

## Lecture 3

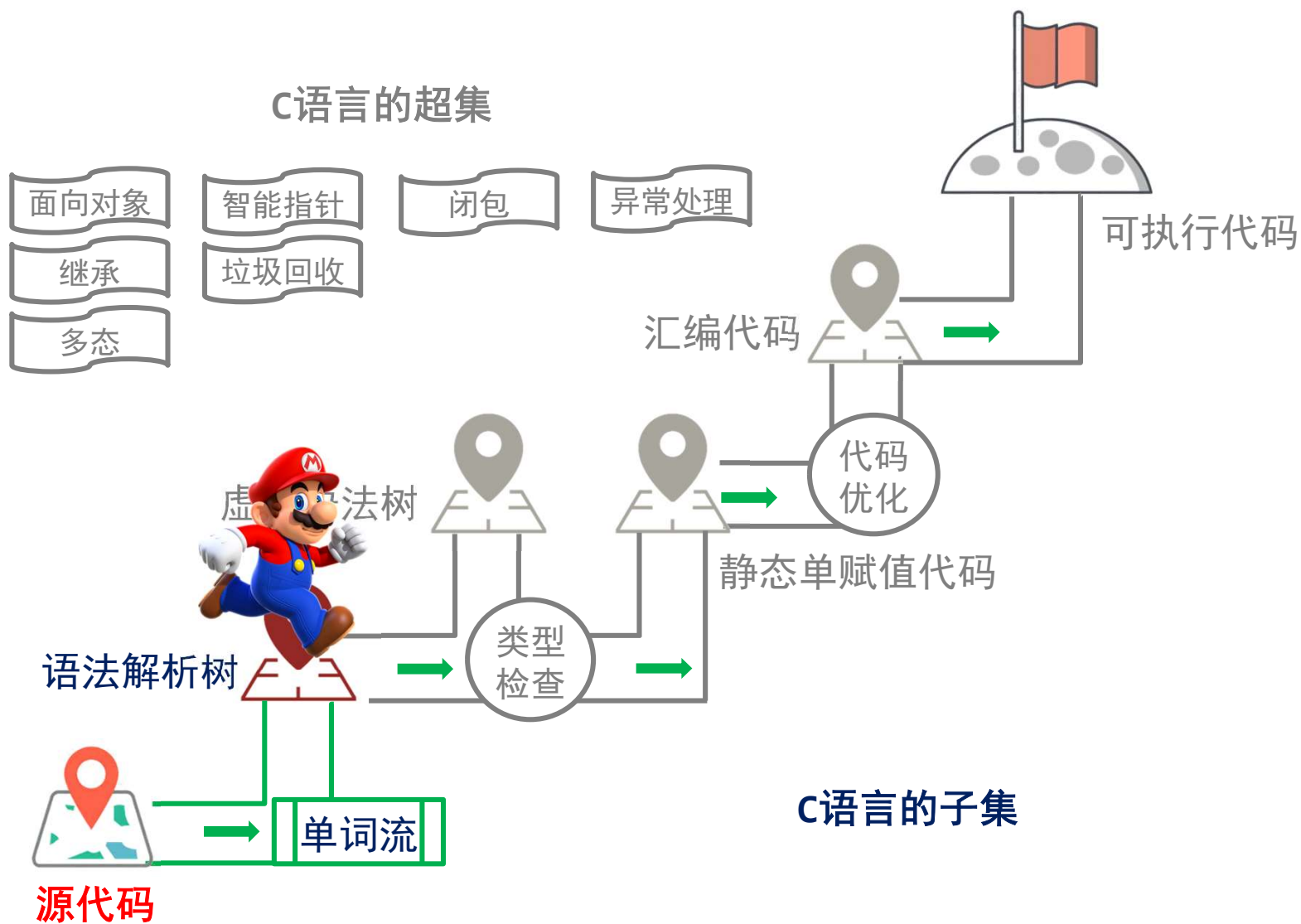
# 句式分析

徐 辉

xuh@fudan.edu.cn



# 学习地图



# 大纲

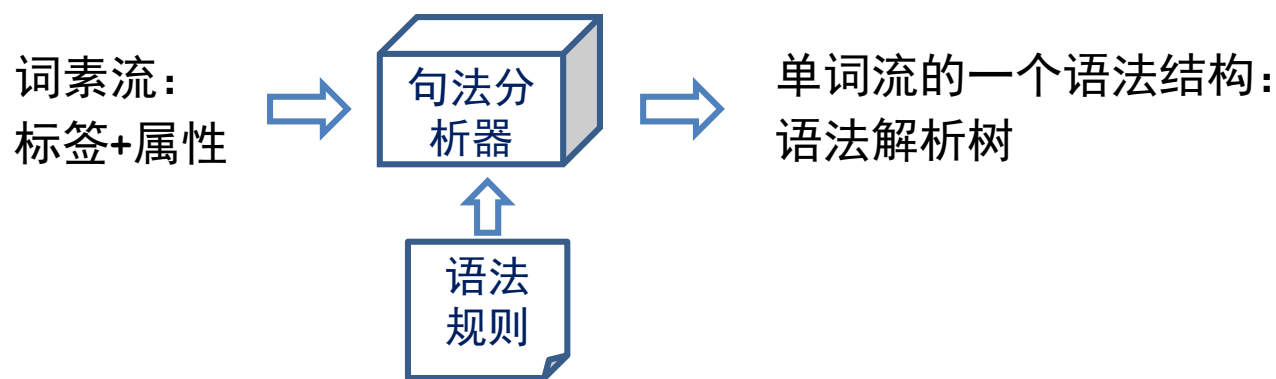
- 一、句式分析的基本概念
- 二、LLVM案例分析
- 三、自顶向下分析
- 四、自底向上分析
- 五、语法分析工具

# 一、句式分析的基本概念

---

# 问题定义

- 给定一个句子和语法规则，找到可生成该句子的一个语法推导。
- 通过词法分析已经将句子转换为了标签流。
- 语法规则（Grammar）定义了：
  - 什么是语法分析器（parser）可接受的标签组成，
  - 及其语法推导方式。



# 基本概念

- 一门语言 (language) 是多个句子 (sentences) 的集合。
- 句子 (sentence) 是由终结符 (terminal symbols) 组成的序列 (sequence)。
- 字符串 (string) 是包含终结符和非终结符的序列。
  - 字符串符号:  $\alpha, \beta, \gamma$
  - 非终结符:  $X, Y, Z$
  - 终结符 (标签):  $a, b, c$
- 一条语法 (grammar) 包括一个开始符号  $S$  和多条推导规则 (productions)
  - $X \rightarrow \beta$ 。

# 语法推导

- 语法 $G$ 的语言 $L(G)$ 是该语法可推导的所有句子的集合。
- 问题：下列语法是否可推导出句子 $aaabbbccc$ ?

语法规则

[1]  $S \rightarrow aBSc$

[2]  $S \rightarrow abc$

[3]  $Ba \rightarrow aB$

[4]  $Bb \rightarrow bb$

推导

[1]  $S \rightarrow aBSc$

[1]  $S \rightarrow aBaBScc$

[3]  $S \rightarrow aaBBScc$

[2]  $S \rightarrow aaBBabccc$

[3]  $S \rightarrow aaBaBbccc$

[3]  $S \rightarrow aaaBBbccc$

[4]  $S \rightarrow aaaBbbccc$

[4]  $S \rightarrow aaabbbccc$

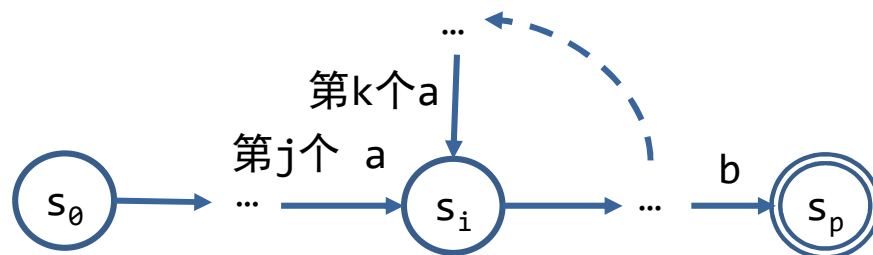
# 语法表示：使用正则表达式？

- 正则表达式是否可识别四则运算？
  - $y = a \times x + b$ 
    - $(var|num)((+|-|\times|\div)(var|num))^*$
  - $y = a \times (x + b)$ 
    - $('(|var|num)((+|-|\times|\div)('(|var|num|')'))^*$
    - 可导致单词流被错误接收：
      - $y = (a \times (x + b)$
      - $y = (a \times (x + (b$
- 正则表达式不能处理括号匹配问题： $(^*)^*$



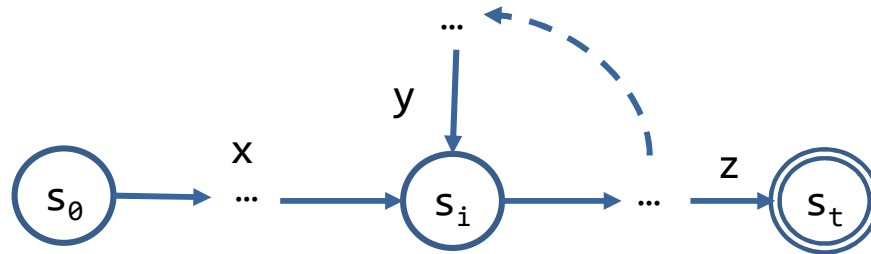
# 非正则语言

- 不能用正则表达式或有穷自动机表示的语言。
  - 正则语言不能计数，如  $L = \{a^n b^n, n > 0\}$
  - 证明：
    - 假设DFA可识别该语言，其包含  $p$  个状态；
    - 假设某词素为  $a^q b^q, q > p$ 。
    - 识别该词素需要经过某状态  $s_i$  至少两次，分别对应第  $j$  和第  $k$  个  $a$ ；
    - 该DFA可同时接受  $a^q b^q$  和  $a^{q+k-j} b^q$ ，推出矛盾。



# 正则语言的泵引理 (Pumping Lemma)

- 词素数量有限的语言一定是正则语言。
- 词素数量无穷多的语言是否为正则语言？
- 某语言 $L(r)$ 是正则语言的必要条件：
  - 任意长度超过 $p$ （泵长）的句子都可以被分解为 $xyz$ 的形式
  - 其中 $x$ 和 $z$ 可为空，
  - 子句 $y$ 被重复任意次（如 $xyyz$ ）后得到的句子仍属于该语言。

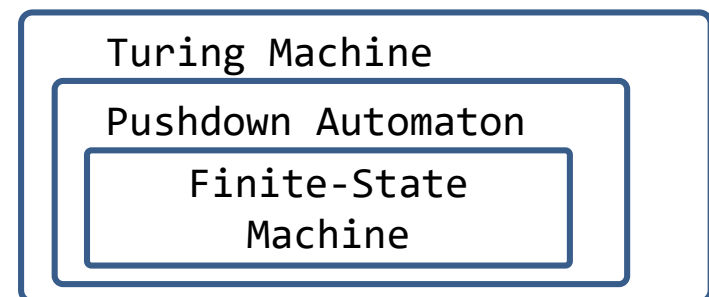


# 语言分析问题难度

- 通常来说，判断一个句子是否属于某个语言  $w \in L(G)$  是不能计算的。

Chomsky Hierarchy

Class	Languages	Automaton	Rules	Word Problem	Example
type-0	recursively enumerable	Turing machine	no restriction	undecidable	Post's corresp. problem
type-1	context sensitive	linear-bounded TM	$\alpha \rightarrow \gamma$ $ \alpha  \leq  \gamma $	PSPACE-complete	$a^n b^n c^n$
type-2	context free	pushdown automaton	$A \rightarrow \gamma$	cubic	$a^n b^n$
type-3	regular	NFA / DFA	$A \rightarrow a$ or $A \rightarrow aB$	linear time	$a^* b^*$



# 上线文无关语法和BNF范式

- 上下文无关语法（CFG/context-free grammar）是一个四元组 $(T, NT, S, P)$ 
  - $T$ : 终结符
  - $NT$ : 非终结符
  - $S$ : 起始符号
  - $P$ : 产生式规则集合 $X \rightarrow \gamma$ ,
    - $X$  是非终结符
    - $\gamma$  是可能包含终结符和非终结符的字符串
- BNF范式（Backus-Naur form）是传统的上下文无关语法表示方法。

$$\langle \textit{SheepNoise} \rangle ::= \textit{baa} \langle \textit{SheepNoise} \rangle$$
$$| \textit{baa}$$

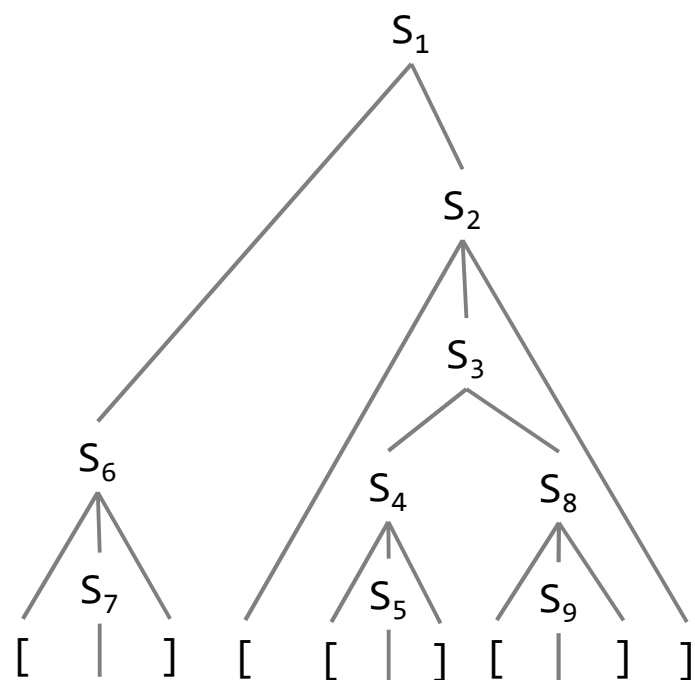
## 上线文无关语法举例

- 给定可生成所有匹配括号对的语法,  $[[[]][[]]]$  是该语法的一个推导吗?

## 语法规则

$$\begin{array}{lcl} [1] & S & \rightarrow \epsilon \\ [2] & & | [S] \\ [3] & & | SS \end{array}$$

## 推导

$$\begin{array}{l} [3] \quad S \rightarrow SS \\ [2] \quad S \rightarrow S[S] \\ [3] \quad S \rightarrow S[SS] \\ [2] \quad S \rightarrow S[S[S]] \\ [1] \quad S \rightarrow S[S[ ]] \\ [2] \quad S \rightarrow S[[S][ ]] \\ [1] \quad S \rightarrow S[[ ][ ]] \\ [2] \quad S \rightarrow [S][ ][ ] \\ [1] \quad S \rightarrow [ ][ ][ ] \end{array}$$


## 语法解析树

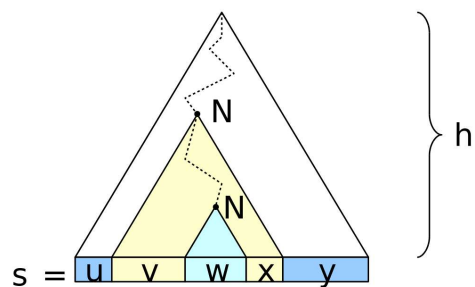
# 非CFG语言：上下文敏感语法

- $L = \{a^n b^n c^n, n > 0\}$  不是CFG语言
- 无法用CFG定义
  - $X \rightarrow \gamma$
- 可以用上下文敏感语法定义
  - $\alpha A \beta \rightarrow \alpha \gamma \beta$
  - $\alpha$ 和 $\beta$ 不变,  $A$ 展开

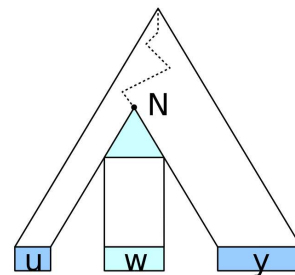
[1]	$S \rightarrow aBC$
[2]	$\quad \mid aSBC$
[3]	$CB \rightarrow CZ$
[4]	$CZ \rightarrow WZ$
[5]	$WZ \rightarrow WC$
[6]	$WC \rightarrow BC$
[7]	$aB \rightarrow ab$
[8]	$bB \rightarrow bb$
[9]	$bC \rightarrow bc$
[10]	$cC \rightarrow cc$

# 非CFG语言的泵引理

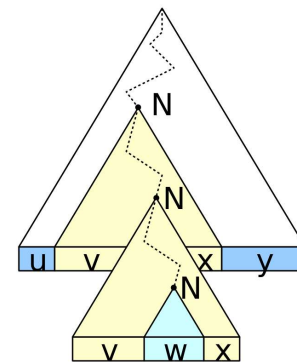
- CFG语言的泵引理（必要条件）：
  - 任意长度超过 $p$ （泵长）的句子可以被拆分为 $uvwxy$ ,
  - 子句 $v$ 和 $x$ 被重复任意次后得到的新句子（如 $uvvwxxy$ ）仍属于该语言。
- 正则属于CFG:  $uv^nwx^n\epsilon$



Sufficiently high  
derivation tree



Generating  $uv^0wx^0y$



Generating  $uv^2wx^2y$

# 澄清几个概念

- Regular language: 用DFA/NFA可以计算
  - Regular expression
    - 狭义: 其表示的都是正则语言, 所有正则语言都可以用正则表达式表示 (Flex工具)
    - 广义 (regex): 字符串匹配工具。
- Context-free language: 需要pushdown automaton计算
  - Context-free grammar
  - 正则语言可以用CFG表示:  $StartSymbol \rightarrow regex$ 
    - 特性: 右侧的非终结符均可替换为终结符
      - $S \rightarrow (0?1)^*$
      - $S \rightarrow 0S1S|1S0S|\epsilon$
- Context-sensitive language: 需要图灵机计算



# 下列语言是否为正则语言？

- 集合表示

- 1)  $L = \{a^n b^n\}$
- 2)  $L = \{a^n b^n | n \leq 100\}$
- 3)  $L = \{a^n | n \geq 1\}$
- 4)  $L = \{a^{2^n} | n \geq 1\}$
- 5)  $L = \{a^p | p \text{ is prime}\}$

- Regex/CFG语法表示

- 1)  $S \rightarrow (0?1)^*$
- 2)  $S \rightarrow aT | \epsilon, T \rightarrow Sb$
- 3)  $S \rightarrow 0S1S | 1S0S | \epsilon$
- 4)  $S \rightarrow A | B$ 
  - $A \rightarrow E1A'E$
  - $A' \rightarrow A | \epsilon$
  - $B \rightarrow E0B'E$
  - $B' \rightarrow B | \epsilon$
  - $E \rightarrow 0S1S | 1S0S | \epsilon$

# 推导 (Derivation) 的优先级

暂不考虑优先级

[1]	$Expr \rightarrow (Expr)$
[2]	$  Expr Op num$
[3]	$  num$
[4]	$Op \rightarrow +$
[5]	$  -$
[6]	$  \times$
[7]	$  \div$

