 且碰到了语法错误，因此设法扫描完包含语法错误的 A 的子串，假装找到了 A 的一个实例



LR语法分析中的错误恢复（2）



■ 短语层次的恢复

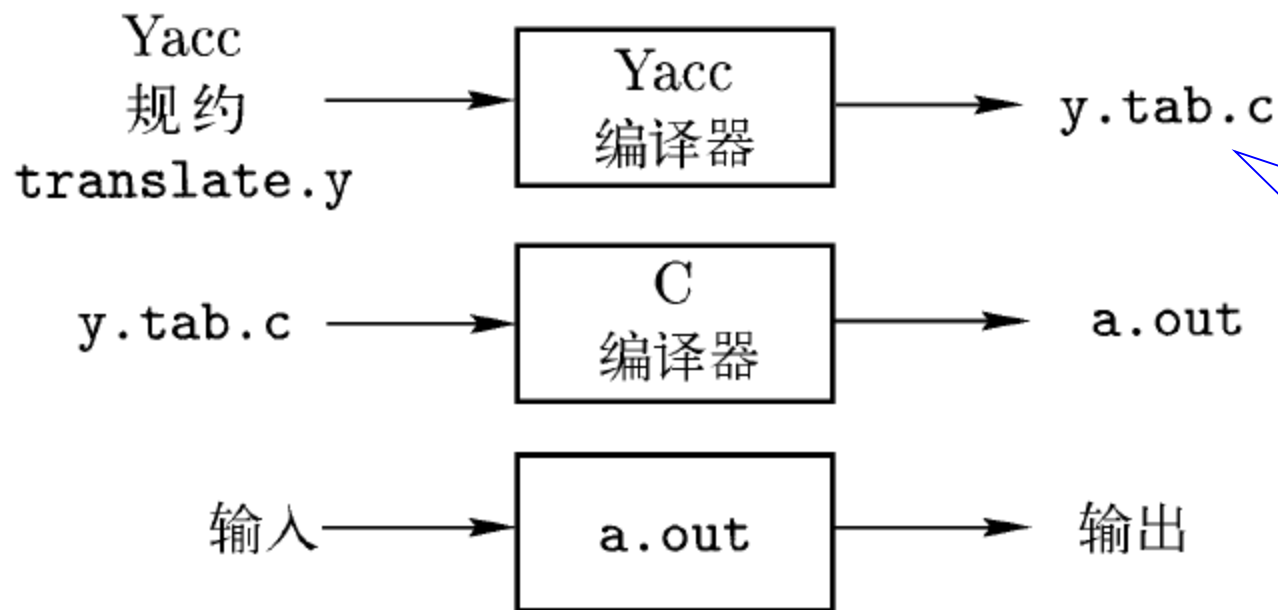
- 检查LR分析表中的每个报错条目，根据语言的特性来确定程序员最可能犯了什么错误，然后构造适当的恢复程序



语法分析器生成工具YACC



■ YACC的使用方法如下：



C语言写的
LALR语法
分析器



YACC源程序的结构



■ 声明

- 分为可选的两节：第一节放置**C**声明，第二节是对词法单元的声明。

声明

%%

翻译规则

%%

辅助性**C**语言程序

■ 翻译规则：

- 指明产生式及相关的语义动作

■ 辅助性**C**语言例程

- 被直接拷贝到生成的**C**语言源程序中，
- 可以在语义动作中调用。
- 其中必须包括**yylex()**。
这个函数返回词法单元，
可以由**LEX**生成



翻译规则的格式



<产生式头> : <产生式体>1 {<语义动作>1}
 | <产生式体>2 {<语义动作>2}

 | <产生式体>n {<语义动作>n}
 ;

- 第一个产生式的头被看作开始符号;
- 语义动作是C语句序列;
- \$\$表示和产生式头相关的属性值, \$i表示产生式体中第i个文法符号的属性值。
- 当我们按照某个产生式归约时, 执行相应的语义动作。通常可以根据\$i来计算\$\$的值。
- 在YACC源程序中, 可以通过定义YYSTYPE来定义\$\$, \$i的类型。



