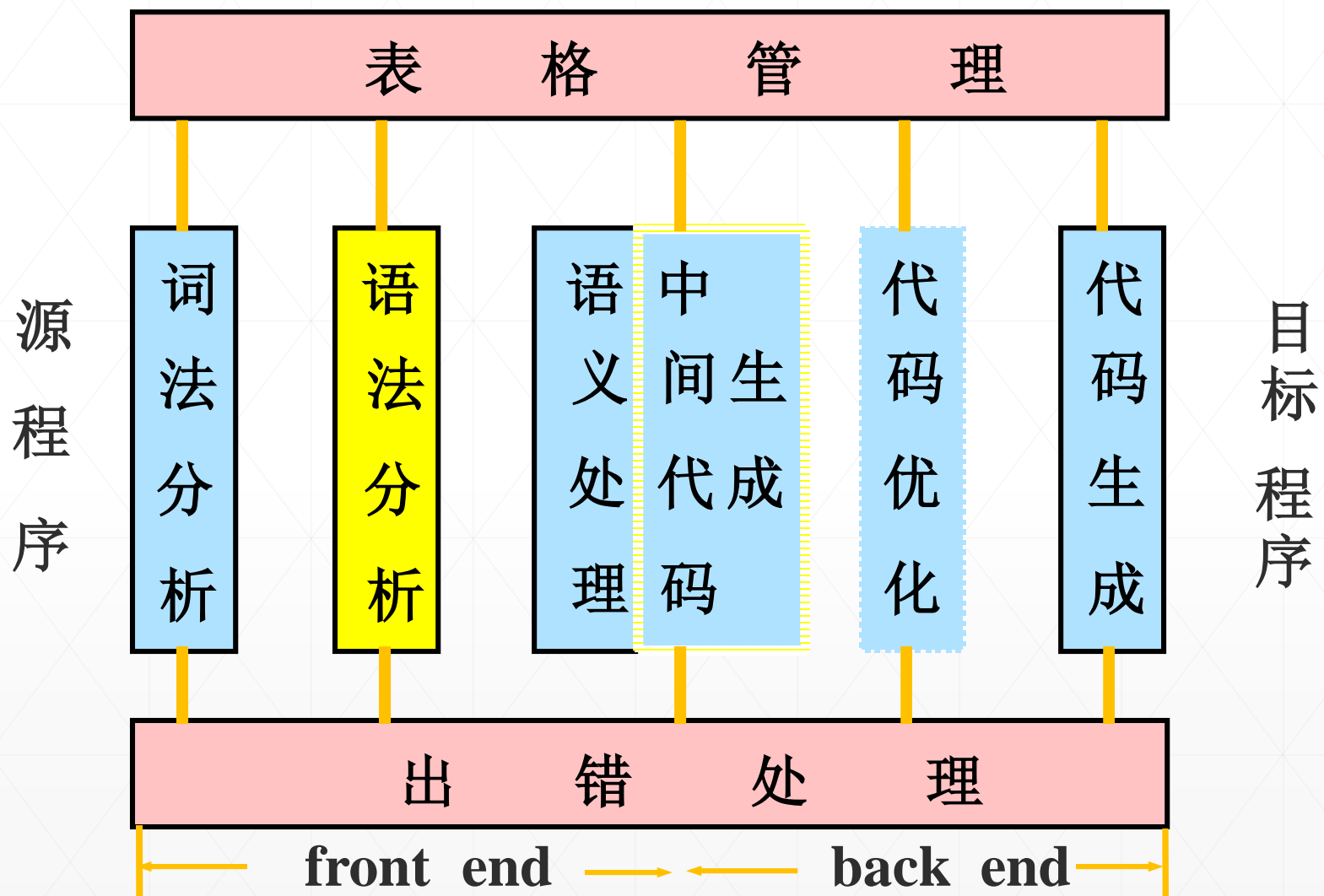




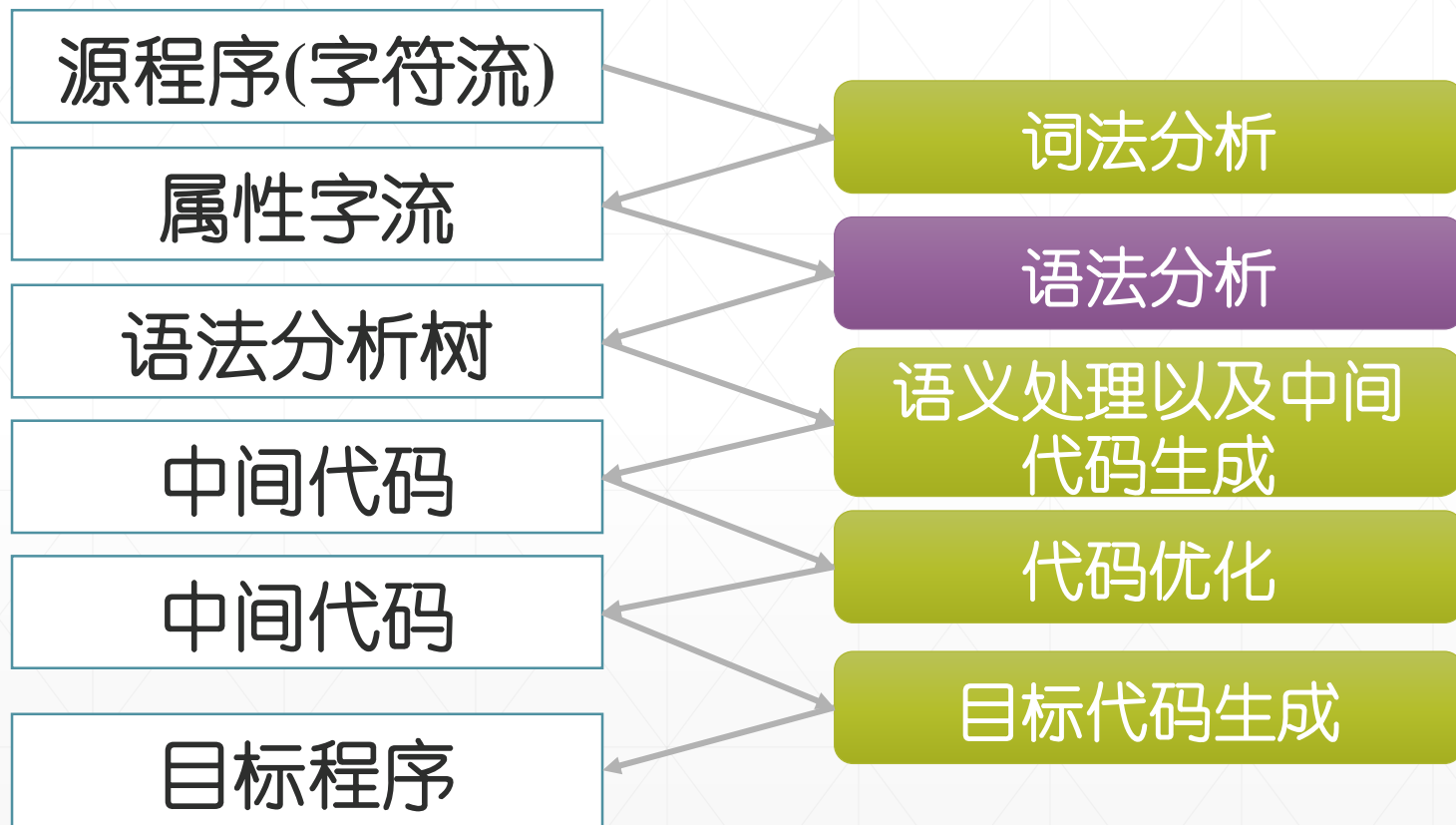
语法分析

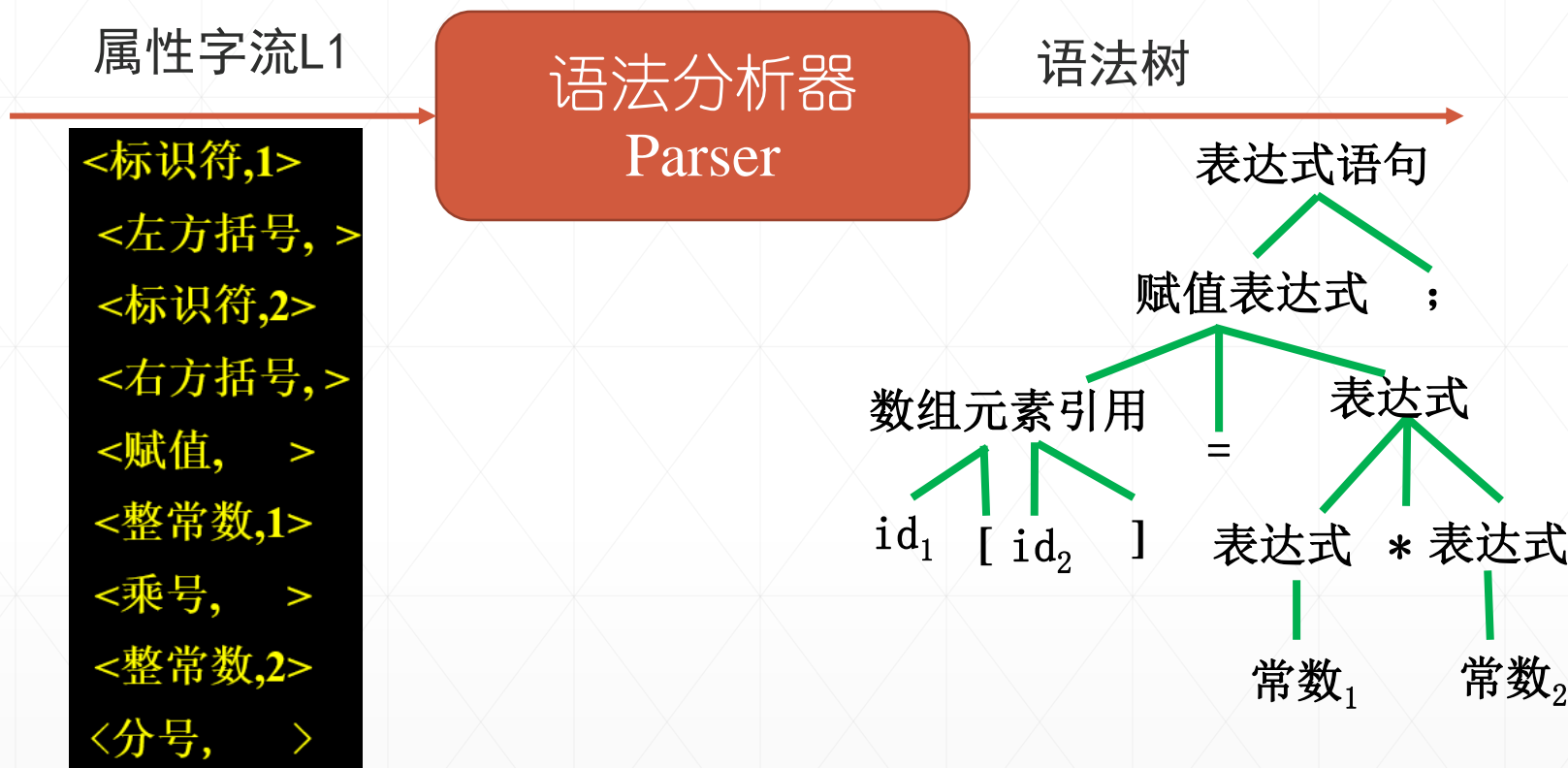
语法分析：概览





■ 基本功能

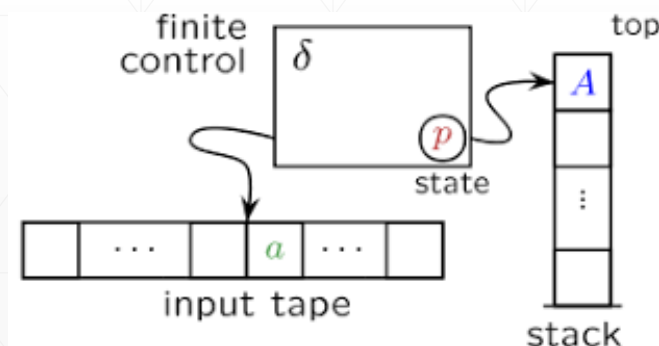




从左至右地扫描Token流，按照语言的语法规则识别语句结构，输出语法树或者语法错误信息。



- 自动生成工具
 - Bison/Yacc, CUP, ANTLR, SableCC, Beaver, JavaCC, ...
- 内部对应一个确定性的下推自动机





语法说明

如何简洁地描述合法程序的结构

上下文无关文法

语法识别

编译器如何判断输入程序是否符合说明给出的结构

LL、LR 分析器



- 讲授内容
 - 文法介绍
 - 自顶向下的分析方法
 - 自底向上的分析方法
 - 二义文法分析与错误处理
 - 自动生成工具简介(自学)



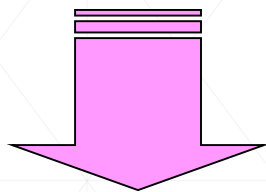
■ 关于语言：

自然语言 —— 人与人交互的工具；

程序设计语言 —— 人机交互的工具。



从文法和语言的直观概念



文法、语言的形式定义

表示方法

类型

二义性问题



2.1 文法和语言

- ➔ 2.1.1 语言的语法和语义
- 2.1.2 文法和语言的定义
- 2.1.3 文法的表示方法
- 2.1.4 语法树与二义性
- 2.1.5 文法和语言的类型

文法：自然语言的问题



dǎ gōng rén
打工
人
职场人的自称



wěi kuǎn rén
尾款
人
付完定金欠下尾款的人



shuāng jié gùn
