



第二章 高级语言及其文法

重点：文法的定义与分类，CFG的语法树及二义性、程序设计语言的定义。

难点：程序设计语言的语义定义。



第2章 高级语言及其文法

2.1 语言概述

2.2 基本定义

2.3 文法的定义

2.4 文法的分类

2.5 CFG的语法树

2.6 CFG的二义性

2.7 本章小结



2.1 语言概述

语言是一定的群体用来进行
什么？
信息交流的工具。



2.1 语言概述

□信息交流的基础是什么？

- 按照共同约定的生成规则和理解规则去生成“句子”和理解“句子”
- 例：
 - “今节日上课始开译第一编”
 - “今日开始上第一节编译课”



2.1 语言概述

□语言的特征

- 自然语言(Natural Language)
 - 是人与人的通讯工具
 - 语义(semantics):环境、背景知识、语气、二义性——难以形式化
- 计算机语言(Computer Language)
 - 计算机系统间、人机间通讯工具
 - 严格的语法(Grammar)、语义(semantics) ——易于形式化：严格



2.1 语言概述

□语言的描述方法——现状

- 自然语言：自然、方便-非形式化
- 数学语言（符号）：严格、准确-形式化
- 形式化描述
 - 高度的抽象，严格的理论基础和方便的计算机表示。



2.1 语言概述

□语言——形式化的内容提取

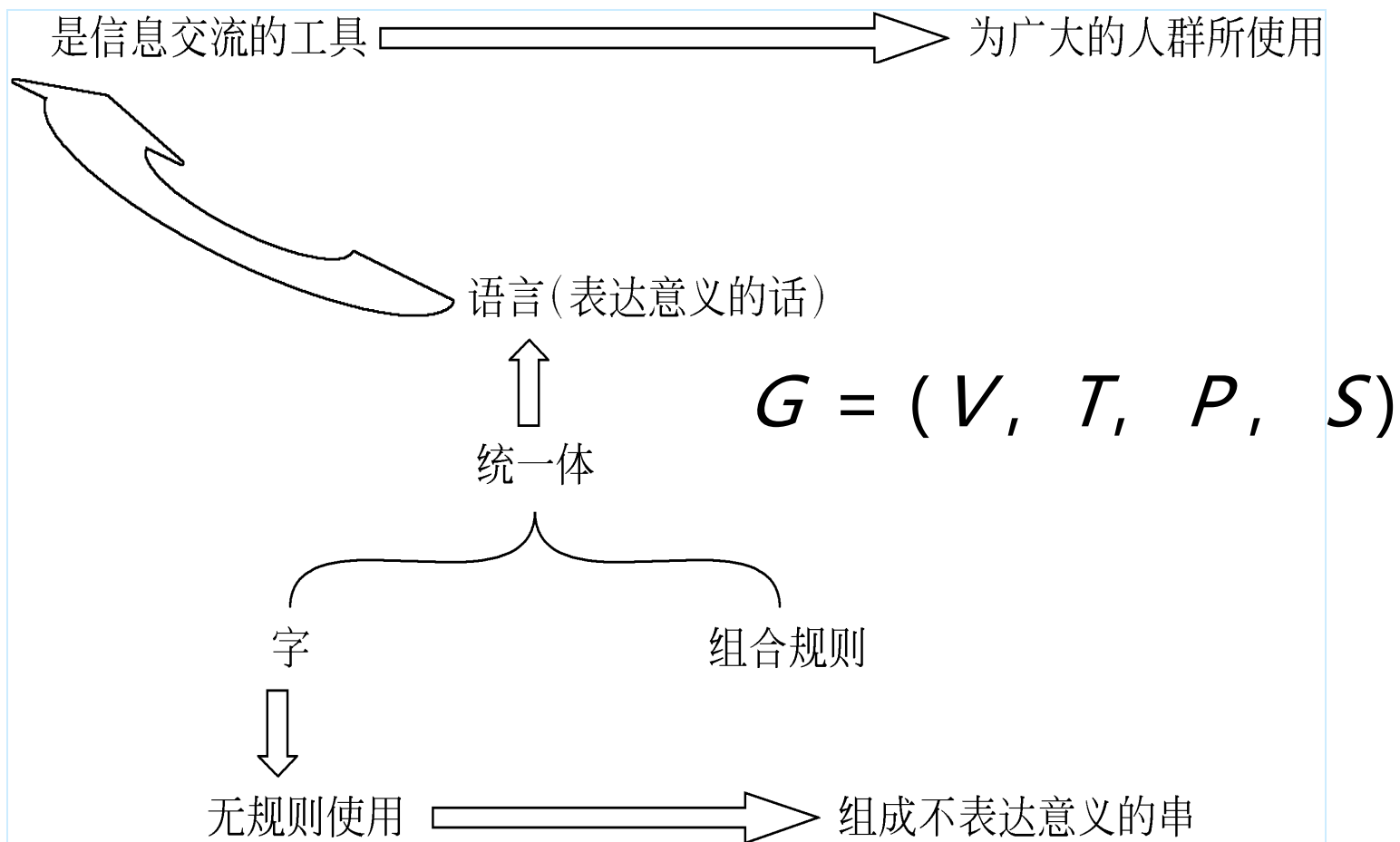
- 语言(Language): 满足一定条件的句子集合
- 句子(Sentence): 满足一定规则的单词序列
- 单词(Token): 满足一定规则的字符(Character)串

□语言是字和组合字的规则

- 例 (自然语言: 第译始二天课今开编上节)
- 今天开始上第二节编译课



2.1 语言概述



语言是字及其组合规则的统一体



2.1 语言概述

□程序设计语言——形式化的内容提取

- 程序设计语言(Programming Language): 组成程序的所有语句的集合。
- 程序(Program): 满足语法规则的语句序列。
- 语句(Sentence): 满足语法规则的单词序列。
- 单词(Token): 满足词法规则的字符串。

□例: 变量:=表达式

- if 条件表达式 then 语句
- while 条件表达式 do 语句
- call 过程名(参数表)



2.1 语言概述

□描述形式——文法

- 语法——语句
 - 语句的组成规则
 - 描述方法：BNF范式、语法(描述)图
- 词法——单词
 - 单词的组成规则
 - 描述方法：BNF范式、正规式



形式语言与自动机理论的产生与作用

□语言学家Chomsky最初从产生语言的角度研究语言。

- 1956年，通过抽象，他将语言形式地定义为是由一个字母表中的字母组成的一些串的集合。可以在字母表上按照一定的规则定义一个文法（Grammar），该文法所能产生的所有句子组成的集合就是该文法产生的语言。

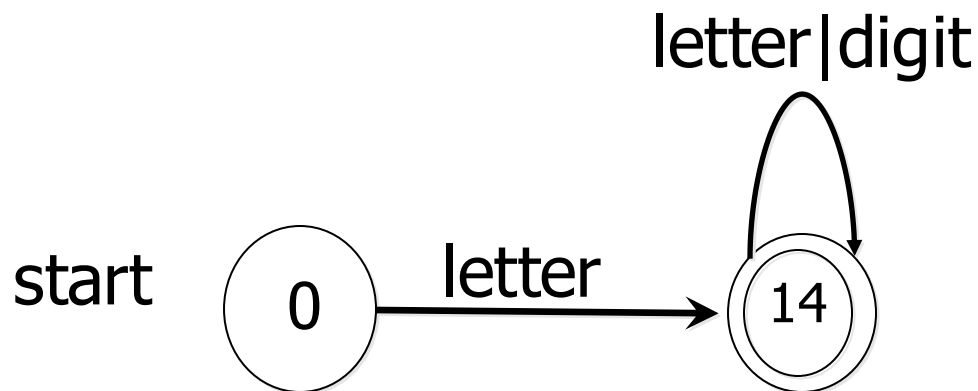
$$G = (V, T, P, S)$$



形式语言与自动机理论的产生与作用

□克林 (Kleene) 在1951年到1956年间，从识别语言的角度研究语言，给出了语言的另一种描述。

- 克林是在研究神经细胞中，建立了自动机，他用这种自动机来识别语言：对于按照一定的规则构造的任一个自动机，该自动机就定义了一个语言，这个语言由该自动机所能识别的所有句子组成。





形式语言与自动机理论的产生与作用

□1959年，Chomsky通过深入研究，将他本人的研究成果与克林的研究成果结合了起来，不仅确定了文法和自动机分别从生成和识别的角度去表达语言，而且证明了文法与自动机的等价性。



形式语言与自动机理论的产生与作用

□20世纪50年代，人们用**巴科斯范式**（Backus Nour Form 或 Backus Normal Form，简记为BNF）成功地对高级语言 ALGOL-60进行了描述。实际上，巴科斯范式就是**上下文无关文法**（Context Free Grammar）的一种表示形式。这一成功，使得形式语言在20世纪60年代得到了大力的发展。



