



# 第4章 语法分析 (1)

---

Syntax Analysis

【对应教材2.2, 4.1-4.2】



# 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
- 语法分析器生成工具YACC



# 程序设计语言构造的描述

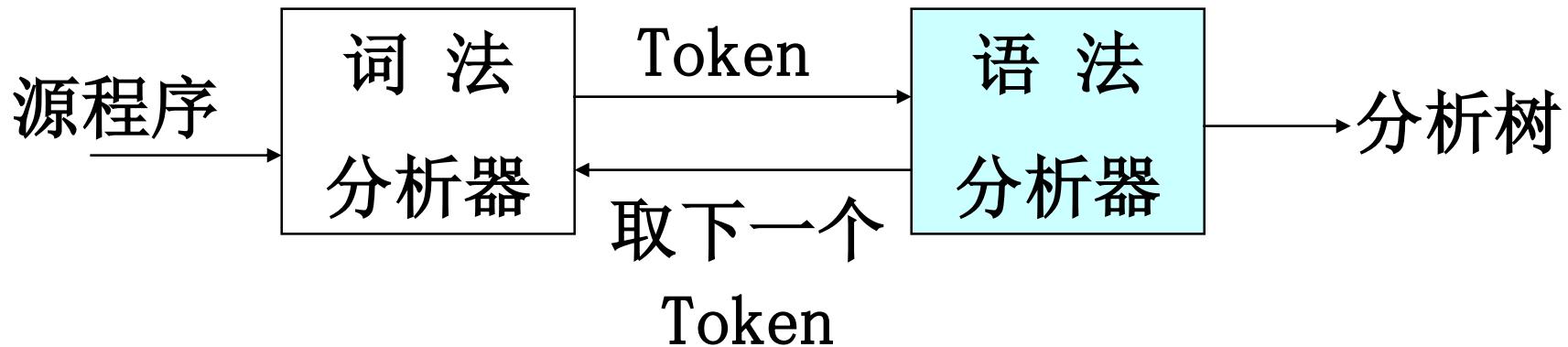
□ 程序设计语言构造的语法可使用上下文无关文法或BNF表示法来描述

- 文法可以给出精确易懂的语法规则
- 可以自动构造出某些类型文法的语法分析器
- 文法指出了语言的结构，有助于进一步的语义处理和代码生成
- 支持语言的演化和迭代

