



# 第 4 章 语法分析 (4)

---

Syntax Analysis (4)

【对应教材 4.5-4.6】



# 内容提要

---

- 语法分析简介
- 上下文无关文法
- 文法的设计方法
- 自顶向下的语法分析
- **自底向上的语法分析**
  - 简单LR分析: LR(0), SLR
  - 更强大的LR分析: LR(1), LALR
  - 二义性文法的使用



# SLR的原理：活前缀（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。
  - 直接把项当作状态，可以构造得到一个NFA/DFA。
- 活前缀 (Viable Prefix, 也译作可行前缀):
  - 可以出现在移进-归约语法分析器栈中的右句型前缀
  - 定义：某个右句型的前缀，且没有越过该句型的句柄的右端。
- 有效项：
  - 如果存在  $S \Rightarrow \alpha A w \Rightarrow \alpha \beta_1 \beta_2 w$ ，那么我们说项  $A \rightarrow \beta_1 \cdot \beta_2$  对  $\alpha \beta_1$  有效。

# SLR的原理：活前缀（2）

- 如果 $A \rightarrow \beta_1 \cdot \beta_2$ 对 $\alpha\beta_1$ 有效，
  - 当 $\beta_2$ 不等于空，表示句柄尚未出现在栈中，应该移进；
  - 如果 $\beta_2$ 等于空，表示句柄出现在栈中，应该归约。
  
- 如果某个时刻存在两个有效项要求执行不同的动作，那么就存在冲突。
  - 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A后面跟着的下一个输入符号在 $FOLLOW(A)$ 中时才可以归约。



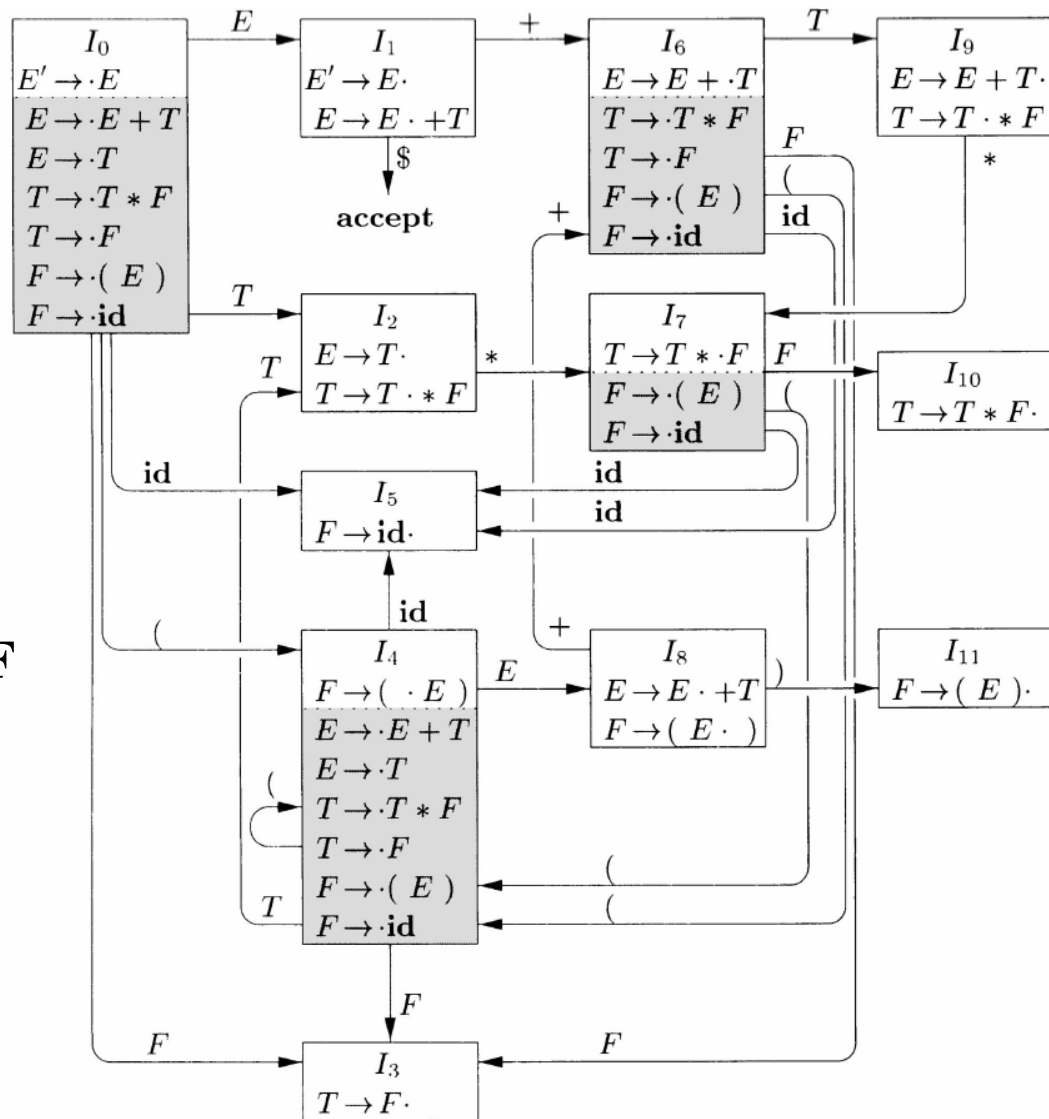
# SLR的原理：活前缀（3）

- 如果  $A \rightarrow \beta_1 \cdot \beta_2$  对  $\alpha\beta_1$  有效，
  - $A ? \alpha\beta_1 ?$
  
- 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  $\gamma$  的路径到达一个状态，那么
  - $\gamma$  是活前缀
  - 这个状态对应的项集就是  $\gamma$  的有效项的集合。



# 活前缀/有效项的例子

- 活前缀  $E+T^*$
- 对应的LR(0)项
  - $T \rightarrow T^* \cdot F$
  - $F \rightarrow \cdot (E)$
  - $F \rightarrow \cdot id$
- 对应的最右推导
  - $E' \Rightarrow E \Rightarrow E+T \Rightarrow E+T^*F$
  - $E' \Rightarrow E \Rightarrow E+T \Rightarrow E+T^*F \Rightarrow E+T^*(E)$
  - $E' \Rightarrow E \Rightarrow E+T \Rightarrow E+T^*F \Rightarrow E+T^*id$





# 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 $[A \rightarrow \alpha \bullet]$ 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 $a$ 还有其它的移进/归约操作, 则出现冲突
- 假设此时栈中的符号串为 $\beta\alpha$ ,
  - 如果 $\beta A a$ 不可能是最右句型的前缀, 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