双节棍
过完两个双十一还是单身的人



jī měi
集美
姐妹的谐音



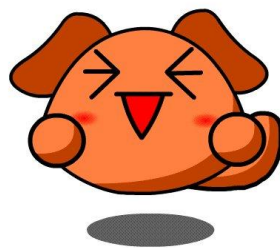
gōng jù rén
工具
人
如同工具一般被使唤的人



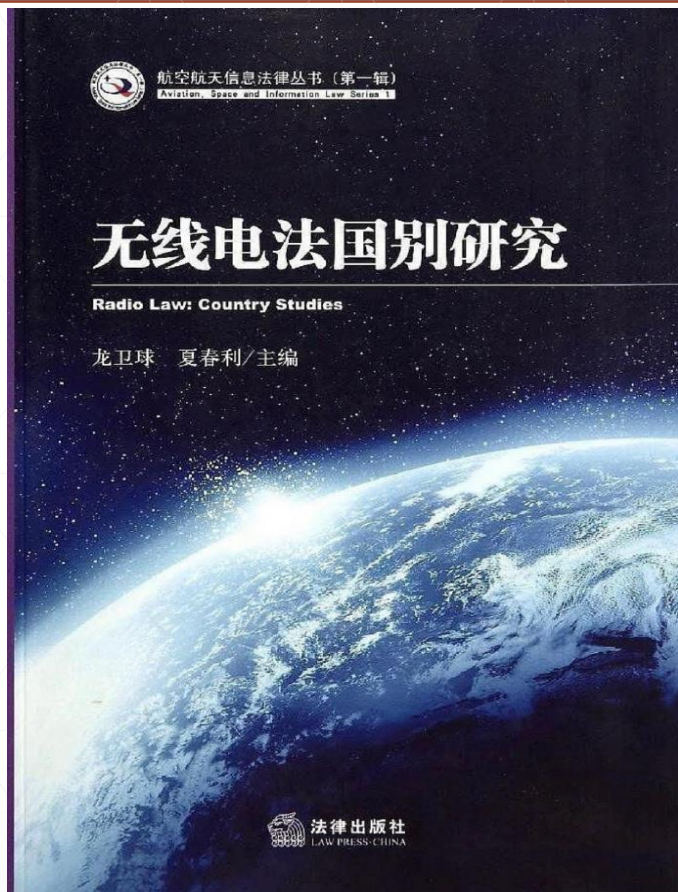
yún jiàn gōng
云监工
看直播监工建医院、送快递



nǐ xíng zhě
逆行
者
逆风前行的抗疫英雄们



2333



中文词性标注 无线电法国别研究

Jieba: 无线电/b 法国/ns 别/r 研究/vn

SnowNLP: 无线电/n 法国/ns 别/d 研究/v

PKUSeg: 无线电/n 法国/ns 别/d 研究/v

Thulac: 无线电/n 法国/ns 别/d 研究/v

HanLP: 无线电/n 法国/nsf 别/d 研究/vn

FoolNLTK: 无线/b 电法/n 国别/n 研究/n

LTP: 无线电/n 法国/ns 别/d 研

- 开会的时候，有人抽烟，老板咬牙道，“抽烟的都掐死。”
- 过几天天天天气不好。
- 我背有点驼，麻麻说“你的背得背背背佳”
- 六十老儿生一子人言非是我子也家产田园尽付与女婿外人不得争执
- 下雨天留客天留我不留



- He is not a grave man until he is a grave man.
- Time flies like an arrow, fruit flies like a banana.
- Flying planes can be dangerous.
- I'm glad I'm a man, and so is Lola.
- John saw the man on the mountain with a telescope.
- ...



