

编译程序设计原理

杜淑敏 王永宁 编著

北京大学出版社

内 容 简 介

本书系统地介绍了编译程序设计的基本原理和方法,较详细地论述了当前各种语法分析技术。书中采用语法制导翻译方法。在文法符号的综合属性和继承属性等概念的基础上,给出了语法制导定义的形式,使得语义分析颇为清晰。在书的编写方式上,着重问题的提出、分析与解决。在内容的选取上,注意深入浅出。本书可作为高等学校计算机系各专业“编译原理”或“编译方法”课程的教材或参考书,也可供其它专业学生及从事计算机工作的有关人员阅读参考。

编译程序设计原理

杜淑敏 王永宁 编著

责任编辑: 邱淑清 彭令华

*

北京大学出版社出版

(北京大学校内)

北京大学印刷厂印刷

新华书店北京发行所发行 各地新华书店经售

*

850×1168 毫米 32 开本 13.5 印张 335 千字

1990 年 11 月第一版 1990 年 11 月第一次印刷

印数: 0001—5,000 册

ISBN 7-301-01210-1/TP·054

定价: 7.30 元

目 录

第一章 语言的基本知识	1
1.1 基本定义	1
1.1.1 文法的讨论	1
1.1.2 字母表	5
1.1.3 符号串	6
1.1.4 文法和语言的形式定义	9
1.2 分析树(parse tree)和二义性	18
1.2.1 分析树	18
1.2.2 子树	21
1.2.3 二义性	23
1.3 形式语言概观	27
练习	29
第二章 编译概述	31
2.1 翻译和解释	31
2.2 编译程序的组成部分	33
2.3 有关编译程序的设计与实现	37
第三章 词法分析	42
3.1 词法分析程序的功能	42
3.1.1 单词符号的种别和属性	43
3.1.2 词法分析程序作为一个独立子程序	47
3.2 输入缓冲	47
3.3 正规表达式与正规集	51
3.4 正规表达式与正规文法	55
3.5 词法与正规文法	58
3.5.1 文法的分解	58

3.5.2	基本符号的文法	58
3.6	状态转换图与基本符号的识别	63
3.6.1	状态转换图	63
3.6.2	状态转换图的实现	68
3.7	有限自动机	71
3.7.1	确定的有限自动机(DFA)	72
3.7.2	非确定有限自动机(NFA)	75
3.7.3	具有 ϵ -转移的非确定有限自动机	78
3.7.4	正规文法与有限自动机(FA)的等价性	82
3.7.5	有限自动机与正规表达式的等价性	87
3.7.6	确定有限自动机的化简	92
	练习	97
第四章	语法分析	100
4.1	语法分析器的作用	100
4.2	上下文无关文法	101
4.3	书写文法	105
4.3.1	关于消除二义性	105
4.3.2	消除左递归	107
4.3.3	提取左因子	112
4.4	自顶向下分析	113
4.4.1	递归下降分析方法	113
4.4.2	预测分析器	116
4.4.3	预测分析器的转换图	121
4.4.4	非递归的预测分析方法	126
4.4.5	函数 FIRST 和 FOLLOW	129
4.4.6	预测分析表的构造	132
4.4.7	LL(1)文法	133
4.4.8	在预测分析法中的错误处理示例	134
4.5	自底向上分析	137
4.5.1	规范归约	137

4.5.2	“移进-归约”分析法的栈实现	142
4.6	算符优先分析法	146
4.6.1	利用算符优先关系寻找右句型的可归约串	149
4.6.2	算符优先关系表的构造	152
4.6.3	优先函数	154
4.6.4	算符优先分析法的错误处理示例	156
4.7	LR 分析器	159
4.7.1	LR 分析器的逻辑结构及工作过程	159
4.7.2	SLR 分析表的构造	163
4.7.3	LR(1)分析表的构造	178
4.7.4	LALR 分析表的构造	185
4.8	LR 分析法对二义文法的应用	197
	练习	202
第五章	语法制导翻译	208
5.1	语法制导定义	209
5.1.1	语法制导定义的形式	210
5.1.2	综合属性	211
5.1.3	继承属性	212
5.1.4	依赖图	214
5.1.5	计算顺序	217
5.2	语法树(syntax tree)的构造	218
5.2.1	语法树	219
5.2.2	建立表达式的语法树	220
5.2.3	建立语法树的语法制导定义	221
5.2.4	关于表达式的有向非循环图	223
5.3	S-属性定义及其自底向上的计算	225
5.4	L-属性定义	229
5.4.1	L-属性定义	230
5.4.2	翻译模式	231
5.5	自顶向下的翻译	236