2.2 基本定义

□定义2.1 **字母表** (Alphabet) Σ 是一个非空有穷集合, 字母表中的元素称为该字母表的一个字母 (Letter), 也叫字符 (Character)

□例 以下是不同的字母表:

(1) $\{a, b, c, d\}$

(2) $\{a, b, c, \dots, z\}$

(3) $\{0, 1\}$

(4) ASCII字母表



2.2 基本定义

□定义2.2 设 Σ_1 、 Σ_2 是两个字母表， Σ_1 与 Σ_2 的乘积(Product)定义为

$$\Sigma_1\Sigma_2 = \{ab \mid a \in \Sigma_1, b \in \Sigma_2\}$$

$$\text{例: } \Sigma_1 = \{0, 1\}, \Sigma_2 = \{a, b\}, \Sigma_1\Sigma_2 = \{0a, 0b, 1a, 1b\}$$

□定义2.3 设 Σ 是一个字母表， Σ 的 n 次幂(Power)递归地定义为：

- (1) $\Sigma^0 = \{\varepsilon\}$
- (2) $\Sigma^n = \Sigma^{n-1}\Sigma \quad n \geq 1$
- 例： $\Sigma_1^3 = \{000, 001, 010, 011, 100, 101, 110, 111\}$



2.2 基本定义

□ 定义2.4 设 Σ 是一个字母表, Σ 的**正闭包**(Positive Closure)定义为:

- $\Sigma^+ = \Sigma \cup \Sigma^2 \cup \Sigma^3 \cup \Sigma^4 \cup \dots$

□ Σ 的**克林闭包**(Kleene Closure)为:

- $\Sigma^* = \Sigma^0 \cup \Sigma^+$

- $= \Sigma^0 \cup \Sigma \cup \Sigma^2 \cup \Sigma^3 \cup \dots$



2.2 基本定义

□ 例

$$\{0,1\}^+ = \{0, 1, 00, 01, 11, 000, 001, 010, 011, 100, \dots\}$$

$$\{a, b, c, d\}^+ = \{a, b, c, d, aa, ab, ac, ad, ba, bb, bc, bd, \dots, aaa, aab, aac, aad, aba, abb, abc, \dots\}$$



2.2 基本定义

□例

$$\{0,1\}^* = \{\varepsilon, 0, 1, 00, 01, 11, 000, 001, 010, 011, \\ 100, \dots\}$$

$$\{a, b, c, d\}^* = \{\varepsilon, a, b, c, d, aa, ab, ac, ad, ba, \\ bb, bc, bd, \dots, aaa, aab, aac, aad, aba, abb, abc, \dots\}$$



2.2 基本定义

- 定义2.5 设 Σ 是一个字母表, $\forall x \in \Sigma^*$, x 称为字母表 Σ 上的一个句子(sentence), ε 叫做 Σ 上的一个空句子。
- 定义2.6 设 Σ 是一个字母表, 对任意的 $x, y \in \Sigma^*$, $a \in \Sigma$, 句子 xay 中的 a 叫做 a 在该句子中的一个出现(occurrence)。
- 定义2.7 设 Σ 是一个字母表, $\forall x \in \Sigma^*$, 句子 x 中字符出现的总个数叫做该字符串的长度(length), 记作 $|x|$ 。



2.2 基本定义

□ 定义2.8 设 Σ 是一个字母表, $\forall x, y \in \Sigma^*$, x, y 的并置 (concatenation) 是这样一串, 该串是由串 x 直接连接串 y 所组成的。记作 xy 。并置又叫做连结。

- 对于 $n \geq 0$, 串 x 的 n 次幂(power)定义为:
- (1) $x^0 = \varepsilon$;
- (2) $x^n = x^{n-1}x$ 。



2.2 基本定义

□ 定义2.9 设 Σ 是一个字母表, 对 $\forall x, y, z \in \Sigma^*$, 如果 $x=yz$, 则称 y 是 x 的**前缀**(prefix), 如果 $z \neq \varepsilon$, 则称 y 是 x 的**真前缀**(proper prefix); z 是 x 的**后缀**(suffix), 如果 $y \neq \varepsilon$, 则称 z 是 x 的**真后缀**(proper suffix)。

- 字母表 $\Sigma = \{a, b\}$ 上的句子 $abaabb$ 的前缀、后缀、真前缀和真后缀如下:
- 前缀: $\varepsilon, a, ab, aba, abaa, abaab, abaabb$
- 真前缀: $\varepsilon, a, ab, aba, abaa, abaab$
- 后缀: $\varepsilon, b, bb, abb, aabb, baabb, abaabb$
- 真后缀: $\varepsilon, b, bb, abb, aabb, baabb$



2.2 基本定义

□ 定义2.10 设 Σ 是一个字母表, 对 $\forall x, y, z, w, v \in \Sigma^*$, 如果 $x=yz$, $w=yv$, 则称 y 是 x 和 w 的**公共前缀**(common prefix)。如果 x 和 w 的任何公共前缀都是 y 的前缀, 则称 y 是 x 和 w 的**最大公共前缀**(maximal common prefix)。如果 $x=zy$, $w=vy$, 则称 y 是 x 和 w 的**公共后缀**(common suffix)。如果 x 和 w 的任何公共后缀都是 y 的后缀, 则称 y 是 x 和 w 的**最大公共后缀**(maximal common suffix)。



2.2 基

GPT-4

LH

□ 定

$w=yv$,

公共前

commo

(commc

x 和 w 的

定义2.10 设 Σ 是一个字母表, 对 $x, y, z, w, v \in \Sigma^*$, 如果 $x=yz$, $w=yv$, 则称 y 是 x 和 w 的公共前缀(common prefix)。如果 x 和 w 的任何公共前缀都是 y 的前缀, 则称 y 是 x 和 w 的最大公共前缀(maximal common prefix)。如果 $x=zy$, $w=vy$, 则称 y 是 x 和 w 的公共后缀(common suffix)。如果 x 和 w 的任何公共后缀都是 y 的后缀, 则称 y 是 x 和 w 的最大公共后缀(maximal common suffix)。

给一个最大公共前缀的例子。

是 $x=yz$,

任何



首先我们可以选择或构造两个字符串 x 和 w 以及一个字符串 y , 满足条件 y 是 x 和 w 的最大公共前缀。

例如, 我们可以取以下字符串:

- $x = \text{"abcd123"}$
- $w = \text{"abcd456"}$
- $y = \text{"abcd"}$

这里, y (即 "abcd") 是 x 和 w 的公共前缀, 因为 x 可以表示为 yz ($\text{"abcd123"} = \text{"abcd"} + \text{"123"}$), w 可以表示为 yv ($\text{"abcd456"} = \text{"abcd"} + \text{"456"}$)。

同时, y 是 x 和 w 的最大公共前缀, 因为不存在一个更长的字符串 y' 使得 y' 是 x 和 w 的公共前缀并且 y' 包含 y 作为它的前缀。

称 y 是



2.2 基本定义

□ 定义2.11 设 Σ 是一个字母表, 对 $\forall w, x, y, z \in \Sigma^*$, 如果 $w = xyz$, 则称 y 是 w 的**子串**(substring)。

□ 定义2.12 设 Σ 是一个字母表, 对 $\forall t, u, v, w, x, y, z \in \Sigma^*$, 如果 $t = uyv$, $w = xyz$, 则称 y 是 t 和 w 的**公共子串**(common substring)。如果 y_1, y_2, \dots, y_n 是 t 和 w 的公共子串, 且 $|y_j| = \max\{|y_1|, |y_2|, \dots, |y_n|\}$, 则称 y_j 是 t 和 w 的**最大公共子串**(maximal common substring)。



2.2 基



定义2.12 设 Σ 是一个字母表, 对 $t, u, v, w, x, y, z \in \Sigma^*$, 如果 $t=uyv$, $w=xyz$, 则称 y 是 t 和 w 的公共子串(common substring)。如果 y_1, y_2, \dots, y_n 是 t 和 w 的公共子串, 且 $|y| = \max\{|y_1|, |y_2|, \dots, |y_n|\}$, 则称 y 是 t 和 w 的最大公共子串(maximal common substring)。

给个最大公共子串的例子。

□ 定

则称 y 是

□ 定

果 $t=uyv$

果 y_1, y_2, \dots, y_n

$|y_n|\}$, 贝



我们可以选择或构造两个字符串 t 和 w , 以及一个字符串 y , 来满足条件 y 是 t 和 w 的最大公共子串。

例如:

- $t = \text{"abcde123fgh"}$
- $w = \text{"ijkl123mnop"}$
- $y = \text{"123"}$

在这里, y (即 "123") 是 t 和 w 的公共子串, 因为我们可以找到 u, v, x, z 使得 $t=uyv$ ($\text{"abcde123fgh"} = \text{"abcde"} + \text{"123"} + \text{"fgh"}$) 和 $w=xyz$ ($\text{"ijkl123mnop"} = \text{"ijkl"} + \text{"123"} + \text{"mnop"}$)。