- 语言要素
 - 语法：语言的描述规则
 - 语义：语言的含义

语法是一种媒介，**语义**以语法为媒介来表述。

语言是由单词按一定规则（文法）组成来表达特定意思的句子的集合。

对语言的分析集中于对句子的分析。

句子的分析依据：语言的文法规则。



■ 例：设有语句“小八哥吃大花生”。

汉语语法规则中的其中七条规则：

〈句子〉 → 〈主语〉 〈谓语〉

〈主语〉 → 〈形容词〉 〈名词〉

〈谓语〉 → 〈动词〉 〈宾语〉

〈宾语〉 → 〈形容词〉 〈名词〉

〈形容词〉 → 小 | 大

〈名词〉 → 八哥 | 花生

〈动词〉 → 吃



- 巴科斯-诺尔范式表示法，简称BNF。

- $\langle \rangle$: 表示语法成分;
- $\rightarrow / ::=$: 表示“定义为”或“由...组合成”;
- $|$: “或”，具有相同左部的产生规则用 $|$ 分开

元语言符号

元语言：描述另一个语言的语言。



■ 例：设有语句“小八哥吃大花生”。

汉语语法规则中的其中七条规则：

〈句子〉 → 〈主语〉 〈谓语〉

〈主语〉 → 〈形容词〉 〈名词〉

〈谓语〉 → 〈动词〉 〈宾语〉

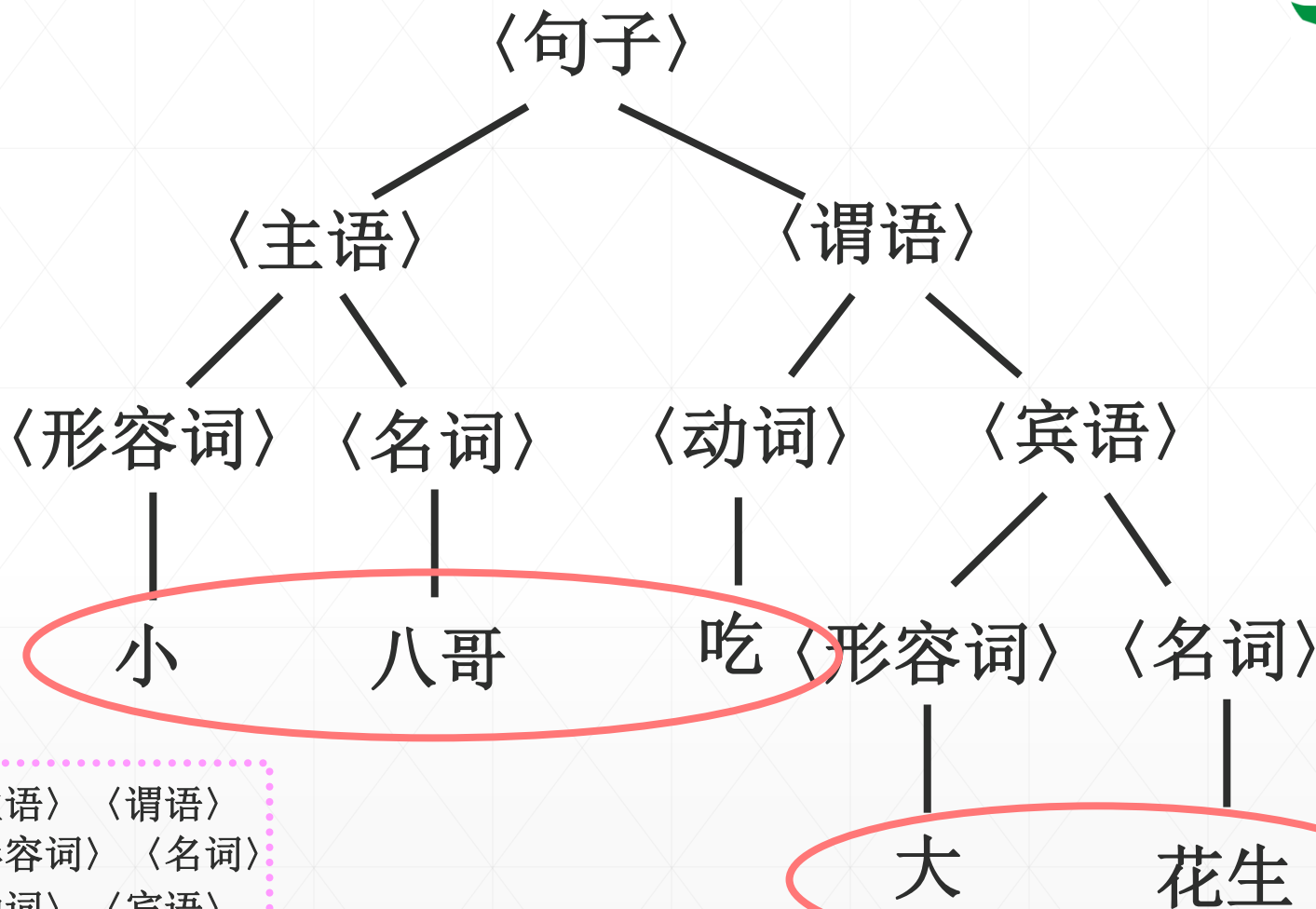
〈宾语〉 → 〈形容词〉 〈名词〉

〈形容词〉 → 小 | 大

〈名词〉 → 八哥 | 花生

〈动词〉 → 吃

语句“小八哥吃大花生”的语法分析树



〈句子〉 → 〈主语〉 〈谓语〉
〈主语〉 → 〈形容词〉 〈名词〉
〈谓语〉 → 〈动词〉 〈宾语〉
〈宾语〉 → 〈形容词〉 〈名词〉
〈形容词〉 → 小 | 大
〈名词〉 → 八哥 | 花生
〈动词〉 → 吃

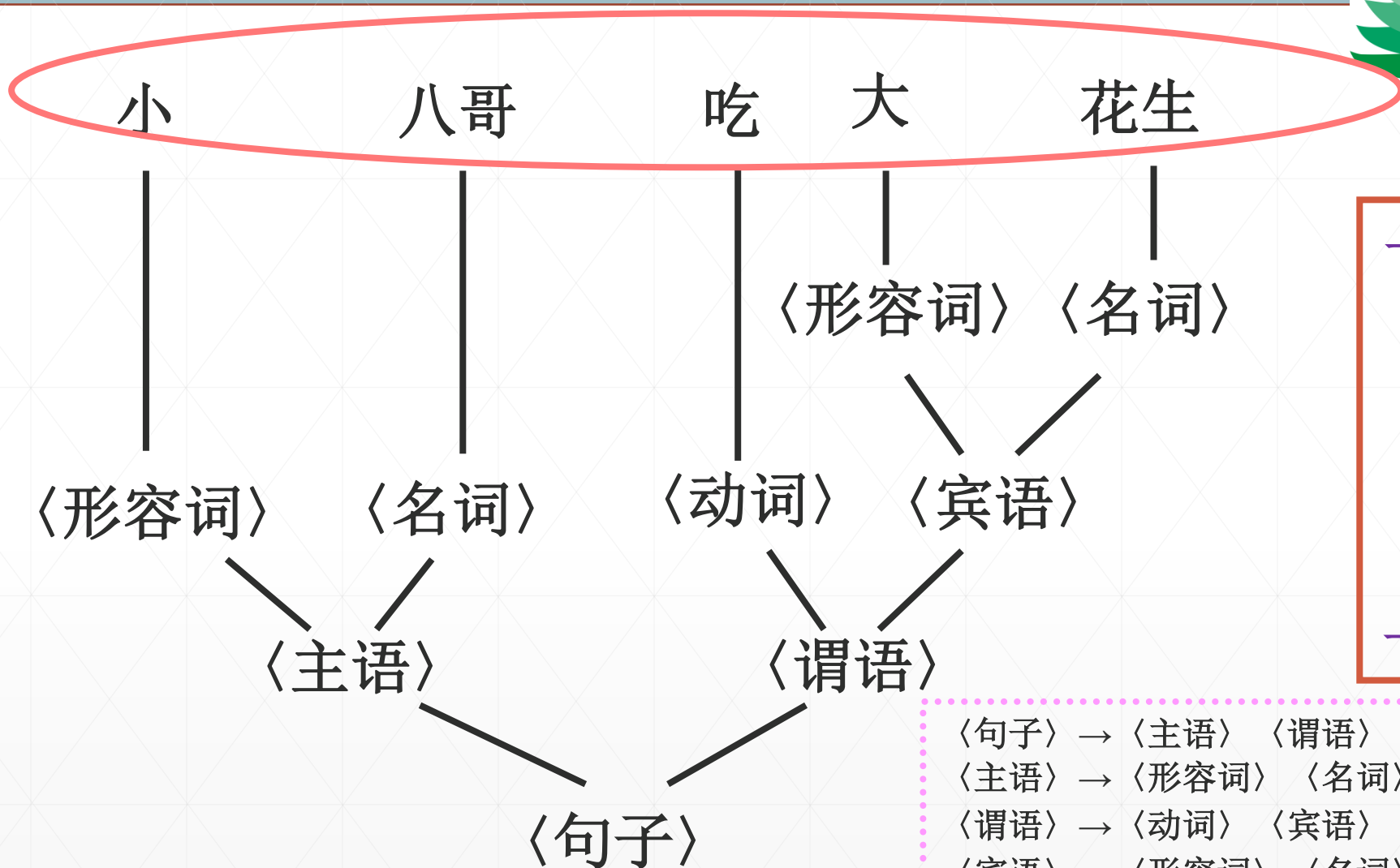


■ 句子的推导

<句子> ⇒ <主语> <谓语>
⇒ <形容词> <名词> <谓语>
⇒ 小 <名词> <谓语>
⇒ 小八哥 <谓语>
⇒ 小八哥 <动词> <宾语>
⇒ 小八哥 吃 <宾语>
⇒ 小八哥 吃 <形容词> <名词>
⇒ 小八哥吃大 <名词>
⇒ 小八哥吃大花生

<句子> → <主语> <谓语>
<主语> → <形容词> <名词>
<谓语> → <动词> <宾语>
<宾语> → <形容词> <名词>
<形容词> → 小 | 大
<名词> → 八哥 | 花生
<动词> → 吃

句子的归约



〈句子〉→ 〈主语〉 〈谓语〉
〈主语〉→ 〈形容词〉 〈名词〉
〈谓语〉→ 〈动词〉 〈宾语〉
〈宾语〉→ 〈形容词〉 〈名词〉
〈形容词〉→ 小 | 大
〈名词〉→ 八哥 | 花生
〈动词〉→ 吃