5.5.1	从翻译模式中消除左递归	236
5.5.2	预测翻译器的设计	242
5.6	自底向上计算继承属性	244
5.6.1	从翻译模式中去掉嵌入的动作	245
5.6.2	分析栈中的继承属性	246
5.6.3	模拟继承属性的计算	248
5.7	递归求值	254
	练习	262
第六章	运行时刻环境	265
6.1	有关源语言中的一些问题的讨论	265
6.1.1	过程	266
6.1.2	活动树	267
6.1.3	控制栈	270
6.1.4	说明的作用域	271
6.1.5	名字的联编	272
6.1.6	提出的问题	273
6.2	存储组织	274
6.2.1	运行时刻内存的划分	274
6.2.2	活动记录	275
6.2.3	编译时刻的局部数据的设计	277
6.3	运行时刻存储分配策略	278
6.3.1	静态存储分配	278
6.3.2	栈式存储分配	282
6.3.3	堆式存储分配	288
6.4	对非局部名字的访问	289
6.4.1	块	290
6.4.2	不含嵌套过程的词法作用域	293
6.4.3	含有嵌套过程的词法作用域	294
6.4.4	动态作用域	301
6.5	参数传递	303

6.5.1	传值调用	304
6.5.2	引用调用	305
6.5.3	复制恢复	307
6.5.4	传名调用	308
6.6	符号表	309
6.6.1	符号表的表项	309
6.6.2	线性表 (linear list)	312
6.6.3	散列表	313
6.6.4	表示作用域的信息	316
练习	319
第七章	中间代码生成	324
7.1	中间语言	324
7.1.1	图表示法	324
7.1.2	三地址代码	327
7.1.3	三地址语句的种类	328
7.1.4	语法制导翻译生成三地址代码	329
7.1.5	三地址代码的具体实现	331
7.2	说明语句	334
7.2.1	过程中的说明语句	334
7.2.2	保留作用域信息	336
7.2.3	记录中的域名	339
7.3	赋值语句	340
7.3.1	符号表中的名字	340
7.3.2	数组元素地址分配	342
7.3.3	访问数组元素的翻译模式	345
7.3.4	访问记录中的域	349
7.4	布尔表达式	349
7.4.1	翻译布尔表达式的方法	350
7.4.2	数值表示法	350
7.4.3	控制流语句	352

7.4.4	控制流语句中的布尔表达式的翻译	355
7.5	CASE 语句	358
7.6	回填	359
7.6.1	使用回填翻译布尔表达式	360
7.6.2	使用回填翻译控制流语句	365
7.6.3	标号和转移语句	371
7.7	过程调用	371
	练习	374
第八章	代码生成	377
8.1	目标机器	378
8.2	运行存储管理	382
8.2.1	静态分配管理	383
8.2.2	栈式分配管理	385
8.2.3	名字的运行地址	389
8.3	基本块和流图	390
8.3.1	基本块	390
8.3.2	流图	392
8.4	下次引用信息	393
8.5	一个简单的代码生成器	395
8.5.1	寄存器描述器和地址描述器	396
8.5.2	代码生成算法	397
8.6	基本块的 dag 表示法	402
8.6.1	dag 的构造	405
8.6.2	dag 的应用	407
8.7	从 dag 生成目标代码	411
	练习	415
	参考文献	417

第一章 语言的基本知识

在本章中我们将简单介绍形式语言的基本知识。这些知识对学习本书其它章节的内容是必不可少的。

1.1 基本定义

语言是人类思考问题和交流思想的工具。它是随着社会的产生而产生、发展而发展的,是人类传达思想、感情和希望的工具。凡是与人类社会生活有关的活动,都离不开语言。可以说,语言是在人们认识世界和改造世界过程中发展起来的。

计算机的出现促进了语言学的研究,特别是形式语言学的研究。自从著名语言学家 Noam Chomsky 于 1956 年建立了形式语言的描述以来,形式语言学理论发展得很快。这种理论对计算机科学有着深刻的影响,对程序语言的设计和编译程序的构造有着重大的作用。这正是本书所要讨论的课题。程序语言是形式化语言,是人和计算机相互传达信息的工具。以下我们将介绍有关形式语言的几个基本概念。这些基本概念来源于自然语言学的研究。为了使读者易于掌握这些内容,在形式地引进定义之前,我们将首先对文法和语言进行一般的讨论。