YACC 源程序 的例子

```
%{  
#include <ctype.h>  
%}  
  
%token DIGIT  
  
%%  
line    : expr '\n'          { printf("%d\n", $1); }  
        ;  
expr    : expr '+' term      { $$ = $1 + $3; }  
        | term  
        ;  
term    : term '*' factor    { $$ = $1 * $3; }  
        | factor  
        ;  
factor  : '(' expr ')'       { $$ = $2; }  
        | DIGIT  
        ;  
  
%%  
yylex() {  
    int c;  
    c = getchar();  
    if (isdigit(c)) {  
        yylval = c-'0';  
        return DIGIT;  
    }  
    return c;  
}
```



YACC对于二义性文法的处理



- 缺省的处理方法
 - 对于归约/归约冲突，选择前面的产生式
 - 对于归约/移入冲突，总是移入（悬空-**else**的解决）
- 运行选项-**v**可以在文件**y.output**中看到冲突的描述及其解决方法；
- 可以通过一些命令来确定终结符号的优先级/结合性，解决移入/归约冲突。
 - 结合性：**%left** **%right** **%nonassoc**
 - 终结符号的优先级通过它们在声明部分的出现顺序而定。
 - 产生式的优先级设定为它的最右终结符号的优先级。也可以加标记**%prec<终结符号>**，指明产生式的优先级等同于该终结符号
 - 移入符号**a**/按照 **$A \rightarrow \alpha$** 归约：当 **$A \rightarrow \alpha$** 的优先级高于**a**，或者两者优先级相同但产生式是左结合时，选择归约，否则移入。



YACC的错误恢复



- 使用错误产生式的方式来完成语法错误恢复
 - 错误产生式 $A \rightarrow \mathbf{error} \alpha$
 - 例如: $stmt \rightarrow \mathbf{error} ;$
- 首先定义哪些“主要”非终结符号有相关的错误恢复动作;
 - 比如: 表达式, 语句, 块, 函数定义等对应的非终结符号
- 当语法分析器碰到错误时
 - 不断弹出栈中状态, 直到有一个状态包含 $A \rightarrow \mathbf{error} \alpha;$
 - 分析器将 **error** 移入栈中
 - 如果 α 为空, 分析器直接执行归约, 并调用相关的语义动作; 否则向前跳过一些符号, 直到找到可以归约为 α 的串



实例



- LR(0)项集族构造
- SLR分析表
- SLR分析
- LR(1)项集族构造
- LR(1)分析表
- LR(1)分析
- LALR项集族构造
- LALR分析表
- LALR分析

$$S' \rightarrow S$$

$$S \rightarrow C \ C$$

$$C \rightarrow c \ C \mid d$$



LR(0)

$$S' \rightarrow S$$

$$S \rightarrow C C$$

$$C \rightarrow c C \mid d$$

$I_0:$

$S' \rightarrow \cdot S$
 $S \rightarrow \cdot CC$
 $C \rightarrow \cdot cC$
 $C \rightarrow \cdot d$



$I_1:$

$S' \rightarrow S \cdot$



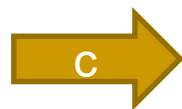
$I_2:$

$S \rightarrow C \cdot C$
 $C \rightarrow \cdot cC$
 $C \rightarrow \cdot d$



$I_5:$

$S \rightarrow CC \cdot$



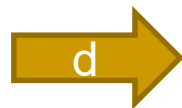
$I_3:$

$C \rightarrow c \cdot C$
 $C \rightarrow \cdot cC$
 $C \rightarrow \cdot d$



$I_6:$

$C \rightarrow cC \cdot$



$I_4:$

$C \rightarrow d \cdot$





SLR分析表



状态	Action			Goto	
	c	d	\$	S	C
0	s3	s4		1	2
1			accept		
2	s3				5
3	s3	s4			6
4	r	r	r		
5			r		
6	r	r	r		

$\text{Follow}(C) = \{c, d, \$\}$

$\text{Follow}(S) = \{\$\}$