同时, y 是 t 和 w 的最大公共子串, 因为在这两个字符串中没有一个更长的公共子串。

xyz ,

*, 如

如

.....,

)。



2.2 基本定义

□ 定义2.13 设 Σ 是一个字母表, $\forall L \subseteq \Sigma^*$, L 称为字母表 Σ 上的一个语言 (Language) , $\forall x \in L$, x 叫做 L 的一个句子。

□ 例: 字母表 $\{0, 1\}$ 上的语言

$\{0, 1\}$

$\{00, 11\}$

$\{0, 1, 00, 11\}$

$\{0, 1, 00, 11, 01, 10\}$

$\{00, 11\}^*$

$\{01, 10\}^*$



2.2 基本定义

□ 2.14 设 Σ_1, Σ_2 是字母表, $L_1 \subseteq \Sigma_1^*, L_2 \subseteq \Sigma_2^*$, 语言 L_1 与 L_2 的乘积(product)是字母表 $\Sigma_1 \cup \Sigma_2$ 上的一个语言, 该语言定义为:

$$L_1 L_2 = \{xy \mid x \in L_1, y \in L_2\}$$



2.2 基本定义

□ 定义2.15 设 Σ 是一个字母表, $\forall L \subseteq \Sigma^*$, L 的 n 次幂(power)是一个语言, 该语言定义为:

(1) 当 $n=0$ 是, $L^n = \{\varepsilon\}$;

(2) 当 $n \geq 1$ 时, $L^n = L^{n-1}L$ 。

L 的正闭包(positive closure) L^+ 是一个语言, 该语言定义为:

$$L^+ = L \cup L^2 \cup L^3 \cup L^4 \cup \dots$$

L 的克林闭包(Kleene closure) L^* 是一个语言, 该语言定义为:

$$L^* = L^0 \cup L \cup L^2 \cup L^3 \cup L^4 \cup \dots$$



2.3 文法的定义

如何实现语言结构的 形式化描述?

考虑赋值语句的形式:

左部量 = 右部表达式

$a = a + a$

$b = m[3] + b$

$m[1] = a + m[2]$



句子的组成规则

- $\langle \text{赋值语句} \rangle \rightarrow \langle \text{左部量} \rangle = \langle \text{右部表达式} \rangle$
- $\langle \text{左部量} \rangle \rightarrow \langle \text{简单变量} \rangle$
- $\langle \text{左部量} \rangle \rightarrow \langle \text{下标变量} \rangle$
- $\langle \text{简单变量} \rangle \rightarrow a$
- $\langle \text{简单变量} \rangle \rightarrow b$
- $\langle \text{简单变量} \rangle \rightarrow c$
- $\langle \text{下标变量} \rangle \rightarrow m[1]$
- $\langle \text{下标变量} \rangle \rightarrow m[2]$
- $\langle \text{下标变量} \rangle \rightarrow m[3]$
- $\langle \text{右部表达式} \rangle \rightarrow \langle \text{简单变量} \rangle \langle \text{运算符} \rangle \langle \text{简单变量} \rangle$
- $\langle \text{右部表达式} \rangle \rightarrow \langle \text{简单变量} \rangle \langle \text{运算符} \rangle \langle \text{下标变量} \rangle$
- $\langle \text{右部表达式} \rangle \rightarrow \langle \text{下标变量} \rangle \langle \text{运算符} \rangle \langle \text{简单变量} \rangle$
- $\langle \text{右部表达式} \rangle \rightarrow \langle \text{下标变量} \rangle \langle \text{运算符} \rangle \langle \text{下标变量} \rangle$
- $\langle \text{运算符} \rangle \rightarrow +$
- $\langle \text{运算符} \rangle \rightarrow -$

问题：如何用符号来描述？即如何形式化？



定义句子的规则的语法组成

——终结符号集, 非终结符号集, 语法规则, 开始符号

非终结符号集 $V =$

$\{ \langle \text{赋值语句} \rangle, \langle \text{左部量} \rangle, \langle \text{右部表达式} \rangle, \langle \text{简单变量} \rangle, \langle \text{下标变量} \rangle, \langle \text{运算符} \rangle \}$

终结符号集 $T =$

$\{ a, b, c, m[1], m[2], m[3], +, - \}$

语法规则集 $P =$

$\{ \langle \text{赋值语句} \rangle \rightarrow \langle \text{左部量} \rangle = \langle \text{右部表达式} \rangle, \dots \}$

开始符号 $S = \langle \text{赋值语句} \rangle$



文法G的形式定义

定义2.16 文法 G 为一个四元组:

$$G = (V, T, P, S)$$

□ V : 非终结符(Variable)集

- 语法变量 (成分) ——代表某个语言的各种子结构

□ T : 终结符(Terminal)集

- 语言的句子中出现的字符, $V \cap T = \emptyset$

□ S : 开始符号(Start Symbol), $S \in V$

- 代表文法所定义的语言, 至少在产生式左侧出现一次



文法G的形式定义

□ P : 产生式(Product)集合

$\alpha \rightarrow \beta$, 被称为产生式 (定义式), 读作: α 定义为 β 。其中 $\alpha \in (V \cup T)^+$, 且 α 中至少有 V 中元素的一个出现。 $\beta \in (V \cup T)^*$ 。 α 称为产生式 $\alpha \rightarrow \beta$ 的左部 (Left Part), β 称为产生式 $\alpha \rightarrow \beta$ 的右部 (Right Part)。

产生式定义各个语法成分的结构 (组成规则)



例 赋值语句的文法

□ $V = \{ \langle \text{赋值语句} \rangle, \langle \text{左部量} \rangle, \langle \text{右部表达式} \rangle, \langle \text{简单变量} \rangle, \langle \text{下标变量} \rangle, \langle \text{运算符} \rangle \}$

□ $T = \{ a, b, c, m[1], m[2], m[3], +, - \}$

□ $P = \{ \langle \text{赋值语句} \rangle \rightarrow \langle \text{左部量} \rangle = \langle \text{右部表达式} \rangle,$

$\langle \text{左部量} \rangle \rightarrow \langle \text{简单变量} \rangle, \langle \text{左部量} \rangle \rightarrow \langle \text{下标变量} \rangle,$

$\langle \text{简单变量} \rangle \rightarrow a, \langle \text{简单变量} \rangle \rightarrow b, \langle \text{简单变量} \rangle \rightarrow c,$

$\langle \text{下标变量} \rangle \rightarrow m[1], \langle \text{下标变量} \rangle \rightarrow m[2],$

$\langle \text{下标变量} \rangle \rightarrow m[3],$

$\langle \text{右部表达式} \rangle \rightarrow \langle \text{简单变量} \rangle \langle \text{运算符} \rangle \langle \text{简单变量} \rangle,$

$\langle \text{右部表达式} \rangle \rightarrow \langle \text{简单变量} \rangle \langle \text{运算符} \rangle \langle \text{下标变量} \rangle,$

$\langle \text{右部表达式} \rangle \rightarrow \langle \text{下标变量} \rangle \langle \text{运算符} \rangle \langle \text{简单变量} \rangle,$

$\langle \text{右部表达式} \rangle \rightarrow \langle \text{下标变量} \rangle \langle \text{运算符} \rangle \langle \text{下标变量} \rangle,$

$\langle \text{运算符} \rangle \rightarrow +, \langle \text{运算符} \rangle \rightarrow - \}$

□ $S = \langle \text{赋值语句} \rangle$



例 赋值语句的文法 (续)

□ 符号化之后:

□ $G = (\{A, B, E, C, D, P\}, \{a, b, c, m[1], m[2], m[3], +, -\}, P, A),$

其中, $P = \{A \rightarrow B = E, B \rightarrow C, B \rightarrow D, C \rightarrow a, C \rightarrow b, C \rightarrow c, D \rightarrow m[1], D \rightarrow m[2], D \rightarrow m[3], E \rightarrow COC, E \rightarrow COD, E \rightarrow DOC, E \rightarrow DOD, O \rightarrow +, O \rightarrow -\}$



产生式的简写

□ 对一组有相同左部的产生式

$$\alpha \rightarrow \beta_1, \alpha \rightarrow \beta_2, \dots, \alpha \rightarrow \beta_n$$

可以简单地记为：

$$\alpha \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$$

读作： α 定义为或者 β_1 ，或者 β_2 ，...，或者 β_n 。并且称它们为 α 产生式。 $\beta_1, \beta_2, \dots, \beta_n$ 称为**候选式**(Candidate)

□ 对一个文法，只列出该文法的所有产生式，且所列的第一个产生式的左部是该文法的开始符号。

问题：有了定义句子的规则，如何判定某一句子是否属于某语言？



句子的派生(推导)——从产生语言的角度

<赋值语句>

$\Rightarrow \langle \text{左部量} \rangle = \langle \text{右部表达式} \rangle$

$\Rightarrow \langle \text{简单变量} \rangle = \langle \text{右部表达式} \rangle$

$\Rightarrow a = \langle \text{右部表达式} \rangle$

$\Rightarrow a = \langle \text{简单变量} \rangle \langle \text{运算符} \rangle \langle \text{简单变量} \rangle$