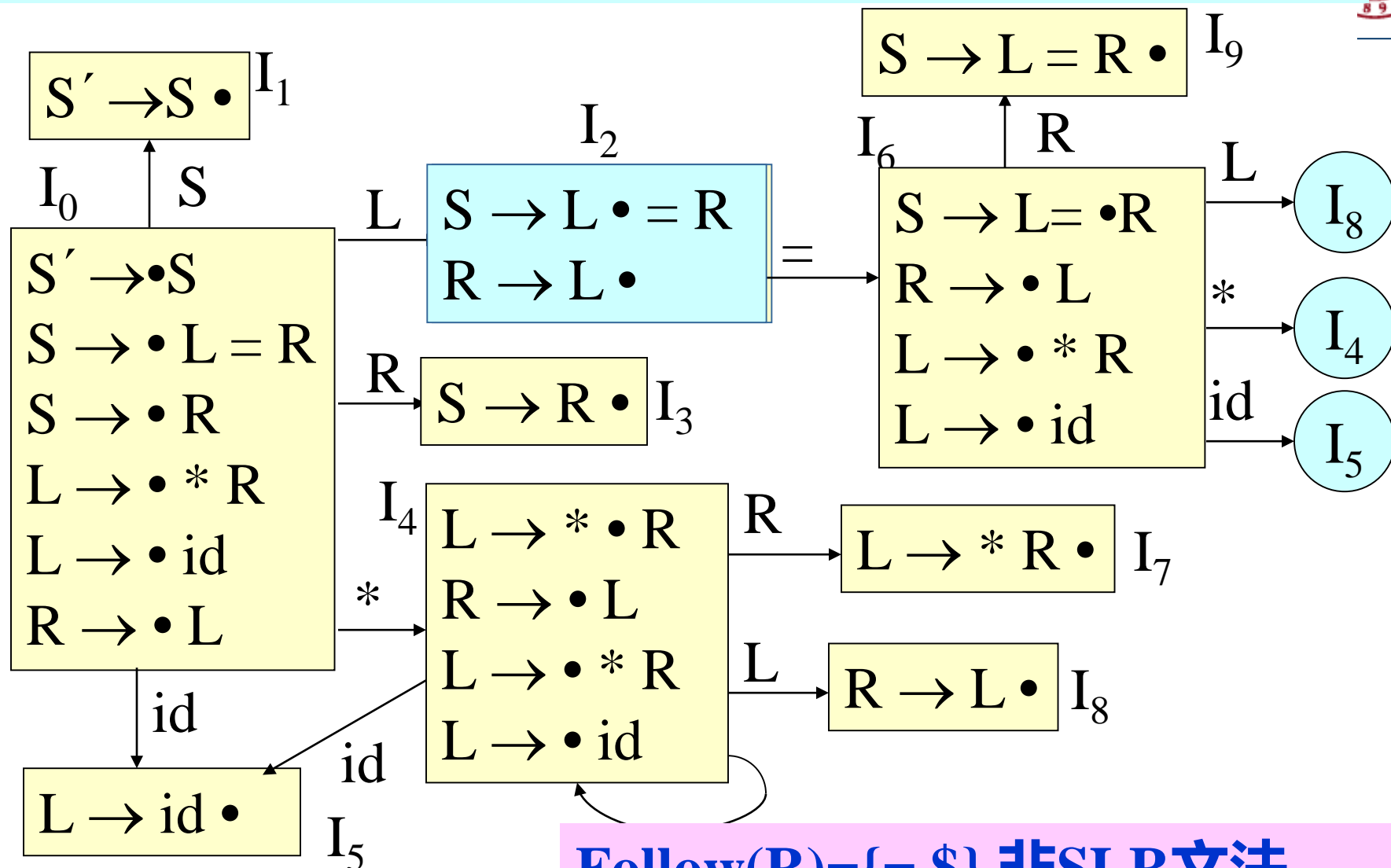


# SLR语法分析器的弱点（2）

- $[A \rightarrow \alpha \bullet]$  出现在项集中的条件：
  - 首先  $[A \rightarrow \bullet \alpha]$  出现在某个项集中，然后逐步读入/归约到  $\alpha$  中的符号，点不断后移，到达末端。
  - 而  $[A \rightarrow \bullet \alpha]$  出现的条件是  $B \rightarrow \beta \bullet A \gamma$  出现在项中
    - 先按照  $A \rightarrow \alpha$  归约，然后将  $B \rightarrow \beta \bullet A \gamma$  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  $\gamma$  的第一个符号。
  - 但是从 LR(0) 项集中不能确定这个信息。



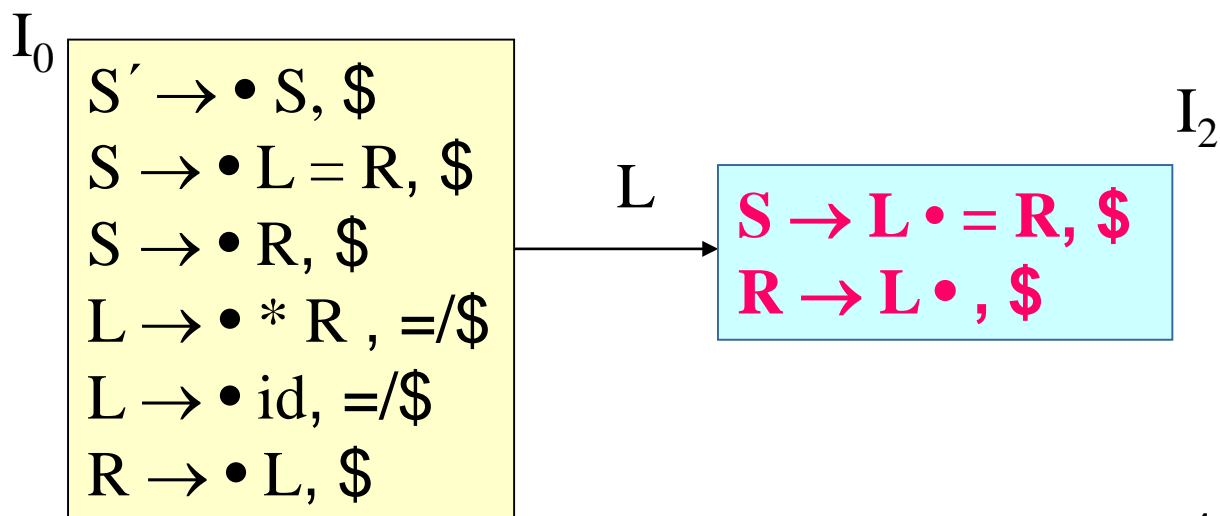
**非SLR文法:** (0)  $S' \rightarrow S$  (1)  $S \rightarrow L = R$  (2)  $S \rightarrow R$   
 (3)  $L \rightarrow^* R$  (4)  $L \rightarrow id$  (5)  $R \rightarrow L$



**Follow(R)={=,\$},非SLR文法。**

# LR(1)分析

- **问题：**栈中已识别出一个L，此时，a是“=”，不能用“ $R \rightarrow L$ ”进行归约，因为不存在“ $R =$ ”的规范句型。
- **解决方法：**添加搜索符号串。
  - 使每个项中含有更多的展望信息，使得能确切知道何时能进行归约。



# LR(k)项

## □ LR(k) 项

形式:  $[A \rightarrow \alpha \bullet \beta, a_1 a_2 \dots a_k]$

- $\beta \neq \epsilon$  时, 是移进项,  $a_1 a_2 \dots a_k$  不起作用。
- 对归约项  $[A \rightarrow \alpha \bullet, a_1 a_2 \dots a_k]$ , 仅当前输入符号串开始的前  $k$  个符号是  $a_1 a_2 \dots a_k$  时, 才能用  $A \rightarrow \alpha$  进行归约。
- $a_1 a_2 \dots a_k$  称为向前搜索符号串。



# LR(1)分析表的构造

## □ LR(1)有效项:

如果存在一个规范推导

$$S \Rightarrow_{rm}^* \delta A w \Rightarrow_{rm} \delta \alpha \beta w$$

则称一个LR(1)项  $[A \rightarrow \alpha \bullet \beta, a]$  对活前缀  $\gamma = \delta \alpha$  是有效的。

- 其中  $w$  的第一个符号为  $a$ , 或  $w = \epsilon$  而  $a$  为  $\$$ 。

例: 考虑文法  $G$ :  $S \rightarrow CC$      $C \rightarrow cC \mid d$   
因为有规范推导

$$S \Rightarrow_{rm}^* ccCcd \Rightarrow_{rm} cccCcd$$

故项  $[C \rightarrow c \bullet C, c]$  对活前缀  $ccc$  是有效的。



# 如何推导LR(1)的有效项

- 若项  $[A \rightarrow \alpha \bullet B \beta, a]$  对活前缀  $\delta \alpha$  是有效的, 则存在一个规范推导

$$S \Rightarrow^*_{rm} \delta A a x \Rightarrow_{rm} \delta \alpha B \beta a x$$

- 假定  $\beta a x \Rightarrow^*_{rm} b y$ , 则对每一个形如  $B \rightarrow \gamma$  的产生式, 我们有规范推导

$$S \Rightarrow^*_{rm} \delta \alpha B b y \Rightarrow_{rm} \delta \alpha \gamma b y$$