- 文法：Grammar
- 上下文无关文法：Context-Free Grammar, CFG



# 语法分析器的作用



- **功能：**根据文法规则，从源程序单词符号串中识别出语法成分，并进行语法检查
- **基本任务：**识别符号串S是否为某个合法的语法单元



# 语法分析器的种类

## □ 通用语法分析器

- 可以对任意文法进行语法分析
- 效率很低，不适合用于编译器

## □ 自顶向下的语法分析器

- 从语法分析树的根部开始构造语法分析树

## □ 自底向上的语法分析器

- 从语法分析树的叶子开始构造语法分析树

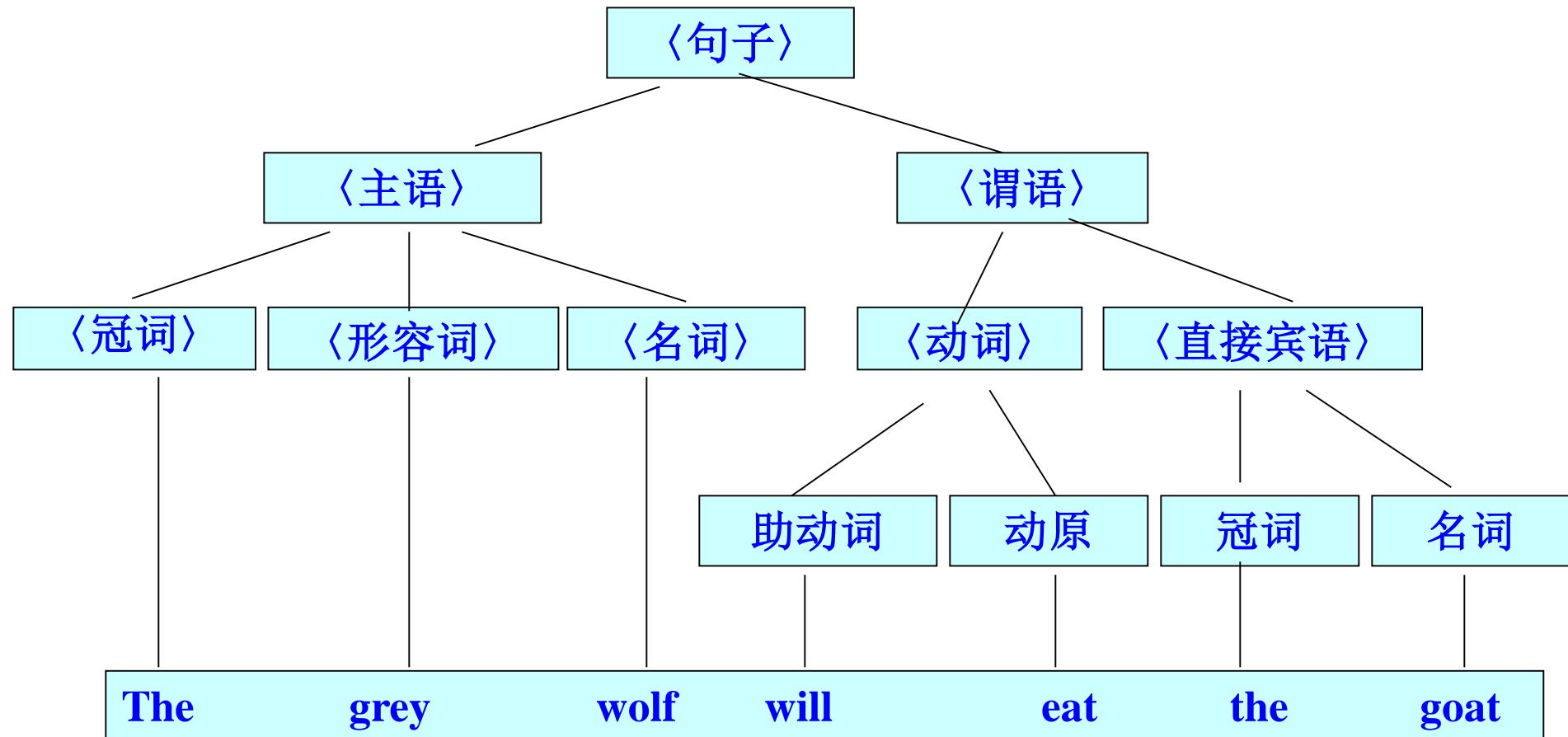
## □ 后两种方法

- 通常总是从左到右、逐个扫描词法单元
- 为了保证效率，只针对特定类型的文法，但是这些文法足以用来描述常见的程序设计语言



# 语言和文法基础：例子

分析：The grey wolf will eat the goat





- 为了进行机器分析，“句子由主语后跟随谓语组成”可以表示为：

<句子> $\rightarrow$  <主语> <谓语> (1)

<主语> $\rightarrow$  <冠词> <形容词> <名词> (2)

<冠词> $\rightarrow$  the (3)

<形容词> $\rightarrow$  grey (4)

<谓语> $\rightarrow$  <动词> <直接宾语> (5)

<动词> $\rightarrow$  <助动词> <动词原形> (6)

<助动词> $\rightarrow$  will (7)

<动词原形> $\rightarrow$  eat (8)

<直接宾语> $\rightarrow$  <冠词> <名词> (9)

<名词> $\rightarrow$  wolf (10)

<名词> $\rightarrow$  goat (11)



# 句子的语法 (Syntax)

- 终结符号集  $V_T = \{\text{the}, \text{grey}, \text{wolf}, \text{will}, \text{eat}, \text{goat}\}$
- 非终结符号集  $V_N = \{\langle\text{句子}\rangle, \langle\text{主语}\rangle, \langle\text{谓语}\rangle, \langle\text{冠词}\rangle, \langle\text{形容词}\rangle, \langle\text{名词}\rangle, \langle\text{动词}\rangle, \langle\text{直接宾语}\rangle, \langle\text{助动词}\rangle, \langle\text{动词原形}\rangle\}$
- 开始符号  $S = \langle\text{句子}\rangle$
- 产生规则集  $P = \{\langle\text{句子}\rangle \rightarrow \langle\text{主语}\rangle \langle\text{谓语}\rangle, \dots\}$



# 上述句子可根据规则得出

<句子>⇒<主语> <谓语>  
⇒ <冠词> <形容词> <名词> <谓语>  
⇒ the <形容词> <名词> <谓语>  
⇒ the grey <名词> <谓语>  
⇒ the grey wolf <谓语>  
⇒ the grey wolf <动词> <直接宾语>  
⇒ .....  
⇒ the grey wolf will eat the goat



# 句子既要符合语法规则又要符合语义规定

<句子>  $\stackrel{+}{\Rightarrow}$  the grey wolf will eat the goat  
the grey wolf will eat the wolf  
the grey goat will eat the goat  
the grey goat will eat the wolf

符合语法规则且符合语义规定的句子仅是：

**the grey wolf will eat the goat**



# 文法 (Grammar) 的正式定义

文法  $G = (V_T, V_N, S, P)$ , 其中:

- $V_T$  是一个非空有穷的终结符号(terminal)集合
- $V_N$  是一个非空有穷的非终结符号(nonterminal)集合,  
且  $V_T \cap V_N = \emptyset$
- $P = \{ \alpha \rightarrow \beta \mid \alpha \in (V_T \cup V_N)^* \text{ 且至少包含一个非终结符} \beta \in (V_T \cup V_N)^* \}$ , 称为产生式(production)集合
- $S \in V_N$ , 称为开始符号(start symbol)
  - $S$ 必须在某个产生式的左部至少出现一次

产生式可以写成  $A ::= \alpha$  或  $A \rightarrow \alpha$

$A \rightarrow \alpha_1 \ A \rightarrow \alpha_2$  可以缩写为:  $A \rightarrow \alpha_1 | \alpha_2$



# 上下文无关文法

- Context-free grammar, 简称CFG
- 所有产生式的左边只有一个非终结符号, 即
  - 产生式的形式为:  $A \rightarrow \beta$
  - 因此不需要任何上下文 (context) 就可以对A进行推导
- 上下文无关文法描述的语言称为上下文无关语言



# BNF范式 (Backus-Naur Form)

- BNF广泛用于描述现代语言语法的形式化表示
- 1959年：John **Backus** 发明的一种用来描述 ALGOL 60 语法规则的标记法
- 1963年：Peter **Naur** 把它称为**Backus Normal Form**, 并进行了化简。
  - 缩小了其中使用的字符集（用`::=`代替了`→`）
- Donald Knuth 最早把它改名为**Backus-Naur Form**
  - Knuth, Donald E. (1964). “Backus Normal Form vs. Backus Naur Form”. *Communications of the ACM* 7 (12): 735–736.