$\Rightarrow a = a \langle \text{运算符} \rangle \langle \text{简单变量} \rangle$

$\Rightarrow a = a + \langle \text{简单变量} \rangle$

$\Rightarrow a = a + a$

句子的识别(归约)——从识别语言的角度

$a = a + a$

$\Leftarrow a = a + \langle \text{简单变量} \rangle$

$\Leftarrow a = a \langle \text{运算符} \rangle \langle \text{简单变量} \rangle$

$\Leftarrow a = \langle \text{简单变量} \rangle \langle \text{运算符} \rangle \langle \text{简单变量} \rangle$

$\Leftarrow a = \langle \text{右部表达式} \rangle$

$\Leftarrow \langle \text{简单变量} \rangle = \langle \text{右部表达式} \rangle$

$\Leftarrow \langle \text{左部量} \rangle = \langle \text{右部表达式} \rangle$

$\Leftarrow \langle \text{赋值语句} \rangle$

□ 派生与识别均根据规则



基于产生式的变换——推导或归约

□定义2.17 设 $G=(V, T, P, S)$ 是一个文法, 如果 $\alpha \rightarrow \beta \in P$, $\gamma, \delta \in (V \cup T)^*$, 则称 $\gamma\alpha\delta$ 在 G 中**直接推导出** $\gamma\beta\delta$, 记作:

$$\gamma\alpha\delta \Rightarrow_G \gamma\beta\delta$$

□读作: $\gamma\alpha\delta$ 在文法 G 中直接推导出 $\gamma\beta\delta$ 。在不特别强调推导的直接性时, “直接推导”可以简称为推导(derivation), 有时我们也称推导为**派生**。

□与之相对应, 也可以称 $\gamma\beta\delta$ 在文法 G 中**直接归约**成 $\gamma\alpha\delta$ 。在不特别强调归约的直接性时, “直接归约”可以简称为**归约**(reduction)。由于这里实际是将 $\gamma\beta\delta$ 中的 β 变成了 α , 而 γ 和 δ 都没有变化, 所以又称将 β 归约成 α 。



(多步)推导

$$\square \alpha_0 \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n$$

- 记为 $\alpha_0 \overset{n}{\Rightarrow} \alpha_n$ (恰用 n 步)
- $\alpha_0 \overset{+}{\Rightarrow} \alpha_n$ (至少一步)
- $\alpha_0 \overset{*}{\Rightarrow} \alpha_n$ (若干步:零步或多步)



文法 G 产生的语言

□定义2.19 设有文法 $G=(V, T, P, S)$, 则称

$$L(G)=\{w \mid S \xRightarrow{+} w \ \& \ w \in T^* \}$$

为文法 G 产生的语言(language)。 $\forall w \in L(G)$, w 称为 G 产生的一个句子(sentence)。

显然, 对于任意一个文法 G , G 产生的语言 $L(G)$ 就是该文法的开始符号 S 所对应的集合 $L(S)$ 。



文法G产生的语言(续)

$$L(G) = \{x \mid S \Rightarrow^* x \text{ and } x \in T^*\}$$

- 文法 G 可以派生出多少个句子?
- 文法 G 的作用——语言的有穷描述
 - 以有限的规则描述无限的语言现象
- 有限:
 - 产生式集合、终结符集合、非终结符集合
- 无限:
 - 可以导出无穷多个句子 (L 也可能是有穷)



推导/归约举例

$A \Rightarrow B=E$ 有产生式 $A \rightarrow B=E$, 所以可以将 A 换成 $B=E$

$\Rightarrow C=E$ 有产生式 $B \rightarrow C$, 所以可以将 $B=E$ 中的 B 换成 C

$\Rightarrow a=E$ 有产生式 $C \rightarrow a$, 所以可以将 $C=E$ 中的 C 换成 a

$\Rightarrow a=COD$ 有产生式 $E \rightarrow COD$, 所以可以将 $a=E$ 中的 E 换成 COD

$\Rightarrow a=bOD$ 有产生式 $C \rightarrow b$, 所以可以将 $a=COD$ 中的 C 换成 b

$\Rightarrow a=b+D$ 有产生式 $O \rightarrow +$, 所以可以将 $a=bOD$ 的 O 换成 $+$

$\Rightarrow a=b+m[1]$ 有产生式 $D \rightarrow m[1]$, 所以可以将 $a=b+D$ 的 D 换成 $m[1]$

$A \rightarrow B=E$
 $B \rightarrow C \mid D$
 $C \rightarrow a \mid b \mid c$
 $D \rightarrow m[1] \mid m[2] \mid m[3]$
 $E \rightarrow COC \mid COD \mid DOC \mid DOD$
 $O \rightarrow + \mid -$



句型与句子

□定义2.20 设文法 $G=(V, T, P, S)$, 对于 $\forall \alpha \in (V \cup T)^*$, 如果 $S \xRightarrow{*} \alpha$, 则称 α 是 G 产生的一个句型(sentential form), 简称为**句型**

□对于任意文法 $G=(V, T, P, S)$, G 产生的句子和句型的区别在于句子满足 $w \in T^*$, 而句型满足 $\alpha \in (V \cup T)^*$



句型与句子

- 句子 w 是从 S 开始, 在 G 中可以推导出来的终结符号行, 它不含语法变量;
- 句型 α 是从 S 开始, 在 G 中可以推导出来的符号行, 它可能含有语法变量;
- 句子一定是句型; 但句型不一定是句子

文法**G**产生的()的全体是该文法描述的语言。

- ☐ A 终结符集
- ☐ B 非终结符集
- ☐ C 句子
- ☐ D 句型

给定文法 $A \rightarrow bA \mid cc$ ，则以下符号串是该文法句子的是()。

- ☐ A cc
- ☐ B bcbc
- ☐ C bbbccc
- ☐ D bccbcc
- ☐ E bbbbcc

文法 $S \rightarrow aaS|abc$ 定义的语言是()。

- A $\{a^{2k-1}bc|k > 0\}$
- B $\{a^k a^k bc|k > 0\}$
- C $\{a^k bc|k > 0\}$
- D $\{a^{2k} bc|k > 0\}$



2.4 文法的分类(Chomsky体系)

□根据语言结构的复杂程度（形式语言）

- 涉及文法的复杂程度、分析方法的选择
- 反映文法描述语言的能力

□如果 G 满足文法定义的要求，则 G 是 **0 型文法**（短语结构文法PSG: Phrase Structure Grammar）。

□ $L(G)$ 为PSL。



(1) 上下文有关文法(CSG)

□如果对于 $\forall \alpha \rightarrow \beta \in P$, 均有 $|\beta| \geq |\alpha|$ 成立($S \rightarrow \varepsilon$ 除外), 则称 G 为
1 型文法

- 即: 上下文有关文法 (CSG——Context Sensitive Grammar)

□ $L(G)$ 为1型/上下文有关/敏感语言(CSL)



(2) 上下文无关文法(CFG)

□如果对于 $\forall \alpha \rightarrow \beta \in P$, 均有 $|\beta| \geq |\alpha|$, 并且 $\alpha \in V$ 成立, 则称 G 为 **2型文法**

- 即: 上下文无关文法(CFG: Context Free Grammar)

□ $L(G)$ 为2型/上下文无关语言 (CFL)

- CFG能描述程序设计语言的多数语法成分



(2) 正规文法(RG)

□ 设 $A, B \in V$, $a \in T$ 或为 ε

右线性(Right Linear)文法: $A \rightarrow aB$ 或 $A \rightarrow a$

左线性(Left Linear)文法: $A \rightarrow Ba$ 或 $A \rightarrow a$

都是 3 型文法(正规文法 Regular Grammar -RG)

□ $L(G)$ 为 3 型/正规集/正则集/正则语言 (RL)

- 能描述程序设计语言的多数单词
- 左、右线性文法不可混用



例 标识符的文法

$$\square S \rightarrow L | LT$$

$$T \rightarrow L | N | TL | TN$$

$$L \rightarrow a | b | c | d$$

$$N \rightarrow 0 | 1 | 2 | 3 | 4 | 5$$

$$\square S \rightarrow a | b | c | d$$

$$S \rightarrow aT | bT | cT | dT$$

$$T \rightarrow a | b | c | d | 0 | 1 | 2 | 3 | 4 | 5$$

$$T \rightarrow aT | bT | cT | dT | 0T$$

$$T \rightarrow 1T | 2T | 3T | 4T | 5T$$



非上下文无关的语言结构

□ 程序设计语言的有些语言结构不能用上下文无关文法描述

□ 例

$$L_1 = \{wcw \mid w \in \{a,b\}^+\}$$

aabcaab 就是 L_1 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



非上下文无关的语言结构

□ 例

$$L_2 = \{a^n b^m c^n d^m \mid n, m \geq 0\}$$

□ 它是检查过程声明的形参个数和过程引用的参数个数是否一致问题的抽象



Chomsky体系——总结

$G = (V, T, P, S)$ 是一个文法, $\alpha \rightarrow \beta \in P$

* G 是0型文法, $L(G)$ 是0型语言;