- 从而项  $[B \rightarrow \bullet \gamma, b]$  对于活前缀  $\delta \alpha$  也是有效的。
- 注意到  $b$  或者是从  $\beta$  推出的第一个终结符号, 或者  $\beta \Rightarrow^*_{rm} \epsilon$  而  $b = a$ 。这两种可能性结合在一起, 则  $b \in \text{first}(\beta a)$ 。

# LR(1)项集的构造

□ 设 $I$ 是 $G$ 的一个LR(1)项集， $\text{closure}(I)$ 是从 $I$ 出发用下面三个规则构造的项集：

1. 每一个 $I$ 中的项都属于 $\text{closure}(I)$ 。
2. 若项  $[A \rightarrow \alpha \bullet B \beta, a]$  属于 $\text{closure}(I)$  且  $B \rightarrow \gamma \in P$ , 则对任何  $b \in \text{FIRST}(\beta a)$ , 把  $[B \rightarrow \bullet \gamma, b]$  加到  $\text{closure}(I)$  中。
3. 重复执行(2)直到 $\text{closure}(I)$ 不再增大为止。

□ 设 $I$ 是 $G$ 的一个LR(1)项集， $X$ 是一个文法符号，定义

$$\text{GOTO}(I, X) = \text{closure}(J),$$

其中  $J = \{[A \rightarrow \alpha X \bullet \beta, a] \mid \text{当 } [A \rightarrow \alpha \bullet X \beta, a] \in I \text{ 时}\}$

# LR(1)项集的CLOSURE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  $I$ ;  
}
```

注意添加 $[B \rightarrow \cdot \gamma, b]$ 时， $b$ 的取值范围。



# LR(1)项集的GOTO算法

```
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( $J$ );  
}
```

和LR(0)项集的GOTO算法基本相同



# LR(1)项集族的主构造算法