$(a + b) \times 2$



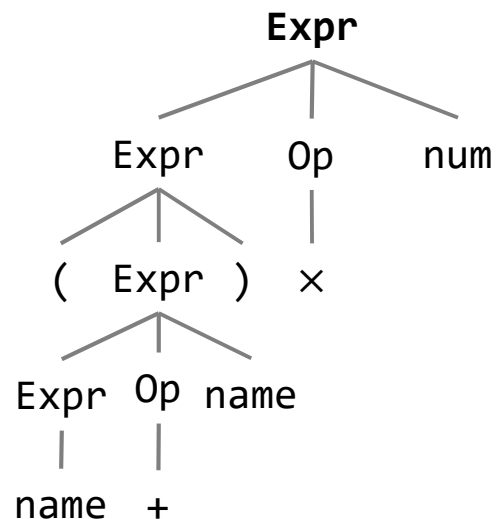
$(a + b) \times 2$

	$Expr$
[2]	$\Rightarrow_{lm} Expr Op num$
[1]	$\Rightarrow_{lm} (Expr) Op num$
[2]	$\Rightarrow_{lm} (Expr Op num) Op num$
[3]	$\Rightarrow_{lm} (num Op num) Op num$
[4]	$\Rightarrow_{lm} (num + num) Op num$
[6]	$\Rightarrow_{lm} (num + int) \times num$

左侧优先推导  
(Leftmost Derivation)

右侧优先推导  
(Rightmost Derivation)

	$Expr$
[2]	$\Rightarrow_{rm} Expr Op num$
[6]	$\Rightarrow_{rm} Expr \times num$
[1]	$\Rightarrow_{rm} (Expr) \times num$
[2]	$\Rightarrow_{rm} (Expr Op num) \times num$
[4]	$\Rightarrow_{rm} (Expr + num) \times num$
[3]	$\Rightarrow_{rm} (num + num) \times num$



语法解析树完全相同

# 练习：语法推导

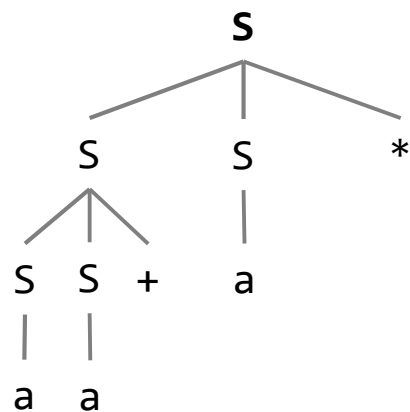
给定下列语法和字符串  $aa + a *$

- 1) 写出左推导
- 2) 写出右推导
- 3) 画出语法推导树

[1]	$S \rightarrow SS +$
[2]	$\quad \mid SS *$
[3]	$\quad \mid a$

$$S \xRightarrow{lm} SS * \xRightarrow{lm} SS + S * \xRightarrow{lm} aS + S * \xRightarrow{lm} aa + S * \xRightarrow{lm} aa + a *$$

$$S \xRightarrow{rm} SS * \xRightarrow{rm} Sa * \xRightarrow{rm} SS + a * \xRightarrow{rm} Sa + a * \xRightarrow{rm} aa + a *$$



## 练习：语法设计

- 为下列语言设计语法规则。
  - 1) 所有0和1组成的字符串，每一个0后面紧跟着若干个1
  - 2) 所有0和1组成的字符串，0和1的个数相同
  - 3) 所有0和1组成的字符串，0和1的个数不相同

$$S \rightarrow (0? 1)^*$$

$$S \rightarrow 0S1S|1S0S|\epsilon$$

$$S \rightarrow A|B$$

$$A \rightarrow E1A'E$$

$$A' \rightarrow A|\epsilon$$

$$B \rightarrow E0B'E$$

$$B' \rightarrow B|\epsilon$$

$$E \rightarrow 0S1S|1S0S|\epsilon$$

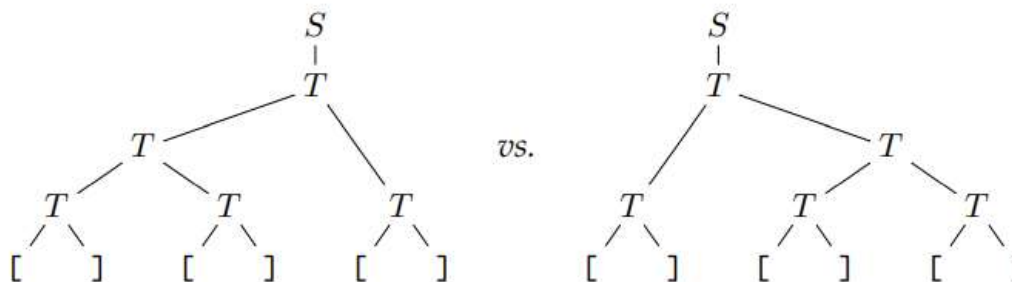
## 二义性 (ambiguity)

- 如果 $L(G)$ 中的某个句子有一个以上的最左（或最右）推导，那么语法 $G$ 就有二义性。
  - 语法解析树不同
- 根据下列语法规则如何推导出 $[][][]$ ?

[1]	$S \rightarrow \epsilon$
[2]	$\quad   T$
[3]	$T \rightarrow []$
[4]	$\quad   [T]$
[5]	$\quad   TT$

$S \rightarrow T \rightarrow TT \rightarrow TTT \rightarrow []TT \rightarrow T[][] \rightarrow [][][]$

$S \rightarrow T \rightarrow TT \rightarrow []T \rightarrow []TT \rightarrow [][]T \rightarrow [][][]$



# 极端情况

- 存在无数棵语法解析树
  - 考虑循环的情况
- 根据下列语法规则如何推导出 $[][[ ][]]$ ?

[1]	$S \rightarrow \epsilon$
[2]	$  [S]$
[3]	$  SS$

$S \rightarrow \textcolor{red}{S}S \rightarrow [\textcolor{red}{S}]S \rightarrow [ ]\textcolor{red}{S} \rightarrow [ ]\textcolor{red}{[S]} \rightarrow [ ][S S] \rightarrow [ ][[S] [S]] \rightarrow [ ][[ ][]]$

$S \rightarrow \textcolor{red}{S}S \rightarrow \textcolor{red}{S} \rightarrow SS \rightarrow \dots$



# 消除二义性

[1]	$S \rightarrow \epsilon$
[2]	[S]
[3]	SS



[1]	$S \rightarrow \epsilon$
[2]	T
[3]	$T \rightarrow []$
[4]	[T]
[5]	TT



[1]	$S \rightarrow \epsilon$
[2]	T
[3]	$T \rightarrow UT$
[4]	U
[5]	$U \rightarrow []$
[6]	[T]

消除循环引起的二义性

$S \rightarrow SS \rightarrow S$

推导[] [] []的例子

$T \rightarrow TT$  左递归容易引起二义性（回溯语法）



# 将语义加入语法中：四则运算的例子

[1]  $Expr \rightarrow Expr \text{ Bop } Expr$   
[2]       |  $num$   
[3]       |  $(Expr)$   
[4]  $Bop \rightarrow +$   
[5]       |  $-$   
[6]       |  $\times$   
[7]       |  $\div$

3 + 4 × 5 的语义?

$Expr$   
 $\rightarrow Expr \text{ Bop } Expr$          $(3 + 4) \times 5$   
 $\rightarrow Expr \text{ Bop } Expr \text{ Bop } Expr$          $3 + (4 \times 5)$   
 $\rightarrow Expr + Expr \times Expr$

将运算符特性加入到语法规则中：

- 优先级：( ) > × / ÷ > + / -
- 结合性：左结合

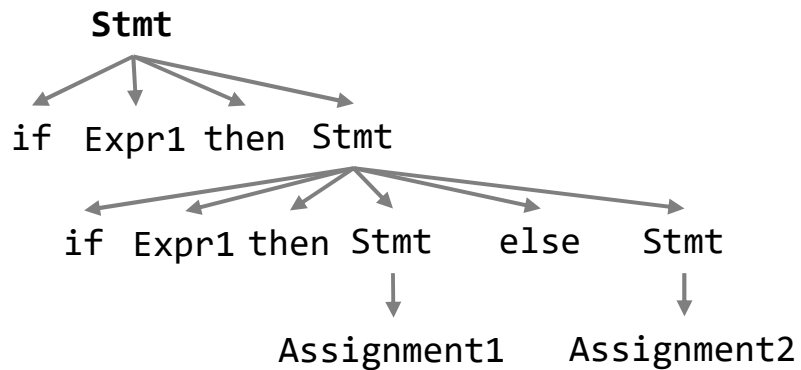
$Expr$   
 $\rightarrow Expr \text{ AMop } Term$   
 $\rightarrow Term \text{ AMop } Term$   
 $\rightarrow Term \text{ AMop } Term \text{ MDop } Factor$   
 $\rightarrow \dots$

[1]  $Expr \rightarrow Expr \text{ AMop } Term$   
[2]       |  $Term$   
[3]  $Term \rightarrow Term \text{ MDop } Factor$   
[4]       |  $Factor$   
[5]  $Factor \rightarrow (Expr)$   
[6]       |  $num$   
[7]  $AMop \rightarrow +$   
[8]       |  $-$   
[9]  $MDop \rightarrow \times$   
[10]       |  $\div$

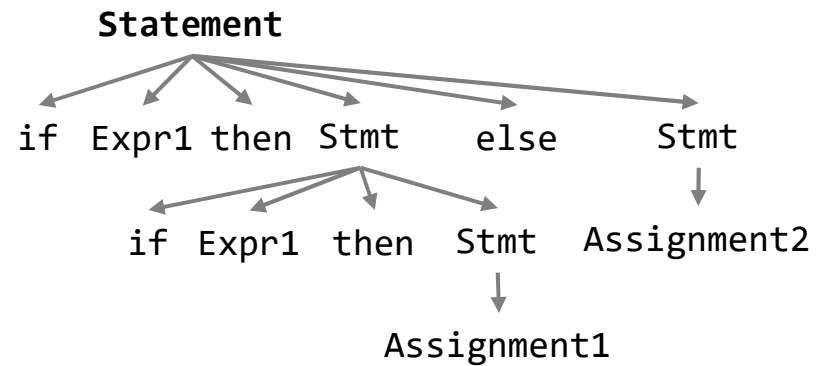
# If-Else嵌套的二义性语法

```
[1] Stmt → if Expr then Stmt else Stmt  
[2]      | if Expr then Stmt  
[3]      | Assignment  
[4]      | ...
```

if Expr1 then if Expr2 then Assignment1 else Assignment2



```
if Expr1 then  
  if Expr2 then  
    Assignment1  
  else  
    Assignment2
```



```
if Expr1 then  
  if Expr2 then  
    Assignment1  
  else  
    Assignment2
```

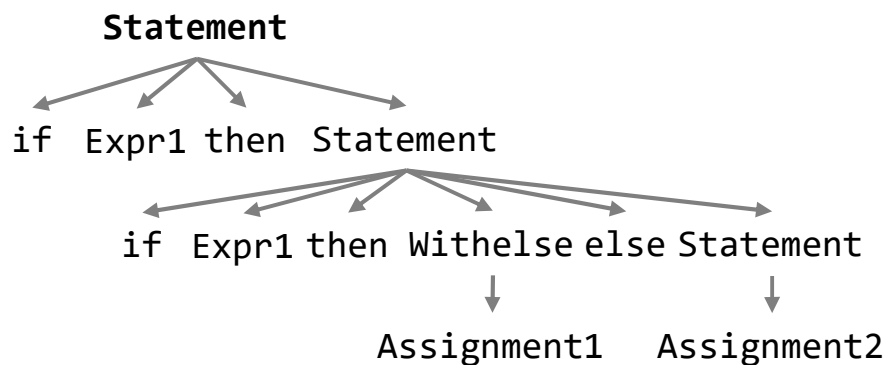
# 消除If-Else语法的二义性

- 将语义编码加入到结构中
  - 要求else优先匹配内层if
    - 如果外层出现else, 则内部嵌套的if语句一定有匹配的else

```
[1] Stmt → if Expr then Withelse else Stmt
[2]       | if Expr then Stmt
[3]       | Assignment
[4] Withelse → if Expr then Withelse else Withelse
[5]       | Assignment
```

```
[1] 有else配套的if语句, 可继续展开[4]
[2] 无配套else的if语句
[3]
[4] 内层有配套else的if语句
[5]
```

if Expr1 then if Expr2 then Assignment1 else Assignment 2



```
if Expr1 then
  if Expr2 then
    Assignment1
  else
    Assignment2
```

不存在其它推导方式

# 练习：二义性分析

- 下列If-Else语法是否存在二义性？

```
[1] Stmt → if Expr then Stmt  
[2]         | matchedStmt  
[3] matchedStmt → if Expr then matchedStmt else Stmt  
[4]         | Assignment
```

if Expr1 then if Expr2 then Assignment1 else Assignment 2



if Expr1 then if Expr2 then Assignment1 else if Expr3 then Assignment2 else Assignment3



```
if Expr1 then  
  if Expr2 then  
    Assignment1  
  else  
    if Expr3 then  
      Assignment2  
    else  
      Assignment3
```

```
if Expr1 then  
  if Expr2 then  
    Assignment1  
  else  
    if Expr3 then  
      Assignment2  
  else  
    Assignment3
```

# 练习

- 为描述正则语言的正则表达式语法设计一种CFG
  - 支持字符[A-Za-z0-9]
  - 支持连接、或|、闭包\*运算
  - 支持()
- 检查语法是否有二义性?

```
[1] < regex > ::= < union > | < concat > | < closure > | < term >
[2] < union > ::= < regex > "|" < regex >
[3] < concat > ::= < regex > < regex >
[4] < closure > ::= < regex > *
[5] < term > ::= < group > | < alphanum >
[6] < group > ::= (< regex >)
[7] < alphanum > ::= A|...|Z|a|...|z|0|...|9
```

- 主要问题:
  - 未考虑运算的优先级: 比如解析ab|c存在歧义。
  - 一般按照运算符优先级由低到高依次展开

# 改写

```
< regex > ::= < union > | < concat > | < closure > | < term >  
< union > ::= < regex > "|" < regex >  
< concat > ::= < regex > < regex >  
< closure > ::= < regex > *  
< term > ::= < group > | < alphanum >  
< group > ::= (< regex >)  
< alphanum > ::= A|...|Z|a|...|z|0|...|9
```



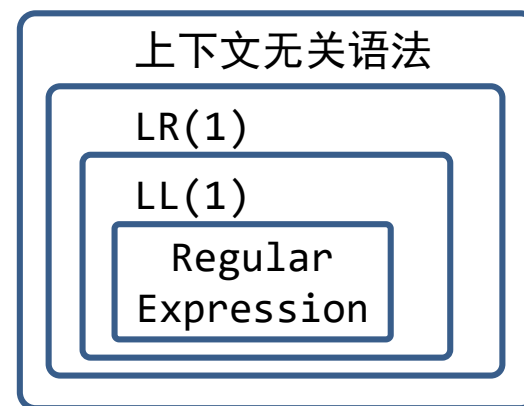
```
< regex > ::= < union > | < concat >  
< union > ::= < regex > "|" < concat >  
< concat > ::= < concat > < term > | < term >  
< term > ::= < element > * | < element >  
< element > ::= (< regex >)| < alphanum >  
< alphanum > ::= A|...|Z|a|...|z|0|...|9
```

## 思考

- 1) 用正则表达式可以定义所有的正则语言吗？
- 2) 用有穷自动机（正则表达式模拟器）可以解析任意正则表达式吗？
- 3) 用CFG可以定义任意正则语言吗？
- 4) 用CFG可以定义任意CFL语言吗？
- 5) 用pushdown automaton（CFG模拟器）可以解析任意正则表达式吗？
- 6) 用pushdown automaton可以解析任意CFG吗？
- 7) 用通用图灵机可以解析任意CFG吗？
- 8) 用通用图灵机可以解析任意程序吗？

# 编译器的任务：找到语法树推导

- 方法：
  - 自顶向下 (top-down parser)
  - 自底向上 (bottom-up parser)
- 语法难度：CFG > LR(1) > LL(1) > RE
  - 任意CFG需要花费更多时间进行语法分析
    - Earley/CYK算法复杂度 $O(n^3)$
  - LL(1)是LR(1)的一个子集
    - Left-to-Right, Leftmost
    - 前瞻单词1个
    - 适合自顶向下分析
  - LR(1)是无歧义CFG的一个子集
    - Left-to-Right, Rightmost
    - 前瞻单词1个
    - 适合自底向上分析





## 二、LLVM案例分析

---

LR( $\emptyset$ )

LR(1): SLR/LALR

# Clang采用自顶向下分析算法

A single unified parser for C, Objective C, C++, and Objective C++

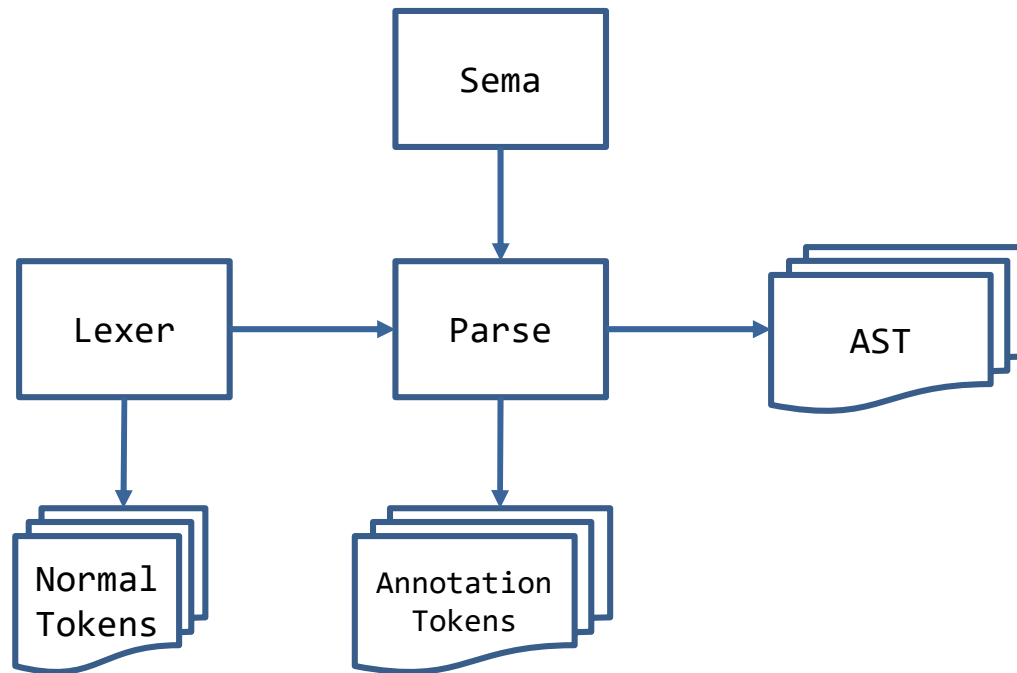
Clang is the "C Language Family Front-end", which means we intend to support the most popular members of the C family. We are convinced that the right parsing technology for this class of languages is a **hand-built recursive-descent parser**. Because it is plain C++ code, recursive descent makes it very easy for new developers to understand the code, it easily supports ad-hoc rules and other strange hacks required by C/C++, and makes it straightforward to implement excellent diagnostics and error recovery.