# 例：用BNF描述的BNF语法

```
<syntax> ::= <rule> | <rule> <syntax>
<rule> ::= <opt-whitespace> "<" <rule-name> ">"
           <opt-whitespace> " ::= " <opt-whitespace>
           <expression> <line-end>
<opt-whitespace> ::= " " <opt-whitespace> | ""
                   <!-- "" is empty string, i.e. no whitespace -->
<expression> ::= <list> | <list> " | " <expression>
<line-end> ::= <opt-whitespace> <EOL> |
               <line-end> <line-end>
<list> ::= <term> | <term> <opt-whitespace> <list>
<term> ::= <literal> | "<" <rule-name> ">"
<literal> ::= ' ' <text> ' ' | ' ' <text> ' '
                  <!-- actually, the original BNF did not use quotes -->
```



# 例：算术表达式的文法G

$$G = (\{a, +, *, (), \}, \{\text{<表达式>}, \text{<项>}, \text{<因子>}\}, \\ \text{<表达式>, P})$$

P: <表达式> $\rightarrow$  <表达式> + <项> | <项>  
<项> $\rightarrow$  <项> \* <因子> | <因子>  
<因子> $\rightarrow$  ( <表达式>) | a

通常可以简写为 (G1[E]):

$E \rightarrow E + T$		T
$T \rightarrow T * F$		F
$F \rightarrow ( E )$		a

Expression  
Term  
Factor



# 关于文法的一些约定

- 通常可以不用将文法G的四元组显式地表示出来，而只需将产生式写出
- 一般约定：
  - 第一条产生式的左部是开始符号
  - 用尖括号括起来的是非终结符号，而不用尖括号的是终结符号，或者
  - 大写字母表示非终结符号，小写字母表示终结符号
  - 小写的希腊字母表示（可能为空的）文法符号串
  - 另外也可以把G表示为G[S]，其中S为开始符号



# 直接推导(Immediate Derivation)

令  $G = (V_T, V_N, S, P)$ , 若  $\alpha \rightarrow \beta \in P$ , 且  
 $\gamma, \delta \in (V_T \cup V_N)^*$ , 则称  $\gamma\alpha\delta$  可以直接推导出  $\gamma\beta\delta$ ,  
表示成  $\frac{\gamma\alpha\delta \Rightarrow \gamma\beta\delta}{}$

若  $\gamma\alpha\delta$  直接推导出  $\gamma\beta\delta$ , 即:

$\frac{\gamma\alpha\delta \Rightarrow \gamma\beta\delta}{\text{则称 } \gamma\beta\delta \text{ 直接归约到 } \gamma\alpha\delta}$

归约: **reduce (vi)**、**reduction(n.)** 是推导的逆过程



# 推导(Derivation)

一个直接推导序列：

$$\alpha_0 \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n \quad (n > 0)$$

可表示成  $\alpha_0 \Rightarrow^+ \alpha_n$

$\alpha_0 \Rightarrow^* \alpha_n$  定义为：或者  $\alpha_0 = \alpha_n$

或者  $\alpha_0 \Rightarrow^+ \alpha_n$



# 例：

$$E \Rightarrow E+T \Rightarrow T+T \Rightarrow F+T \Rightarrow a+T \Rightarrow a+F \Rightarrow a+a$$

$\alpha A\beta$	$\alpha\gamma\beta$	产生式	$\alpha$	$\beta$
E	$\Rightarrow E+T$	$E \rightarrow E+T$	$\epsilon$	$\epsilon$
E+T	$\Rightarrow T+T$	$E \rightarrow T$	$\epsilon$	+T
T+T	$\Rightarrow F+T$	$T \rightarrow F$	$\epsilon$	+T
F+T	$\Rightarrow a+T$	$F \rightarrow a$	$\epsilon$	+T
a+T	$\Rightarrow a+F$	$T \rightarrow F$	a+	$\epsilon$
a+F	$\Rightarrow a+a$	$F \rightarrow a$	a+	$\epsilon$



# 最左推导和最右推导

对于 $w$ 和 $G$ ,  $w \in L(G)$ , 是否存在 $S \xrightarrow{*} w$ ? 如何构造这个推导?

例如,  $G[E]$  (表达式文法) 和  $w = a + a * a$

$$\begin{array}{c} E \xrightarrow{\underline{lm}} E+T \xrightarrow{\underline{lm}} T+T \xrightarrow{\underline{lm}} F+T \xrightarrow{\underline{lm}} a+T \xrightarrow{\underline{lm}} a+T*F \\ \xrightarrow{\underline{lm}} a+F*F \xrightarrow{\underline{lm}} a+a*F \xrightarrow{\underline{lm}} a+a*a \end{array}$$

特点:  $\alpha A \beta \xrightarrow{\underline{lm}} \alpha \gamma \beta$  ( $A \rightarrow \gamma$ ),  $\alpha \in V_T^*$  (最左)

$$\begin{array}{c} E \xrightarrow{\underline{rm}} E+T \xrightarrow{\underline{rm}} E+T*F \xrightarrow{\underline{rm}} E+T*a \xrightarrow{\underline{rm}} E+F*a \\ \xleftarrow{\underline{rm}} E+a*a \xrightarrow{\underline{rm}} T+a*a \xrightarrow{\underline{rm}} F+a*a \xrightarrow{\underline{rm}} a+a*a \end{array}$$

特点:  $\alpha A \beta \Rightarrow \alpha \gamma \beta$  ( $A \rightarrow \gamma$ ),  $\beta \in V_T^*$  (最右)