1.1.1 文法的讨论

在本节中,我们从一个具体英语例句的分析出发,引进有关文法和语言的基本概念。请考虑英语句子“The grey wolf will eat the goat”。根据英语知识,知道这是一个英语句子,并且可以进行图

解,如图 1.1 所示。这是一棵有序有向树(简称树)。这种图解把句子分解成它的各个组成部分。由图 1.1 可以看出,〈句子〉是由〈主语〉后随〈谓语〉组合而成的;〈主语〉又是由〈冠词〉后随〈形容词〉、再随〈名词〉构成的,等等。

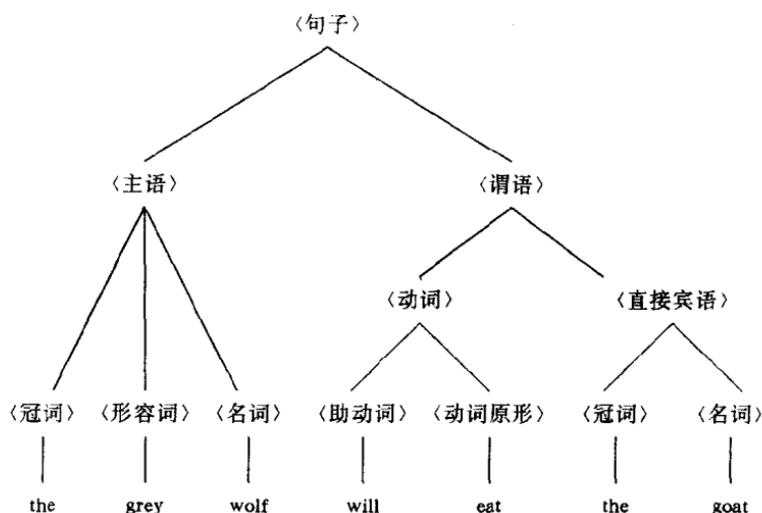


图 1.1 一棵树

为了描述这种结构,在图 1.1 中我们使用了一些新的符号,即所谓语法单位(或称语法实体),如〈句子〉、〈主语〉、〈冠词〉等等。图 1.1 中凡是用尖括号〈和〉括起来的都是语法单位。这样就把所用到的语法单位和语言中的基本字较为明显地区分开了。

如前所述,图 1.1 表明〈句子〉是由〈主语〉后随〈谓语〉组合成的。为了能机械地进行这样的分解,我们必须给出一些形式的精确的规则,用以表明句子的结构。如果我们用符号“→”(或“::=”)表示“定义为”(或“由...组合成的”),用符号“|”表示“或”,那末,

我们可使用这些符号构成规则。在此例中所需要的全部规则如下：

- 〈句子〉→〈主语〉〈谓语〉 (1)
- 〈主语〉→〈冠词〉〈形容词〉〈名词〉 (2)
- 〈冠词〉→ the (3)
- 〈形容词〉→ grey (4)
- 〈谓语〉→〈动词〉〈直接宾语〉 (5)
- 〈动词〉→〈助动词〉〈动词原形〉 (6)
- 〈助动词〉→will (7)
- 〈动词原形〉→eat (8)
- 〈直接宾语〉→〈冠词〉〈名词〉 (9)
- 〈名词〉→wolf (10)
- 〈名词〉→goat (11)

而后面的两个规则可以写在一起：

〈名词〉→wolf | goat

规则(1)读作：〈句子〉是由〈主语〉后随〈谓语〉组合而成的，或读作：〈句子〉定义为〈主语〉后随〈谓语〉。规则(2)则读作：〈主语〉是由〈冠词〉后随〈形容词〉再后随〈名词〉组合而成，等等。

我们把描述语言的语法结构的形式规则称为文法。上述十一个规则，可以说是构成了一个小型文法。让我们仔细分析一下这个文法。此文法包括四个组成部分：

1. 一组终结符号。这里包括 the, grey, wolf, will, eat 和 goat 等六个字。

2. 一组非终结符号。非终结符号用来代表语法单位。这里包括〈句子〉、〈主语〉、〈冠词〉、〈形容词〉、〈谓词〉、〈动词〉、〈助动词〉、〈动词原形〉、〈直接宾语〉和〈名词〉等十个语法单位。