We believe that implementing C/C++/ObjC in a single unified parser makes the end result easier to maintain and evolve than maintaining a separate C and C++ parser which must be bugfixed and maintained independently of each other.

# 作业回顾

$\langle \text{program} \rangle$	$::=$	$\langle \text{gdecl} \rangle^* \langle \text{function} \rangle^*$
$\langle \text{gdecl} \rangle$	$::=$	$\text{extern } \langle \text{prototype} \rangle;$
$\langle \text{function} \rangle$	$::=$	$\langle \text{prototype} \rangle \langle \text{body} \rangle$
$\langle \text{prototype} \rangle$	$::=$	$\langle \text{type} \rangle \langle \text{ident} \rangle (\langle \text{paramlist} \rangle)$
$\langle \text{paramlist} \rangle$	$::=$	$\epsilon   \langle \text{type} \rangle \langle \text{ident} \rangle [, \langle \text{type} \rangle \langle \text{ident} \rangle]^*$
$\langle \text{body} \rangle$	$::=$	$\{ \langle \text{stmt} \rangle \}$
$\langle \text{stmt} \rangle$	$::=$	$\langle \text{exp} \rangle$
$\langle \text{exp} \rangle$	$::=$	$(\langle \text{exp} \rangle)   \langle \text{const} \rangle   \langle \text{ident} \rangle  $ $\langle \text{exp} \rangle \langle \text{binop} \rangle \langle \text{exp} \rangle   \langle \text{callee} \rangle$
$\langle \text{callee} \rangle$	$::=$	$\langle \text{ident} \rangle (\epsilon   \langle \text{exp} \rangle [, \langle \text{exp} \rangle]^*)$
$\langle \text{ident} \rangle$	$::=$	$[A - Z\_a - z][0 - 9A - Z\_a - z]^*$
$\langle \text{const} \rangle$	$::=$	$\langle \text{intconst} \rangle   \langle \text{doubleconst} \rangle$
$\langle \text{binop} \rangle$	$::=$	$+   -   *   <$
$\langle \text{intconst} \rangle$	$::=$	$[0 - 9][0 - 9]^*$
$\langle \text{doubleconst} \rangle$	$::=$	$\langle \text{intconst} \rangle . \langle \text{intconst} \rangle$
$\langle \text{type} \rangle$	$::=$	$\text{int}   \text{double}$

# LLVM前端架构



# Lexer::LexTokenInternal()

```
bool Lexer::LexTokenInternal(Token &Result, bool TokAtPhysicalStartOfLine) {
LexNextToken:
    // New token, can't need cleaning yet.
    Result.clearFlag(Token::NeedsCleaning);
    Result.setIdentifierInfo(nullptr);
    const char *CurPtr = BufferPtr;
    if (isHorizontalWhitespace(*CurPtr)) {
        do {
            ++CurPtr;
        } while (isHorizontalWhitespace(*CurPtr));
        BufferPtr = CurPtr;
        Result.setFlag(Token::LeadingSpace);
    }
    char Char = getAndAdvanceChar(CurPtr, Result);
    tok::TokenKind Kind;
    switch (Char) {
        case '0': case '1': case '2': case '3': case '4':
        case '5': case '6': case '7': case '8': case '9':
            MIOpt.ReadToken();
            return LexNumericConstant(Result, CurPtr);
        case ...
```

# 根据首字符大致分类

`0-9` → NumericConstant

`A-Za-z_` → Identifier

特殊情况:

- 1) 保留字
- 2) u/U/L

`'\'` → CharConstant

`'\"'` → StringLiteral

`'?'` → tok::question

`'['` → tok::l\_square

`']'` → tok::r\_square

`'('` → tok::l\_paren

`')'` → tok::r\_paren

`'{'` → tok::l\_brace

`'}'` → tok::r\_brace

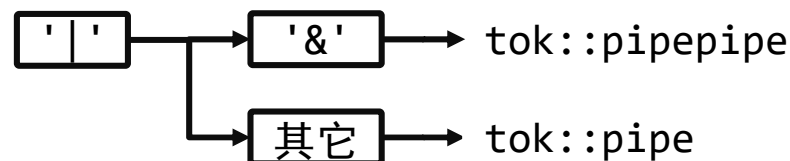
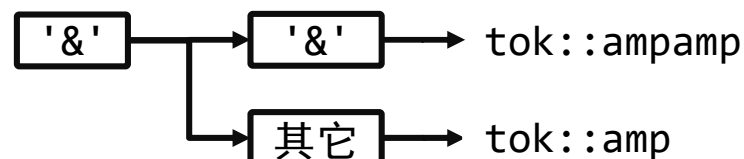
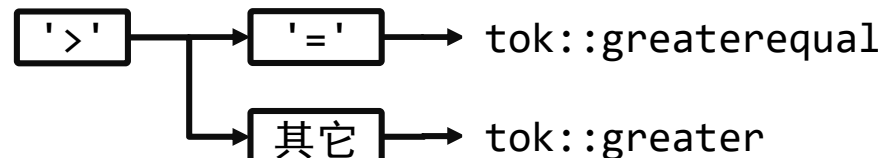
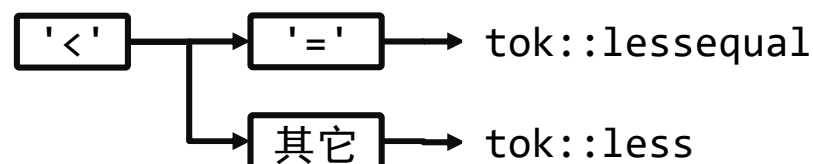
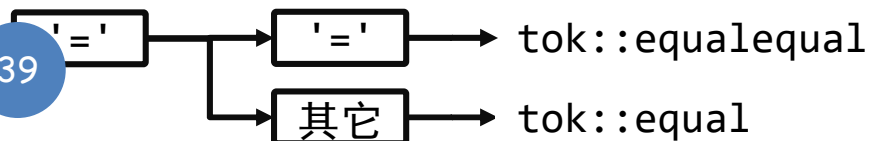
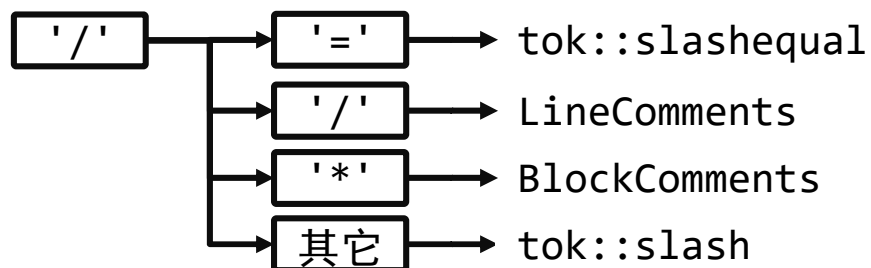
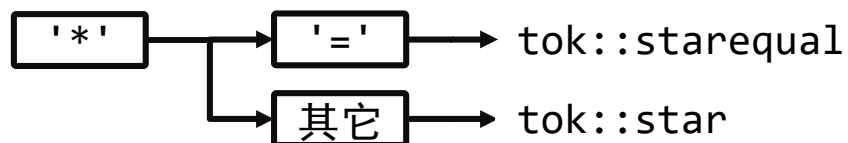
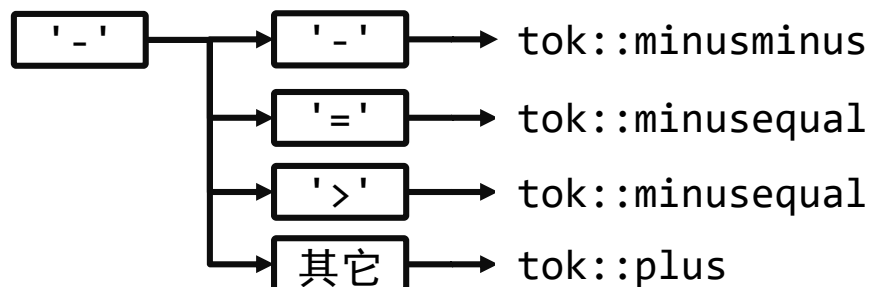
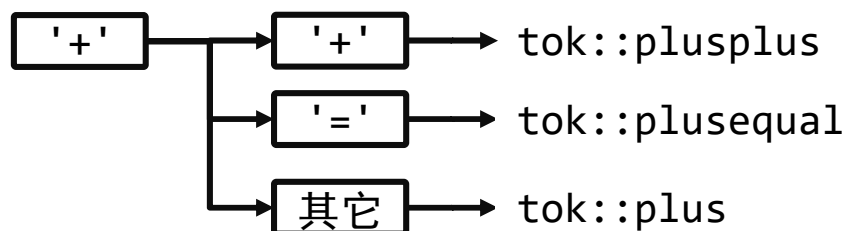
`';'` → tok::semi

`','` → tok::comma

`':'` → ?

`'.'` → ?

# 根据前两个字符大致分类



# Token类

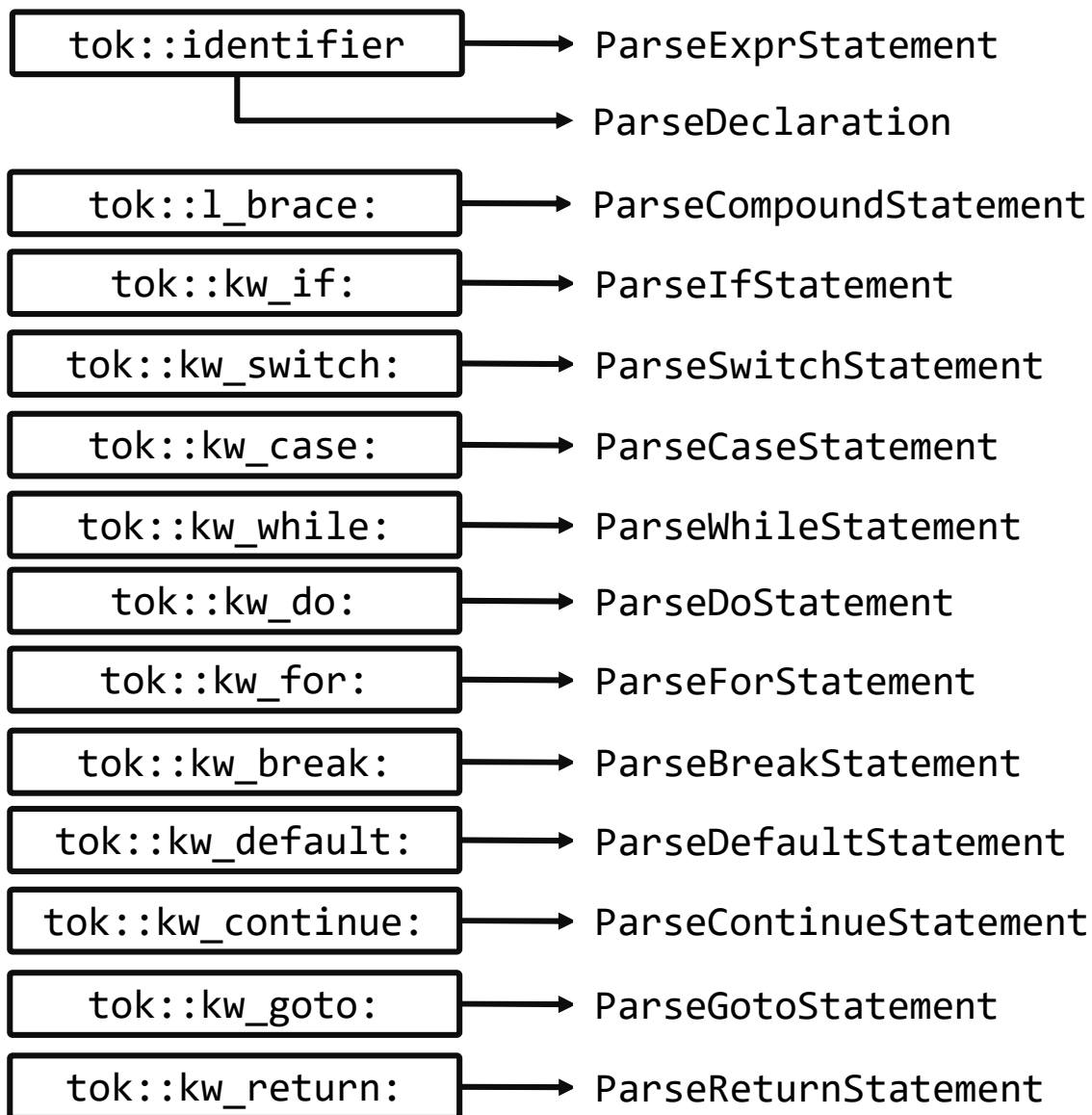
- 单个token，供lexer和parser使用，AST中不存在。
- `Sizeof(Token)`是16字节
- 分为两类：
  - Normal tokens: lexer返回的结果
    - 标识符信息（指向用于检索的哈希值）
    - Token类型（参照TokenKinds.def中的定义）
    - 位置、长度、Flags
    - 缺少了什么？
      - 数值信息
  - Annotation tokens: parser处理后的结果，添加了语义信息



## Parser::ParseStatementOrDeclarationAfterAttributes()

```
Parser::ParseStatementOrDeclarationAfterAttributes(  
    StmtVector &Stmts, ParsedStmtContext StmtCtx,  
    SourceLocation *TrailingElseLoc, ParsedAttributesWithRange &Attrs) {  
    tok::TokenKind Kind = Tok.getKind();  
    switch (Kind) {  
        case tok::at: // May be a @try or @throw statement  
        {  
            AtLoc = ConsumeToken(); // consume @  
            return ParseObjCAtStatement(AtLoc, StmtCtx);  
        }  
        case tok::identifier: {  
            Token Next = NextToken();  
            if (Next.is(tok::colon)) { // C99 6.8.1: labeled-statement  
                return ParseLabeledStatement(Attrs, StmtCtx);  
            }  
            if (Next.isNot(tok::coloncolon)) { ... }  
            default: {  
                ...  
            }  
        }  
        ...  
    }  
}
```

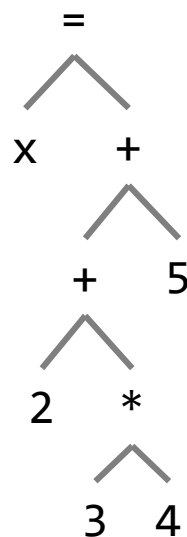
# 根据首字符大致分类



# 如何解析表达式？

$x = 2 + 3 * 4 + 5$

- 根据运算符优先级,
- 假设  $* > + > =$



# Pratt Parsing

优先级

=: [1,2]

+: [3,4]

\*: [5,6]

```
Parse(token, precedence) {  
  left = token.next();  
  if left.type != tok::num  
    return -1;  
  while true:  
    op = token.peek();  
    if op.tokenype != tok::binop  
      return -1;  
    lp, rp = Precedence[op];  
    if lp < precedence  
      break;  
    token.next();  
    right = Parse(token, rp)  
    left = (op, left, right)  
  return left  
}
```

0 1 2 3 4 5 6 3 4 0

x = 2 + 3 \* 4 + 5

Parse(start,0)

op: =,

Parse(=,2)

op: +

Parse(+,4)

op: \*

Parse(\*,+)

Return left = 4

left = \*(3, 4)

0 1 2 3 4 5 6 3 4 0

x = 2 + 3 \* 4 + 5

Return left = +(2, \*(3, 4))

0 1 2 3 4 5 6 3 4 0

x = 2 + 3 \* 4 + 5

Parse(pre5,+)

Return left = +(+(2, \*(3, 4)), 5)

0 1 2 3 4 5 6 3 4 0

x = 2 + 3 \* 4 + 5

Return left = = (x, +(+(2, \*(3, 4)), 5))

## 三、自顶向下分析

---

Top-down

Recursive-descent

# 自顶向下构建语法解析树

假设每次都能选对规则

如何解析  $(num+num) \times num$ ?

```
[1] Expr → Expr + Term
[2]      | Expr - Term
[3]      | Term
[4] Term → Term × Factor
[5]      | Term ÷ Factor
[6]      | Factor
[7] Factor → (Expr)
[8]      | num
[9]      | id
```

word	cur	Rule	Stack
(	Expr	[3]	Term
(	Term	[4]	Term, ×, Factor
(	Term	[6]	Factor, ×, Factor
(	Factor	[7]	(, Expr, ), ×, Factor
(	(	-	Expr, ), ×, Factor
num	Expr	[1]	Expr, +, Term, ), ×, Factor
num	Expr	[3]	Term, +, Term, ), ×, Factor
num	Term	[6]	Factor, +, Term, ), ×, Factor
num	Term	[8]	num, +, Term, ), ×, Factor
num	num	-	+, Term, ), ×, Factor
+	+	-	Term, ), ×, Factor
num	Term	[6]	Factor, ), ×, Factor
num	Factor	[8]	num, ), ×, Factor
num	num	-	), ×, Factor
)	)	-	×, Factor
×	×	-	Factor
num	Factor	[8]	num
num	num	-	null
eof			

# 自动搜索语法树的算法...

```
输入：程序单词流 seq;  
      CFG语法 rules;  
Output: accept: 语法解析树 ptree,  
        reject;  
初始化:  
let ptree = start symbol;  
let ptr = root;  
let st = stack();  
st.push(null);  
  
开始:  
let word = seq.NextWord();  
While (true) do:  
    if (!ptr.nodeType().isTerminal())  
        For each rule in  $\{A \rightarrow \beta_1, \dots, \beta_n; A \rightarrow \dots\}$   
            ptr.children =  $(\beta_1, \dots, \beta_n)$ ;  
            For  $1 < i < n+1$   
                st.push( $\beta_{n+2-i}$ );  
            ptr =  $\beta_1$ ;  
        //ptr.nodeType()=terminal  
    else if (word == ptr) //单词匹配成功  
        word = seq.NextWord(); //下一个单词  
        ptr = st.pop();  
    else if (word == eof && cur == null)  
        accept and return ptree;  
    else  
        backtrack(); //回溯
```

- 不考虑递归的情况:
  - $A \rightarrow \dots \rightarrow A$
- 不考虑左递归的情况:
  - $A \rightarrow \dots \rightarrow AX$
- 复杂度高:
  - 单词个数
  - 规则个数
  - 规则生产的符号数
  - 基于栈的回溯

# 左递归问题

- 对CFG的一个规则来说，其右侧的第一个符号与左侧符号相同或者能够推导出左侧符号。
- 主要问题：可使搜索算法无限递归下去，不终止。

[1]	$Expr \rightarrow Expr + Term$
[2]	$  Expr - Term$
[3]	$  Term$

word	cur	Rule	Stack
(	Expr	[1]	<i>Expr</i> , +, <i>Term</i>
(	Expr	[1]	<i>Expr</i> , +, <i>Term</i> , +, <i>Term</i>
(	Expr	[1]	<i>Expr</i> , +, <i>Term</i> , +, <i>Term</i> , +, <i>Term</i>
			...