```
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;  
}
```

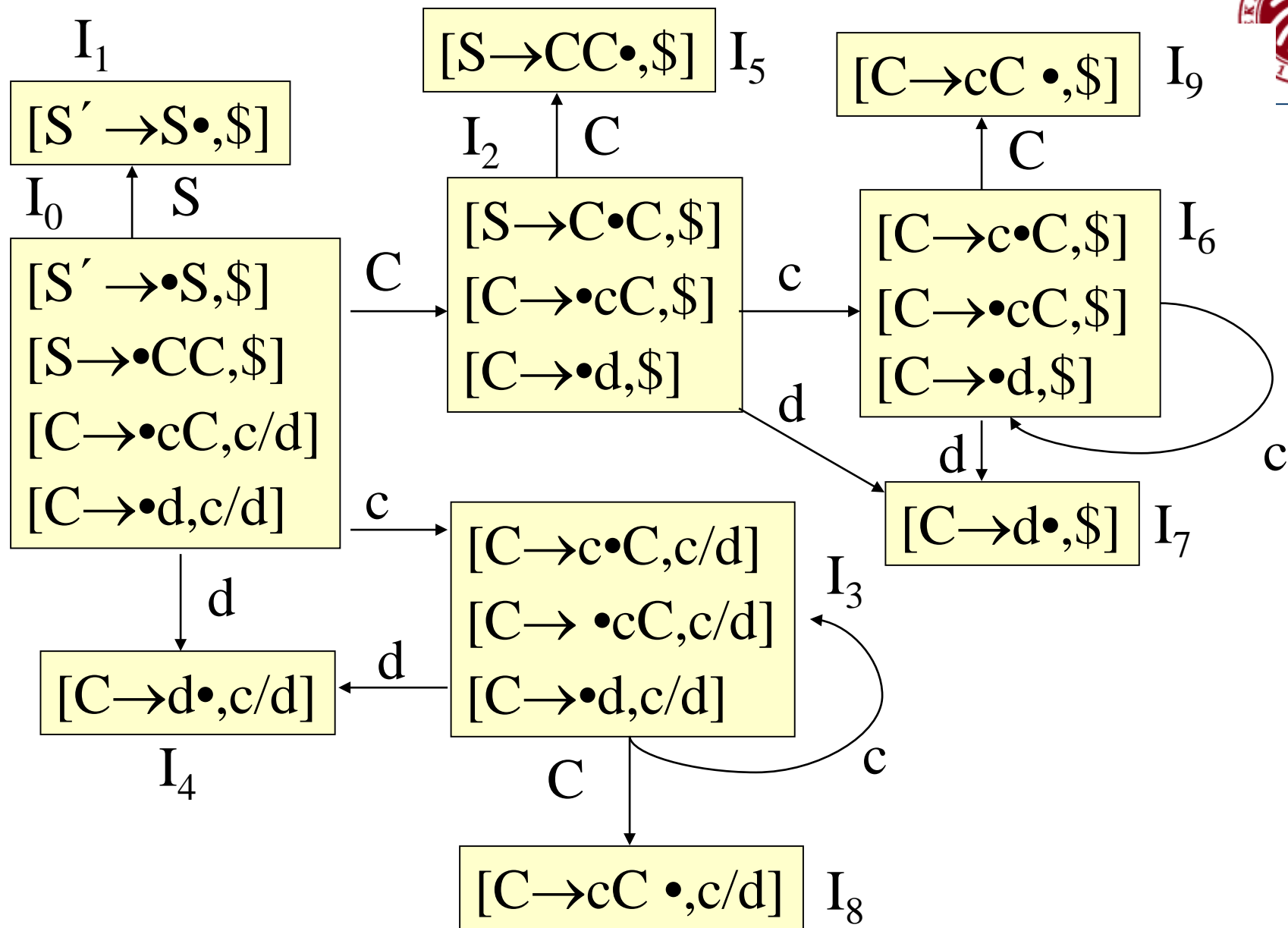
和LR(0)项集族的构造算法相同

# LR(1)项集及规范族的构造过程

例4\_9  $G(S')$ : (0)  $S' \rightarrow S$  (1)  $S \rightarrow CC$   
(2)  $C \rightarrow cC$  (3)  $C \rightarrow d$

若  $[A \rightarrow \alpha \bullet B \beta, a] \in I$   
且  $B \rightarrow \gamma \in P$   
则  $[B \rightarrow \bullet \gamma, b] \in I$   
其中  $b \in \text{FIRST}(\beta a)$

$I_0$   
 $[S' \rightarrow \bullet S, \$]$   
 $[S \rightarrow \bullet CC, \$]$   
 $[C \rightarrow \bullet cC, c/d]$   
 $[C \rightarrow \bullet d, c/d]$



# 构造LR(1)分析表

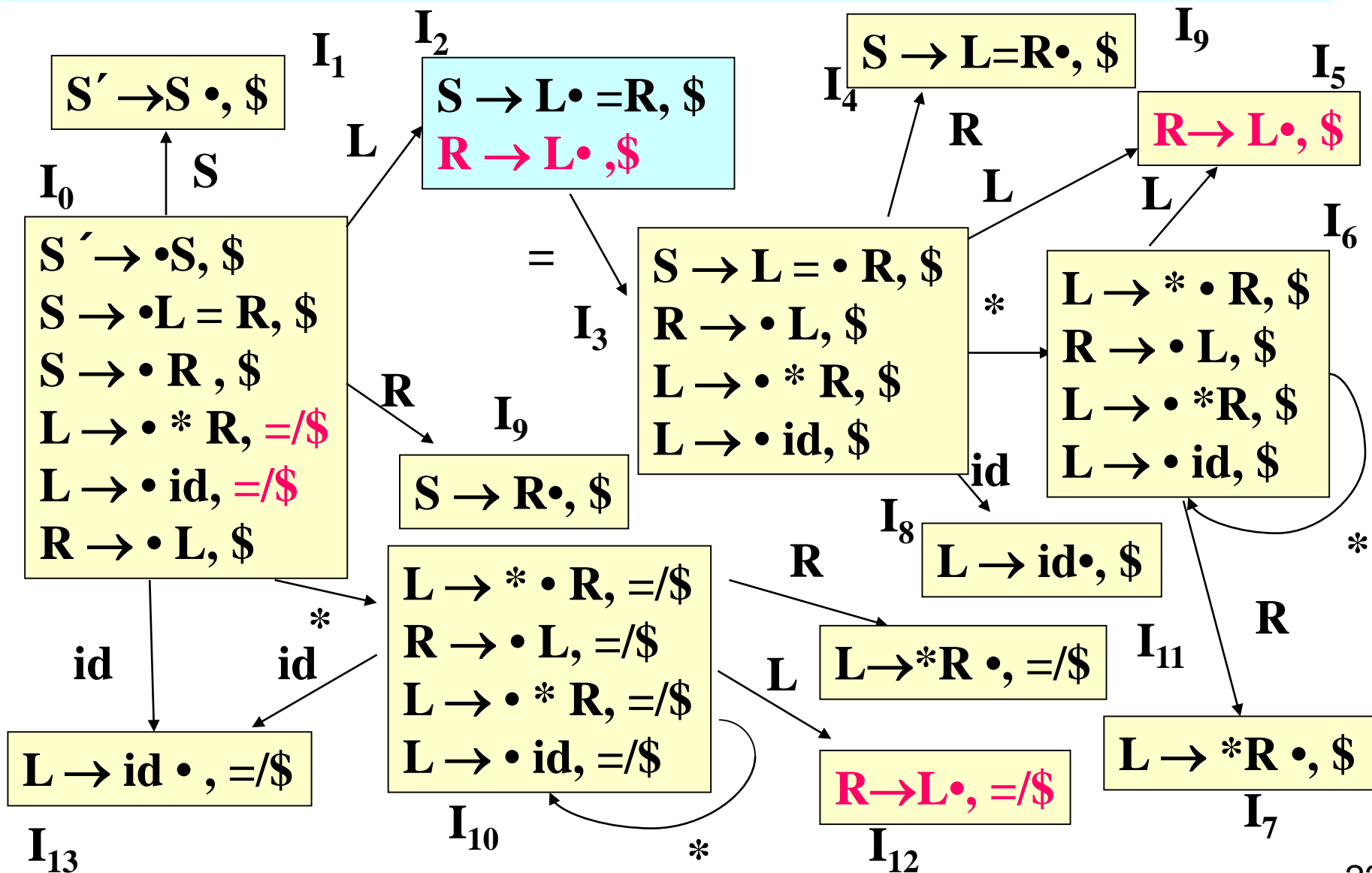
- 根据DFA构造LR(1)分析表M:
  - DFA中的每个状态对应分析表中的一行
  - 对于DFA中的每一个从状态*i*到状态*j*的转移
    - 如果转移符号为终结符*a*: 在相应的表项  $M[i, a]$  中填写移进动作  $S_j$
    - 如果转移符号为非终结符*A*: 在相应的表项  $M[i, A]$  中填写要转移到的状态  $j$
  - 对于包含归约项  $[A \rightarrow \alpha \bullet, a]$  的状态  $i$ 
    - 在相应的表项  $M[i, a]$  中填写归约动作  $r_k$ , 其中  $k$  是产生式  $A \rightarrow \alpha$  的编号。

如果按照上面的步骤填写的分析表中没有冲突（即每个单元格中只包含一个动作），那么得到的就是合法的LR(1)分析表