# 句型/句子/语言

## □ 句型 (sentential form) :

- 如果  $S \Rightarrow^* a$ , 那么  $a$  是文法的句型
- 句型可能既包含非终结符号, 又包含终结符号
- 句型也可以是空串

## □ 句子 (sentence)

- 文法的句子是不包含非终结符号的句型

## □ 语言

- 文法  $G$  的语言是  $G$  的所有句子的集合, 记为  $L(G)$
- $w$  在  $L(G)$  中当且仅当  $w$  是  $G$  的句子, 即  $S \Rightarrow^* w$

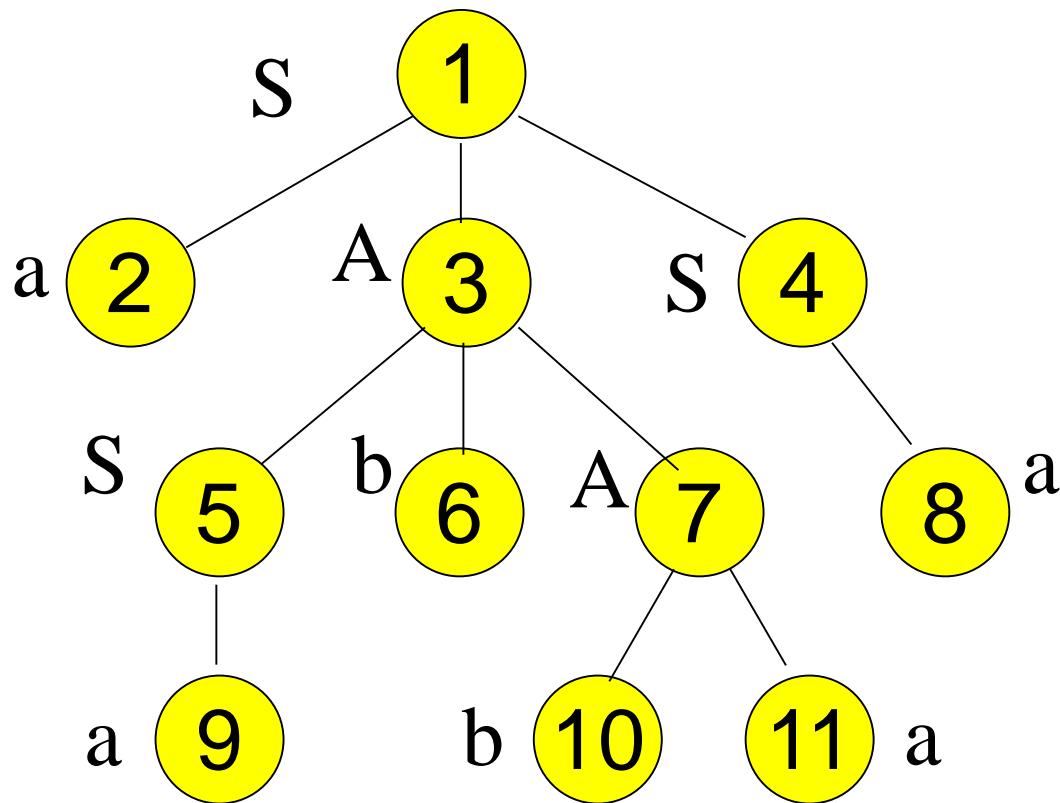


# 语法分析树 (Parse Tree)

- 语法分析树是推导的一种图形表示形式
  - 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 $\epsilon$
  - 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，该结点的子结点从左到右对应产生式的右部
- 有时允许树根不是开始符号（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