# 消除左递归

- 引入新的非终结符，基本规则：

$$\begin{array}{|l} E \rightarrow E \alpha \\ | \beta \end{array} \Rightarrow \begin{array}{|l} E \rightarrow \beta E' \\ E' \rightarrow \alpha E' \\ | \epsilon \end{array} \qquad \begin{array}{|l} E \rightarrow E \alpha \\ | \beta \\ | \gamma \end{array} \Rightarrow \begin{array}{|l} E \rightarrow \beta E' \mid \gamma E' \\ E' \rightarrow \alpha E' \\ | \epsilon \end{array}$$

举例：

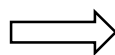
$$\begin{array}{|l} [1] \textit{Expr} \rightarrow \textit{Expr} + \textit{Term} \\ [2] \quad \quad | \textit{Expr} - \textit{Term} \\ [3] \quad \quad | \textit{Term} \end{array} \Rightarrow \begin{array}{|l} [1] \textit{Expr} \rightarrow \textit{Term} \textit{Expr}' \\ [2] \textit{Expr}' \rightarrow + \textit{Term} \textit{Expr}' \\ [3] \quad \quad | - \textit{Term} \textit{Expr}' \\ [4] \quad \quad | \epsilon \end{array}$$

$$\begin{array}{|l} [4] \textit{Term} \rightarrow \textit{Term} \times \textit{Factor} \\ [5] \quad \quad | \textit{Term} \div \textit{Factor} \\ [6] \quad \quad | \textit{Factor} \end{array} \Rightarrow \begin{array}{|l} [5] \textit{Term} \rightarrow \textit{Factor} \textit{Term}' \\ [6] \textit{Term}' \rightarrow \times \textit{Factor} \textit{Term}' \\ [7] \quad \quad | \div \textit{Factor} \textit{Term}' \\ [8] \quad \quad | \epsilon \end{array}$$

$$\begin{array}{|l} [7] \textit{Facor} \rightarrow (\textit{Expr}) \\ [8] \quad \quad | \textit{num} \\ [9] \quad \quad | \textit{name} \end{array} \Rightarrow \begin{array}{|l} [9] \textit{Facor} \rightarrow (\textit{Expr}) \\ [10] \quad \quad | \textit{num} \\ [11] \quad \quad | \textit{name} \end{array}$$

## 更多例子

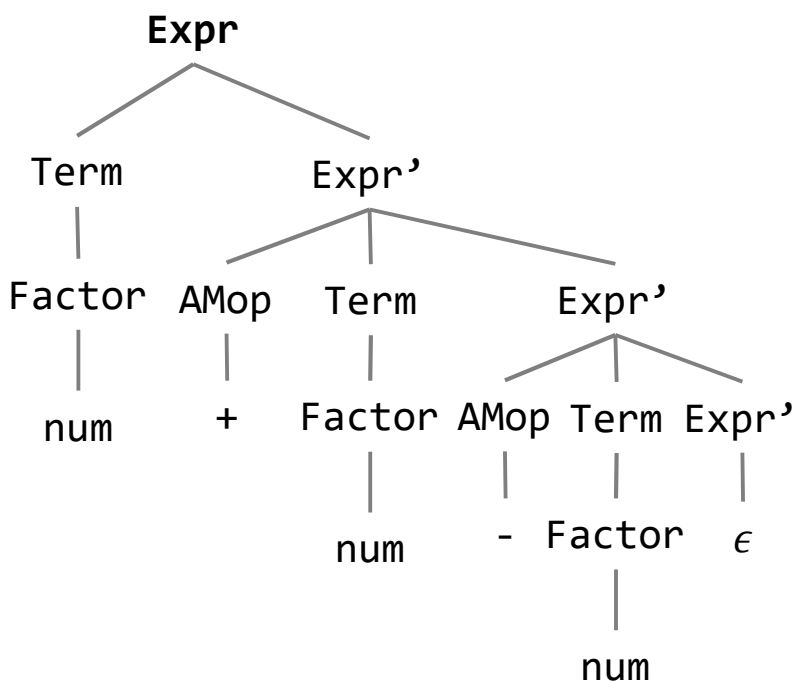
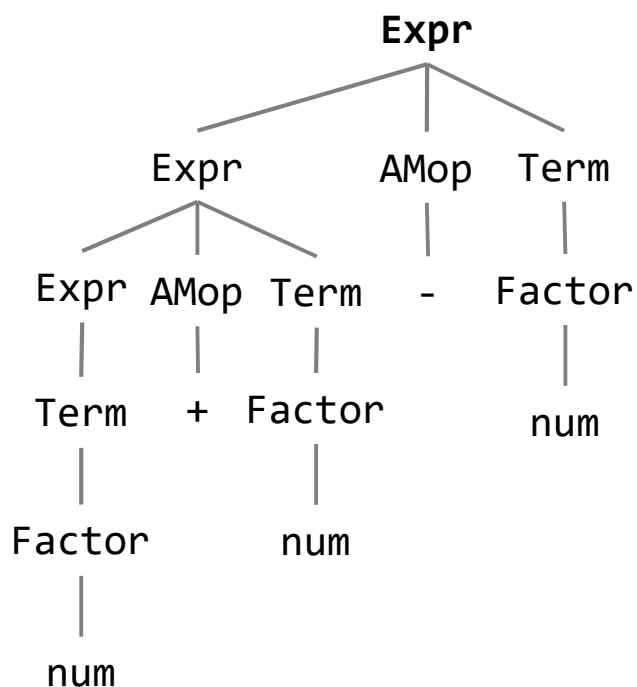
[1]  $Expr \rightarrow Expr \text{ AMop } Term$   
[2]             $| Term$   
[3]  $Term \rightarrow Term \text{ MDop } Factor$   
[4]             $| Factor$   
[5]  $Factor \rightarrow (Expr)$   
[6]             $| num$   
[7]  $AMop \rightarrow +$   
[8]             $| -$   
[9]  $MDop \rightarrow \times$   
[10]            $| \div$



[1]  $Expr \rightarrow Term \text{ Expr}'$   
[2]  $Expr' \rightarrow \text{AMop } Term \text{ Expr}'$   
[3]             $| \epsilon$   
[4]  $Term \rightarrow Factor \text{ Term}'$   
[5]  $Term' \rightarrow \text{MDop } Factor \text{ Term}'$   
[6]             $| \epsilon$   
[7]  $Factor \rightarrow (Expr)$   
[8]             $| num$   
[9]  $AMop \rightarrow +$   
[10]            $| -$   
[11]  $MDop \rightarrow \times$   
[12]            $| \div$

# 思考

- 是否会改变运算符结合性？如 $3+4-5$ 。
  - 语法解析树和虚拟语法树的区别



# 间接左递归问题

$$\begin{array}{l} E \rightarrow \alpha \\ \alpha \rightarrow \beta + \\ \beta \rightarrow E \end{array} \quad \Longrightarrow \quad E \rightarrow E +$$

展开所有非终结符NT检测和消除间接左递归

输入: Grammar{T,NT}

开始:

for i=1 to n

  for j=1 to i-1

    if  $\exists NT_i \rightarrow NT_j \gamma$

      展开  $NT_i \rightarrow NT_j \gamma$  中的非终结符  $NT_j$

      重写会造成  $NT_i$  左递归的规则

# 无回溯语法

- 目的：消除语法生成规则选择时的不确定性，避免回溯。
- 思路：如果对于每个非终结符的任意两个生成式，其产生的首个终结符号不同，则在前瞻一个单词的情况下，总能够选择正确的生成式规则。
  - [1]  $NT_1 \rightarrow NT_i \rightarrow \dots \rightarrow \text{term}_1 NT_p$
  - [2]  $NT_1 \rightarrow NT_j \rightarrow \dots \rightarrow \text{term}_2 NT_q$
- 预测解析 (Predictive Parsing) : LL(1)语法
  - Left-to-Right, Leftmost, 前瞻一个字符

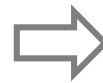
# 消除回溯：提取左因子

- 对一组产生式提取并隔离共同前缀

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_n | \gamma_1 | \dots | \gamma_j \quad \Rightarrow \quad \begin{aligned} A &\rightarrow \alpha B | \gamma_1 | \dots | \gamma_j \\ B &\rightarrow \beta_1 | \beta_2 | \dots | \beta_n \end{aligned}$$

应用举例：

```
[11] Factor → name
[12]         | name [ArgList]
[13]         | name (ArgList)
[14] ArgList → Expr MoreArgs
[15] MoreArgs → , Expr MoreArgs
[16]           | ε
```



```
[11] Factor → name Arguments
[12] Arguments → [ArgList]
[13]             | (ArgList)
[14]             | ε
[15] ArgList → Expr MoreArgs
[16] MoreArgs → , Expr MoreArgs
[17]           | ε
```

# 无回溯语法的必要性质

$$First^+(A \rightarrow \beta) = \begin{cases} First(\beta), & \text{if } \epsilon \notin First(\beta) \\ First(\beta) \cup Follow(A), & \text{otherwise} \end{cases}$$

$$\forall 1 \leq i, j \leq n, First^+(A \rightarrow \beta_i) \cap First^+(A \rightarrow \beta_j) = \emptyset$$

- 同一非终结符 $A$ 的任意两个语法推导 $(A \rightarrow \beta_i)$ 和 $(A \rightarrow \beta_j)$ 所产生的的首个终结符不能相通。
- $First(\beta)$ 是从语法符号 $\beta$ 推导出的每个子句的第一个终结符的集合，其值域是 $T \cup \{\epsilon, eof\}$ 。
- 如果 $First(\beta)$ 是 $\{\epsilon\}$ ，则计算紧随 $A$ 之后出现的终结符的集合 $Follow(A)$ 。

# First集合计算

- 对于生成式  $A \rightarrow \beta_1 \beta_2 \dots \beta_n$  来说：
  - 如果  $\epsilon \notin First(\beta_1)$ , 则  $First(A) = First(\beta_1)$
  - 如果  $\epsilon \in First(\beta_1) \& \dots \& \epsilon \in First(\beta_i)$ , 则  $First(A) = First(\beta_1) \cup \dots \cup First(\beta_{i+1})$

[1]  $Expr \rightarrow Term Expr'$   
 [2]  $Expr' \rightarrow + Term Expr'$   
 [3]       |  $- Term Expr'$   
 [4]       |  $\epsilon$

[5]  $Term \rightarrow Factor Term'$   
 [6]  $Term' \rightarrow \times Factor Term'$   
 [7]       |  $\div Factor Term'$   
 [8]       |  $\epsilon$

[9]  $Factor \rightarrow (Expr)$   
 [10]       |  $num$   
 [11]       |  $id$

	num	id	+	-	$\times$	$\div$	(	)	$\epsilon$
<i>Expr</i>	✓	✓					✓		
<i>Expr'</i>			✓	✓					✓
<i>Term</i>	✓	✓					✓		
<i>Term'</i>					✓	✓			✓
<i>Factor</i>	✓	✓					✓		



# Follow集合计算

- 紧随非终结符之后出现的所有可能的终结符

[1]  $Expr \rightarrow Term\ Expr'$   
 [2]  $Expr' \rightarrow +\ Term\ Expr'$   
 [3]       |  $-\ Term\ Expr'$   
 [4]       |  $\epsilon$

[5]  $Term \rightarrow Factor\ Term'$   
 [6]  $Term' \rightarrow \times\ Factor\ Term'$   
 [7]       |  $\div\ Factor\ Term'$   
 [8]       |  $\epsilon$

[9]  $Factor \rightarrow (Expr)$   
 [10]       |  $num$   
 [11]       |  $id$

	num	id	+	-	$\times$	$\div$	(	)	$\epsilon$	eof
$Expr$	✓	✓					✓	⊙		⊙
$Expr'$			✓	✓				⊙	✓	⊙
$Term$	✓	✓	⊙	⊙			✓	⊙		⊙
$Term'$			⊙	⊙	✓	✓		⊙	✓	⊙
$Factor$	✓	✓	⊙	⊙	⊙	⊙	✓	⊙		⊙

# First+集合计算

$$First^+(A \rightarrow \beta) = \begin{cases} First(\beta), & \text{if } \epsilon \notin First(\beta) \\ First(\beta) \cup Follow(A), & \text{otherwise} \end{cases}$$

	num	id	+	-	×	÷	(	)	ε	eof
<i>Expr</i>	✓	✓					✓	$\odot$		$\odot$
<i>Expr'</i>			✓	✓				$\odot$	✓	$\odot$
<i>Term</i>	✓	✓	$\odot$	$\odot$			✓	$\odot$		$\odot$
<i>Term'</i>			$\odot$	$\odot$	✓	✓		$\odot$	✓	$\odot$
<i>Facor</i>	✓	✓	$\odot$	$\odot$	$\odot$	$\odot$	✓	$\odot$		$\odot$

# 解析表构造：应用哪条规则可得到目标终结符？

	num	id	+	-	×	÷	(	)	ε	eof
<i>Expr</i>	✓	✓					✓			
<i>Expr'</i>			✓	✓				⊙	✓	⊙
<i>Term</i>	✓	✓					✓			
<i>Term'</i>			⊙	⊙	✓	✓		⊙	✓	⊙
<i>Factor</i>	✓	✓					✓			

[1]  $Expr \rightarrow Term\ Expr'$   
 [2]  $Expr' \rightarrow +\ Term\ Expr'$   
 [3]       |  $-\ Term\ Expr'$   
 [4]       |  $\epsilon$

[5]  $Term \rightarrow Factor\ Term'$   
 [6]  $Term' \rightarrow \times\ Factor\ Term'$   
 [7]       |  $\div\ Factor\ Term'$   
 [8]       |  $\epsilon$

[9]  $Factor \rightarrow (Expr)$   
 [10]       | *num*  
 [11]       | *id*

	num	id	+	-	×	÷	(	)	ε	eof
<i>Expr</i>	1	1					1			
<i>Expr'</i>			2	3				4		4
<i>Term</i>	5	5					5			
<i>Term'</i>			8	8	6	7		8		8
<i>Factor</i>	10	11					9			

# 应用解析表

(num+num)×num

[1]  $Expr \rightarrow Term\ Expr'$   
 [2]  $Expr' \rightarrow +\ Term\ Expr'$   
 [3]       |  $-\ Term\ Expr'$   
 [4]       |  $\epsilon$

[5]  $Term \rightarrow Factor\ Term'$   
 [6]  $Term' \rightarrow \times\ Factor\ Term'$   
 [7]       |  $\div\ Factor\ Term'$   
 [8]       |  $\epsilon$

[9]  $Facor \rightarrow (Expr)$   
 [10]       |  $num$   
 [11]       |  $id$

	num	id	+	-	×	÷	(	)	<i>eof</i>
<i>Expr</i>	1	1					1		
<i>Expr'</i>			2	3				4	4
<i>Term</i>	5	5					5		
<i>Term'</i>			8	8	6	7		8	8
<i>Facor</i>	10	11					9		

word	cur	Rule	Stack
(	Expr	[1]	Term, Expr'
(	Term	[5]	Factor, Term', Expr'
(	Factor	[9]	(, Expr, ), Term', Expr'
(	(	-	Expr, ), Term', Expr'
num	Expr	[1]	Term, Expr', ), Term', Expr'
num	Term	[5]	Factor, Term', Expr', ), Term', Expr'
num	Factor	[10]	num, Term', Expr', ), Term', Expr'
num	num	-	Term', Expr', ), Term', Expr'
+	Term'	[8]	Expr', ), Term', Expr'
+	Expr'	[2]	+, Term, Expr', ), Term', Expr'
+	+	-	Term, Expr', ), Term', Expr'
num	Term	[5]	Factor, Term', Expr', ), Term', Expr'
num	Factor	[10]	num, Term', Expr', ), Term', Expr'
num	num	-	Term', Expr', ), Term', Expr'
)	Term'	[8]	Expr', ), Term', Expr'
)	Expr'	[4]	), Term', Expr'
)	)	-	Term', Expr'
+	Term'	[8]	Expr'
...	...	...	...

## 练习：

- 将正则表达式CFG改写为LL(1)语法并写出应用解析表

```
< regex > ::= < union > | < concat >  
< union > ::= < regex > "|" < concat >  
< concat > ::= < concat > < term > | < term >  
< term > ::= < element > * | < element >  
< element > ::= (< regex >)| < alphanum >
```

# 解答：消除左递归

- 左递归问题一：
  - $\langle concat \rangle ::= \langle concat \rangle \langle term \rangle \mid \langle term \rangle$
- 改写结果：
  - $\langle concat \rangle ::= \langle term \rangle \langle concat' \rangle$
  - $\langle concat' \rangle ::= \langle term \rangle \langle concat' \rangle \mid \epsilon$
- 和下列形式等价（但转换AST时要注意结合性）：
  - $\langle concat \rangle ::= \langle term \rangle \langle concat \rangle \mid \langle term \rangle$
- 左递归问题二（间接左递归）：
  - $\langle regex \rangle ::= \langle union \rangle \mid \langle concat \rangle$
  - $\langle union \rangle ::= \langle regex \rangle \mid \langle concat \rangle$
  - $\Rightarrow \langle regex \rangle ::= \langle regex \rangle \mid \langle concat \rangle \mid \langle concat \rangle$
- 改写结果：
  - $\langle regex \rangle ::= \langle concat \rangle \langle regex' \rangle$
  - $\langle regex' \rangle ::= \mid \langle concat \rangle \langle regex' \rangle \mid \epsilon$

```
[1] < regex > ::= < concat > < regex' >
[2] < regex' > ::= " | " < concat > < regex' >
[3]                | ε
[4] < concat > ::= < term > < concat' >
[5] < concat' > ::= < term > < concat' >
[6]                | ε
[7] < term > ::= < element > *
[8]                | < element >
[9] < element > ::= ( < regex > )
[10]                | < alphanum >
```

# 解答：消除回溯

- 回溯问题语法
  - $\langle term \rangle ::= \langle element \rangle^*$
  - $\quad \quad \quad | \langle element \rangle$
- 改写结果：
  - $\langle term \rangle ::= \langle element \rangle \langle follow \rangle$
  - $\langle follow \rangle ::= * | \epsilon$

```
[1]  $\langle regex \rangle ::= \langle concat \rangle \langle regex' \rangle$ 
[2]  $\langle regex' \rangle ::= "|" \langle concat \rangle \langle regex' \rangle$ 
[3]  $\quad \quad \quad | \epsilon$ 
[4]  $\langle concat \rangle ::= \langle term \rangle \langle concat' \rangle$ 
[5]  $\langle concat' \rangle ::= \langle term \rangle \langle concat' \rangle$ 
[6]  $\quad \quad \quad | \epsilon$ 
[7]  $\langle term \rangle ::= \langle element \rangle \langle follow \rangle$ 
[8]  $\langle follow \rangle ::= *$ 
[9]  $\quad \quad \quad | \epsilon$ 
[10]  $\langle element \rangle ::= (\langle regex \rangle)$ 
[11]  $\quad \quad \quad | \langle alphanum \rangle$ 
```

## 解答：构建LL(1)解析表

	(	)	A – Z, a – z, 0 – 9		*	ε	eof
<i>regex</i>	√[1]		√[1]				
<i>regex'</i>				√[1]		√	⊙[3]
<i>concat</i>	√[4]		√[4]				
<i>concat'</i>	√[5]		√[5]			√	⊙[6]
<i>term</i>	√[7]		√[7]				
<i>element</i>	√[10]		√[11]				
<i>follow</i>	⊙[9]		⊙[9]		√[8]	√	
<i>alphanum</i>			√[12]				

√:first()

⊙:follow()

[ ]:语法规则条目



# 通用自顶向下语法分析算法：Earley算法

- 三种基本操作：
  - 预测 (Prediction) : 对于每个状态  $X \rightarrow \alpha \circ Y \beta$ , 根据语法规则预测  $Y \rightarrow \circ \gamma$ 。
  - 扫描 (Scanning) : 如果下一个待处理的符号是  $a$ , 并且存在状态  $X \rightarrow \alpha \circ a \beta$ , 则扫描该字符并且将状态变更为  $X \rightarrow \alpha a \circ \beta$ 。
  - 完成 (Completion) :  $Y \rightarrow \gamma \circ$  完成了对  $Y$  的分析, 进而更新  $X \rightarrow \alpha \circ Y \beta$  为  $X \rightarrow \alpha Y \circ \beta$ 。

# Earley

## 算法参考

```
DECLARE ARRAY S;

function INIT(words) {
    S ← CREATE-ARRAY(LENGTH(words) + 1)
    for k ← from 0 to LENGTH(words)
        S[k] ← EMPTY-ORDERED-SET
}
function EARLEY-PARSE(words, grammar) {
    INIT(words)
    ADD-TO-SET(( $\gamma \rightarrow \bullet S$ , 0), S[0])
    for k ← from 0 to LENGTH(words)
        for each state in S[k] { // S[k] can expand during this loop
            if not FINISHED(state) {
                if NEXT-ELEMENT-OF(state) is a nonterminal then
                    PREDICTOR(state, k, grammar) // non-terminal
                else
                    SCANNER(state, k, words) // terminal
            }
            else
                COMPLETER(state, k)
        }
    }
    return chart
}
procedure PREDICTOR(( $A \rightarrow \alpha \bullet B \beta$ , j), k, grammar) {
    for each ( $B \rightarrow \gamma$ ) in GRAMMAR-RULES-FOR(B, grammar)
        ADD-TO-SET(( $B \rightarrow \bullet \gamma$ , k), S[k])
}
procedure SCANNER(( $A \rightarrow \alpha \bullet a \beta$ , j), k, words)
    if  $a \in \text{PARTS-OF-SPEECH}(\text{words}[k])$ 
        ADD-TO-SET(( $A \rightarrow \alpha a \bullet \beta$ , j), S[k+1])

procedure COMPLETER(( $B \rightarrow \gamma \bullet$ , x), k)
    for each ( $A \rightarrow \alpha \bullet B \beta$ , j) in S[x]
        ADD-TO-SET(( $A \rightarrow \alpha B \bullet \beta$ , j), S[k])
```

# Earley算法

如何根据下列语法规则  
解析  $(3+4) \times 5$ ?