3. 一个开始符号。开始符号是一个特殊的非终结符号。这里，〈句子〉是开始符号。

4. 一组规则(也称产生式或产生规则)。这里共包括十一个规

则。

综上所述,如果我们用 G 表示这一文法,则可写成:

$$G = (\{ \text{the, grey, wolf, will, eat, goat}, \\ \langle \text{句子} \rangle, \langle \text{主语} \rangle, \langle \text{冠词} \rangle, \langle \text{形容词} \rangle, \langle \text{谓语} \rangle, \\ \langle \text{动词} \rangle, \langle \text{助动词} \rangle, \langle \text{动词原形} \rangle, \langle \text{直接宾语} \rangle, \\ \langle \text{名词} \rangle\}, \langle \text{句子} \rangle, \mathcal{P}),$$

其中 \mathcal{P} 表示由上述十一个规则所组成的集合。

当我们有了一组规则之后,就可以按照下述方式用它们去**推导或产生**句子。我们从开始符号 $\langle \text{句子} \rangle$ 开始,在 \rightarrow 的左端找带有 $\langle \text{句子} \rangle$ 的规则,并把 $\langle \text{句子} \rangle$ 改写成 \rightarrow 的右端的符号串:

$$\langle \text{句子} \rangle \Rightarrow \langle \text{主语} \rangle \langle \text{谓语} \rangle$$

重复此过程,在符号串 $\langle \text{主语} \rangle \langle \text{谓语} \rangle$ 中取一语法单位,比如说 $\langle \text{主语} \rangle$,我们再找 \rightarrow 的左端有 $\langle \text{主语} \rangle$ 的规则,并用相应的右端的符号串代替原来符号串中的 $\langle \text{主语} \rangle$,这就产生了:

$$\langle \text{主语} \rangle \langle \text{谓语} \rangle \Rightarrow \langle \text{冠词} \rangle \langle \text{形容词} \rangle \langle \text{名词} \rangle \langle \text{谓语} \rangle$$

符号“ \Rightarrow ”意味着:使用文法中的一条规则,代替 \Rightarrow 左端的某个符号,产生 \Rightarrow 的右端的符号串。句子“The grey wolf will eat the goat”的全部推导过程将是:

$$\langle \text{句子} \rangle \Rightarrow \langle \text{主语} \rangle \langle \text{谓语} \rangle$$

$$\Rightarrow \langle \text{冠词} \rangle \langle \text{形容词} \rangle \langle \text{名词} \rangle \langle \text{谓语} \rangle$$

$$\Rightarrow \text{the} \langle \text{形容词} \rangle \langle \text{名词} \rangle \langle \text{谓语} \rangle$$

$$\Rightarrow \text{the grey} \langle \text{名词} \rangle \langle \text{谓语} \rangle$$

$$\Rightarrow \text{the grey wolf} \langle \text{谓语} \rangle$$

$$\Rightarrow \text{the grey wolf} \langle \text{动词} \rangle \langle \text{直接宾语} \rangle$$

$$\Rightarrow \text{the grey wolf} \langle \text{助动词} \rangle \langle \text{动词原形} \rangle \langle \text{直接宾语} \rangle$$

$$\Rightarrow \text{the grey wolf will} \langle \text{动词原形} \rangle \langle \text{直接宾语} \rangle$$

$$\Rightarrow \text{the grey wolf will eat} \langle \text{直接宾语} \rangle$$

$$\Rightarrow \text{the grey wolf will eat} \langle \text{冠词} \rangle \langle \text{名词} \rangle$$

⇒the grey wolf will eat the 〈名词〉

⇒the grey wolf will eat the goat

我们用一个新符号“ $\overset{\pm}{\Rightarrow}$ ”，把该句子的上述推导缩写为：

〈句子〉 $\overset{\pm}{\Rightarrow}$ the grey wolf will eat the goat

通常在每一步推导中，可以去代替任意一个语法单位。只是在上述推导中，总是去代替最左面的那一个语法单位。

文法的目的之一就是用适当条数的规则把语言的全部句子描述出来。对上述文法 G 而言，除上述句子外它还可以描述以下几个句子：

“the grey wolf will eat the wolf”