# 例：LR(1)分析表

	c	d	\$	S	C
0	S <sub>3</sub>	S <sub>4</sub>		1	2
1			acc		
2	S <sub>6</sub>	S <sub>7</sub>			5
3	S <sub>3</sub>	S <sub>4</sub>			8
4	r <sub>3</sub>	r <sub>3</sub>			
5			r <sub>1</sub>		
6	S <sub>6</sub>	S <sub>7</sub>			9
7			r <sub>3</sub>		
8	r <sub>2</sub>	r <sub>2</sub>			
9			r <sub>2</sub>		

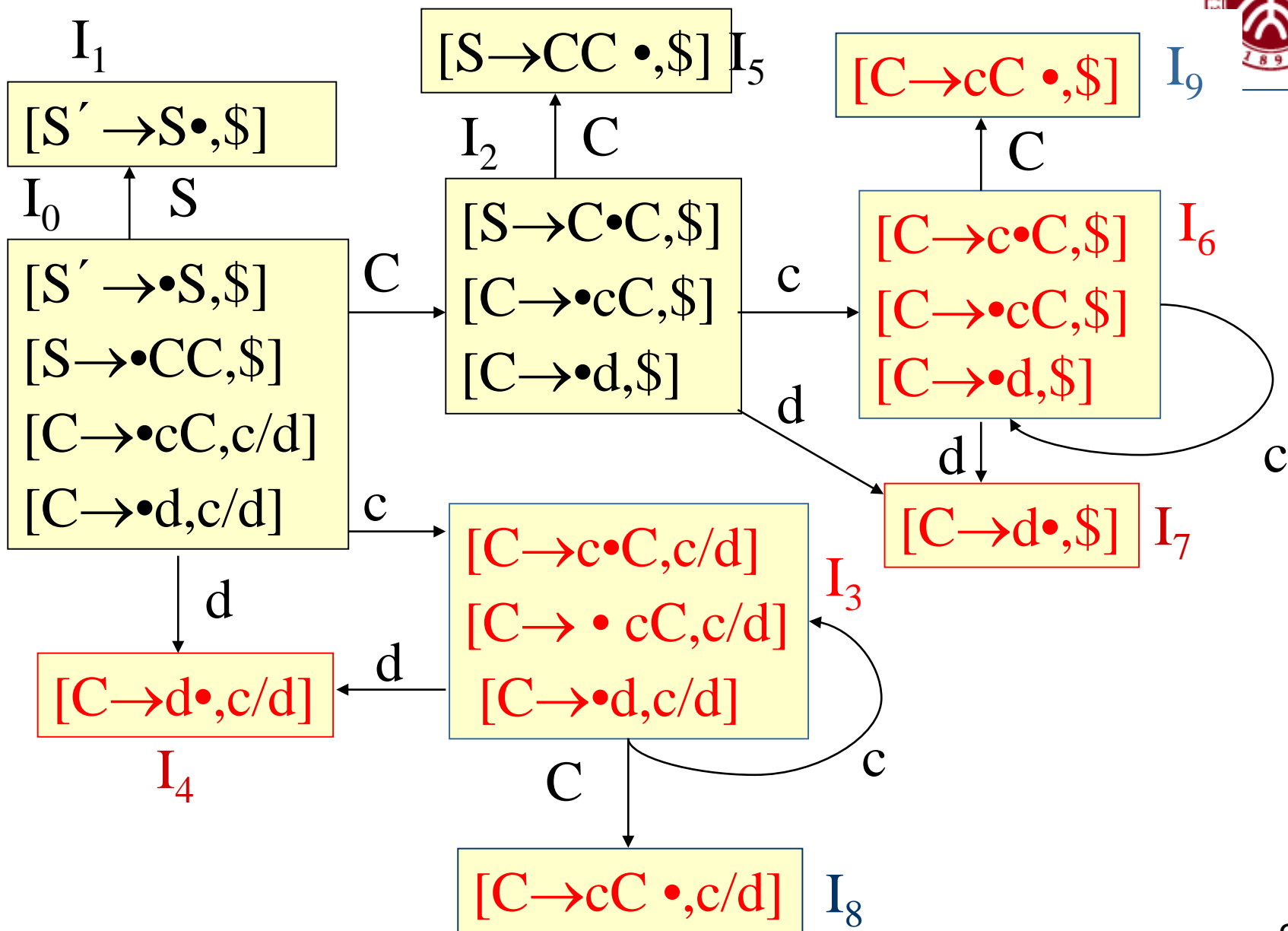
回顾: (0)  $S' \rightarrow S$  (1)  $S \rightarrow L = R$  (2)  $S \rightarrow R$   
 (3)  $L \rightarrow^* R$  (4)  $L \rightarrow id$  (5)  $R \rightarrow L$



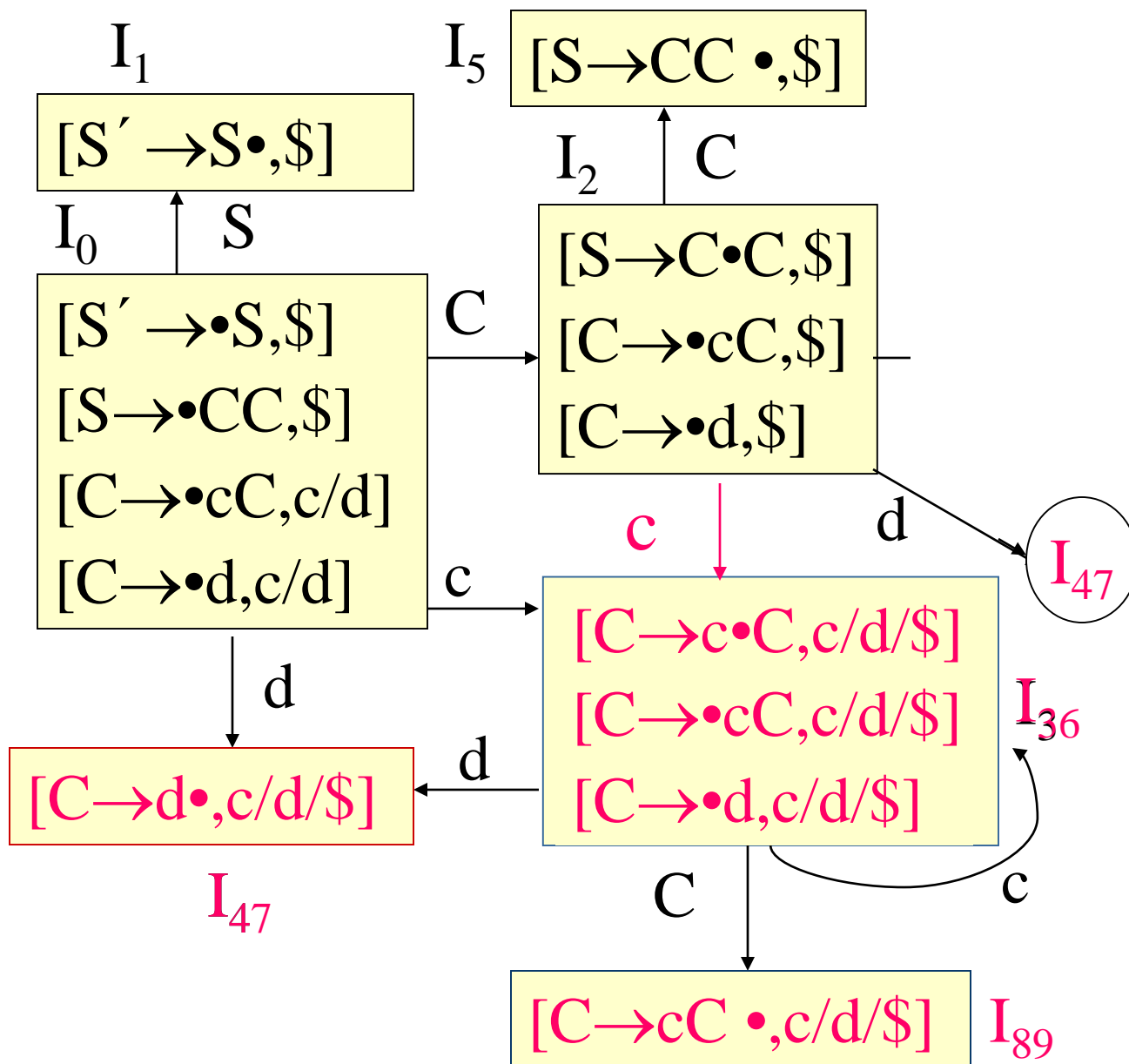


# LALR文法

- LR(1)分析表常常会含有比较多的状态，在实际使用中不是很常用。
- **同心集**：如果两个LR(1)项集去掉搜索符之后是相同的，则称这两个项集具有相同的核心(core)。
- **合并同心集**（合并搜索字符串）后，像构造LR(1)分析表一样构造出的LR分析表称作是**LALR(1)分析表**。







# 对应的LALR 分析表

	c	d	\$	S	C
0	S <sub>3</sub>	S <sub>4</sub>		1	2
1			acc		
2	S <sub>6</sub>	S <sub>7</sub>			5
3(6)	S <sub>3</sub>	S <sub>4</sub>			8
4(7)	r <sub>3</sub>	r <sub>3</sub>	r <sub>3</sub>		
5			r <sub>1</sub>		
6	S <sub>6</sub>	S <sub>7</sub>			9
7			r <sub>3</sub>		