$(\text{num}+\text{num}) \times \text{num}$

[1]	$Expr \rightarrow Expr + Term$
[2]	$  Expr - Term$
[3]	$  Term$
[4]	$Term \rightarrow Term \times Factor$
[5]	$  Term \div Factor$
[6]	$  Factor$
[7]	$Factor \rightarrow (Expr)$
[8]	$  num$
[9]	$  id$

no	production	origin	comment
$s(0) = \circ(\text{num}+\text{num}) \times \text{num}$			
1	$Expr \rightarrow \circ Expr + Term$	$s(0)$	start rule
2	$Expr \rightarrow \circ Expr - Term$	$s(0)$	start rule
3	$Expr \rightarrow \circ Term$	$s(0)$	start rule
4	$Term \rightarrow \circ Term \times Factor$	$s(0)$	Predict from [0][3]
5	$Term \rightarrow \circ Term \div Factor$	$s(0)$	Predict from [0][3]
6	$Term \rightarrow \circ Factor$	$s(0)$	Predict from [0][3]
7	$Factor \rightarrow \circ (Expr)$	$s(0)$	Predict from [0][6]
8	$Factor \rightarrow \circ num$	$s(0)$	Predict from [0][6]
9	$Factor \rightarrow \circ id$	$s(0)$	Predict from [0][6]
$s(1) = (\circ \text{num}+\text{num}) \times \text{num}$			
1	$Factor \rightarrow (\circ Expr)$	$s(0)$	Scan from [0][7]
2	$Expr \rightarrow \circ Expr + Term$	$s(1)$	Predict from [1][1]
3	$Expr \rightarrow \circ Expr - Term$	$s(1)$	Predict from [1][1]
4	$Expr \rightarrow \circ Term$	$s(1)$	Predict from [1][1]
5	$Term \rightarrow \circ Term \times Factor$	$s(1)$	Predict from [1][4]
6	$Term \rightarrow \circ Term \div Factor$	$s(1)$	Predict from [1][4]
7	$Term \rightarrow \circ Factor$	$s(1)$	Predict from [1][4]
8	$Factor \rightarrow \circ (Expr)$	$s(1)$	Predict from [1][7]
9	$Factor \rightarrow \circ num$	$s(1)$	Predict from [1][7]
10	$Factor \rightarrow \circ id$	$s(1)$	Predict from [1][7]

# Earley算法

[1]  $Expr \rightarrow Expr + Term$   
 [2]       |  $Expr - Term$   
 [3]       |  $Term$   
 [4]  $Term \rightarrow Term \times Factor$   
 [5]       |  $Term \div Factor$   
 [6]       |  $Factor$   
 [7]  $Factor \rightarrow (Expr)$   
 [8]       |  $num$   
 [9]       |  $id$

no	production	origin	comment
$s(2) = (num \circ + num) \times num$			
1	$Factor \rightarrow num \circ$	$s(1)$	Scan from [1][9]
2	$Term \rightarrow Factor \circ$	$s(1)$	Complete [1][7]
3	$Term \rightarrow Term \circ \times Factor$	$s(1)$	Complete [1][5]
4	$Term \rightarrow Term \circ \div Factor$	$s(1)$	Complete [1][6]
5	$Expr \rightarrow Term \circ$	$s(1)$	Complete [1][4]
6	$Expr \rightarrow Expr \circ + Term$	$s(1)$	Complete [1][2]
7	$Expr \rightarrow Expr \circ - Term$	$s(1)$	Complete [1][3]
8	$Factor \rightarrow (Expr \circ)$	$s(0)$	Complete [1][1]
$s(3) = (num + \circ num) \times num$			
1	$Expr \rightarrow Expr + \circ Term$	$s(1)$	Scan from [2][6]
2	$Term \rightarrow \circ Term \times Factor$	$s(3)$	Predict from [3][1]
3	$Term \rightarrow \circ Term \div Factor$	$s(3)$	Predict from [3][1]
4	$Term \rightarrow \circ Factor$	$s(3)$	Predict from [3][1]
5	$Factor \rightarrow \circ (Expr)$	$s(3)$	Predict from [3][4]
6	$Factor \rightarrow \circ num$	$s(3)$	Predict from [3][4]
7	$Factor \rightarrow \circ id$	$s(3)$	Predict from [3][4]

# Earley算法

[1]  $Expr \rightarrow Expr + Term$   
 [2]     |  $Expr - Term$   
 [3]     |  $Term$   
 [4]  $Term \rightarrow Term \times Factor$   
 [5]     |  $Term \div Factor$   
 [6]     |  $Factor$   
 [7]  $Factor \rightarrow (Expr)$   
 [8]     |  $num$   
 [9]     |  $id$

no	production	origin	comment
$s(4) = (num+num \circ) \times num$			
1	$Factor \rightarrow num \circ$	$s(3)$	Scan from [3][6]
2	$Term \rightarrow Factor \circ$	$s(3)$	Complete [3][7]
3	$Term \rightarrow Term \circ \times Factor$	$s(3)$	Complete [3][5]
4	$Term \rightarrow Term \circ \div Factor$	$s(3)$	Complete [3][6]
5	$Expr \rightarrow Term \circ$	$s(3)$	Complete [3][4]
6	$Expr \rightarrow Expr \circ + Term$	$s(3)$	Complete [3][2]
7	$Expr \rightarrow Expr \circ - Term$	$s(3)$	Complete [3][3]
8	$Expr \rightarrow Expr + Term \circ$	$s(1)$	Complete [3][1]
9	$Factor \rightarrow (Expr \circ)$	$s(0)$	Complete [1][1]
$s(5) = (num+num) \circ \times num$			
1	$Factor \rightarrow (Expr) \circ$	$s(0)$	Scan from [4][9]
2	$Term \rightarrow Factor \circ$	$s(0)$	Complete [0][6]
3	$Term \rightarrow Term \circ \times Factor$	$s(0)$	Complete [0][4]
4	$Term \rightarrow Term \circ \div Factor$	$s(0)$	Complete [0][5]
5			
6			
7			

# Earley算法

```
[1] Expr → Expr + Term
[2]      | Expr - Term
[3]      | Term
[4] Term → Term × Factor
[5]      | Term ÷ Factor
[6]      | Factor
[7] Facor → (Expr)
[8]      | num
[9]      | id
```

no	production	origin	comment
s(6) = (num+num)×°num			
1	$Term \rightarrow Term \times \circ Factor$	s(0)	Scan from [6][3]
2	$Facor \rightarrow \circ (Expr)$	s(0)	Predict from [6][1]
3	$Facor \rightarrow \circ num$	s(6)	Predict from [6][1]
4	$Facor \rightarrow \circ id$	s(6)	Predict from [6][1]
s(7) = (num+num)×num°			
1	$Facor \rightarrow num \circ$	s(6)	Scan from [6][4]
2	$Term \rightarrow Term \times Factor \circ$	s(0)	Complete [6][1]
3	$Expr \rightarrow Term \circ$	s(0)	Complete [0][3]
4			
5			
6			
7			

Earley算法能否解析非CFG语法?

## 练习：

- 应用Earley算法解析正则 $\text{num} + \text{num} - \text{num}$
- 应用Earley算法解析正则  $(a|b)^*$

```
[1] Expr → Expr + Term
[2]      | Expr - Term
[3]      | Term
[4] Term → Term × Factor
[5]      | Term ÷ Factor
[6]      | Factor
[7] Factor → (Expr)
[8]      | num
[9]      | id
```

```
[1] < regex > ::= < union >
[2]      | < concat >
[3] < union > ::= < regex > "|" < concat >
[4] < concat > ::= < concat > < term >
[5]      | < term >
[6] < term > ::= < element > *
[7]      | < element >
[8] < element > ::= (< regex >)
[9]      | < alphanum >
```

# 自顶向下解析算法小结

- Earley算法
  - 支持所有CFG
  - 复杂度 $O(n^3)$ 
    - 无歧义语法:  $O(n^2)$
- 能否有更快的算法?
  - 要求CFG为LL(1)
    - 不存在左递归
    - 不存在回溯



## 四、自底向上分析

---

SLR/LALR/LR(1)

# 基本思路：基于规约的方法

- 如果在句法分析栈的上边缘找到 $\beta$ 且 $A \rightarrow \beta$ ，则将其规约为 $A$ ；
- 否则移进

[1]	$S \rightarrow \epsilon$
[2]	[S]S

	当前栈	待输入
第1步:		[ ]
第2步:	[	]
第3步:	[S	]
第4步:	[S]S	
第5步:	S	

# 移进和规约 (Shift-Reduce)

Shift  
移进表示:

$$\frac{w:\beta}{wa:\beta a} \text{shift}$$

Reduce  
规约表示:

$$\frac{[r] X \rightarrow \alpha \quad w:\beta \alpha}{w:\beta X} \text{reduce}(r)$$

示例:

$$\frac{\epsilon:\beta_0}{[:\beta_0[} \text{shift}$$

$$\frac{[:\beta_1}{[]:\beta_1]} \text{shift}$$

...

$$\frac{[] [] []:\beta_7}{[] [] []:\beta_7]} \text{shift}$$

$$[1] S \rightarrow \epsilon$$

$$\frac{w:\beta}{w:\beta S} \text{reduce}([1])$$

$$[2] S \rightarrow [S]S$$

$$\frac{w:\beta [S]S}{w:\beta S} \text{reduce}([2])$$

# 移进-规约应用

## 最右推导

	[ ] [ ] [ ] [ ]
[	[ ] [ ] [ ] [ ]
[S	] [ ] [ ] [ ]
[S]	[ ] [ ] [ ] [ ]
[S][	[ ] [ ] [ ] [ ]
[S][[	] [ ] [ ] [ ]
[S][[S	] [ ] [ ] [ ]
[S][[S]	[ ] [ ] [ ] [ ]
[S][[S][	] [ ] [ ] [ ]
[S][[S][S	] [ ] [ ] [ ]
[S][[S][S]	] [ ] [ ] [ ]
[S][[S][S]S	] [ ] [ ] [ ]
[S][[S]S	] [ ] [ ] [ ]
[S][S	] [ ] [ ] [ ]
[S][S]S	] [ ] [ ] [ ]
[S]S	] [ ] [ ] [ ]
S	] [ ] [ ] [ ]

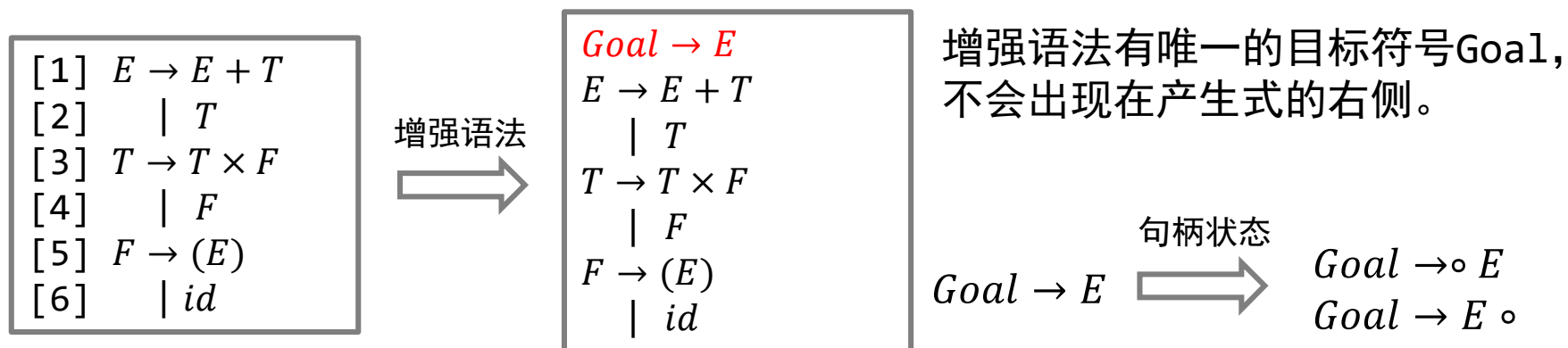
[illegible]
$$\begin{array}{l} [1] \quad S \rightarrow \epsilon \\ [2] \quad \quad | [S]S \end{array}$$

```

shift [
reduce([1])
shift
shift
shift
reduce([1])
shift
shift
reduce([1])
shift
reduce([1])
reduce([2])
reduce([2])
shift
reduce([1])
reduce([2])
reduce([2])

```

# LR(0)句柄分析



可应用  $E \rightarrow E + T$   
规约的句柄状态  $\xRightarrow{\hspace{1cm}}$

$E \rightarrow \circ E + T$   
 $E \rightarrow E \circ + T$   
 $E \rightarrow E + \circ T$   
 $E \rightarrow E + T \circ$

可应用  $E \rightarrow T$   
规约的句柄状态  $\xRightarrow{\hspace{1cm}}$

$E \rightarrow \circ T$   
 $E \rightarrow T \circ$

可应用  $T \rightarrow T \times F$   
规约的句柄状态  $\xRightarrow{\hspace{1cm}}$

$T \rightarrow \circ T \times F$   
 $T \rightarrow T \circ \times F$   
 $T \rightarrow T \times \circ F$   
 $T \rightarrow T \times F \circ$

可应用  $T \rightarrow F$   
规约的句柄状态  $\xRightarrow{\hspace{1cm}}$

$T \rightarrow \circ F$   
 $T \rightarrow F \circ$

可应用  $F \rightarrow (E)$   
规约的句柄状态  $\xRightarrow{\hspace{1cm}}$

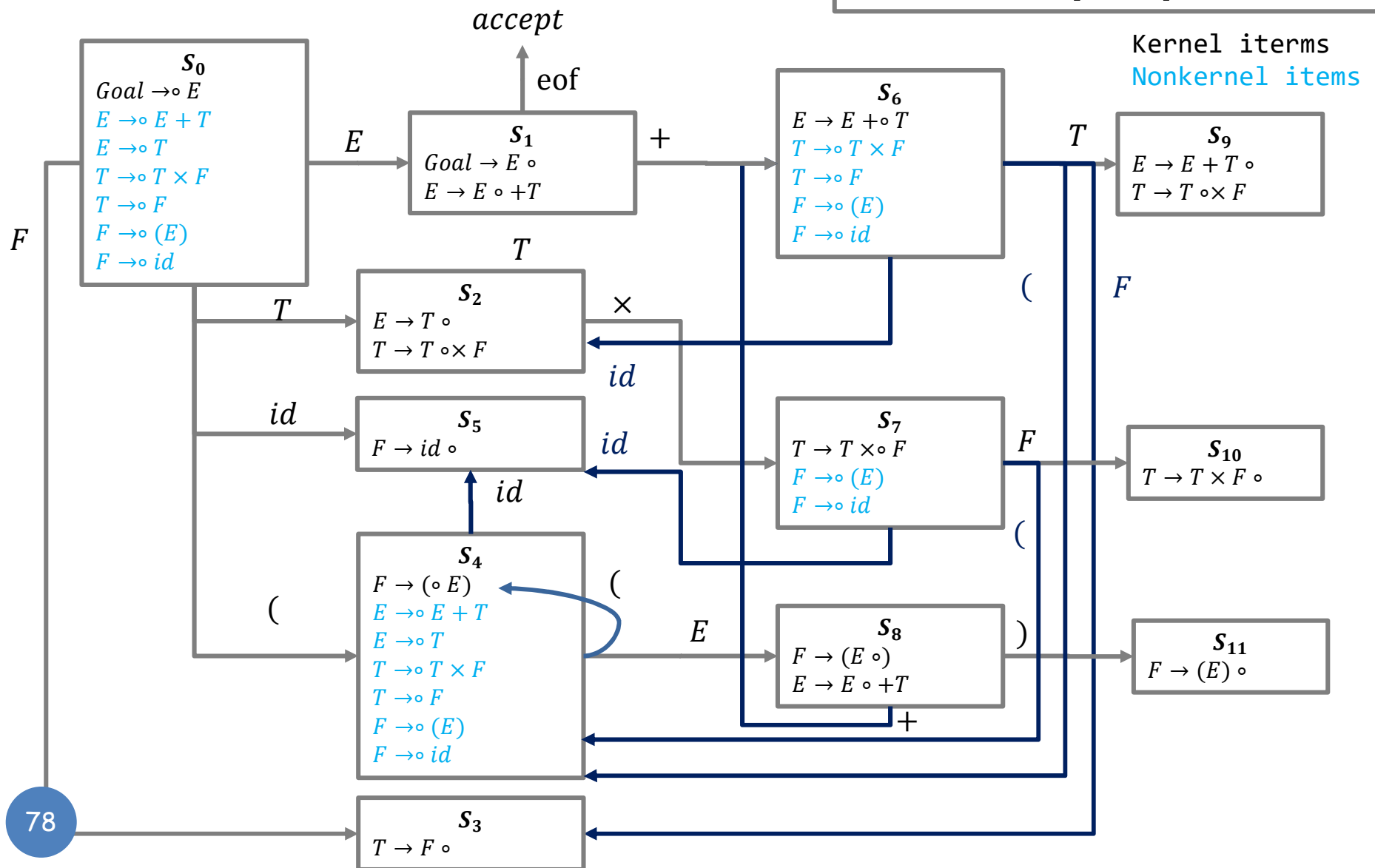
$F \rightarrow \circ (E)$   
 $F \rightarrow (\circ E)$   
 $F \rightarrow (E \circ)$   
 $F \rightarrow (E) \circ$