例:  $G=(V_T, V_N, S, P)$ , 其中

P:  $S \rightarrow aAS \mid a \quad A \rightarrow SbA \mid SS \mid ba$

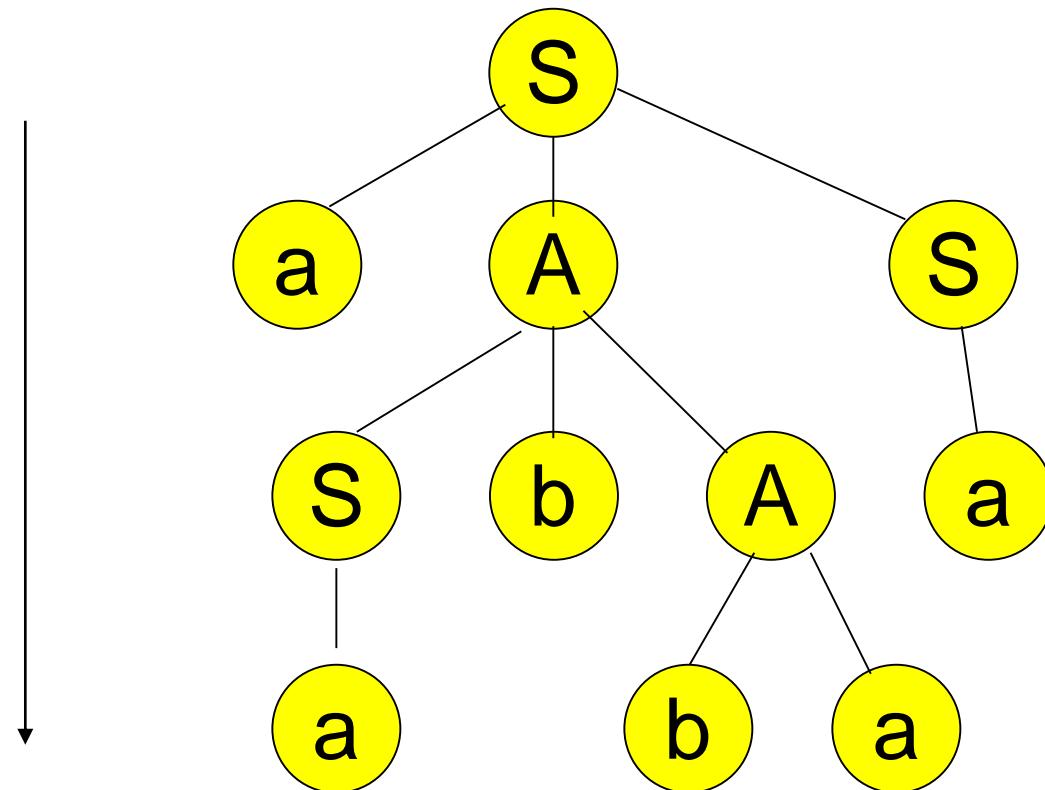




# 如何画出分析树 (1. 自顶向下)

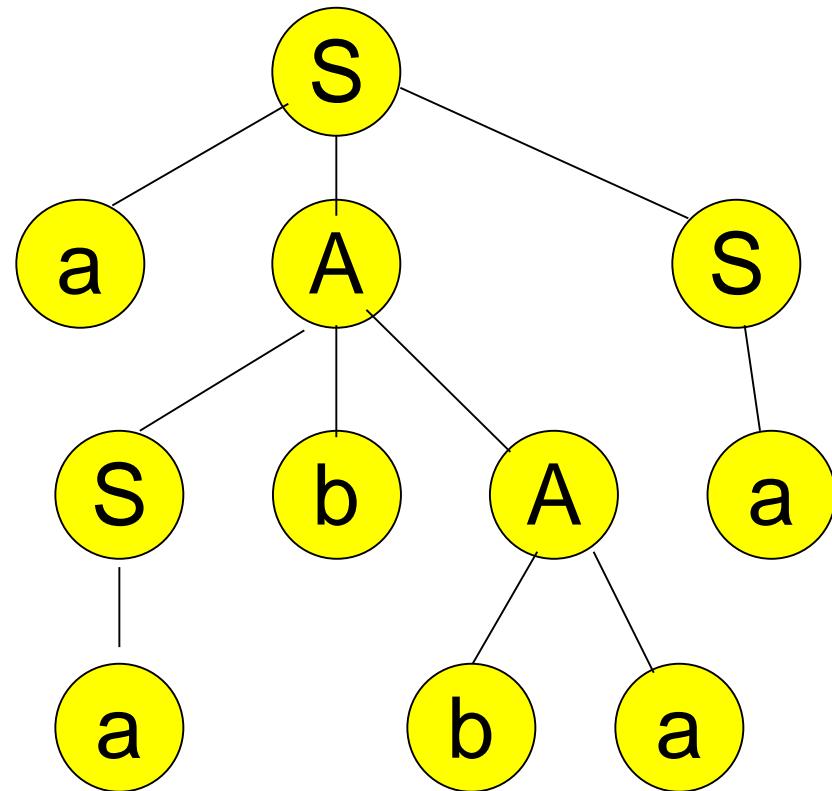
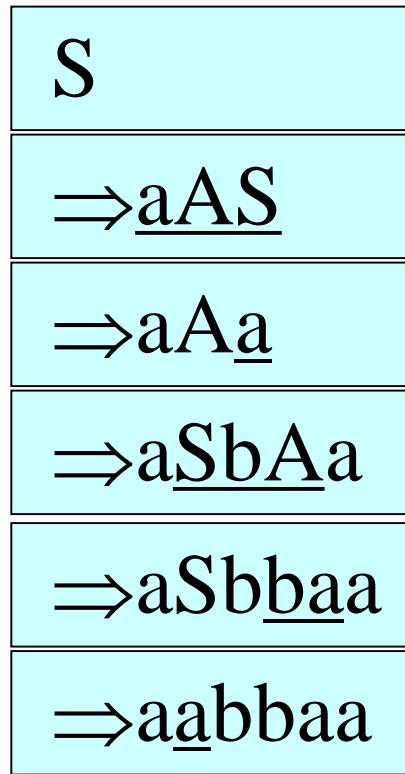
根据推导序列，对每步推导画相应分枝

S
$\Rightarrow aAS$
$\Rightarrow aSbAS$
$\Rightarrow aabAS$
$\Rightarrow aabbAS$
$\Rightarrow aabbaa$



# 如何画出分析树（2. 自底向上）

根据归约序列，对每步归约画相应分枝





# 文法的二义性 (Ambiguity)

例：卡塔尔世界杯中日韩都进了十六强。

1. 一个句子的结构可能不唯一
2. 一个句子的对应的分析树可能不唯一

考虑下面的表达式文法  $G2[E]$ , 其产生式如下:

$$E \rightarrow E+E \mid E^*E \mid (E) \mid a$$

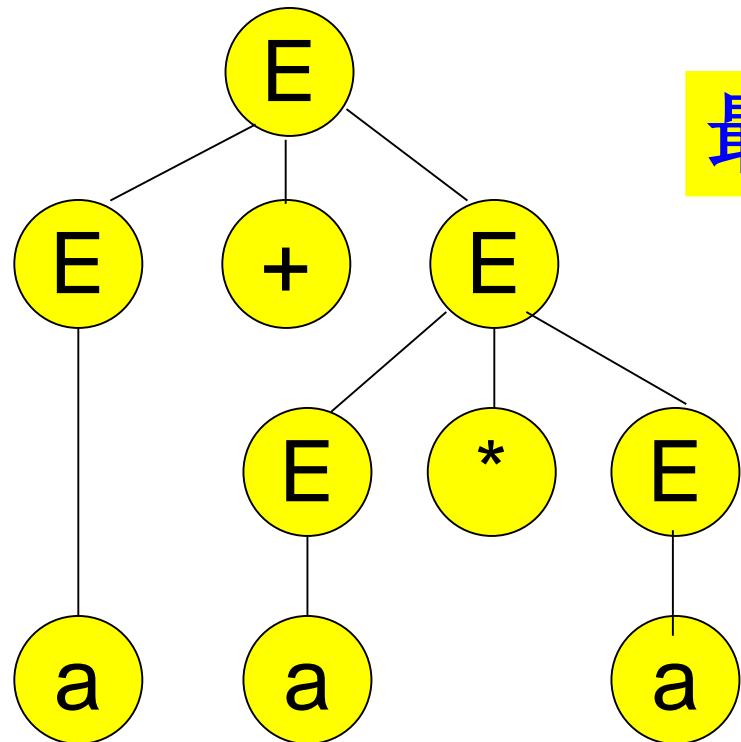
对于句子  $a+a^*a$ , 有如下两个最左推导:

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

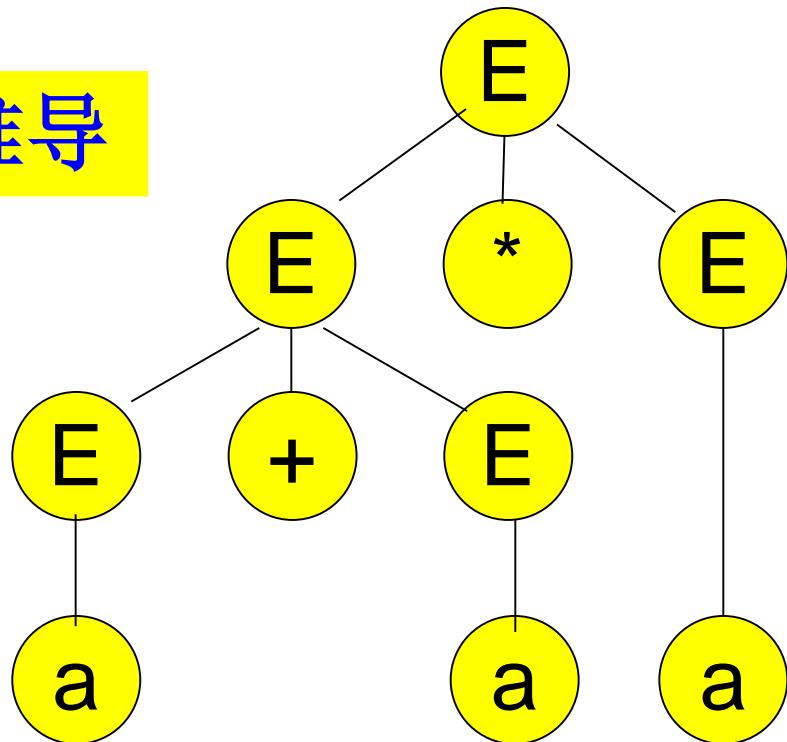
$$E \Rightarrow E^*E \Rightarrow E+E^*E \Rightarrow a+E^*E \Rightarrow a+a^*E \Rightarrow a+a^*a$$

$E \Rightarrow E+E \Rightarrow a+E$   
 $\Rightarrow a+E*E \Rightarrow a+a*E$   
 $\Rightarrow a+a*a$

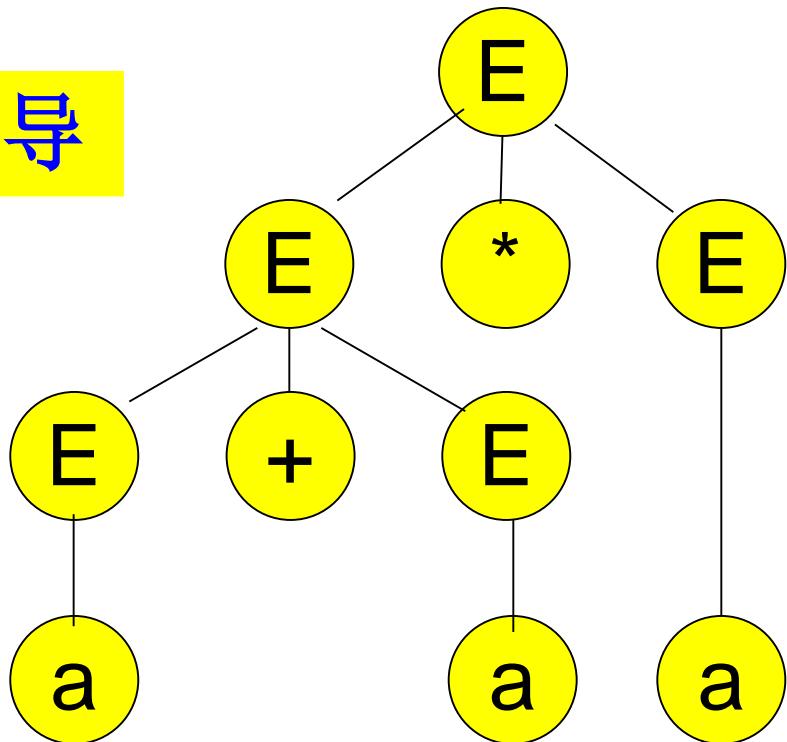
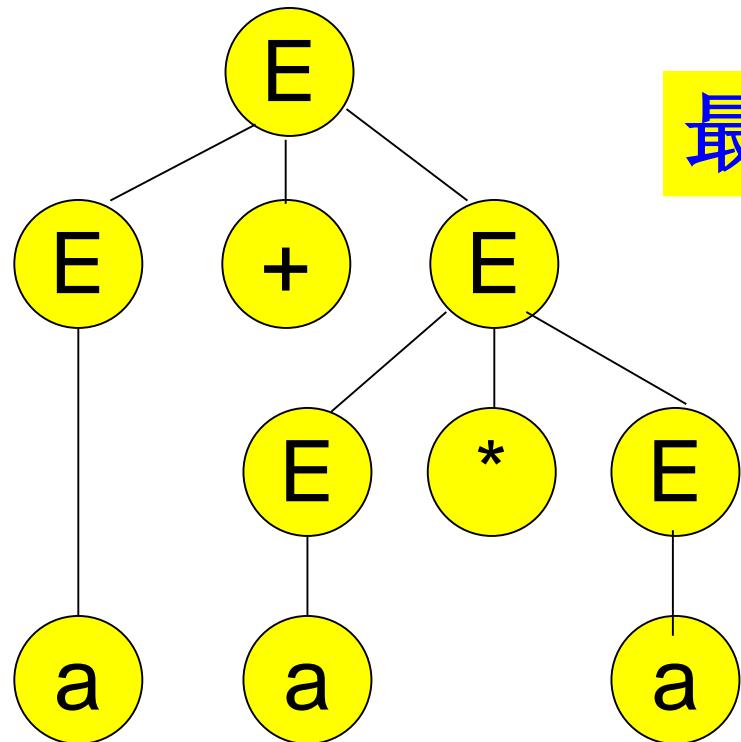
$E \Rightarrow E*E \Rightarrow E+E*E$   
 $\Rightarrow a+E*E \Rightarrow a+a*E$   
 $\Rightarrow a+a*a$