2.1 文法和语言

2.1.1 语言的语法和语义

→ 2.1.2 文法和语言的定义

2.1.3 文法的表示方法

2.1.4 语法树与二义性

2.1.5 文法和语言的类型



- 字符、字符串

任何一种语言，都是由该语言的基本字符所组成的字符串的集合。

例如，程序设计语言的基本字符集是由字母、数字、运算符等其它符号组成，则任何程序都是由这些基本字符组成的序列。

- 字母表

字母表是元素的非空有穷集合。字母表中的元素称为符号，因此字母表也称为符号集。



- 字母表

- 用希腊字母 Σ 或大写英文字母等表示字母表
- 用集合的列举法表示枚举出字母表中的符号

- 例如

- **汉语**：字母表中包括汉字、数字和标点符号等
- **机器语言**：字母表是 $\Sigma = \{0, 1\}$
- **C语言**：字母表是一切可打印字符的集合



■ 字符串/符号串

由字母表中的符号组成的任何有穷序列

已知 Σ 是字母表, Σ 上的字符串的集合 Σ^* 可递归定义如下:

- i) ϵ (由 Σ 中的0个字符组成的符号串, 称为空串) $\in \Sigma^*$;
- ii) 如果 $\alpha \in \Sigma^*$ 且 $a \in \Sigma$, 那么 $\alpha a \in \Sigma^*$;
- iii) $\alpha \in \Sigma^*$ 当且仅当 α 由有限步 i) 和 ii) 产生。

Σ^* 中的元素称为字符串。

📖 通常使用小写希腊字母或小写字母表示符号串,

如 $\alpha = STR$ 表示 α 是由符号 **S**、**T** 和 **R**, 并按此顺序组成的符号串。



■ 符号串长度

如果某符号串 α 中有 m 个 **符号**，则称其长度为 m ，记为 $|\alpha| = m$ 。

空串 ε 的长度定义为0，记为 $|\varepsilon| = 0$ 。

例：1) 定义在字母表 $\Sigma = \{0, 1\}$ 上的符号串

$$|001110| = 6$$

2) 定义在字母表 $\Sigma = \{x, y, z, =, 0, 1, ++, ;\}$ 上的符号串

$$|x=y++; z=0;| = 9$$



■ 符号串的前缀、后缀和子符号串

设 α 是一个符号串，

从 α 的尾部删去0个或若干个符号之后剩余的部分称为 α 的前缀。

从 α 的首部删去0个或若干个符号之后剩余的部分称为 α 的后缀。

若 α 的前缀(后缀)不是 α 自身，则将其称为 α 的真前缀(真后缀)。

从一个符号串中删去它的一个前缀和一个后缀之后剩余的部分称为该符号串的子符号串或子串。



- 例：设 $\alpha = abc$
 - ε, a, ab, abc 都是 α 的前缀，且除 abc 外都为真前缀
 - ε, c, bc, abc 都是 α 的后缀，且除 abc 外都为真后缀
 - abc 是 α 的前缀或后缀，但既不是真前缀也不是真后缀。
- 例：设 $\alpha = abcd$
 - $\varepsilon, a, b, c, ab, bc, cd, abc, bcd$ 及 $abcd$ 都是 α 的子串
 - ac, ad, cb, bd, ba 等都不是 α 的子串



■ 符号串的连接

设 α 和 β 是两个符号串，如果将符号串 β 直接拼接在符号串 α 之后，则称此操作为符号串 α 和 β 的连接，记作 $\alpha\beta$ 。

例如，设有字符串 $\alpha = \textcolor{red}{abc}$ ， $\beta = \textcolor{violet}{xyz}$

则 $\alpha\beta = \textcolor{red}{abc}\textcolor{violet}{xyz}$ ， $\beta\alpha = \textcolor{violet}{xyz}\textcolor{red}{abc}$ 。

 连接运算是有序的。

一般 $\alpha\beta \neq \beta\alpha$ 。



■ 符号串的方幂

设 α 是某字母表上符号串，把 α 自身连接 n 次得到符号串 β ，即 $\beta = \alpha \alpha \dots \alpha$ (n 个 α)，称 β 是符号串 α 的 n 次幂，记作 $\beta = \alpha^n$ 。

设 α 是符号串，则有

$$\alpha^0 = \varepsilon$$

$$\alpha^1 = \alpha$$

$$\alpha^2 = \alpha \alpha$$

$$\alpha^3 = \alpha^2 \alpha = \alpha \alpha^2 = \alpha \alpha \alpha$$

.....

$$\alpha^n = \alpha^{n-1} \alpha = \alpha \alpha^{n-1} = \underbrace{\alpha \alpha \dots \alpha}_{n \text{ 个}}$$

注意



■ 符号串集合的乘积

设 A 、 B 是两个符号串集合， AB 表示 A 与 B 的乘积，定义

$$AB = \{ \alpha\beta \mid (\alpha \in A) \wedge (\beta \in B) \}$$

例如，设 $A = \{ab, c\}$, $B = \{d, ef\}$,

则 $AB = \{abd, abef, cd, cef\}$

 注意：

有 $\{\varepsilon\}A = A\{\varepsilon\} = A$, $\emptyset A = A\emptyset = \emptyset$ ，其中 \emptyset 为空集。

$\emptyset = \{ \} \neq \{\varepsilon\}$ 。



■ 符号串集合的方幂

设 A 是符号串集合， A 自身的乘积可以用方幂表示。

$$A^0 = \{\varepsilon\}$$

$$A^1 = A$$

$$A^2 = AA$$

$$A^3 = A^2A = AAA$$

.....

$$A^n = A^{n-1}A = AA \dots A$$

注意

显然有：

$$A^{i+j} = A^i A^j$$

例： $A = \{ ab, x, aby \}$,

则 $A^2 = AA$

$$= \{ abab, abx, ababy, xab, xx, xaby, abyab, abyx, abyaby \}$$



■ 符号串集合的并

设 P 、 Q 为字符串集，集合 $P \cup Q$ 为 P 和 Q 的并，它的元素是 P 或 Q 中的元素。

例如， $P=\{0, 1, 01\}$ $Q=\{0, 10, 11, 00\}$,

则 $P \cup Q = \{0, 1, 01, 10, 11, 00\}$



■ 符号串集合的闭包

设 A 为符号串集,

A 的正闭包记作 A^+ , 定义为,

$$A^+ = A^1 \cup A^2 \cup \dots \cup A^n \cup \dots$$

A 的自反闭包记作 A^* , 定义为,

$$A^* = A^0 \cup A^1 \cup A^2 \cup \dots \cup A^n \cup \dots$$

由定义知,

$$\begin{cases} \underline{A^+ = AA^*} \\ \underline{A^* = A^0 \cup A^+} \end{cases}$$

例如, 设有 $A = \{01, 10\}$

则 $A^* = \{\epsilon, 01, 10, 0101, 0110, 1001, 1010, 010101, 010110, \dots\}$

$A^+ = \{01, 10, 0101, 0110, 1001, 1010, 010101, 010110, \dots\}$



▪ 串集合运算应用

- 设有 $L = \{ A..Z, a..z \}$, $D = \{ 0, 1, 2, \dots, 9 \}$
 - $L \cup D = \{ \text{由字母和数字构成的集合} \}$
 - $LD = \{ \text{所有一个字母后跟随一个数字组成的字符串的集合} \}$
 - $L^* = \{ \text{由所有字母按任意顺序组成的字符串且包含} \varepsilon \text{ 所构成的集合} \}$
 - $L(L \cup D)^* = \{ \text{由所有一个字母开头后跟随字母或数字组成的字符串或} \varepsilon \text{ 的集合} \}$

文法：文法形式定义



一部文法 G 是一个四元组 $G = (V_N, V_T, S, P)$

其中：

V_N ：非空有限的非终结符号集(一般用大写字母表示)。

其中的元素称为非终结符，或语法变量。

代表了一个语法范畴，表示某种语法结构。

V_T ：非空有限的终结符号集（一般用小写字母表示）。

V_N 、 V_T 合称为文法 G 的符号集 V ， $V = V_T \cup V_N$ 。 $V_T \cap V_N = \emptyset$ 。

S ：文法的开始符号或识别符号，亦称公理， $S \in V_N$ 。

S 代表语言最终要得到的语法范畴。

P ：有限产生式集。

按一定格式书写的定义语法范畴的文法规则，文法的实体。



产生式的形式(BNF):

$$\alpha \rightarrow \beta \text{ 或 } \alpha ::= \beta$$

其中： α 称为产生式的左部，

β 称为产生式的右部或称为 α 的候选式。

$\alpha \in V^+$ ，且 α 中至少包含 V_N 中的一个元素， $\beta \in V^*$ 。

注意，开始符号 S 至少且必须在文法某个产生式的左部出现一次。

以 S 为开始符号的文法 G 可记为 $G(S)$ 。



例：简单的算术表达式文法 G_1 定义为

$\{\{E\}, \{i, +, *, (,)\}, E, \{E \rightarrow i \mid E + E \mid E * E \mid (E)\}\}$

四元式形式

文法的简化表示：

例：数字文法

$\langle \text{NUMBER} \rangle \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid \dots \mid 9$

其中： $V_N = \{ \text{NUMBER} \}$, $V_T = \{0, 1, 2, \dots, 9\}$,

$S = \text{NUMBER}$, P 为定义式本身。



汉语语法规则中的其中七条规则：

$S \bullet \bullet \bullet$

〈句子〉 → 〈主语〉 〈谓语〉

〈主语〉 → 〈形容词〉 〈名词〉

〈谓语〉 → 〈动词〉 〈宾语〉

〈宾语〉 → 〈形容词〉 〈名词〉

〈形容词〉 → 小 | 大

〈名词〉 → 八哥 | 花生

〈动词〉 → 吃

P

V_T

V_N



1. 语言的非形式化定义

给定一部文法 G , 从 G 的开始符号 S 出发, 反复使用产生式对非终结符进行替换, 最后所能得到的终结符号串的全体, 即为文法 G 所描述的语言 $L(G)$ 。

例：设有文法 G

$$S \rightarrow P \mid aPb$$

$$P \rightarrow ba \mid bQa$$

$$Q \rightarrow ab$$

写出该文法所描述的语言。



■ 直接推导 “ \Rightarrow ”

$\lambda = \alpha A \beta$, $\mu = \alpha \gamma \beta$, ($\alpha, \beta, \gamma \in V^*$).