“the grey goat will eat the wolf”

“the grey goat will eat the goat”

尽管它们之中在语义上(含义上)有的可能是不合适的，但在语法上是正确的。小型文法 G 描述了由四个句子构成的语言。以后我们会看到，常常是用有限个规则把具有无穷多个句子的语言描述出来，因而这样做是有意义的。

通过以上介绍，我们大致上做好了描述文法和语言的形式化概念的准备作。

1.1.2 字母表

字母表是符号的非空有穷集合。“符号”是一个抽象实体，我们将不去形式地定义它，就如同几何学中的“点”和“线”一样往往不加定义。字母和数字是经常使用的符号。例如，由符号“ a ”组成的字母表记作 $\{a\}$ ；由符号“ a ”和“ b ”组成的字母表记作 $\{a, b\}$ ；由符号“ a ”，“ b ”和“ c ”组成的字母表记作 $\{a, b, c\}$ ；由符号“ 0 ”和“ 1 ”组成的字母表记作 $\{0, 1\}$ ；而 ASCII 字符集是一个常用的计算机字母表的例子。

1.1.3 符号串

由字母表中的符号所组成的任何有穷序列被称之为该字母表上的符号串。符号串也称作“字”。例如,设有字母表 $\Sigma = \{a, b, c\}$, 那么,序列 ab 是 Σ 上的一个符号串;同样序列 ba , 序列 abc , 序列 $bc-ca$ 等都是 Σ 上的符号串。在语言的理论中,术语“句子”和“字”常常用作术语“符号串”的同义语。符号串 s 的长度记作 $|s|$, 是组成该符号串的符号的个数。例如,上述 Σ 上的符号串 ab 的长度是 2, 记作 $|ab| = 2$ 。符号串 abc 的长度是 3, 记作 $|abc| = 3$ 。空符号串记作 ε , 它由零个符号组成, 于是 $|\varepsilon| = 0$ 。在本书中常用的与符号串有关的几个术语有:

1. 符号串 s 的前缀: 移走符号串 s 的尾部的零个或多个符号所得到的一个符号串。例如 ban 是符号串 $banana$ 的一个前缀。

2. 符号串 s 的后缀: 删去符号串 s 的头部的零个或多个符号所得到的一个符号串。例如 $nana$ 是符号串 $banana$ 的一个后缀。

3. 符号串 s 的子串: 从 s 中删去一个前缀和一个后缀而得到的符号串。例如 nan 是 $banana$ 的一个子串。 s 的每一个前缀和每一个后缀是 s 的一个子串, 但是 s 的每一个子串不一定是 s 的前缀或后缀。对于每一个符号串 s , s 和 ε 两者都是 s 的前缀、后缀和子串。

4. 符号串 s 的真前缀、真后缀、真子串: 任何非空符号串 x , 相应地, 是 s 的前缀、后缀或子串, 并且 $s \neq x$ 。

5. 符号串 s 的子序列: 从符号串 s 中删去零个或多个符号(这些符号不要求是连续的)而得到的符号串。例如 $baaa$ 是 $banana$ 的一个子序列。

术语“语言”表示某个确定的字母表上的符号串的任何集合。

这个定义是非常广泛的。在此意义下,不含任何元素的空集合 \emptyset ,即集合 $\{\}$,以及只含空符号串 ε 的集合 $\{\varepsilon\}$ 都是语言。大家更为熟悉的例子有如:所有合乎 Pascal 语法的程序所组成的集合以及所有合乎英文文法的句子所组成的集合都是语言。当然这后两者难于给出详细说明。另外要注意的是,这个定义没有描述在一个语言中的符号串的任何含义。有关描述符号串的含义的方法将在第五章进行讨论。

下面让我们介绍符号串之间的运算。

符号串 α, β 的连接 $\alpha\beta$ 是把符号串 β 写在符号串 α 之后得到的符号串。例如,若 $\alpha=ab, \beta=bc$,则有, $\alpha\beta=abbc, \beta\alpha=bcab$ 。按此约定,因 ε 是不包含任何符号的符号串,所以,对于任意符号串 α 而言,可以得到

$$\varepsilon\alpha = \alpha\varepsilon = \alpha$$

我们还可以定义符号串的方幂。假设 α 是符号串, α^n 定义为:

$$\underbrace{\alpha\alpha\cdots\alpha}_n \uparrow \alpha$$

当 $n=0$ 时, α^0 是空符号串 ε 。例如,若 $\alpha=a$,则有

$$\begin{aligned} \alpha^0 &= \varepsilon \\ \alpha^1 &= a \\ \alpha^2 &= aa \\ &\dots\dots \\ \alpha^n &= \underbrace{aa\cdots aa}_n \uparrow a \end{aligned}$$