最左推导



$$\begin{aligned}
 E &\Rightarrow E+E \Rightarrow E+E^*E \\
 &\Rightarrow E+E^*a \Rightarrow E+a^*a \\
 &\Rightarrow a+a^*a
 \end{aligned}$$

$$\begin{aligned}
 E &\Rightarrow E^*E \Rightarrow E^*a \\
 &\Rightarrow E+E^*a \Rightarrow E+a^*a \\
 &\Rightarrow a+a^*a
 \end{aligned}$$




# 二义性（或歧义性， Ambiguity）

## 定义

- 如果一个文法中存在某个句子有两棵分析树，那么该句子是二义性的
- 如果一个文法产生二义性的句子，则称这个文法是二义性的
- 否则，该文法是无二义性的



# 关于二义性的几点说明-1

1. 一般来说，程序语言存在无二义性文法
  - 对于表达式来说，文法 G1[E] 是无二义性的
2. 在能驾驭的情况下，经常使用二义性文法  
对于条件语句，经常使用二义性文法描述它：

$$\begin{array}{l} S \rightarrow \text{if expr then } S \\ \quad | \quad \text{if expr then } S \text{ else } S \\ \quad | \quad \text{other} \end{array}$$

二义性的句子：

**if e1 then if e2 then s1 else s2**



# 关于二义性的几点说明-2

3. 对于任意一个上下文无关文法，不存在一个算法，判定它是无二义性的；但能给出一组充分条件，满足这组充分条件的文法是无二义性的。

4. 存在先天二义性的语言。例如，

$$\{a^i b^i c^j \mid i, j \geq 1\} \cup \{a^i b^j c^j \mid i, j \geq 1\}$$

存在一个二义性的句子  $a^k b^k c^k$ 。



# 证明文法生成的语言

- 证明文法  $G$  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  $L(G) \subseteq L$
  - 然后证明  $L \subseteq L(G)$
  - 一般可以使用数学归纳法
- 证明  $L(G) \subseteq L$ ：
  - 可以按照推导序列长度进行数学归纳
- 证明  $L \subseteq L(G)$ ：
  - 通常可按照符号串的长度来构造推导序列



# 文法生成语言的例子（1）

- 文法G:  $S \rightarrow (S)S \mid \epsilon$
- 语言L: 所有具有对称括号对的串
- $L(G) \subseteq L$  的证明:
  - 归纳基础: 推导长度为  $n = 1$ ,  $S \Rightarrow \epsilon$ , 括号对称
  - 归纳步骤: 假设长度小于  $n$  的推导都能得到括号对称的句子。考虑推导步骤为  $n$  的最左推导:

$$S \xrightarrow{lm} (S)S \xrightarrow{lm} (x)S \xrightarrow{lm} (x)y$$

- 其中  $x$  和  $y$  的推导步骤都小于  $n$ , 因此  $x$  和  $y$  都是括号对称的句子, 因此  $(x)y$  也是括号对称的句子
- $\Rightarrow G$  推导出的所有句子都是括号对称的



# 文法生成语言的例子（2）

## □ $L \subseteq L(G)$ 的证明：

- 注意：括号对称的串的长度必然是偶数
- 归纳基础：如果括号对称的串的长度为0，显然它可以从S推导得到
- 归纳步骤：假设长度小于 $2n$ 的括号对称的串都能够由S推导得到。假设w是括号对称且长度为 $2n$ 的串
  - w必然以左括号开头，且w可以写成 $(x)y$ 的形式，其中x也是括号对称的。因为x、y的长度都小于 $2n$ ，根据归纳假设，x和y都可以从S推导得到
  - 因此  $S \Rightarrow^* (S)S \Rightarrow^* (x)y$