---其能力相当于图灵机

* $|\alpha| \leq |\beta|$: G 是1型文法, $L(G)$ 是1型语言(除 $S \rightarrow \varepsilon$);

---其识别系统是线性界限自动机

* $\alpha \in V$: G 是2型文法, $L(G)$ 是2型语言;

---其识别系统是不确定的下推自动机

* $A \rightarrow aB$ 或 $A \rightarrow a$: G 是右线性文法, $L(G)$ 是3型语言

$A \rightarrow Ba$ 或 $A \rightarrow a$: G 是左线性文法, $L(G)$ 是3型语言

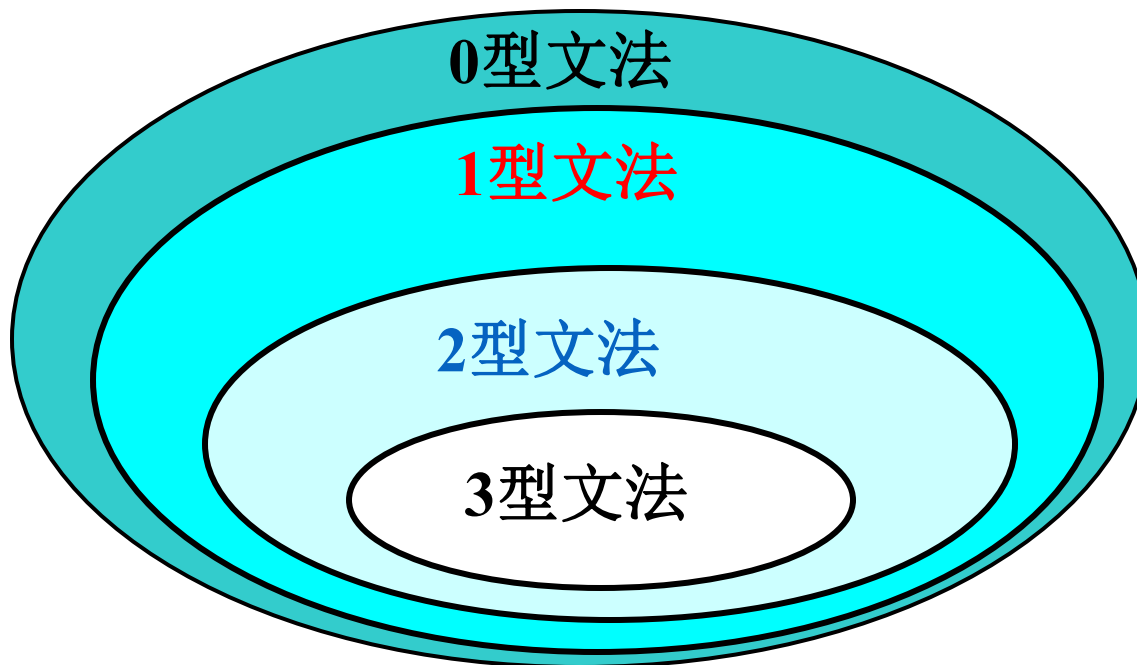
---其识别系统是有穷自动机



文法的类型

四种文法之间的关系是将产生式作进一步限制而定义的。

四种文法之间的逐级“包含”关系如下：





2.5 CFG的语法树

Parse Tree

□用树的形式表示句型的生成

- 树根： 开始符号
- 中间结点： 非终结符
- 叶结点： 终结符

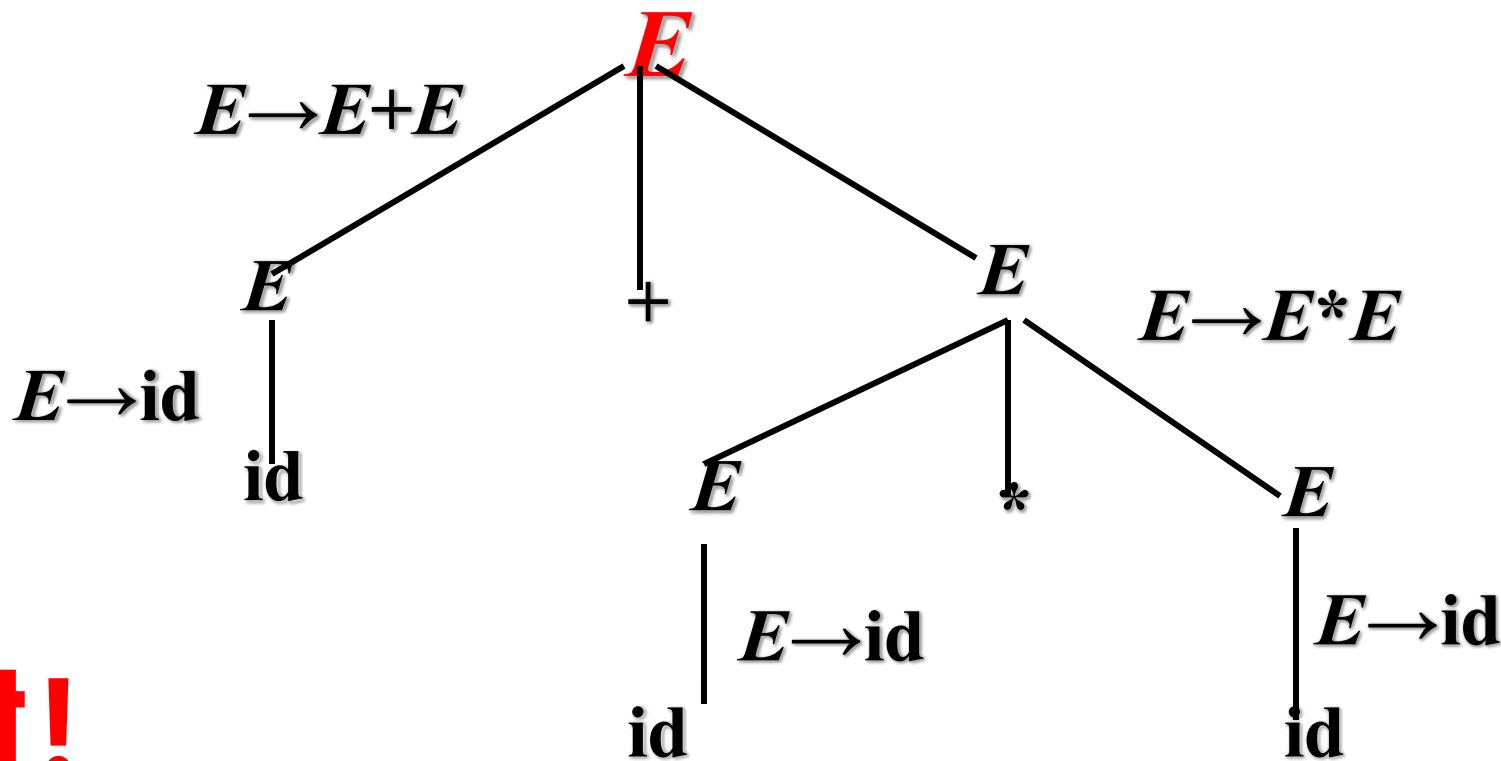
□每个推导对应一个中间结点及其儿子——一个二级子树-直接短语

□又称为分析树(parse tree)、推导树(derivation tree)、派生树(derivation tree)



例 句子结构的表示 (文法 $E \rightarrow E + E \mid E * E \mid (E) \mid id$)

$E \Rightarrow E + E \Rightarrow id + E \Rightarrow id + E * E \Rightarrow id + id * E \Rightarrow id + id * id$



一棵树!