P 中存在一条规则 $A \rightarrow \gamma$,

称 λ 直接推导出 μ (或 μ 直接归约到 λ)

记作: $\lambda \Rightarrow \mu$ 。

■ 直接推导序列

如果存在 $\lambda = \alpha_0 \Rightarrow \alpha_1, \alpha_1 \Rightarrow \alpha_2, \dots, \alpha_{n-1} \Rightarrow \alpha_n = \mu$

或 $\alpha_0 \Rightarrow \alpha_1 \Rightarrow \alpha_2 \Rightarrow \alpha_3 \Rightarrow \dots \Rightarrow \alpha_{n-1} \Rightarrow \alpha_n$,

则 λ 经过 n 步 ($n > 0$) 可以推导出 μ , 记作: $\lambda \overset{\pm}{\Rightarrow} \mu$ 。当

$\lambda \overset{\pm}{\Rightarrow} \mu$ 或 $\lambda = \mu$, 记作: $\lambda \overset{*}{\Rightarrow} \mu$ 。



■ 句型

对文法 $G[S]$ ，若 $S \xRightarrow{*} \alpha$ ($\alpha \in V^*$)，则称 α 为 $G[S]$ 的句型。

■ 句子

对文法 $G[S]$ ，若 $S \xRightarrow{*} \alpha$ ($\alpha \in V_T^*$)，则称 α 为 $G[S]$ 的句子。

■ 最左(右)推导

推导过程中，总是对句型中的最左(右)边的非终结符进行替换，称为最左(右)推导。

■ 规范推导/规范句型/ 规范归约

最右推导也称**规范推导**。仅用规范推导得到的句型称为**规范句型**。规范推导的逆序为**规范归约**。

文法： 语言的定义



例：设有文法 $G[E]$: $E \rightarrow E * E \mid E + E \mid (E) \mid i$

判断\$1: $i*i+i$ 是该文法的句子

$$\underline{E} \xRightarrow{L} E * E \xRightarrow{L} i * E \xRightarrow{L} i * E + E \xRightarrow{L} i * i + E \xRightarrow{L} \underline{i * i + i}$$

最左推导序列

句子

$$E \xRightarrow{R} E * E \xRightarrow{R} \underline{E * E + E} \xRightarrow{R} E * E + i \xRightarrow{R} E * i + i \xRightarrow{R} \underline{i * i + i}$$

句型

最右推导/规范推导序列

规范句型



2. 语言的形式定义

■ 语言

文法 G 所产生(描述)的语言 $L(G)$:

$$L(G) = \{ \alpha \mid \alpha \in V_T^* \wedge S \Rightarrow \alpha, \text{ } S \text{ 是文法 } G \text{ 的开始符号} \}$$

文法： 语言的定义



文法 \Rightarrow 语言

语言:句子的集合。

由给定文法构造句子的思想:

从文法的开始符号出发,
利用直接推导替换非终结符,
直至最终符号串全由终结符号组成。

例: 设有文法 G

$$S \rightarrow P \mid aPb$$

$$P \rightarrow ba \mid bQa$$

$$Q \rightarrow ab$$

写出该文法所描述的语言。



$$S \Rightarrow P \Rightarrow ba$$

$$S \Rightarrow P \Rightarrow bQa \Rightarrow baba$$

$$S \Rightarrow aPb \Rightarrow abab$$

$$S \Rightarrow aPb \Rightarrow abQab \Rightarrow ababab$$

则： $L(G) = \{\underline{ba}, \underline{baba}, \underline{abab}, \underline{ababab}\}$

文法G:

$$S \rightarrow P \mid aPb$$

$$P \rightarrow ba \mid bQa$$

$$Q \rightarrow ab$$



■ 文法的递归

设有文法 G , $A \rightarrow \gamma$ 是 G 的产生式, 若 γ 具有 $\alpha A \beta$ 的形式, 或 $\gamma \xrightarrow{+} \alpha A \beta$, 则称 G 是递归文法。

间接递归文法

直接递归文法

若 $\alpha = \varepsilon$, 则 G 为左递归文法。

若 $\beta = \varepsilon$, 则 G 为右递归文法。

递归文法 { 直接递归
 间接递归
 左（右）递归



例： 设有文法 G_1 ： $E \rightarrow E+E \mid E * E \mid (E) \mid i$

其中有 $E \rightarrow E \dots$ 这样的产生式，所以文法 G_1 是**直接左递归文法**。

例： 设有文法 G_2 ：

$$T \rightarrow Qc \mid c$$

$$Q \rightarrow Rb \mid b$$

$$R \rightarrow Ta \mid a$$

其中有 $T \rightarrow Qc$ $Qc \Rightarrow Rbc \Rightarrow Tabc$,

或 $T \overset{+}{\Rightarrow} Tabc$ ，则文法 G_2 是**间接左递归文法**。



文法和语言之间的相互转换举例



文法和语言之间的相互转换举例



例：设有文法 G : $S \rightarrow S0 \mid 0$ 求 $L(G)$?

$$L(G) = \{ 00^n \mid n \geq 0 \}$$

$$= \{ 0^m \mid m \geq 1 \}$$

$$G': S \rightarrow 0S \mid 0$$

$$L(G') = L(G)$$

后面为 0^n ($n \geq 1$)

S 替换为 0

■ 文法等价

若 $L(G_1) = L(G_2)$, 则称文法 G_1 和 G_2 是等价的。

文法和语言之间的相互转换举例



例：设有文法 G ：

$$S \rightarrow 0S1 \mid \varepsilon$$

S 替换为空串

求 $L(G)$?

$$\begin{aligned} L(G) &= \{ 0^n \varepsilon 1^n \mid n \geq 0 \} \\ &= \{ 0^n 1^n \mid n \geq 0 \} \end{aligned}$$

前面为 0^n , 后面为 $1^n (n \geq 1)$

例：设有语言 $L(G1) = \{ a b^n a \mid n \geq 0 \}$,

给出文法 $G1$?

定义一个语法成分

$$G1(S) : S \rightarrow aRa$$

$$R \rightarrow Rb \mid \varepsilon$$

$$G1(S) : S \rightarrow aT \quad T \rightarrow Ra$$

$$R \rightarrow Rb \mid \varepsilon$$

$$G1(S) : S \rightarrow aa \mid aRa$$

$$R \rightarrow b \mid Rb$$



例： 设有字母表 $\{a,b,(,)\}$ 上的语言 L ：

$$L = \{ a(b^n)a \mid n \geq 0 \}$$

写出描述语言 L 的文法。

$$G: S \rightarrow a(B)a$$

$$B \rightarrow Bb \mid \varepsilon$$

$$G': S \rightarrow a()a \mid a(B)a$$

$$B \rightarrow Bb \mid b$$

文法和语言之间的相互转换举例



例：写出文法 G $S \rightarrow Sb \mid R$

$R \rightarrow aRb \mid ab$ 描述的语言 $L(G)$

后面为 b^m
($m \geq 1$)

产生串 $a^n b^n$ ($n \geq 1$)

$$\begin{aligned} \text{所以 } L(G) &= \{a^n b^n b^m \mid n \geq 1, m \geq 0\} \\ &= \{a^n b^{n+m} \mid n \geq 1, m \geq 0\} \\ &= \{a^n b^j \mid j \geq n \geq 1\} \end{aligned}$$

文法和语言之间的相互转换举例



例： 设 定义在字母表 $\{a,b,(,)\}$ 上的语言

$$L = \{ (a^n)(b^n) \mid n=1,2,3, \dots \}$$

写出该语言的文法。

$$G: \quad S \rightarrow (A)(B)$$

$$A \rightarrow aA \mid a$$

$$B \rightarrow bB \mid b$$

✗

一个语法成分 (非终结符) 描述 a 、 b 个数相等的性质

a, b 个数可以不相同

$$B \rightarrow aBb \mid a)(b$$

$$G: \quad S \rightarrow (B)$$

$$B \rightarrow a)(b \mid aBb$$

文法和语言之间的相互转换举例



例： 设 定义在字母表 $\{a,b\}$ 上的语言

$$L = \{ (ab)^n \mid n=1,2,3, \dots \}$$

写出该语言的文法。

$$G: \quad S \rightarrow abS \mid ab$$



他

你

我



定义在字母表 $\{a,b\}$ 上的语言

$$L = \{ a^n b^n \mid n=1,2,3, \dots \}$$

$$G: \quad S \rightarrow aSb \mid ab$$

文法和语言之间的相互转换举例



例:

设 $L(G) = \{ 1^n 0^m 1^m 0^n \mid n, m \geq 0 \}$,

求 $G(S)$?