# 上下文无关文法和正则表达式 (1)

- 上下文无关文法比正则表达式的能力更强：
  - 所有的正则语言都可以使用上下文无关文法描述
  - 但是一些用上下文无关文法描述的语言不能用正则文法描述
- 证明：
  - 首先证明：存在上下文无关文法  $S \rightarrow aSb \mid ab$  描述了语言  $\{a^n b^n | n > 0\}$ ，但是它无法用 DFA 识别
  - 反证法：假设有 DFA 识别该语言，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。



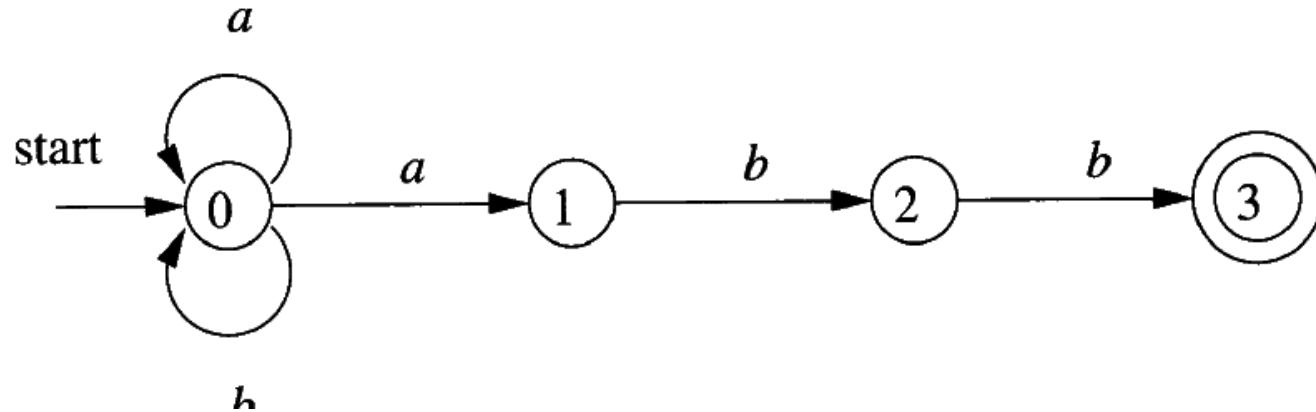
# 上下文无关文法和正则表达式（2）

## □ 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等价的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 $\epsilon$ 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



# NFA构造文法的例子



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

$$A_1 \rightarrow bA_2$$

$$A_2 \rightarrow bA_3$$

$$A_3 \rightarrow \epsilon$$

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。（可以考虑baabb的推导和接受过程）



# 非上下文无关的语言结构-1

- 在我们使用的程序语言中,有些语言结构并不是总能用上下文无关文法描述

**例1:**  $L1 = \{ wcw \mid w \in \{a,b\}^+ \}$ 。例如, `aabcaab` 就是  $L1$  的一个句子。这个语言是检查程序中标识符的声明应先于引用的抽象。

**例2:**  $L2 = \{ a^n b^m c^n d^m \mid n, m \geq 0 \}$ , 它是检查过程声明的形参数个数和过程调用的实参数个数一致问题的抽象。



# 非上下文无关的语言结构-2

- 语言中过程定义和引用的语法并不涉及到参数个数，例如，C语言的函数语句可描述为

**s-call → id (r-list)**  
**r-list → r-list, r**

| r

- 实参和形参个数的一致性检查也是放在语义分析阶段完成的



# 文法分类 (Chomsky)

**0型** (任意文法) :  $G=(V_T, V_N, S, P)$

规则形式:  $\alpha \rightarrow \beta$      $\alpha, \beta \in (V_T \cup V_N)^*$ ,  $\alpha \neq \epsilon$

推导:  $\gamma\alpha\delta \Rightarrow \gamma\beta\delta$

**1型** (上下文有关, Context-Sensitive Grammar)

规则形式:  $\alpha A \beta \rightarrow \alpha \gamma \beta$

$A \in V_N$ ,  $\alpha, \gamma, \beta \in (V_T \cup V_N)^*$ ,  $\gamma \neq \epsilon$

(注: 可以包含  $S \rightarrow \epsilon$ , 但此时不允许  $S$  出现在产生式右边)

**2型** (上下文无关, Context-Free Grammar, CFG)

规则形式:  $A \rightarrow \beta$ ,  $A \in V_N$ ,  $\beta \in (V_T \cup V_N)^*$

**3型** (正则文法, Regular Grammar)

(右线性):  $A \rightarrow aB$     $A \rightarrow a$

(左线性):  $A \rightarrow Ba$     $A \rightarrow a$      $a \in V_T \cup \{\epsilon\}$

每一类逐渐对产生式施加限制, 表示范围逐步缩小。



# 在程序语言中的实际应用

- 与词法有关的规则属于正则文法
- 与局部语法有关的规则属于上下文无关文法
- 而与全局语法和语义有关的部分往往要用上下文有关文法来描述
  - 实际上很少使用
- 为简化分析过程，会把描述词法的正则文法从描述语法的上下文无关文法中分离出来
  - 在分离出正则文法后的上下文无关文法中，这些单词符号是属于终结符号 $V_T$ 中的符号
  - 例：表达式文法G1[E]中，a是终结符号



# 作业

- 10月9日交
- 文法分析
  - Ex. 2.2.1, Ex. 4.2.1
- 上下文无关文法
  - Ex. 4.2.3