8(9)	r <sub>2</sub>	r <sub>2</sub>	r <sub>2</sub>		
9			r <sub>2</sub>		

# LALR分析表的高效构造算法

- LALR的项集可以看作是在LR(0)项集上添加了适当的向前看符号。
- 基本方法
  - 只使用内核项来表示LR(0)或LR(1)项集。只使用 $[S' \rightarrow \bullet S]$ 或 $[S' \rightarrow \bullet S, \$]$ ，以及点不在最左端的项
  - 使用传播和自发生成的过程来生成向前看符号，得到LALR(1)内核项
  - 使用CLOSURE函数，求出内核的闭包。然后得到LALR分析表

详细算法参见龙书第4.7.5节



# LALR(1)的讨论

1. goto(I, X) 仅仅依赖于I的核心，因此 LR(1)项集合并后的转换函数 goto(I, X) 随自身的合并而得到。
2. 动作action应当进行修改，反映各被合并集合的既定动作。
3. 把同心的项集合并的时候，有可能导致冲突，这种冲突不会是移进-归约冲突；但可能引起归约-归约冲突。

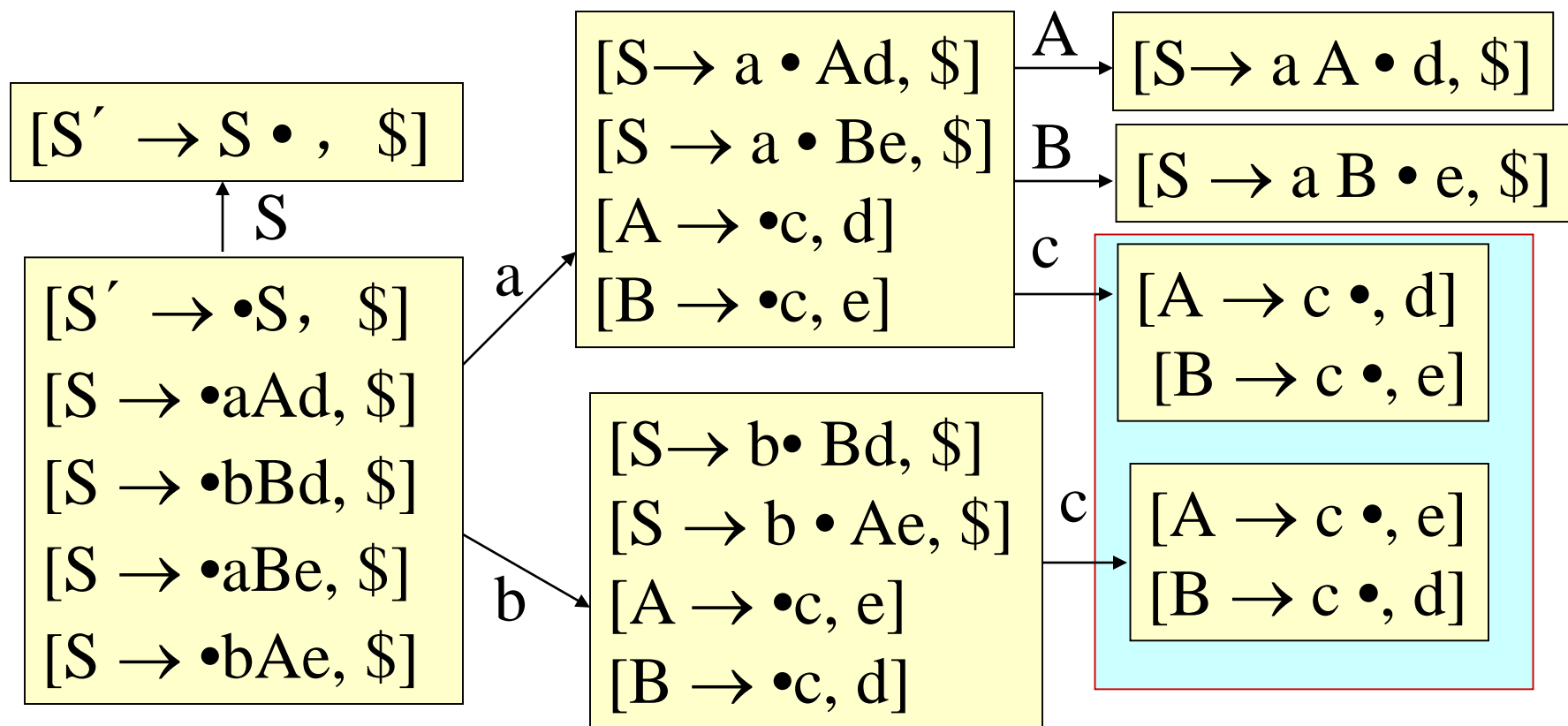
$$I_k: \{ [A \rightarrow \alpha \bullet, u_1] \quad [B \rightarrow \beta \bullet a \gamma, b] \} \quad a \cap u_1 = \emptyset$$

$$I_j: \{ [A \rightarrow \alpha \bullet, u_2] \quad [B \rightarrow \beta \bullet a \gamma, c] \} \quad a \cap u_2 = \emptyset$$

$$I_{kj}: \{ [A \rightarrow \alpha \bullet, u_1/u_2] \quad [B \rightarrow \beta \bullet a \gamma, b/c] \} \quad a \cap \{u_1, u_2\} = \emptyset$$

# 非LALR(1)的LR(1)文法

例:  $S' \rightarrow S$        $S \rightarrow aAd \mid bBd \mid aBe \mid bAe$   