短语(Phrase)

□定义2.27 设有CFG $G=(V, T, P, S)$, $\exists \alpha, \gamma, \beta \in (V \cup T)^*$, $S \xRightarrow{*} \gamma A \beta$, $A \xRightarrow{+} \alpha$, 则称 α 是句型 $\gamma \alpha \beta$ 的相对于变量 A 的短语(phrase);

如果此时有 $A \Rightarrow \alpha$, 则称 α 是句型 $\gamma \alpha \beta$ 的相对于变量 A 的直接短语(immediate phrase)

在无意义冲突时, 简称为句型 $\gamma \alpha \beta$ 的直接短语。直接短语也叫做简单短语(simple phrase)。

□定义2.28 设有CFG $G=(V, T, P, S)$, G 的句型的最左直接短语叫做句柄(handle)。



例：(直接)短语

$$\begin{aligned} E &\Rightarrow E+T && \Rightarrow (E+T)+T \\ &\Rightarrow T+T && \Rightarrow (T+T)+T \\ &\Rightarrow F+T && \Rightarrow (F+T)+T \\ &\Rightarrow (E)+T && \Rightarrow (\text{id}+T)+T \\ &\Rightarrow (E+T)+T && \Rightarrow (a+T^*F)+T \\ & && \Rightarrow (a+F^*F)+T \\ & && \Rightarrow (a+b^*F)+T \\ & && \Rightarrow (a+b^*c)+T \\ & && \Rightarrow (a+b^*c)+F \\ & && \Rightarrow (a+b^*c)+d \end{aligned}$$

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

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

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



句柄(Handle): 最左直接短语

$\square T \rightarrow T^* F \mid F$

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

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

$E \Rightarrow E+T$

$\Rightarrow T+T$

$\Rightarrow F+T$

$\Rightarrow (E)+T$

$\Rightarrow (E+T)+T$

$\Rightarrow (E+T)+T$

$\Rightarrow (T+T)+T$

$\Rightarrow (F+T)+T$

$\Rightarrow (\text{id}+T)+T$

$\Rightarrow (a+T^*F)+T$

$\Rightarrow (a+F^*F)+T$

$\Rightarrow (a+b^*F)+T$

$\Rightarrow (a+b^*c)+T$

$\Rightarrow (a+b^*c)+F$

$\Rightarrow (a+b^*c)+d$

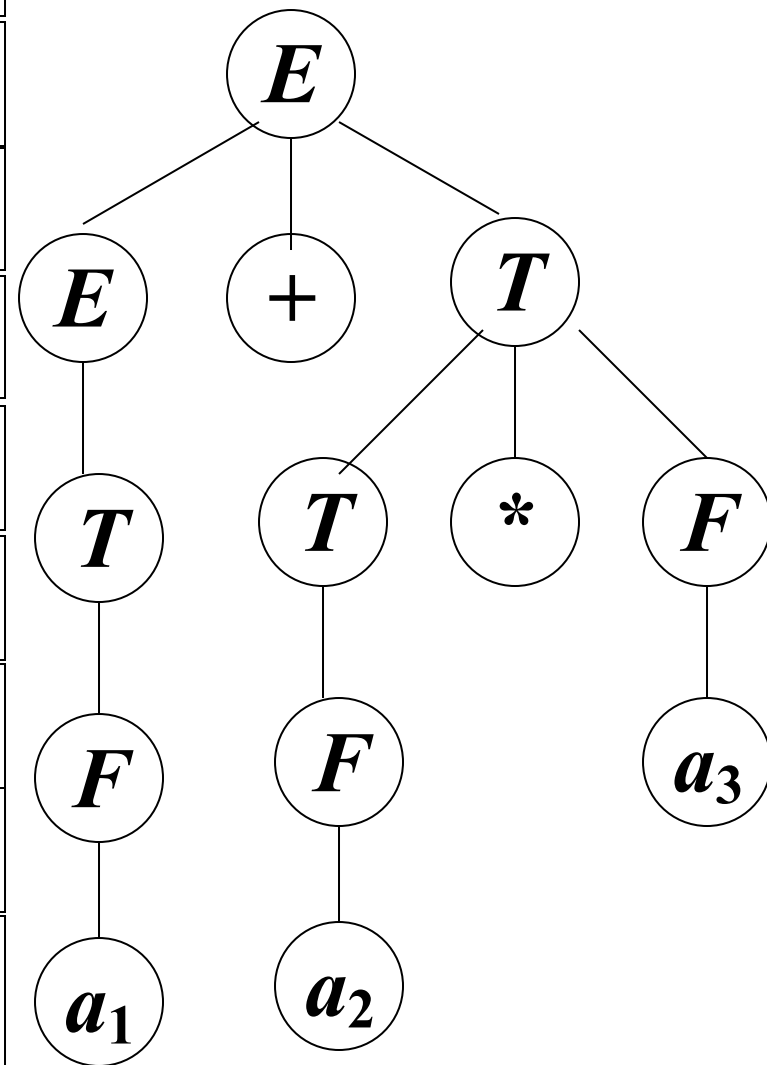


用子树解释短语，直接短语，句柄

- **短语**：一棵子树的所有叶子自左至右排列起来形成一个相对于子树根的短语。
- **直接短语**：仅有父子两代的一棵子树，它的所有叶子自左至右排列起来所形成的符号串。
- **句柄**：一个句型的分析树中最左那棵只有父子两代的子树的所有叶子的自左至右排列。
- 例如，对表达式文法 GE 和句子 $a_1 + a_2^* a_3$ ，挑选出推导过程中产生的句型中的短语，直接短语，句柄



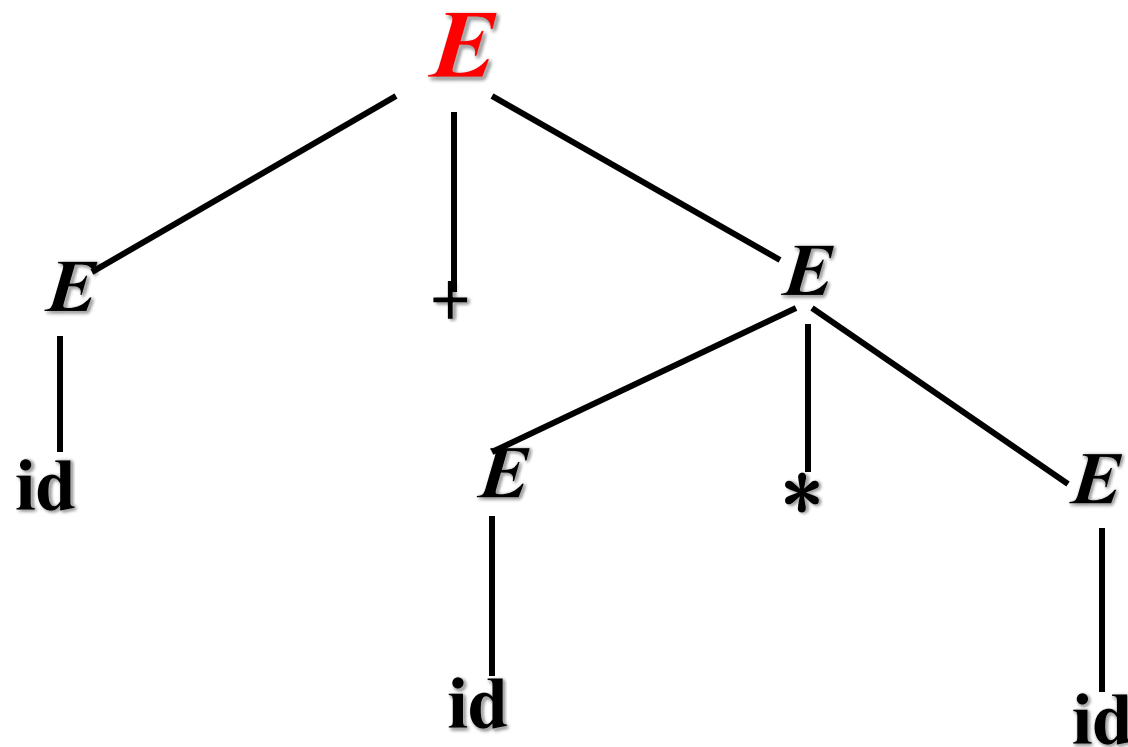
E	短语
$\Rightarrow \underline{E}+T$	$\underline{E+T}$
$\Rightarrow \underline{T}+T$	$T, \underline{T+T}$
$\Rightarrow \underline{F}+T$	$\underline{F}, \underline{F+T}$
$\Rightarrow \underline{a_1}+T$	$\underline{a_1}, \underline{a_1+T}$
$\Rightarrow \underline{a_1}+T*F$	$\underline{a_1}, \underline{T*F}, \underline{a_1+T*F}$
$\Rightarrow \underline{a_1}+F*F$	$\underline{a_1}, \underline{F}, \underline{F*F},$
$\Rightarrow \underline{a_1}+a_2 * F$	$\underline{a_1}, \underline{a_2}, \underline{a_1+a_2 * F}, \underline{a_2 * F}$
$\Rightarrow \underline{a_1}+a_2 * a_3$	$\underline{a_1}, \underline{a_2}, \underline{a_3}, \underline{a_2 * a_3}$ $\underline{a_1+a_2 * a_3}$





例 短语与分析树

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



一棵子树的叶子!



id+id*id的不同推导 $E \rightarrow E+E \mid E^*E \mid (E) \mid id$

$E \Rightarrow E^*E$
 $\Rightarrow E+E^*E$
 $\Rightarrow E+id^*E$
 $\Rightarrow id+id^*E$
 $\Rightarrow id+id^*id$

不做限制

句型 (sentential Form)

(归约)

$E \Rightarrow^* id+id^*id$

$E \Rightarrow E+E$
 $\Rightarrow id+E$
 $\Rightarrow id+E^*E$
 $\Rightarrow id+id^*E$
 $\Rightarrow id+id^*id$