一个语法成分 (非终结符) 描述

$S \rightarrow 1S0 \mid B$

一个语法成分 (非终结符) 描述

$B \rightarrow 0B1 \mid \varepsilon$

$G(S):$ $S \rightarrow 1S0 \mid B$
 $B \rightarrow 0B1 \mid \varepsilon$

文法和语言之间的相互转换举例



例: 设语言 $S = \{a^i b^j \mid 0 \leq i \leq j\}$, 满足 $L(G) = S$ 的文

法 G 为 【 $T \rightarrow AB \quad A \rightarrow aAb \mid \varepsilon \quad B \rightarrow Bb \mid \varepsilon$ 】。

$$T \rightarrow Tb \mid A \quad A \rightarrow aAb \mid \varepsilon$$

$$T \rightarrow aTb \mid B \quad B \rightarrow Bb \mid \varepsilon$$

$$T \rightarrow aTb \mid Tb \mid \varepsilon$$

$$j = i + m \quad m \geq 0$$

$$a^i b^j = a^i b^{i+m} = a^i b^i b^m = a^i b^m b^i$$

文法和语言之间的相互转换举例



例: 设语言 $S = \{a^i b^j \mid 0 \leq i \leq j \leq 2i\}$, 满足 $L(G) = S$ 的文

法 G 为 【 $T \rightarrow aTbb \mid A \quad A \rightarrow aAb \mid \varepsilon$
 $T \rightarrow aTb \mid aTbb \mid \varepsilon$ 】。

$$j = i + m \quad 0 \leq m \leq i$$

$$i = m + n \quad n \geq 0$$

$$j = 2m + n \quad n \geq 0 \quad m \geq 0$$

$$a^i b^j = a^{m+n} b^{2m+n} = a^m a^n b^n b^{2m}$$



$$\text{令: } j = i + m + n$$

$$0 \leq m \leq i \quad 0 \leq n \leq m$$



设语言 $S = \{a^i b^j \mid 0 \leq i \leq j \leq 3i\}$, 满足 $L(G) = S$ 的文法 G 为 【 _____ 】。

文法和语言之间的相互转换举例



例: 写一个文法, 使其语言是奇整数的集合, 每个奇整数不以0为前导。

解: 语言集合 $\{1, -1, 3, -3, 5, -5, 7, -7, 9, -9, 11, -11, \dots\}$

$V_T = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, -\}$

$V_N = \{S, E, D, C, B, A\}$

开始符号: S

$S \rightarrow E | -E$

$E \rightarrow A | BA | BDA$

$D \rightarrow CD | C$

$C \rightarrow 0 | B$

$B \rightarrow 2 | 4 | 6 | 8 | A$

$A \rightarrow 1 | 3 | 5 | 7 | 9$

文法和语言之间的相互转换举例



例: 设集合 A 为字母表 $\Sigma=\{0,1\}$ 上的有相同个数的 0 和 1 组成的符号串集合，给出正确描述集合 A 的文法。

解: 语言集合 $\{\varepsilon, 01, 10, 0011, 0101, 0110, 1001, 1010, 1100, \dots\}$

$$V_T = \{0, 1\}$$

开始符号: S

$$V_N = \{S, P, Q\}$$

P : 1 的个数 = 0 的个数 + 1,

Q : 0 的个数 = 1 的个数 + 1,

$$S \rightarrow \varepsilon \mid 0P \mid 1Q$$

$$P \rightarrow 0PP \mid 1S$$

$$Q \rightarrow 0S \mid 1QQ$$



2.1 文法和语言

2.1.1 语言的语法和语义

2.1.2 文法和语言的定义

 2.1.3 文法的表示方法

2.1.4 语法树与二义性

2.1.5 文法和语言的类型



1. BNF表示法

元语言符号集: $\{ \rightarrow (::=), < >, | \}$

2. 扩充BNF表示法 (EBNF)

元语言符号集: $\{ \rightarrow (::=), < >, |, \{ \}, (), [] \}$

3. 语法图(上下文无关文法)



2. 扩充BNF表示法 (EBNF)

1) 花括号 — $\{t\}_m^n$

串 t 重复的最大次数

串 t 重复的最小次数

省略 m, n : t 可重复0到任意多次。

2) 圆括号 — “(” “)”

提取产生式中的公共因子，简化产生式的表示。

3) 方括号—— $[t]$
 t 可有可无。

文法的表示方法



例： FORTRAN语言中标识符的定义：

长度 ≤ 8 的字母开头后跟字母或数字的字符串

$\langle \text{FORTRAN标识符} \rangle \rightarrow \langle \text{字母} \rangle \{ \langle \text{字母} \rangle | \langle \text{数字} \rangle \}_0^7$

例： 有文法规则

$$U \rightarrow xa \mid xb \mid \dots \mid xz$$

等价于 $U \rightarrow x (a \mid b \mid \dots \mid z)$



注意：

与终结符区分开。

例： 设有条件语句的文法

$\langle \text{条件语句} \rangle \rightarrow \langle \text{如果子句} \rangle | \langle \text{如果子句} \rangle \text{else} \langle \text{语句} \rangle$

$\langle \text{如果子句} \rangle \rightarrow \text{if} \langle \text{布尔表达式} \rangle \text{ then} \langle \text{语句} \rangle$

$\langle \text{条件语句} \rangle \rightarrow \langle \text{如果子句} \rangle [\text{else} \langle \text{语句} \rangle]$



3. 语法图(上下文无关文法)