 $A \rightarrow c$            $B \rightarrow c$





# LR(1)和LALR(1)在分析上的差别

例4\_9  $G(S')$ : (0)  $S' \rightarrow S$  (1)  $S \rightarrow CC$   
(2)  $C \rightarrow cC$  (3)  $C \rightarrow d$

输入: ccd\$

LR: ① c ③ c ③ d ④

报错

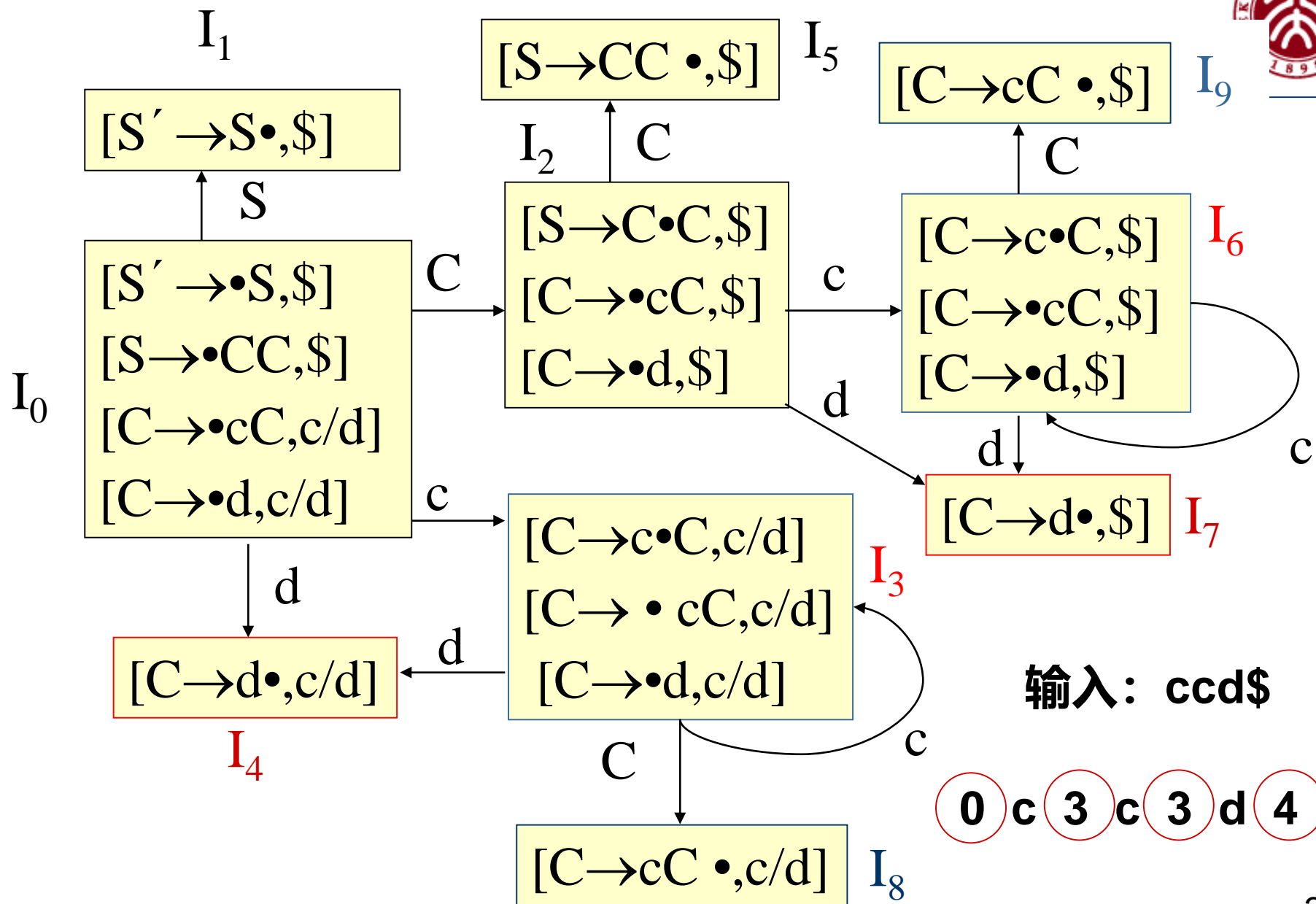
LALR: ① c ③⑥ c ③⑥ d ④⑦

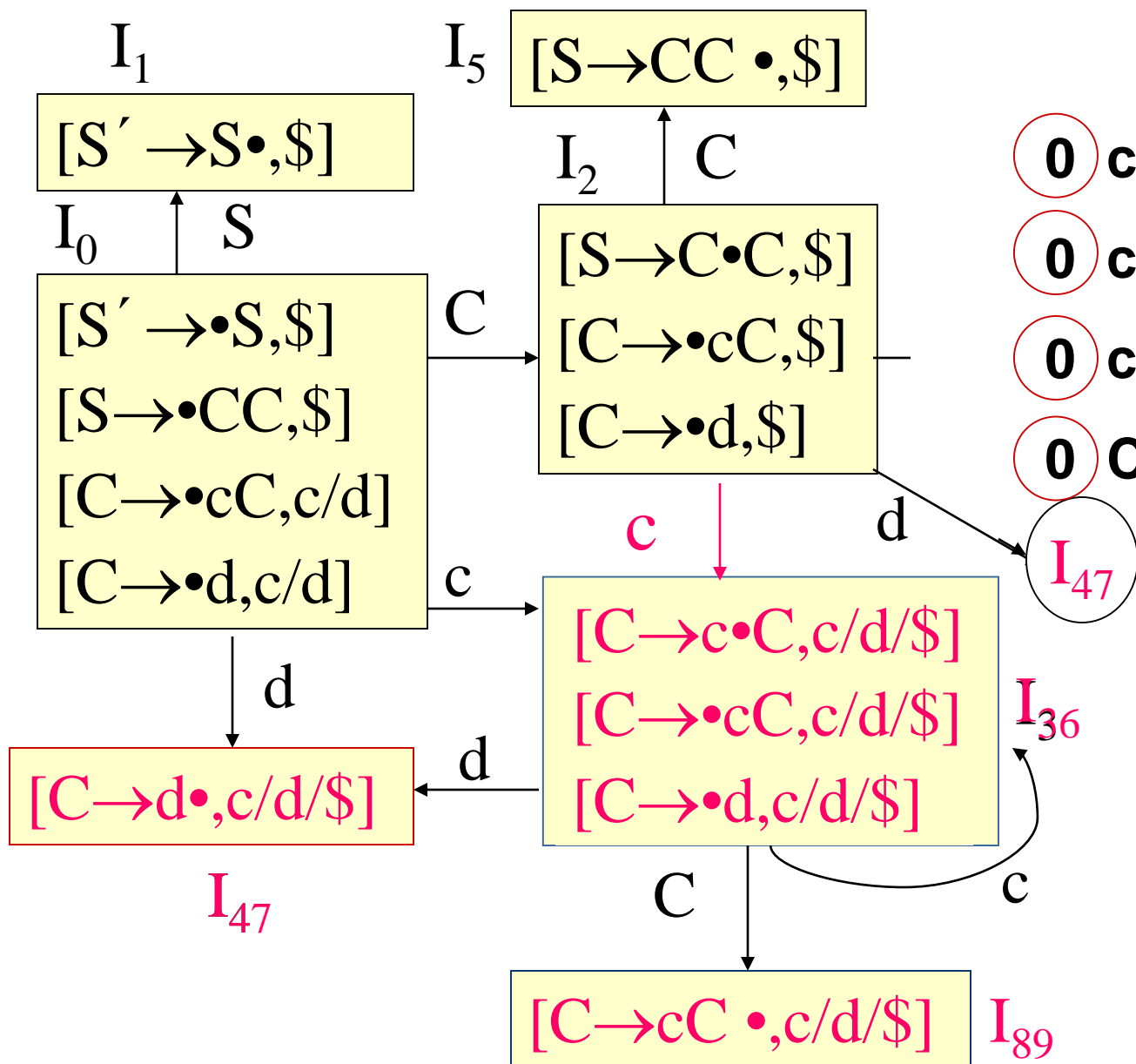
① c ③⑥ c ③⑥ C ⑧⑨

① c ③⑥ C ⑧⑨

① C ②

报错





输入: ccd\$

0	c	36	c	36	d	47
0	c	36	c	36	C	89
0	c	36	C	89		
0	C	2				





# 二义性文法的使用

- 二义性文法都不是LR的。
- 某些二义性文法是有用的
  - 可以简洁地描述某些结构；
  - 隔离某些语法结构，对其进行特殊处理。
- 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。



# 利用优先级/结合性消除冲突

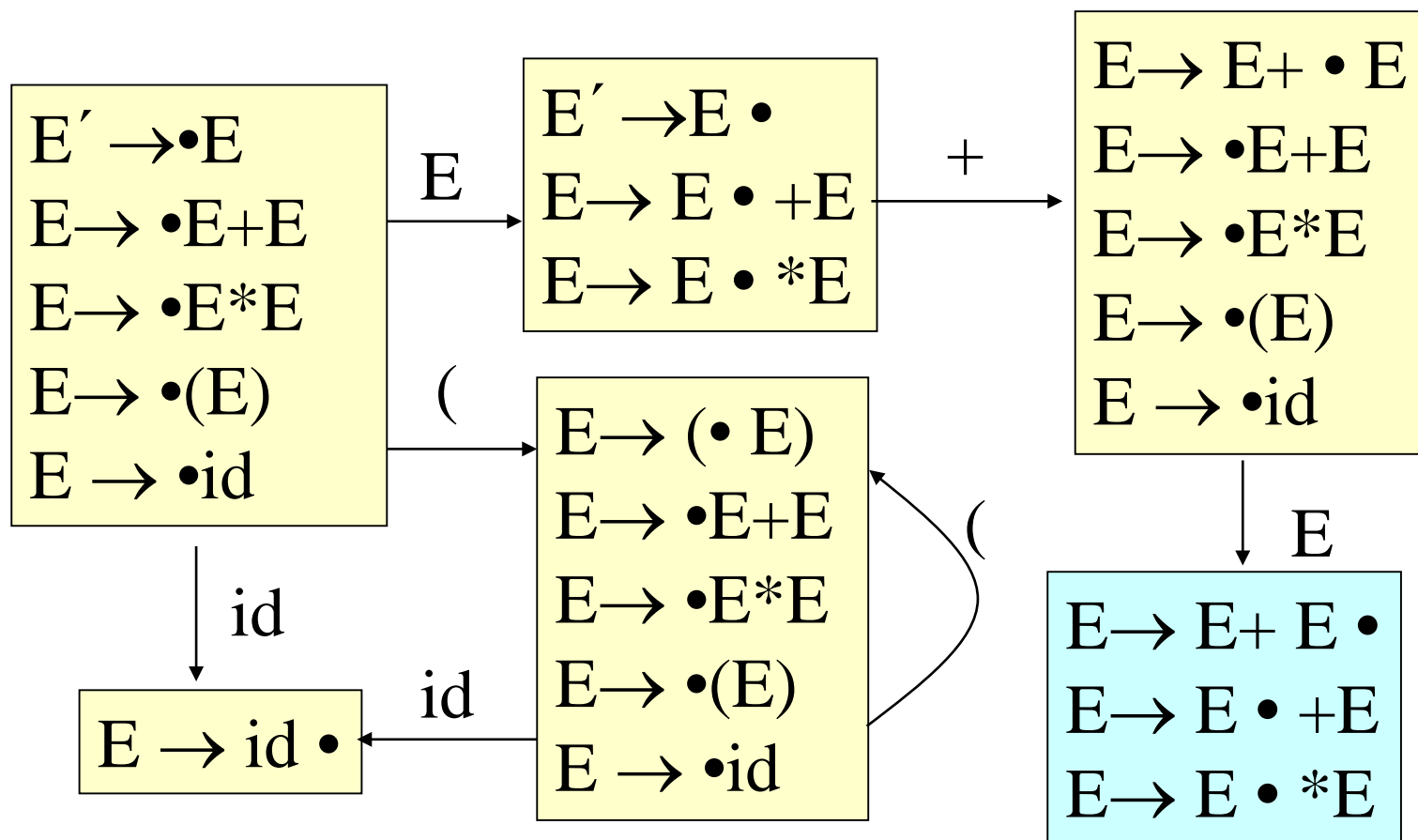
- 二义性文法:  $E \rightarrow E + E \mid E * E \mid (E) \mid id$
- 等价于:  $E \rightarrow E + T \mid T \quad T \rightarrow T * F \mid F$   
 $F \rightarrow (E) \mid id$
- 二义性文法的优点:
  - 容易修改算符的优先级和结合性。
  - 简洁: 如果有多个优先级, 那么无二义性文法讲引入太多的非终结符号
  - 高效: 不需要处理  $E \rightarrow T$  这样的归约。



# LR分析方法对二义文法的应用

龙书4.8节

表达式文法:  $E \rightarrow E+E \mid E * E \mid (E) \mid id$