可应用  $F \rightarrow id$   
规约的句柄状态  $\xRightarrow{\hspace{1cm}}$

$F \rightarrow \circ id$   
 $F \rightarrow id \circ$

# LR(0)自动机构建

While (S has changed)  
 for each item  $[A \rightarrow \beta \circ C \delta, a] \in S$   
 for each production  $[C \rightarrow \lambda] \in G$   
 if  $[C \rightarrow \circ \lambda] \notin S$   
 $S \leftarrow S \cup [C \rightarrow \circ \lambda]$



# LR(0)自动机的状态转移关系表

迭代	规范项	id	+	×	(	)	eof	E	T	F
0	$S_0$	$S_5$	$\emptyset$	$\emptyset$	$S_4$	$\emptyset$	$\emptyset$	$S_1$	$S_2$	$S_3$
1	$S_1$	$\emptyset$	$S_6$	$\emptyset$	$\emptyset$	$\emptyset$	<i>accept</i>	$\emptyset$	$\emptyset$	$\emptyset$
	$S_2$	$\emptyset$	$\emptyset$	$S_7$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
	$S_3$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
	$S_4$	$S_5$	$\emptyset$	$\emptyset$	$S_4$	$\emptyset$	$\emptyset$	$S_8$	$S_2$	$S_3$
	$S_5$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
2	$S_6$	$S_5$	$\emptyset$	$\emptyset$	$S_4$	$\emptyset$	$\emptyset$	$\emptyset$	$S_9$	$S_3$
	$S_7$	$S_5$	$\emptyset$	$\emptyset$	$S_4$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$S_{10}$
	$S_8$	$\emptyset$	$S_6$	$\emptyset$	$\emptyset$	$S_{11}$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
3	$S_9$	$\emptyset$	$\emptyset$	$S_7$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
	$S_{10}$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$
	$S_{11}$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$

# 构建SLR解析表

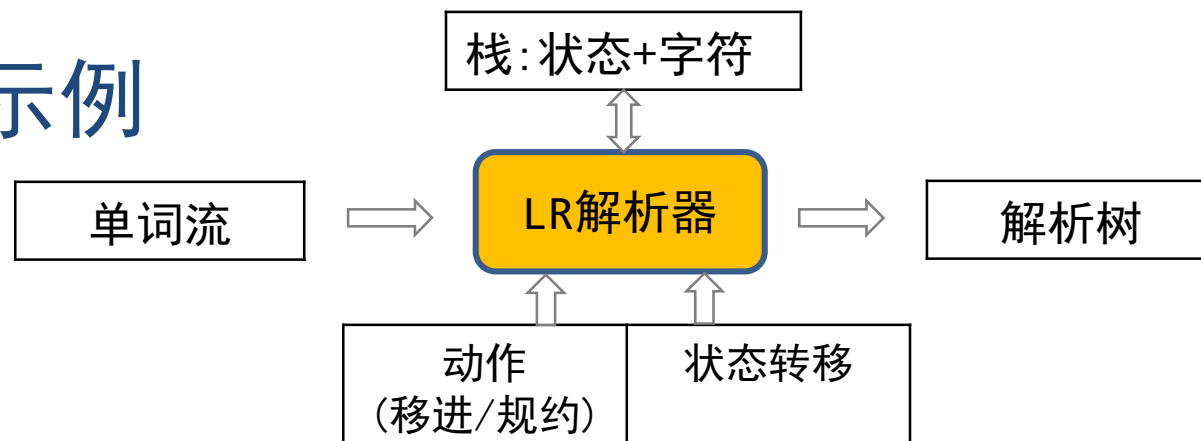
移进条件：如果 $A \rightarrow \alpha \circ a\beta \in S_i$ ，并且 $Goto(S_i, a) = S_j$ ，设置 $Action(S_i, a) = \text{"shift } j\text{"}$

规约条件：如果 $A \rightarrow \alpha \circ \in S_i$ ， $\forall a \in Follow(A)$ ，设置 $Action(S_i, a) = \text{"reduce } A \rightarrow \alpha\text{"}$

规范项	Action						Goto		
	id	+	×	(	)	eof	E	T	F
S <sub>0</sub>	shift S <sub>5</sub>			shift S <sub>4</sub>			S <sub>1</sub>	S <sub>2</sub>	S <sub>3</sub>
S <sub>1</sub>		shift S <sub>6</sub>				accept			
S <sub>2</sub>		reduce [2]	shift S <sub>7</sub>		reduce [2]	reduce [2]			
S <sub>3</sub>									
S <sub>4</sub>	shift S <sub>5</sub>			shift S <sub>4</sub>			S <sub>8</sub>	S <sub>2</sub>	S <sub>3</sub>
S <sub>5</sub>		reduce [6]	reduce [6]		reduce [6]	reduce [6]			
S <sub>6</sub>	shift S <sub>5</sub>			shift S <sub>4</sub>				S <sub>9</sub>	S <sub>3</sub>
S <sub>7</sub>	shift S <sub>5</sub>			shift S <sub>4</sub>					S <sub>10</sub>
S <sub>8</sub>		shift S <sub>6</sub>			shift S <sub>11</sub>				
S <sub>9</sub>		reduce [1]	shift S <sub>7</sub>		reduce [1]	reduce [1]			
S <sub>10</sub>		reduce [3]	reduce [3]		reduce [3]	reduce [3]			
S <sub>11</sub>		reduce [5]	reduce [5]		reduce [5]	reduce [5]			



# SLR应用示例



Stack	Symbols	Input	Action
$S_0$		id×id \$	shift id, goto $S_5$
$S_0S_5$	id	×id \$	reduce by $F \rightarrow id$ , back to $S_0$ , goto $S_3$
$S_0S_3$	F	×id \$	reduce by $T \rightarrow F$ , back to $S_0$ , goto $S_2$
$S_0S_2$	T	×id \$	shift ×, goto $S_7$
$S_0S_2S_7$	T ×	id \$	shift id, goto $S_5$
$S_0S_2S_7S_5$	T × id	\$	reduce by $F \rightarrow id$ , back to $S_7$ , goto $S_{10}$
$S_0S_2S_7S_{10}$	T × F	\$	reduce by $T \rightarrow T \times F$ , back to $S_7S_2S_0$ , goto $S_2$
$S_0S_2$	T	\$	reduce by $E \rightarrow T$ , back to $S_0$ , goto $S_1$
$S_0S_1$	E	\$	

## 练习

- 下面的语法是否是LL(1)? 是否是SLR(1)

[1]  $S \rightarrow AaAb$   
[2]  $\quad \quad |BbBa$   
[3]  $A \rightarrow \epsilon$   
[4]  $B \rightarrow \epsilon$

是LL(1)

$First^+(S \rightarrow AaAb) = \{a\}$

$First^+(S \rightarrow BbBa) = \{b\}$

$S_0$   
 $S' \rightarrow \circ S$   
 $S \rightarrow \circ AaAb$   
 $S \rightarrow \circ BbAa$   
 $A \rightarrow \circ$   
 $B \rightarrow \circ$

不是SLR(1)

$Follow(A) = Follow(B) = \{a, b\}$

$Action(S_0, a) = reduce[3] \text{ 或 } reduce[4]$

## 练习

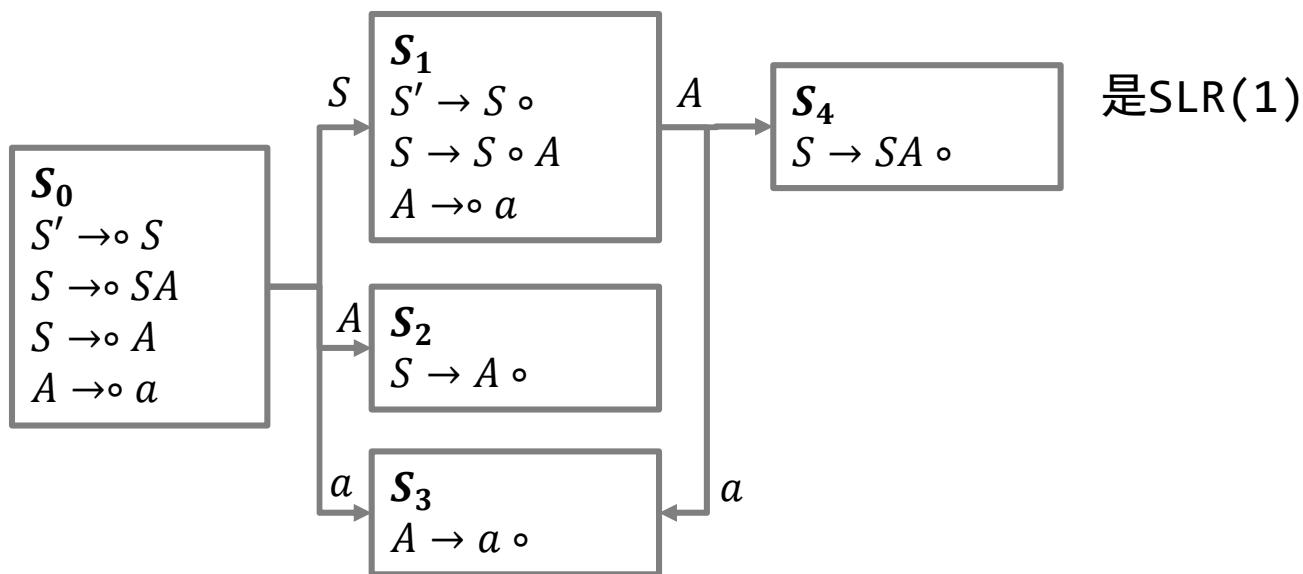
- 下面的语法是否是LL(1)? 是否是SLR(1)

[1]  $S \rightarrow SA$   
[2]      $|A$   
[3]  $A \rightarrow a$

不是LL(1)

$First^+(S \rightarrow SA) = \{a\}$

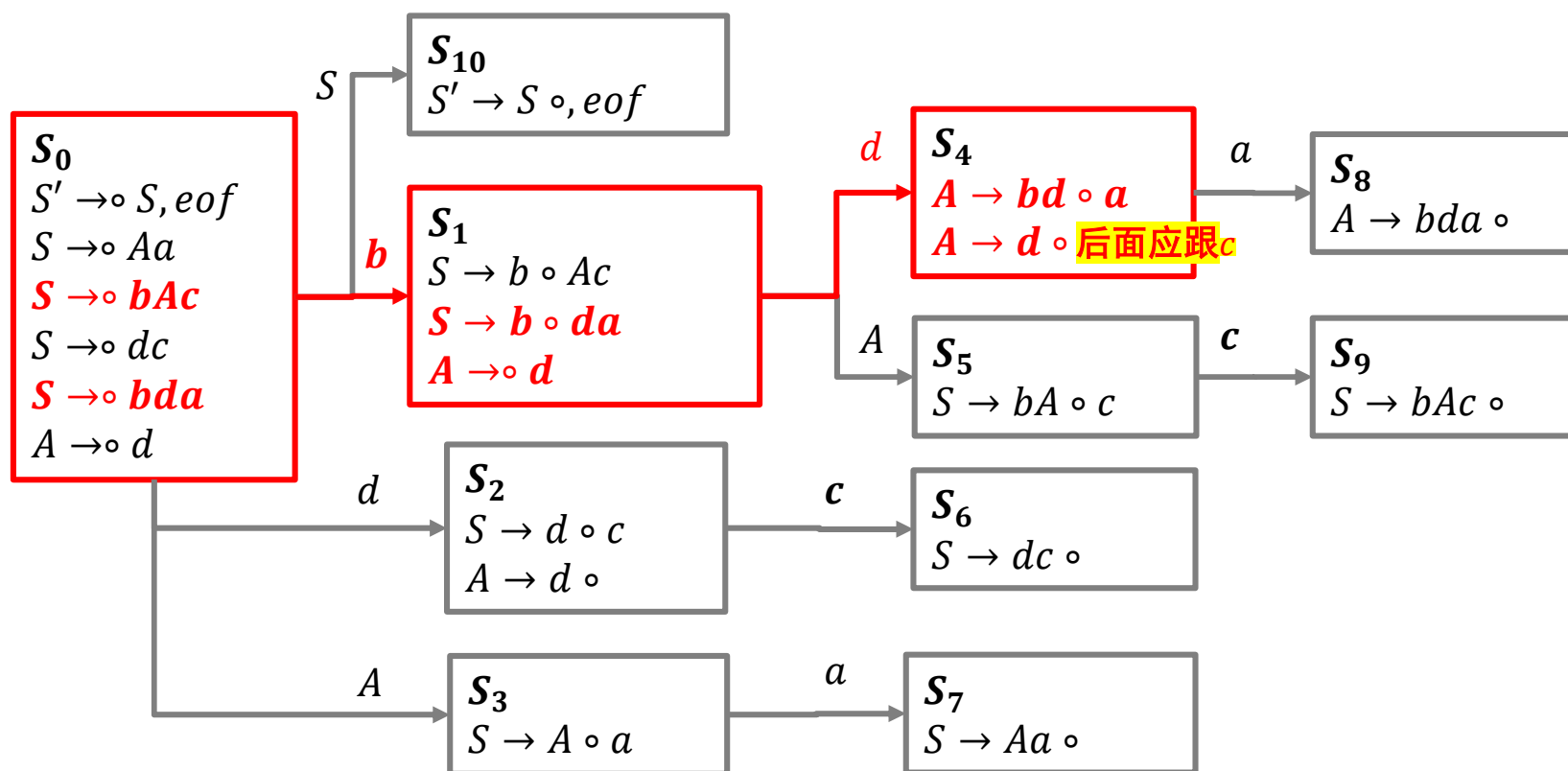
$First^+(S \rightarrow A) = \{a\}$



## 二义性语法：移进-规约冲突

- 构造SLR解析表则解析bda时存在移进-规约冲突
- $S_4$ 下一个字符为 $a$ ，可移进
- $a \in \text{Follow}(A)$ ，可规约

[1]	$S \rightarrow Aa$
[2]	$\quad  bAc$
[3]	$\quad  dc$
[4]	$\quad  bda$
[5]	$A \rightarrow d$



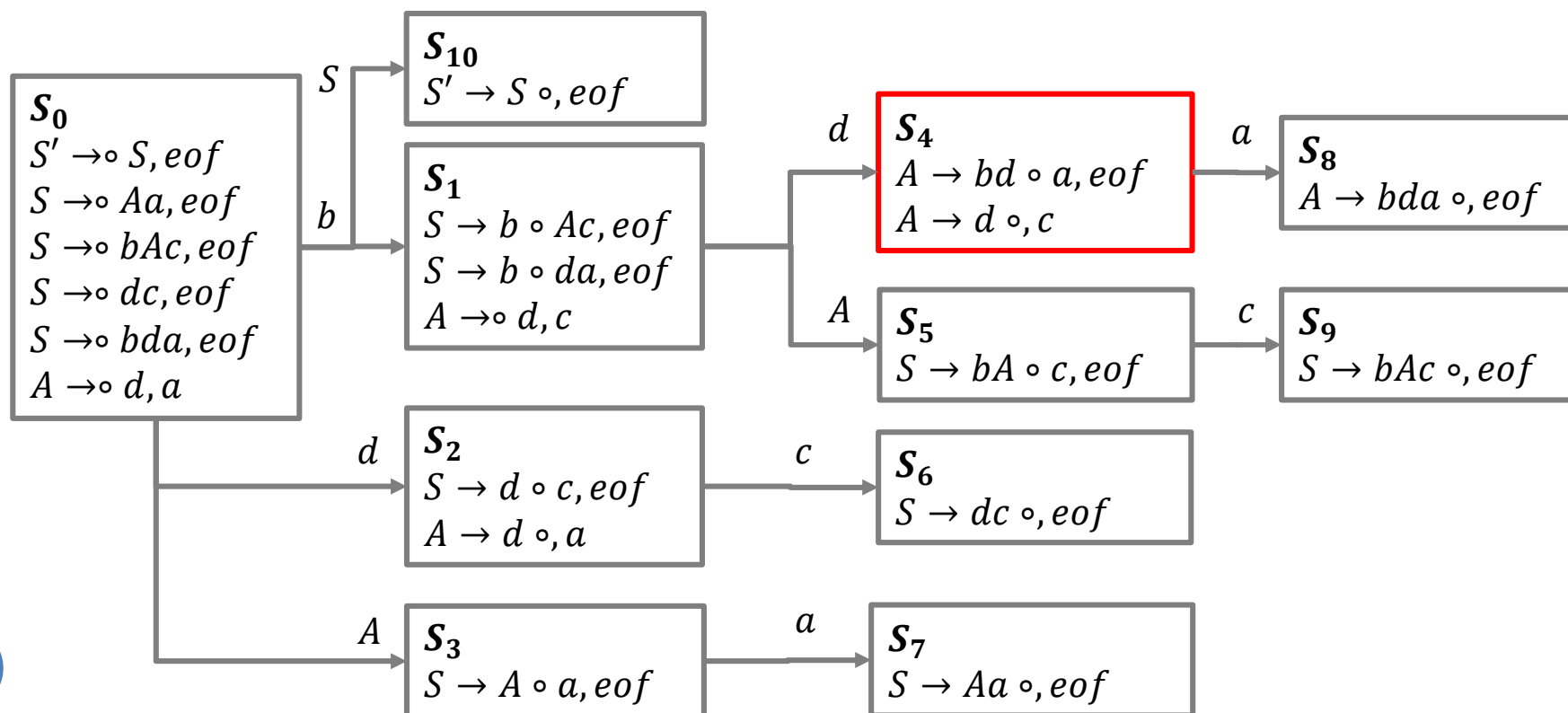
# LR(1)自动机构造

[0]  $S' \rightarrow S$   
 [1]  $S \rightarrow Aa$   
 [2]  $\quad |bAc$   
 [3]  $\quad |dc$   
 [4]  $\quad |bda$   
 [5]  $A \rightarrow d$