施于最左变量

左句型(left-~)

(最右归约)

$E \Rightarrow id+id^*id$

$E \Rightarrow E+E$
 $\Rightarrow E+E^*E$
 $\Rightarrow E+E^*id$
 $\Rightarrow E+id^*id$
 $\Rightarrow id+id^*id$

施于最右变量

右句型/规范句型

(canonical ~)

(最左/规范归约)

$E \Rightarrow id+id^*id$



最左推导与最右推导

□最左推导(Left-most Derivation)

- 每次推导都施加在句型的最左边的语法变量上。——与最右归约对应

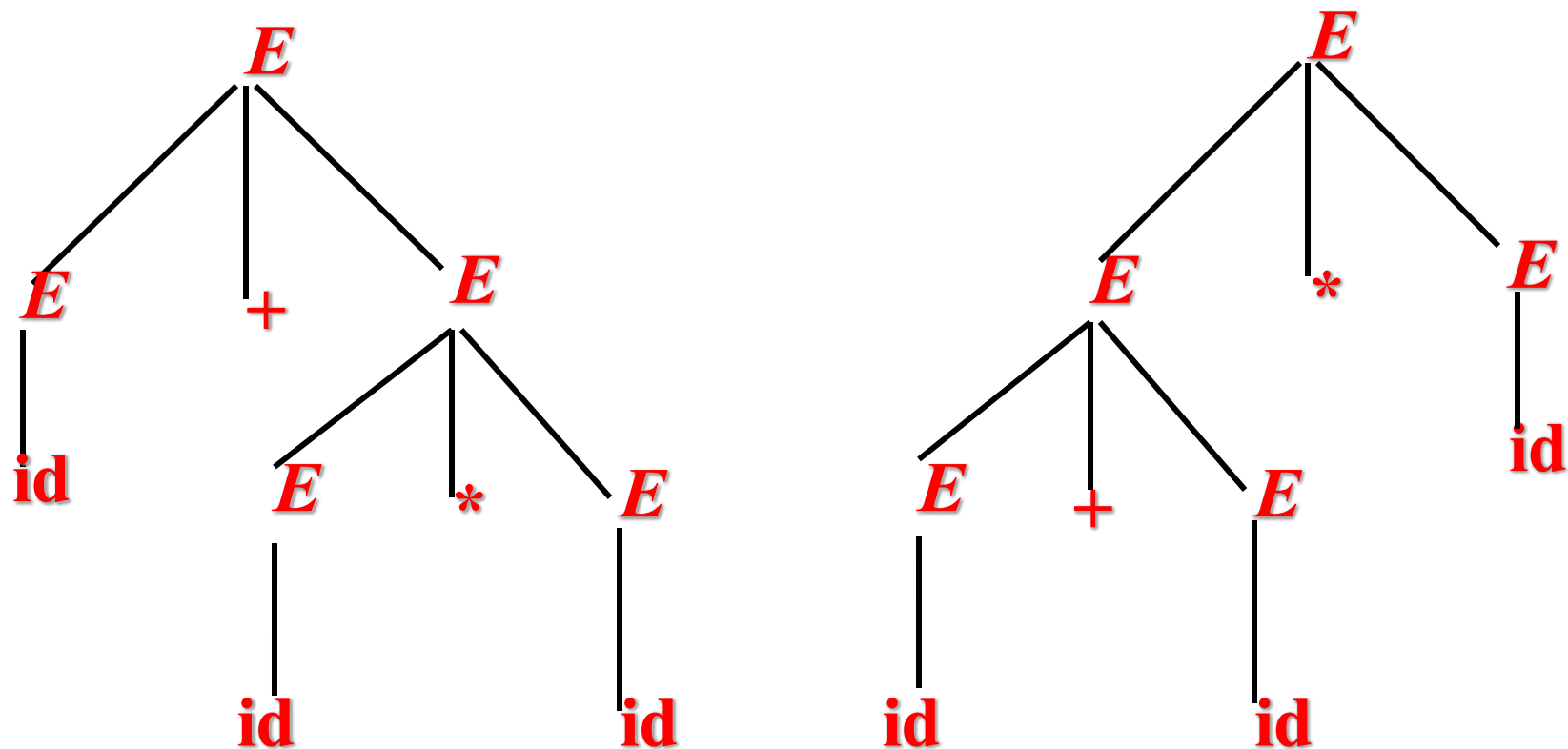
□最右推导(Right-most Derivation)

- 每次推导都施加在句型的最右边的语法变量上。——与最左归约（规范规约）
对应的规范(Canonical)句型



2.6 CFG的二义性

- 对同一句子存在两棵语法分析树
 - 在理论上不可判定





文法的二义性

1. 描述一个句子的文法不是唯一的;
2. 对于一个句子的分析应是唯一的。

考虑表达式下面的文法 $G[E]$, 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:

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

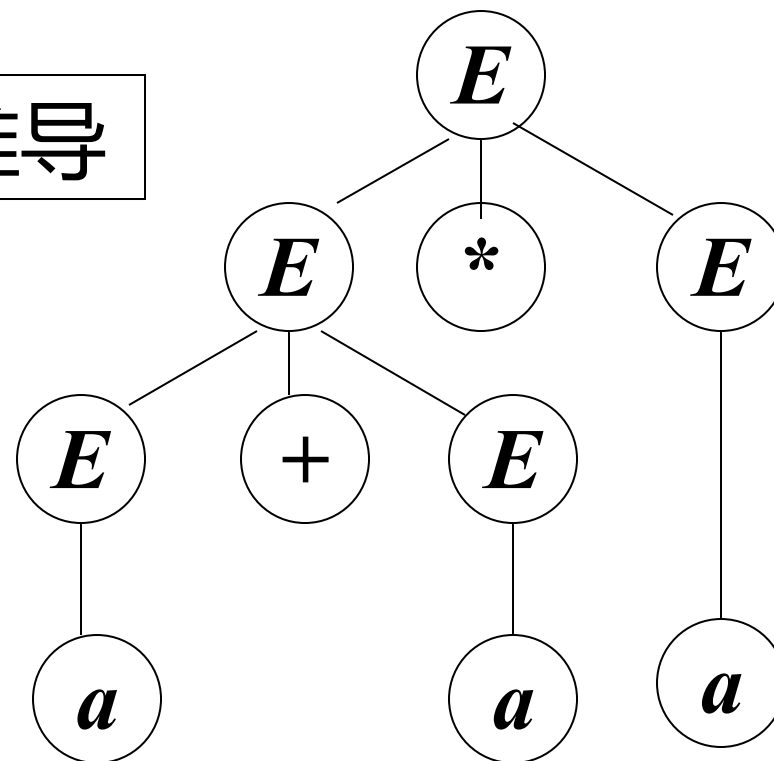
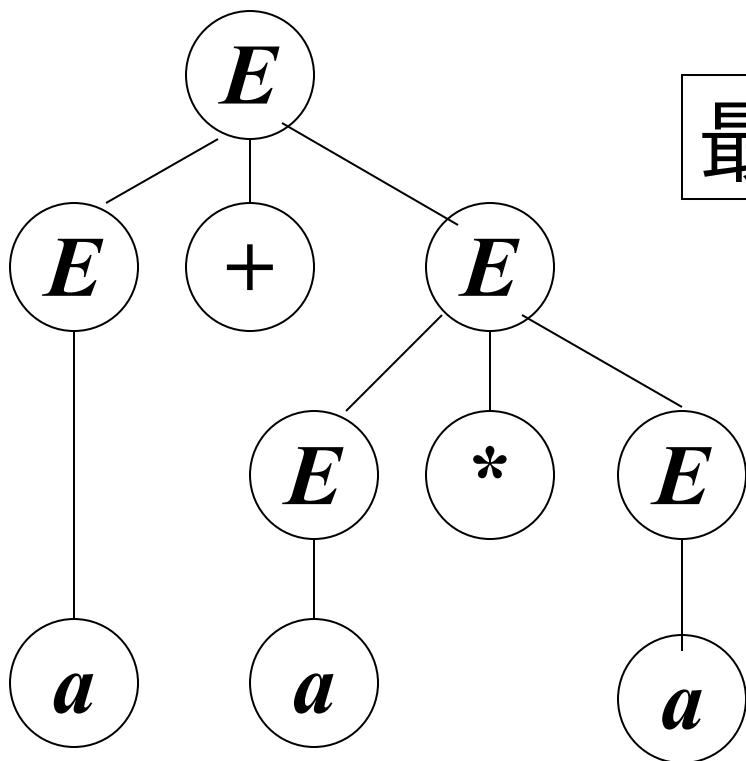
$$E \Rightarrow E * E \Rightarrow E + E * E \Rightarrow a + E * E \Rightarrow a + a * E \Rightarrow a + a * a$$