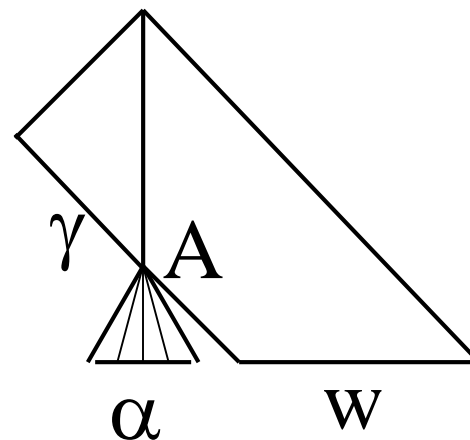
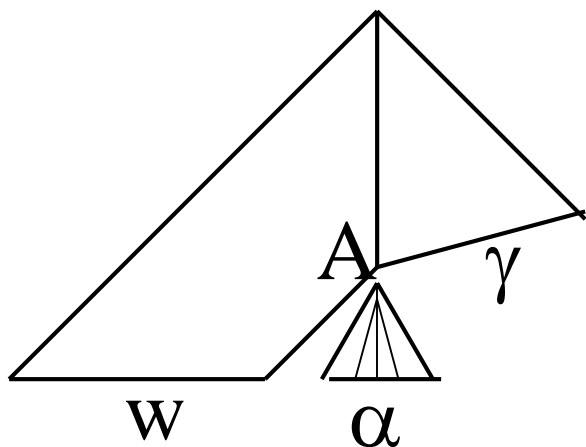
# LR分析总结：四种LR分析的对比

- 如果构造LR(0)的DFA
  - 没有归约冲突就是LR(0)文法
  - 有冲突但是可以通过Follow集合解决冲突就是SLR文法
  - 否则不是SLR文法。
- 如果构造LR(1)的DFA
  - 没有冲突就是LR(1)文法
  - 如果合并同心集之后也没有冲突，那么就是LALR(1)文法。
- 包含关系：LR(0) < SLR < LALR < LR(1)

# LR(k)和LL(k)的比较 (1)

1.  $A \rightarrow \alpha_1 \mid \alpha_2$

LL(k)根据 $\text{FIRST}(\alpha_i)$ 确定使用哪条产生式；而LR(k)是在识别出整个 $\alpha$ 后, 再往后看k个符号, 然后确定使用哪条产生式。





# LR(k)和LL(k)的比较 (2)

2. LL(k)文法都是LR(k)文法。都能用形式化方法实现。

3. 存在非LR的结构

例:  $L = \{ww^R \mid w \in \{a,b\}^*\}$

$G[S]: S \rightarrow aSa \mid bSb \mid \varepsilon$

4. LR(k)分析用手工构造是不可能的

- 类Pascal语言的LR(1)分析表, 需要数千个状态;
- 由于存在自动生成工具 (比如YACC), LR分析受到广泛重视。



# 本章小结

---

## □ 上下文无关文法

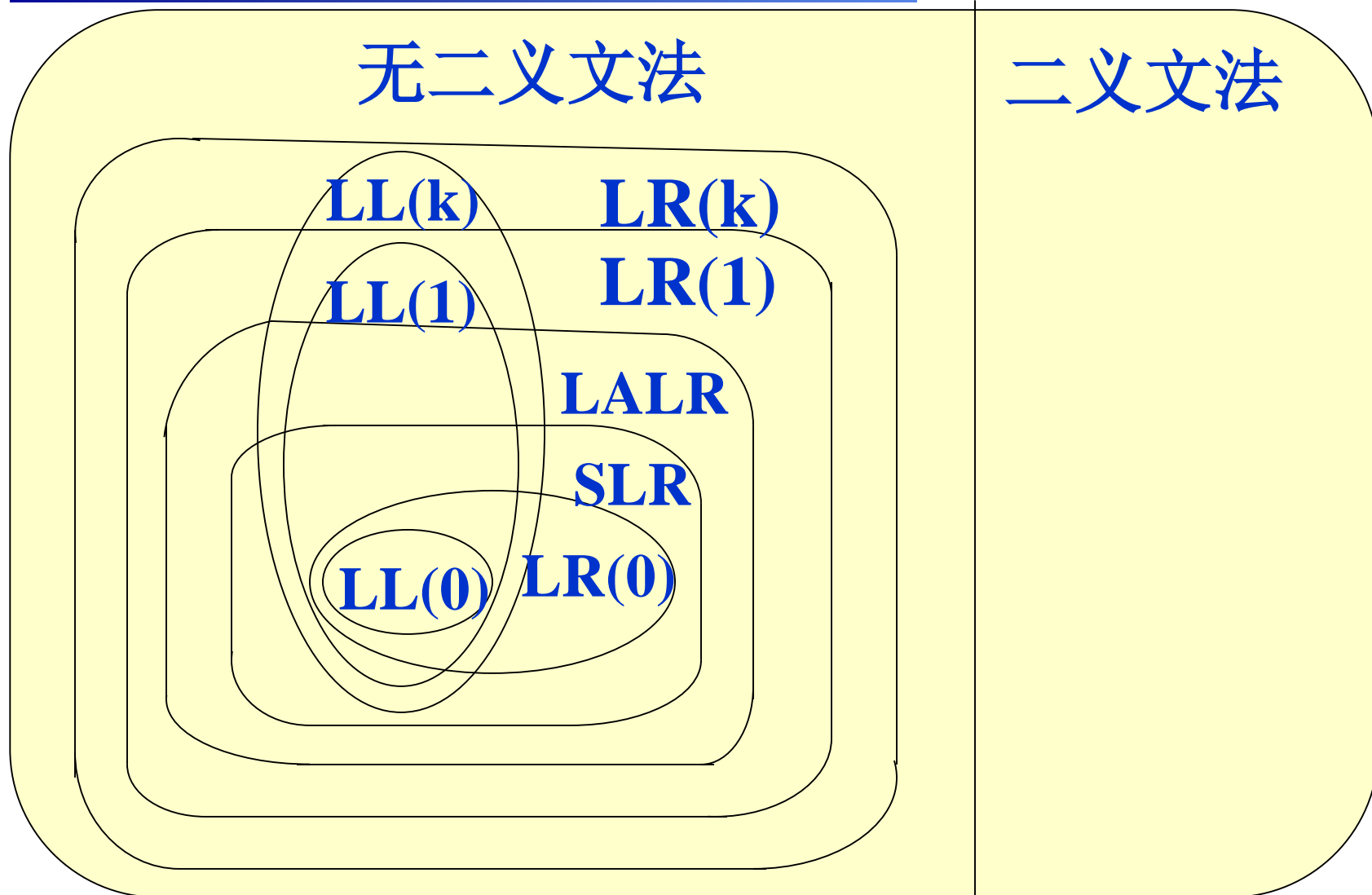
## □ 自顶向下分析

- 递归预测分析（递归子程序法）
- 非递归预测分析——LL(1)
- **注意：**首先消除左递归和提取左公因子。

## □ 自底向上分析

- LR分析: SLR(1), LR(1), LALR(1)

# 文法类的谱系







# 作业

---

- LR分析
  - 4.5.1, 4.5.2
- LR(0), SLR
  - 4.6.2, 4.6.3, 4.6.5
- LR(1), LALR(1)
  - 4.7.1, 4.7.4, 4.7.5