1) 构造其LR(1)项的全集

2) 迭代过程

- 通过闭包找到规范族
- 分析规范族之间的状态转移关系



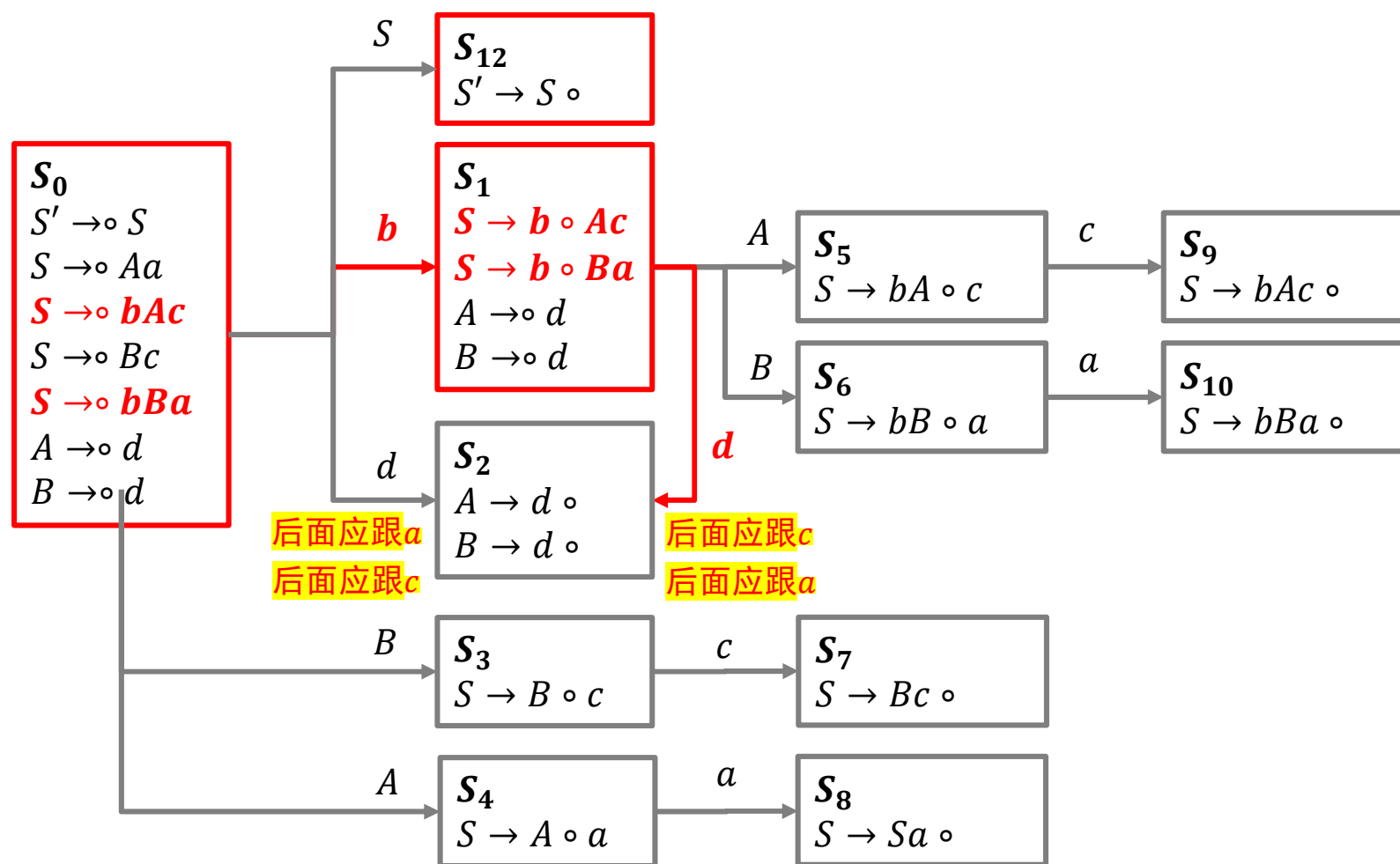
# 得到LR(1)解析表

规范项	Action					Goto	
	a	b	c	d	eof	S	A
S <sub>0</sub>		shift S <sub>1</sub>		shift S <sub>2</sub>		S <sub>10</sub>	S <sub>3</sub>
S <sub>1</sub>				shift S <sub>4</sub>			S <sub>5</sub>
S <sub>2</sub>	reduce [5]		shift S <sub>6</sub>				
S <sub>3</sub>	shift S <sub>7</sub>						
S <sub>4</sub>	shift S <sub>8</sub>		reduce [5]				
S <sub>5</sub>			shift S <sub>9</sub>				
S <sub>6</sub>					reduce [3]		
S <sub>7</sub>					reduce [1]		
S <sub>8</sub>					reduce [4]		
S <sub>9</sub>					reduce [2]		
S <sub>10</sub>					accept		

# 二义性语法：规约-规约冲突

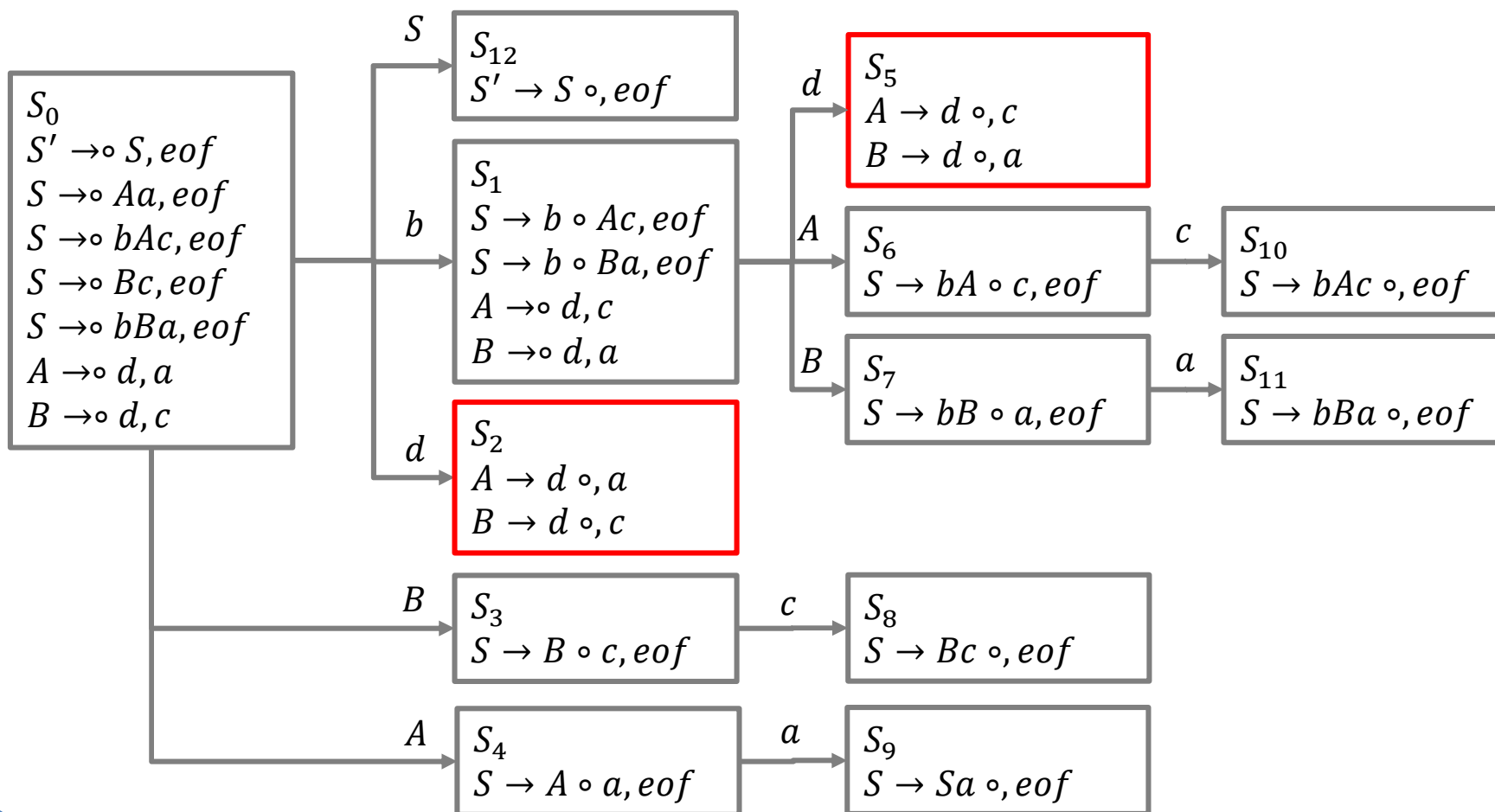
- 解析bdc时存在规约( $A \rightarrow d$ )-规约( $B \rightarrow d$ )冲突
- $a \in \text{Follow}(A)$  且  $a \in \text{Follow}(B)$

[1]	$S \rightarrow Aa$
[2]	$ bAc$
[3]	$ Bc$
[4]	$ bBa$
[5]	$A \rightarrow d$
[6]	$B \rightarrow d$



# LR(1)自动机构造

[1]  $S \rightarrow Aa$   
 [2]  $\quad \quad |bAc$   
 [3]  $\quad \quad |Bc$   
 [4]  $\quad \quad |bBa$   
 [5]  $A \rightarrow d$   
 [6]  $B \rightarrow d$





# 如何选取移进-规约操作？

- 根据当前的栈顶句柄信息：SLR (Simple LR)
  - 通过构造LR( $\theta$ )自动机和下一个字符判断是否可以移进
  - 需要规约时根据Follow判断是否可行
- 自动机构造时考虑Follow信息：经典LR(1):
  - 问题一：LR(1) 支持语法有限，用GLR (Generalized LR)
    - 如果遇到冲突则分叉：复制栈状态，尝试不同规则或移进操作
  - 问题二：LR(1) 规范集太多，用LALR (Lookahead LR)
    - 自动机构造时考虑Follow信息
    - 同时精简规范集

# LALR

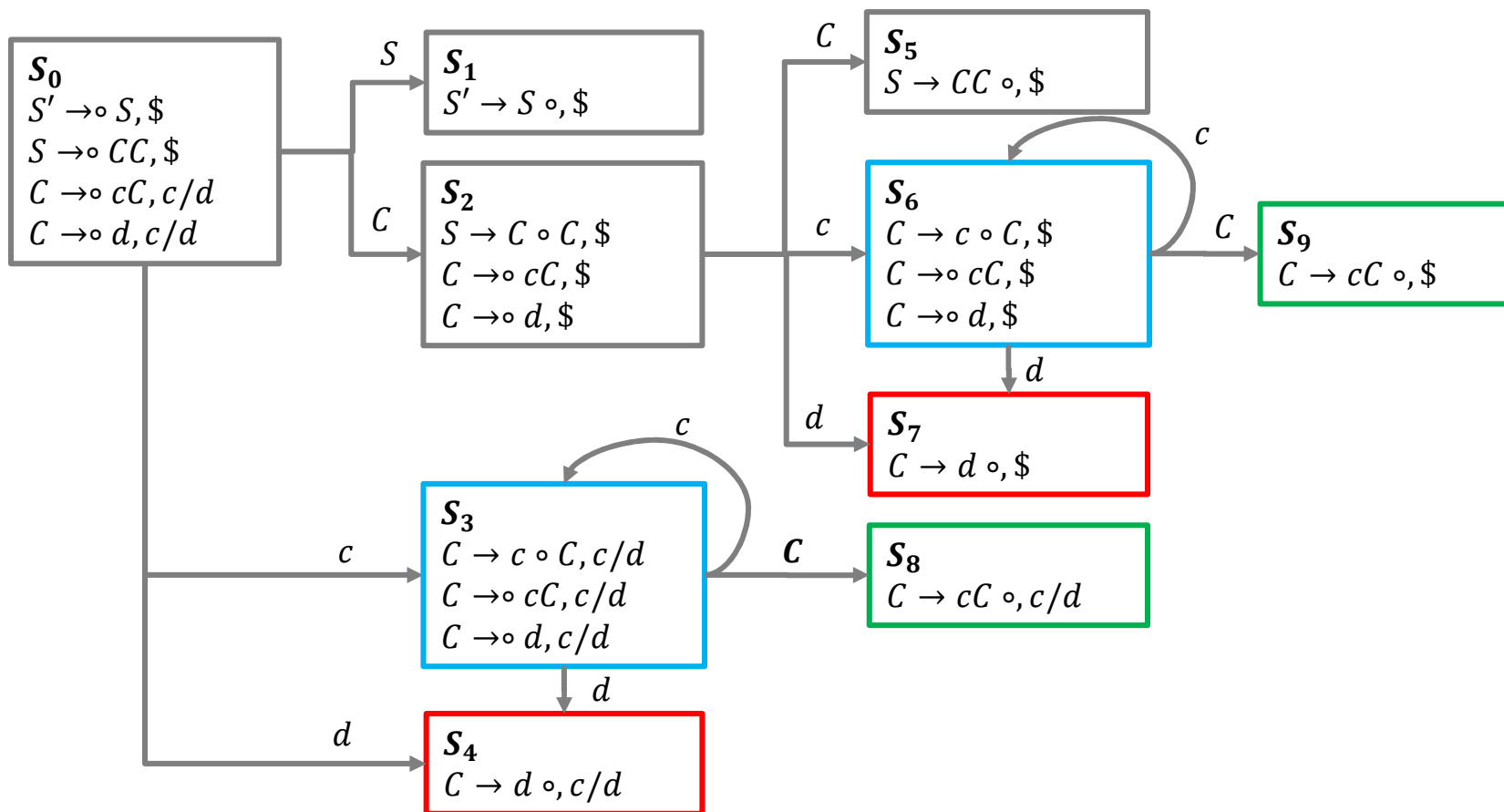
- SLR(1)存在移进-规约、规约-规约冲突问题，支持的语法范围小；
- LR(1)在规范集构造时融合了Follow信息，可以避免很多冲突问题，但解析表可能会比较大；
- LALR是一种折中方法，解析表大小和SLR相同；
- LALR构造思路：合并句柄状态完全相同的状态集

# LALR语法举例

- [1]  $S' \rightarrow S$
- [2]  $S \rightarrow CC$
- [3]  $C \rightarrow cC$
- [4]  $\quad \quad | d$

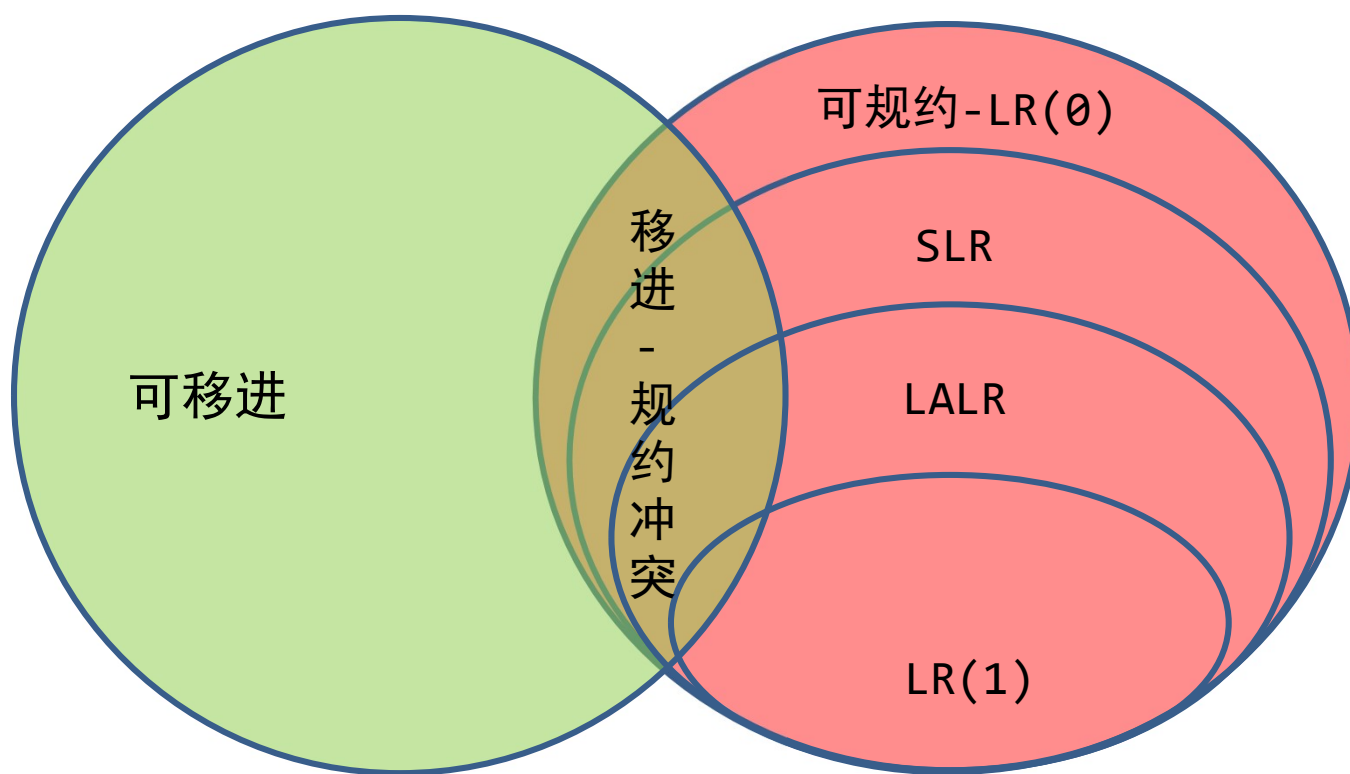
• 可以合并的规范集

- S3和S6、S4和S7、S8和S9。



# 几种语法的关系

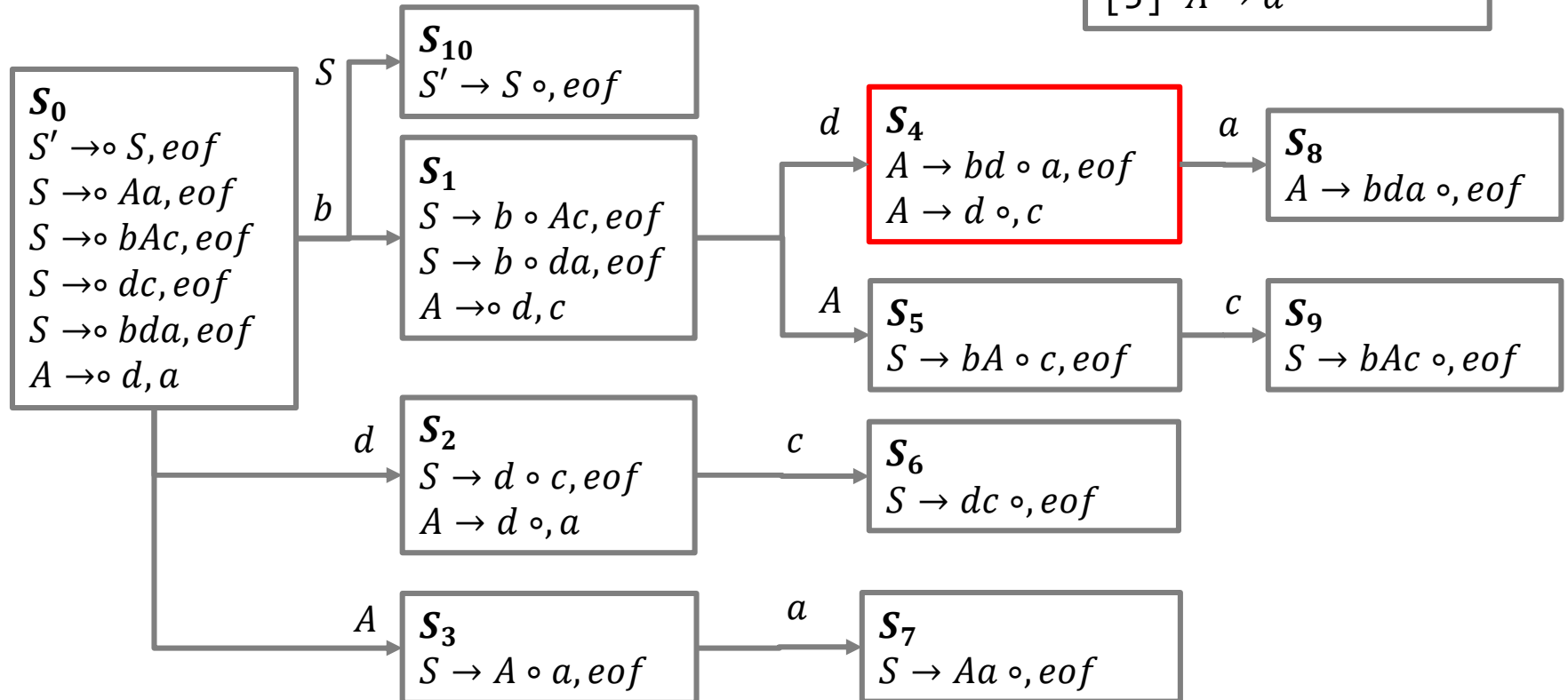
- 语法表达能力:  $LR(1) > LALR(1) > SLR$ 
  - 规约条件严苛:  $LR(1) > LALR(1) > SLR$
  - 移进条件同 $LR(0)$ ?