❖ 用图描述产生式规则

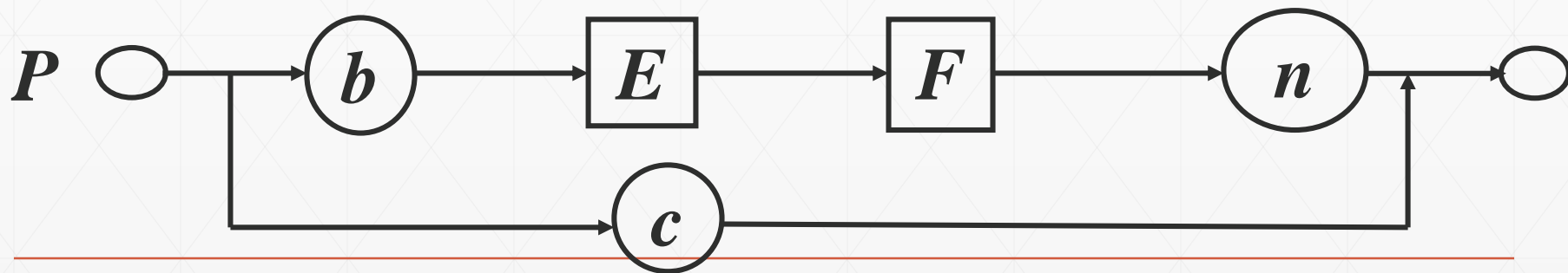
❖ 每个图都有一个起始结点和一个终止结点，其他的结点标记为文法符号。

❖ **终结符**结点用圆形表示。

❖ **非终结符**结点用方形表示。

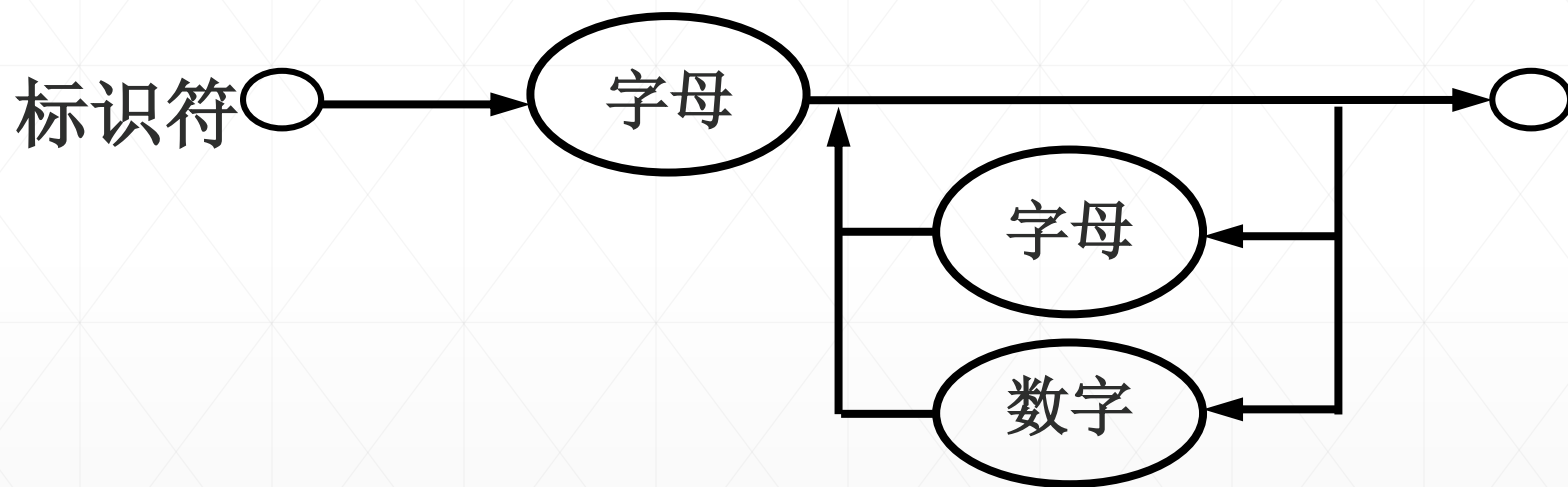
❖ 从起始结点到终止结点的所有路径(标记为结点序列)定义**候选式**。

例：产生式 $P \rightarrow bEFn \mid c$ 。表示为





例： PASCAL语言中标识符的定义如下图。





2.1 文法和语言

2.1.1 语言的语法和语义

2.1.2 文法和语言的定义

2.1.3 文法的表示方法

 2.1.4 语法树与二义性

2.1.5 文法和语言的类型



语法分析树是句子结构的图形表示，展示出句子的推导过程，也利于理解句子语法结构的层次。

■ 表示：

分析树的根结点



G 的 S

分析树的中间结点



G 的产生式左部



父子结点间关系



G 的产生式规则

叶结点从左到右连接的符号串



句型

语法分析树与二义性



例： 设有无符号整数的文法

$\langle \text{无符号整数} \rangle \rightarrow \langle \text{数字串} \rangle$

$\langle \text{数字串} \rangle \rightarrow \langle \text{数字串} \rangle \langle \text{数字} \rangle | \langle \text{数字} \rangle$

$\langle \text{数字} \rangle \rightarrow 0 | 1 | 2 | \dots | 9$

句子25的最左推导：

$\langle \text{无符号整数} \rangle \Rightarrow \langle \text{数字串} \rangle$

$\Rightarrow \langle \text{数字串} \rangle \langle \text{数字} \rangle$

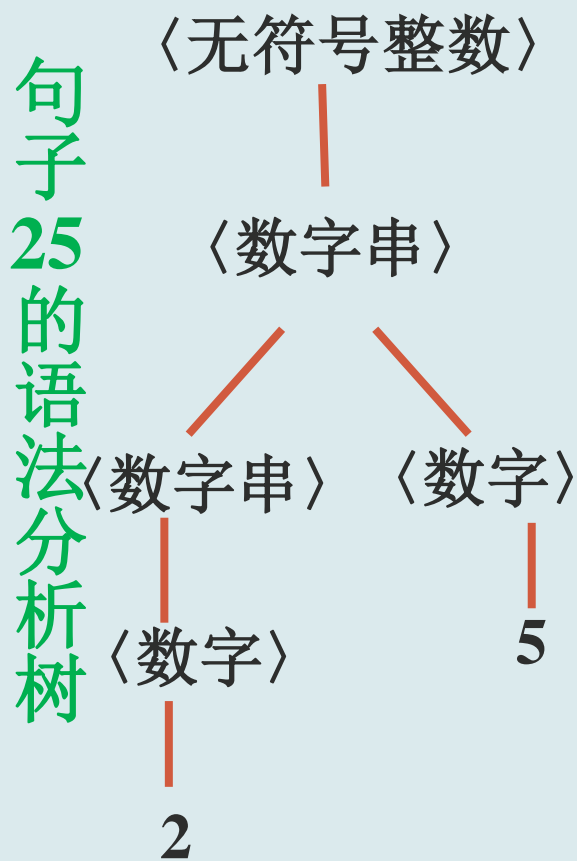
$\Rightarrow \langle \text{数字} \rangle \langle \text{数字} \rangle \Rightarrow 2 \langle \text{数字} \rangle \Rightarrow 25$

句子25的最右推导：

$\langle \text{无符号整数} \rangle \Rightarrow \langle \text{数字串} \rangle$

$\Rightarrow \langle \text{数字串} \rangle \langle \text{数字} \rangle$

$\Rightarrow \langle \text{数字串} \rangle 5 \Rightarrow \langle \text{数字} \rangle 5 \Rightarrow 25$

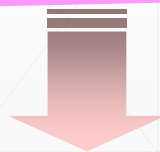




任何一个句型的一棵分析树包括了这个句型的所有可能的推导过程？

一个句型是否只对应一棵分析树？

或者是否只有惟一的最左(右)推导？



文法二义性问题



■ 二义文法

一部文法 G ，如果至少存在一个句子（或句型），
有两棵(或两棵以上)不同的分析树（或最左推导或最
右推导），

称该句子（或句型）是二义性的。

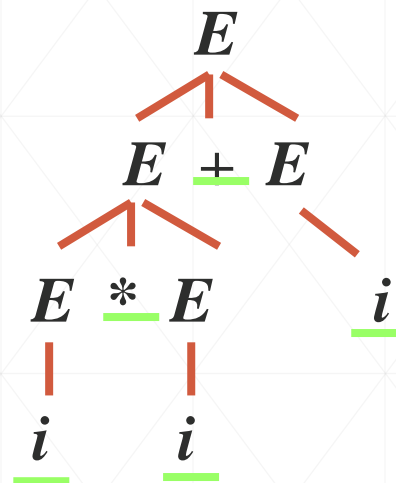
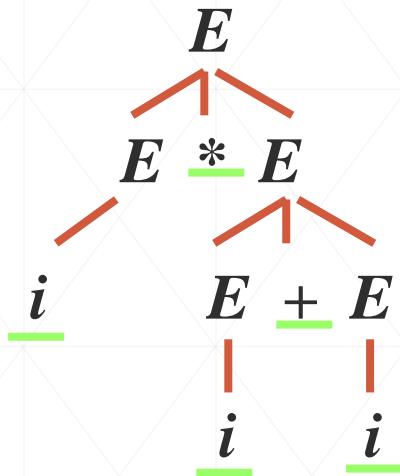
包含有二义性句子（或句型）的文法称为二义文法。
否则，该文法是无二义性的。

定义提供了对给定文法在某一范围内判定是否是
二义性文法的充分条件。

语法分析树与二义性



例：设有文法 G_1 ： $E \rightarrow E+E \mid E * E \mid (E) \mid i$



文法 G_1 中的句子 $i*i+i$ 存在两个不同的最左推导，所以文法 G_1 为二义文法。

$$E \underset{\text{L}}{=} E * E \underset{\text{L}}{=} i * E \underset{\text{L}}{=} i * E + E \underset{\text{L}}{=} i * i + E \underset{\text{L}}{=} i * i + i$$

$$E \underset{\text{L}}{=} E + E \underset{\text{L}}{=} E * E + E \underset{\text{L}}{=} i * E + E \underset{\text{L}}{=} i * i + E \underset{\text{L}}{=} i * i + i$$

文法和语言之间的相互转换举例



例:设集合A为字母表 $\Sigma=\{0,1\}$ 上的有相同个数的0和1组成的符号串集合，给出正确描述集合A的文法。

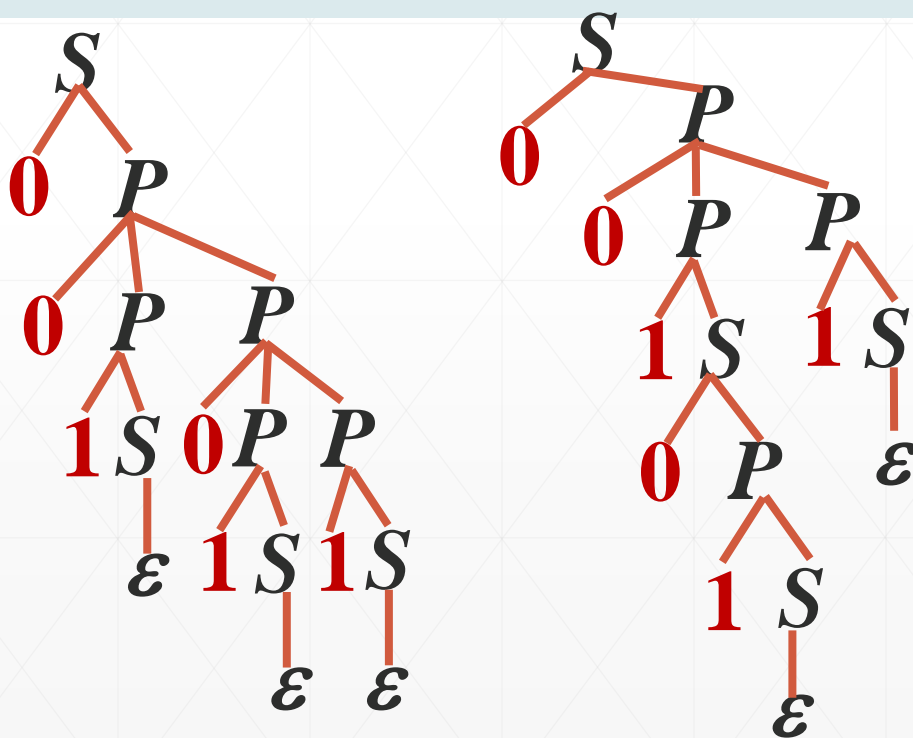
$$S \rightarrow \varepsilon \mid 0P \mid 1Q$$

$$P \rightarrow 0PP \mid 1S$$

$$Q \rightarrow 0S \mid 1QQ$$

句子**001011**

存在两棵不同的分析树，所以文法为二义文法。





🔔 注意:

文法的二义性与语义的二义性是完全不同的概念。并非文法是二义的，语义就二义。

Time Flies. 既有语法又有语义的二义性

i+i+i 有语法的二义性但无语义的二义性

你真可以。 有语义的二义性但无语法的二义性



文法二义性的消除

例如， 对有文法 G_1 ： $E \rightarrow E+E \mid E * E \mid (E) \mid i$

1) 分析二义性原因

- a) 运算符 “+” 和 “*” 未体现优先级；
- b) “+” 和 “*” 自身结合规则不明确；

2) 构造 G_1' ，使 $L(G_1) = L(G_1')$

$$G_1' : \quad E \rightarrow T \mid E+T$$

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

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



2.1 文法和语言

2.1.1 语言的语法和语义

2.1.2 文法和语言的定义

2.1.3 文法的表示方法

2.1.4 语法树与二义性

 2.1.5 文法和语言的类型



■ N.Chomsky 分类：

四类文法：

0型文法（短语文法）

1型文法（上下文有关文法）

2型文法（上下文无关文法）

3型文法（线性文法、正则文法）



1. 0型文法（短语文法）

文法 G 中的规则 $\alpha \rightarrow \beta$ 不加任何限制。

可用0型文法描述的语言为0型语言 L_0 ,

0型语言可由图灵机（Turing）来识别。

例：设有文法 G

$S \rightarrow 0|AC0B$

$C0 \rightarrow 00C$

$CB \rightarrow DB/E$

$0D \rightarrow D0$

$AD \rightarrow AC$

$0E \rightarrow E0$

$AE \rightarrow \varepsilon$

G 是0型文法

$L(G) = \{0^n \mid n \text{ 为 } 2 \text{ 的非负整数次幂}\}$



2. 1型文法（上下文有关文法）

- 每个产生式限制为：

$$\alpha A \beta \rightarrow \alpha \gamma \beta$$

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

可用1型文法描述的语言为1型语言 L_1 ，

1型语言可由线性有界自动机来识别。

文法和语言的类型



例：设有文法 G_1

$$S \rightarrow aSBC \mid abC$$

$$\left. \begin{array}{l} CB \rightarrow CD \\ CD \rightarrow BD \\ BD \rightarrow BC \end{array} \right\} \Rightarrow CB \Rightarrow BC$$

$$bB \rightarrow bb$$

$$bC \rightarrow bc$$

$$cC \rightarrow cc$$



$$S \xRightarrow{*} a^{n+1}bC(BC)^n$$

$$= a^{n+1}b(CB)^nC$$

$$\xRightarrow{*} a^{n+1}b(BC)^nC$$

$$= a^{n+1}bB(CB)^{n-1}C^2$$

$$\Rightarrow a^{n+1}bb(CB)^{n-1}C^2$$

$$\xRightarrow{*} a^{n+1}bb(BC)^{n-1}C^2$$

$$\xRightarrow{*} a^{n+1}b^{n+1}C^{n+1}$$

G_1 是1型文法， G_1 产生的语言：

$$L(G_1) = \{a^n b^n c^n \mid n \geq 1\}$$

$$\xRightarrow{*} a^{n+1}b^{n+1}c^{n+1}$$



3. 2型文法（上下文无关文法）

每个产生式限制形如：

$$A \rightarrow \alpha$$

其中， $A \in V_N$ ， $\alpha \in (V_T \cup V_N)^*$

能用2型文法描述的语言为2型语言 L_2 ，

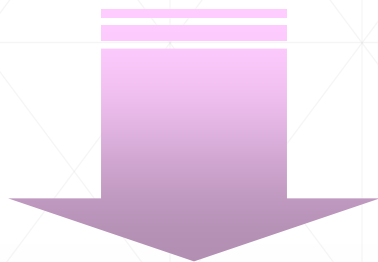
2型语言可由非确定的下推自动机来识别。



例：设有文法 G_2

$$S \rightarrow Ac \mid Sc$$

$$A \rightarrow ab \mid aAb$$



G_2 是2型文法

G_2 产生的语言：

$$L(G_2) = \{a^n b^n c^m \mid n, m \geq 1\}$$



4. 3型文法（正则文法、线性文法）

右线性文法：每个产生式形如

$$A \rightarrow \alpha B \quad \text{或} \quad A \rightarrow \alpha$$

左线性文法：每个产生式形如：

$$A \rightarrow B \alpha \quad \text{或} \quad A \rightarrow \alpha$$

其中， $A, B \in V_N$ ， $\alpha \in V_T^*$

右线性文法和左线性文法统称为3型文法。

也叫正则文法或线性文法。

能用3型文法描述的语言为3型语言 L_3 ，

3型语言可由确定的有限状态自动机来识别。

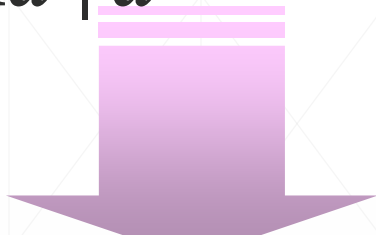


例：设有文法 G_3

$$S \rightarrow Bc \mid Sc$$

$$B \rightarrow Ab \mid Bb$$

$$A \rightarrow Aa \mid a$$



G_3 是左线性文法

G_3 是3型文法

G_3 产生的语言：

$$L(G_3) = \{a^n b^m c^k \mid n, m, k \geq 1\}$$

文法和语言的类型



例：设有文法 G_3

$$S \rightarrow aB \mid c$$

$$B \rightarrow Sb \mid \underline{b}$$



G_3 不是右线性文法，也不是左线性文法

G_3 不是3型文法

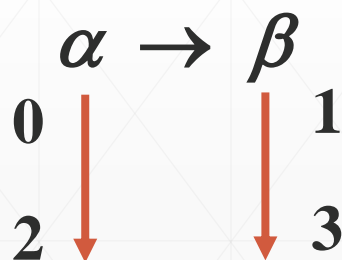
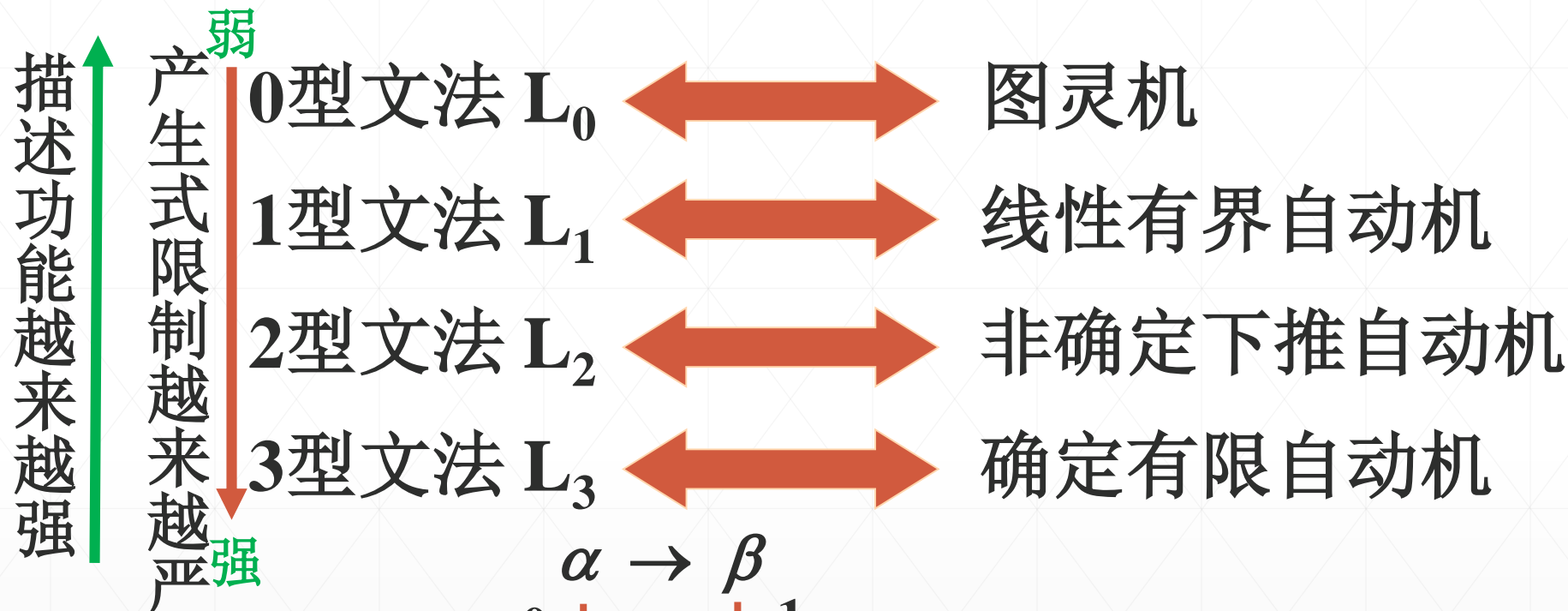
G_3 是2型文法

G_3 产生的语言：

$$L(G) = \{a^n cb^n \mid n \geq 0\} \cup \{a^n b^n \mid n \geq 1\}$$



5. 四类文法和语言小结及关系



0型语言 \supset 1型语言 \supset 2型语言 \supset 3型语言

$(L_0 \supset L_1 \supset L_2 \supset L_3)$