又如,若 $\alpha=ab$,则有

$$\begin{aligned} \alpha^0 &= \varepsilon \\ \alpha^1 &= ab \\ \alpha^2 &= abab \end{aligned}$$

$$\dots\dots\dots$$

$$\alpha^n = \underbrace{abab\dots abab}_n$$

我们可以进一步考虑定义在语言之上的几个重要运算。假设 L 和 M 表示两个语言。

1. 语言 L 和 M 的合并(union), 记作 $L \cup M$, 定义为:

$$L \cup M = \{s \mid s \text{ is in } L \text{ or } s \text{ is in } M\}$$

此集合读作“所有属于 L 或属于 M 的符号串 s 所组成的集合”。并注意不要把运算符号 \cup 与字母 U 相混。

2. 语言 L 和 M 的连接(concatenation), 记作 LM , 定义为:

$$LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$$

此集合读作“满足 s 属于 L 和 t 属于 M 的所有符号串 st 所组成的集合”。

因为对于任意符号串 x 总有 $\epsilon x = x\epsilon = x$, 所以有:

$$\{\epsilon\}L = L\{\epsilon\} = L$$

注意, 一般而言, $LM \neq ML$, 但 $(LM)W = L(MW)$ 。 L 自身的 n 次连接记为:

$$L^n = \underbrace{LL\dots L}_n$$

并规定 $L^0 = \{\epsilon\}$ 。

3. 语言 L 的 Kleene 闭包, 记作 L^* , 定义为:

$$L^* = \bigcup_{i=0}^{\infty} L^i = L^0 \cup L^1 \cup L^2 \cup L^3 \cup \dots$$

4. 语言 L 的正闭包, 记作 L^+ , 定义为:

$$L^+ = \bigcup_{i=1}^{\infty} L^i = L^1 \cup L^2 \cup L^3 \cup L^4 \cup \dots$$

以上运算是在语言上定义的, 即它们对于某个字母表上的任何符号串集合都是适用的。

例 1.1 令 L 是集合 $\{A, B, C, \dots, Z, a, b, c, \dots, z\}$, 并且 D 是集

合 $\{0, 1, 2, \dots, 9\}$ 。我们可以用两种方式来考虑 L 和 D 。其一, 我们可以把 L 当作包括 26 个大写字母和 26 个小写字母的字母表, 把 D 当作包括 10 个数字的字母表; 其二, 我们可以把 L 和 D 都看作是由长度为 1 的符号串所组成的集合, 这样, 集合 L 和 D 每一个都是有限的语言。下面就是由语言 L 和 D 应用上述运算后而创建的一些新的语言:

1. $L \cup D$ 是由所有上述字母和数字所构成的集合。
2. LD 是由所有用一个字母后跟一个数字组成的符号串所构成的集合。
3. L^4 是由所有的四个字母的符号串构成的集合。
4. L^* 是由所有的字母组成的符号串, 包括空符号串 ϵ , 所构成的集合。
5. $L(L \cup D)^*$ 是由所有的字母打头的字母和数字符号串所构成的集合。
6. D^+ 是由所有的长度大于等于 1 的数字串所构成的集合。

我们知道, 一种语言是某个确定的字母表上的某些符号串的集合。此集合是字母表上所有符号串, 包括 ϵ , 所构成的集合的一个子集。设字母表为 Σ , 则集合

$$\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots = \bigcup_{i=0}^{\infty} \Sigma^i$$

是 Σ 上的所有符号串, 包括 ϵ , 所构成的集合。为了详细刻画语言所表示的集合, 需要我们对语言(形式语言)给出更为具体而形式化的定义。这是下面将要讨论的内容。

1.1.4 文法和语言的形式定义

本节我们将介绍上下文无关文法和上下文无关语言。所谓上下文无关文法是这样一种文法, 它所定义的语法单位(或称语法范畴, 或称语法实体)是完全独立于这种语法单位可能出现的上下文