# 举例：LALR，非SLR(1)语法

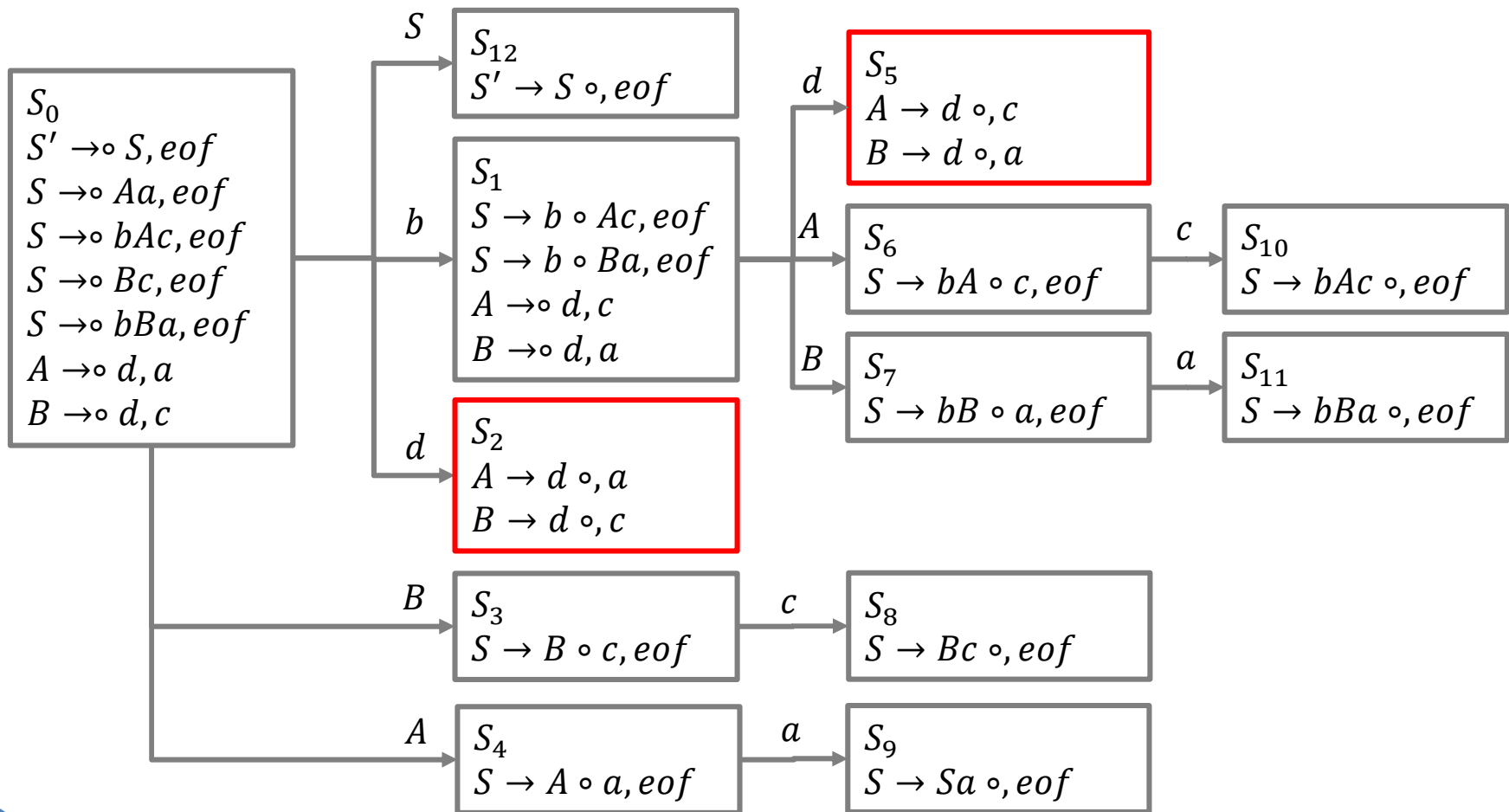
- 构造SLR解析表则解析bda时存在移进-规约冲突
- LALR解析方法可以避免冲突

[1]	$S \rightarrow Aa$
[2]	$ bAc$
[3]	$ dc$
[4]	$ bda$
[5]	$A \rightarrow d$



# 举例：LR(1)，非LALR语法

[1]  $S \rightarrow Aa$   
 [2]  $\quad \quad |bAc$   
 [3]  $\quad \quad |Bc$   
 [4]  $\quad \quad |bBa$   
 [5]  $A \rightarrow d$   
 [6]  $B \rightarrow d$



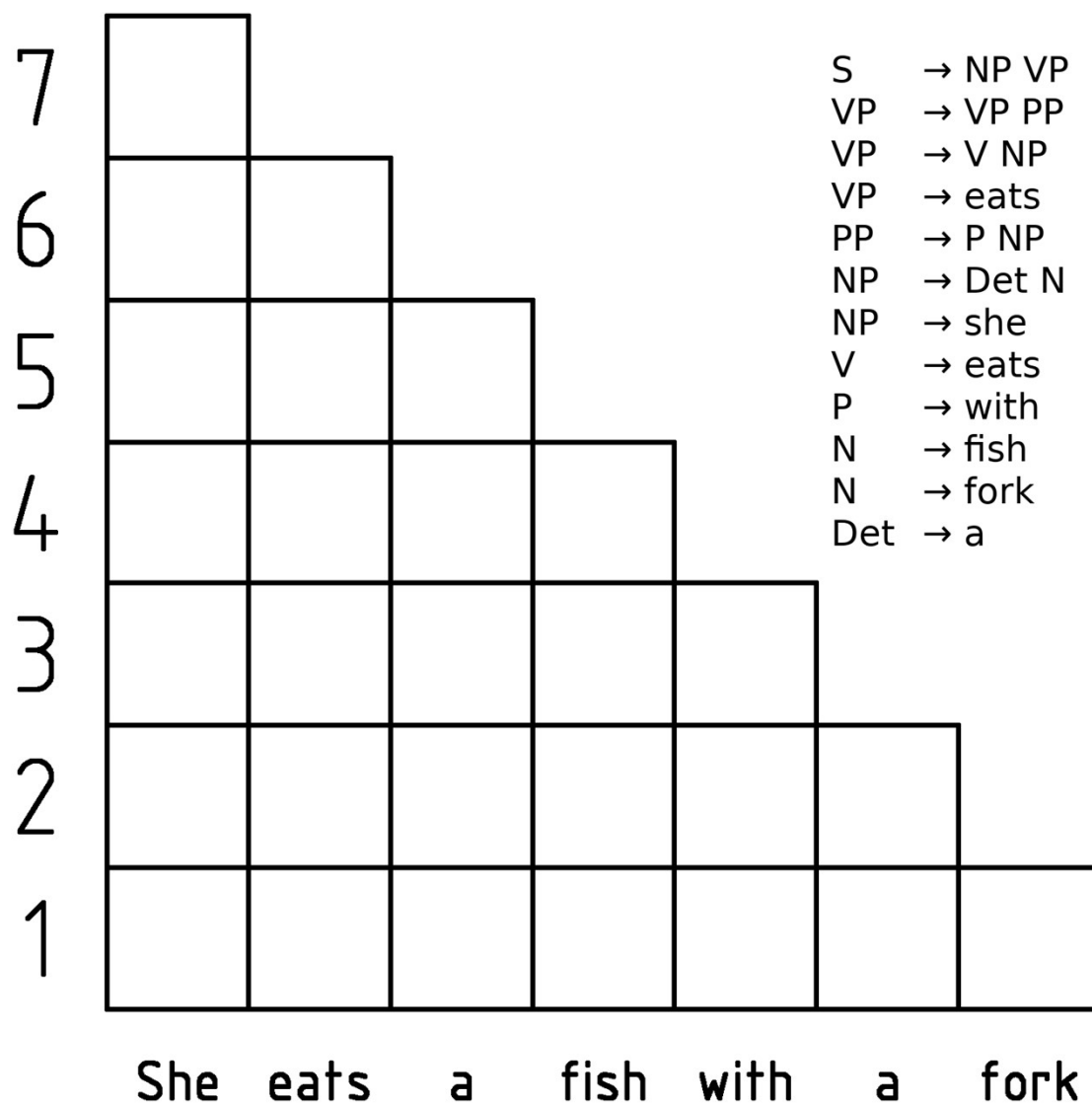
# 思考

- 下列语法是否是SLR或LR(1) ?

```
< regex > ::= < union > | < concat >  
< union > ::= < regex > "|" < concat >  
< concat > ::= < concat > < term > | < term >  
< term > ::= < element > * | < element >  
< element > ::= (< regex > ) | < alphanum >
```

```
[1] < regex > ::= < concat > < regex' >  
[2] < regex' > ::= "|" < concat > < regex' >  
[3]           | $\epsilon$   
[4] < concat > ::= < term > < concat' >  
[5] < concat' > ::= < term > < concat' >  
[6]           | $\epsilon$   
[7] < term > ::= < element > < follow >  
[8] < follow > ::= *  
[9]           | $\epsilon$   
[10] < element > ::= (< regex > )  
[11]           | < alphanum >
```

# 通用自底向上CFG分析： CYK算法





# CYK解析算法伪代码参考

INIT:

Grammar:  $R_1, R_2, \dots, R_r$

String to parse:  $w = w_1, w_2, \dots, w_l$

$P[n, n, r]$  initied with false

Foreach  $i = 1$  to  $n$ :

    Foreach  $R_r \rightarrow a_i$

$P[1, i, r] = \text{True}$

Foreach  $l = 2$  to  $n$ :

    Foreach  $i = 1$  to  $n-l+1$ :

        Foreach  $j = 1$  to  $l-1$ :

            Foreach  $R_r \rightarrow R_a R_b$

                If  $P[j, i, a]$  and  $P[l-j, i+j, b]$ :

$P[l, i, r] = \text{True}$

If  $P[n, 1, 1]$ :

$w$  is a string of the language

Else:

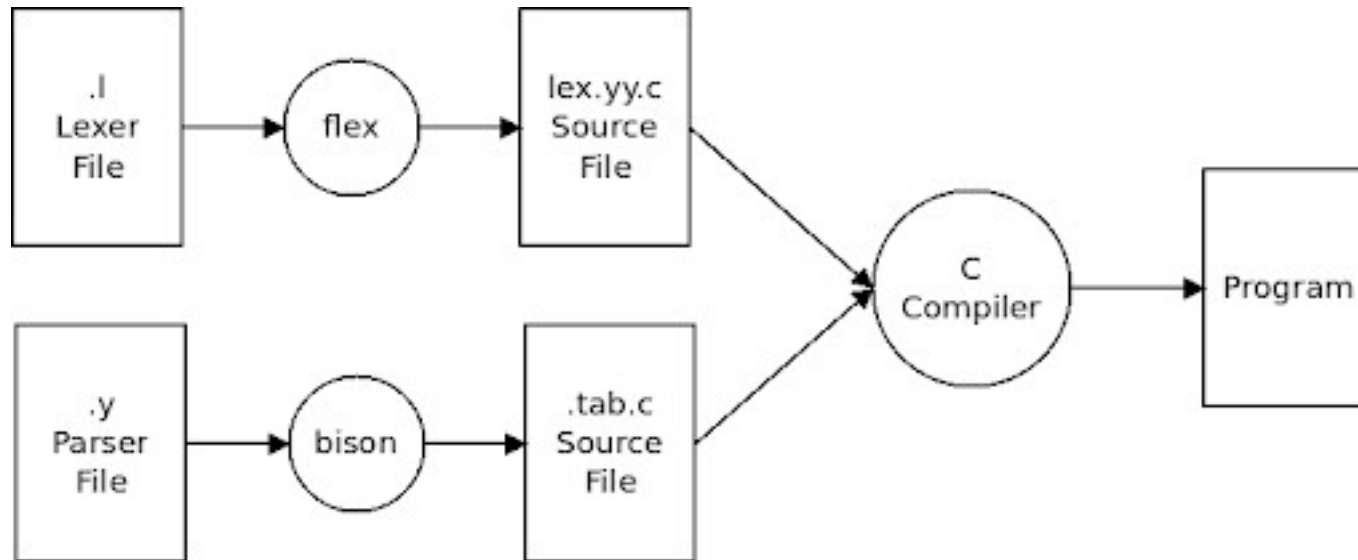
$w$  is not a string of the language

## 五、语法分析工具

---

# Bison

- 语法分析工具YACC(POSIX)/Bison (GNU)
  - 默认采用LALR(1)解析
  - 支持LR(1)等方法



# 计算器程序示例

Main.c 计算器文件  
Lexer.l 词法定义  
Parser.y 语法定义  
Expression.h 文件  
Expression.c 功能函数

# Lex文件

```
%{ /* Lexer.l file */
#include "Expression.h"
#include "Parser.h"
#include <stdio.h> %}

%option outfile="Lexer.c"
header-file="Lexer.h"
%option warn nodefault
%option reentrant noyywrap never-interactive nounistd
%option bison-bridge

%%

[ \r\n\t]* { continue; /* Skip blanks. */ }
[0-9]+ { sscanf(yytext, "%d", &yyval->value); return TOKEN_NUMBER; }
"*" { return TOKEN_STAR; }
"+" { return TOKEN_PLUS; }
"(" { return TOKEN_LPAREN; }
")" { return TOKEN_RPAREN; }
. { continue; /* Ignore unexpected characters. */}

%%

int yyerror(const char *msg) {
    fprintf(stderr, "Error: %s\n", msg);
    return 0;
}
```

```

%{ /* * Parser.y file * */
#include "Expression.h"
#include "Parser.h"
#include "Lexer.h"
int yyerror(SExpression **expression, yyscan_t scanner, const char *msg) { /* Add error
handling routine as needed */ }
%}

%code requires { typedef void* yyscan_t; }
%output "Parser.c"
%defines "Parser.h"
%define api.pure
%lex-param { yyscan_t scanner }
%parse-param { SExpression **expression }
%parse-param { yyscan_t scanner }
%union { int value; SExpression *expression; }
%token TOKEN_LPAREN "("
%token TOKEN_RPAREN ")"
%token TOKEN_PLUS "+"
%token TOKEN_STAR "*"
%token <value> TOKEN_NUMBER "number"
%type <expression> expr

/* Precedence (increasing) and associativity */
%left "+"
%left "*" %

%%
input : expr { *expression = $1; } ;
expr : expr[L] "+" expr[R] { $$ = createOperation( eADD, $L, $R ); }
    | expr[L] "*" expr[R] { $$ = createOperation( eMULTIPLY, $L, $R ); }
    | "(" expr[E] ")" { $$ = $E; }
    | "number" { $$ = createNumber($1); }
} ; %%

```

# main.c

```
int yyparse(SExpression **expression, yyscan_t scanner);
SExpression *getAST(const char *expr) {
    SExpression *expression;
    yyscan_t scanner;
    YY_BUFFER_STATE state;
    if (yylex_init(&scanner)) return NULL;
    state = yy_scan_string(expr, scanner);
    if (yyparse(&expression, scanner)) return NULL;
    yy_delete_buffer(state, scanner);
    yylex_destroy(scanner);
    return expression;
}

int evaluate(SExpression *e) {
    switch (e->type) {
        case eVALUE: return e->value;
        case eMULTIPLY: return evaluate(e->left) * evaluate(e->right);
        case eADD: return evaluate(e->left) + evaluate(e->right);
        default: /* should not be here */ return 0;
    }
}

int main(void) {
    char expr[256];
    scanf("%s", expr);
    SExpression *e = getAST(test);
    int result = evaluate(e);
    printf("Result of '%s' is %d\n", test, result);
    deleteExpression(e);
    return 0;
}
```

# Expression.h

```
#ifndef __EXPRESSION_H__
#define __EXPRESSION_H__

typedef enum tagEOperationType {
    eVALUE,
    eMULTIPLY,
    eADD
} EOperationType;

typedef struct tagSExpression {
    EOperationType type; /* /< type of operation */
    int value; /* /< valid only when type is eVALUE */
    struct tagSExpression *left; /* /< left side of the tree */
    struct tagSExpression *right; /* /< right side of the tree */
} SExpression;

SExpression *createNumber(int value);
SExpression *createOperation(EOperationType type, SExpression *left, SExpression *right);
void deleteExpression(SExpression *b);

#endif
```



# Expression.c

```
static SExpression *allocateExpression() {
    SExpression *b = (SExpression *)malloc(sizeof(SExpression));
    if (b == NULL) return NULL;
    b->type = eVALUE;
    b->value = 0;
    b->left = NULL;
    b->right = NULL;
    return b;
}

SExpression *createNumber(int value) {
    SExpression *b = allocateExpression();
    if (b == NULL) return NULL;
    b->type = eVALUE;
    b->value = value;
    return b;
}

SExpression *createOperation(EOperationType type, SExpression *left, SExpression *right) {
    SExpression *b = allocateExpression();
    if (b == NULL) return NULL;
    b->type = type;
    b->left = left;
    b->right = right;
    return b;
}

void deleteExpression(SExpression *b) {
    if (b == NULL) return;
    deleteExpression(b->left);
    deleteExpression(b->right);
    free(b);
}
```

# 总结

一、句式分析的基本概念

二、LLVM案例分析

三、自顶向下分析

- LL(1)语言
- 通用算法：Earley算法

四、自底向上分析

- SLR、LALR、LR(1)语言
- 通用算法：GLR算法、CYK算法

五、语法分析工具