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

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

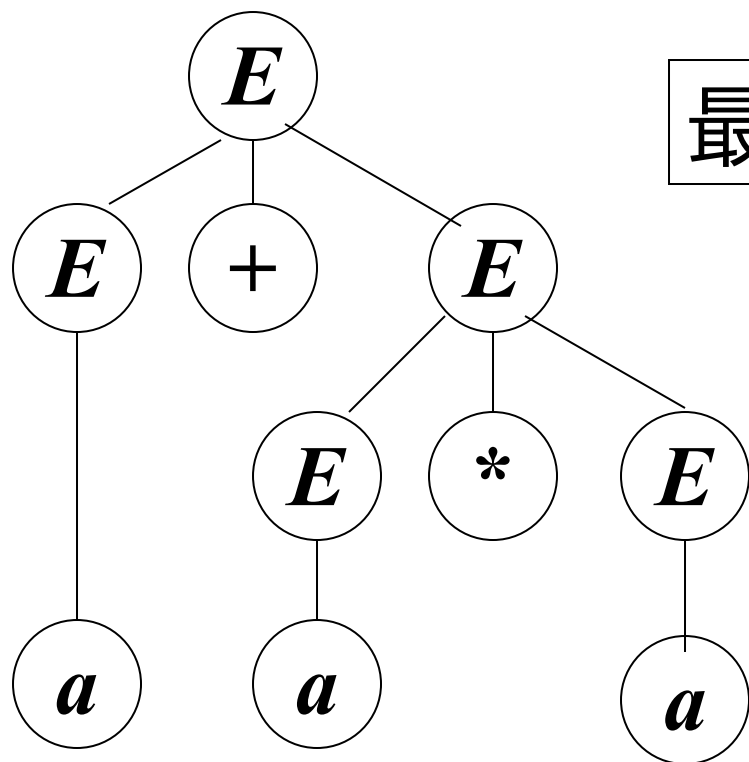
最左推导



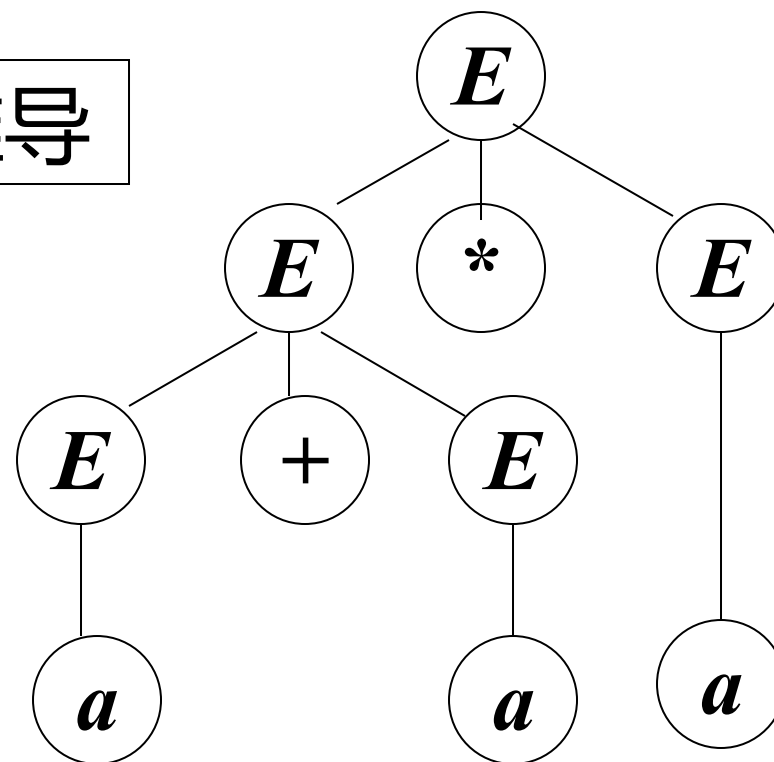


$$\begin{aligned} E &\Rightarrow E + E \Rightarrow E + E^* E \\ &\Rightarrow E + E^* a \Rightarrow E + a^* a \\ &\Rightarrow a + a^* a \end{aligned}$$

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



最右推导





二义性 (ambiguity)的定义

□如果一个文法的句子存在两棵分析树,那么,该句子是二义性的。如果一个文法包含二义性的句子,则称这个文法是二义性的; 否则, 该文法是无二义性的。几点说明:

□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 (2 .1) 是二义性的。对于条件语句, 经常使用二义性文法描述它:

□ $S \rightarrow \text{if } expr \text{ then } S$
 | $\text{if } expr \text{ then } S \text{ else } S$
 | other

□二义性的句子: $\text{if } e_1 \text{ then if } e_2 \text{ then } s_1 \text{ else } s_2$

描述if语句的无二义性文法的产生式

下面是

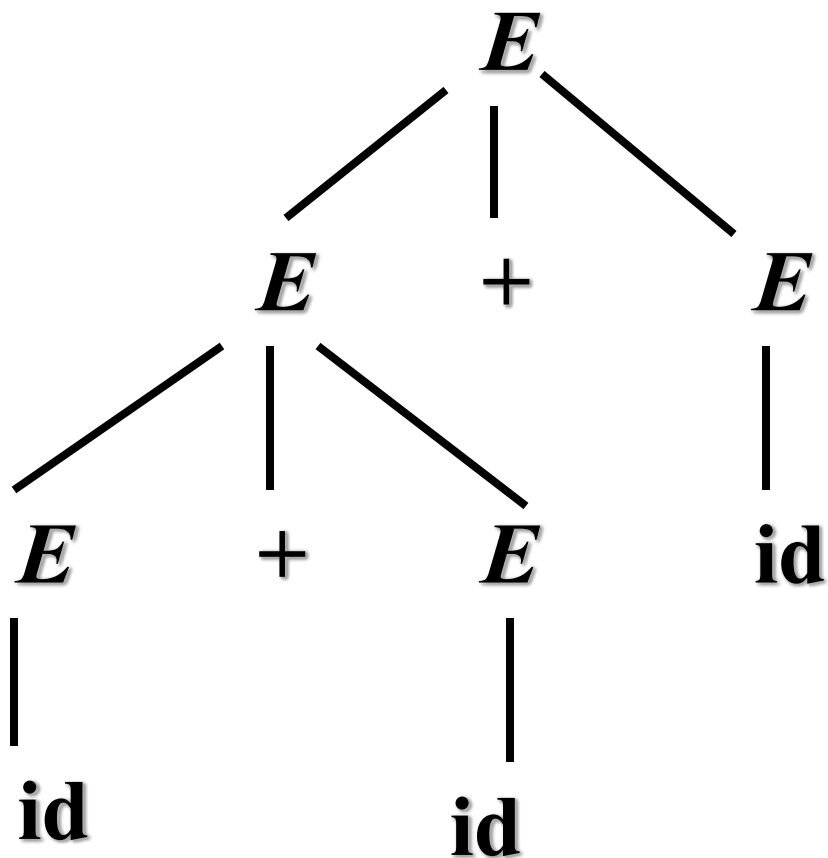
$$\begin{aligned} S &\rightarrow matched_s \mid unmatched_s \\ matched_s &\rightarrow \text{if } expr \text{ then } matched_s \\ &\quad \text{else } matched_s \mid \text{other} \\ unmatched_s &\rightarrow \text{if } expr \text{ then } S \\ &\quad \mid \text{if } expr \text{ then } matched_s \\ &\quad \text{else } unmatched_s \end{aligned}$$

它显然比较复杂，因此：

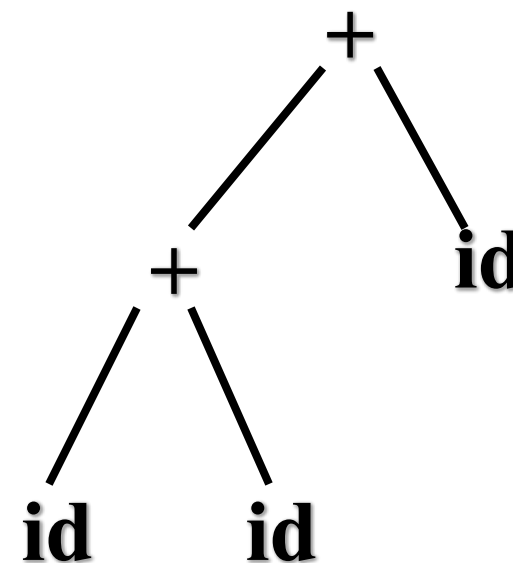
2. 在能驾驭的情况下，使用二义性文法。



分析树与(抽象)语法树不同



(语法)分析树



(抽象)语法树

乔姆斯基(Chomsky)把文法分为四种类型，即0型、1型、2型、3型。其中3型文法是()。

- ☐ A 正则文法
- ☐ B 上下文无关文法
- ☐ C 上下文有关文法
- ☐ D 非限制文法



2.7 本章小结

- 语言及其描述
- 文法的基本概念
- Chomsky体系
- CFG的语法树
- 文法的二义性
- 语言的固有二义性