

形式语言
及其
句法分析

〔美〕 A. V. 阿霍 J. D. 厄尔曼 著

科学出版社

形式语言及其句法分析

〔美〕 A. V. 阿霍 J. D. 厄尔曼 著

石青云 译

科学出版社

1987

内 容 简 介

本书系统介绍形式语言及其句法分析方法。全书共分七章，其中包括数学预备知识、编译导引、语言理论基础、翻译理论、一般句法分析方法、单路无回溯剖析和回溯量有限制的剖析算法。书中所述内容是计算机编译理论的重要组成部分，也是句法模式识别的理论基础。

本书论述严谨而富于启发性，各章节末配有由易到难的大量习题，便于读者学习和研讨。

本书可供从事计算机理论、软件工程、模式识别、人工智能等方面工作的科研人员参考，也可作为高等院校有关专业的教材。

Alfred V. Aho Jeffrey D. Ullman
THE THEORY OF PARSING, TRANSLATION, AND COMPILING
VOLUME 1: PARSING
Prentice-Hall, Inc., 1972

形式语言及其句法分析

(美) A. V. 阿霍 J. D. 厄尔曼 著

石青云 译

责任编辑 鞠丽娜

科学出版社出版

北京朝阳门内大街 137 号

中国科学院印刷厂印刷

新华书店北京发行所发行 各地新华书店经售

*

1987年9月第 一 版 开本：850×1168 1/32

1987年9月第一次印刷 印张：18 1/2

印数：0001—3,750 字数：484,000

统一书号：15031·840

本社书号：4951·15—8

定 价：5.50 元

译 者 的 话

形式语言理论自二十世纪五十年代问世以来，已逐渐成为计算机科学的一个重要领域。其原因不仅在于它提供了计算机程序设计语言的数学模型及编译程序设计的基本理论和方法，而且在于它奠定了模式识别的语言结构法的基础。

这本书的原书名是 A. V. 阿霍和 J. D. 厄尔曼合著的《剖析、翻译和编译理论(卷 1: 剖析)》。为使本书书名更能体现该书内容，翻译时，我们将该书书名改为《形式语言及其句法分析》。书中除了讲述形式语言与自动机的关系、语言的闭包性质和可判定性问题等理论性较强的内容以外，一个突出的特点是系统介绍了各种句法分析(剖析)算法，并分析了它们的计算复杂度。这一点对于模式识别尤为重要。当一类模式所对应的形式语言由产生它的文法来描述时，相应的句法分析算法就可以作为该类模式的识别程序。为方便读者系统了解这方面的知识，我们选择了这本书。如果从事计算机软件工作的科研人员和高校师生希望了解作者对编译理论的进一步阐述，请参阅 A. V. 阿霍和 J. D. 厄尔曼合著的《剖析、翻译和编译理论(卷 2: 编译)》。

本书结构严谨完整，概念清晰，富于启发性，特别是着重于交待处理各种论题的基本思想。书中选用了很多例子来帮助读者理解基本概念和理论推导，在各章节末尾配有由易到难的大量习题，并给出了进一步研究的问题和有关参考文献的注释。长期以来美国一些大学一直把它作为教科书使用。

本书的译稿曾请中国科学院自动化所戴汝为同志进行了校对，在此一并致谢。

由于译者水平所限，加之时间仓促，书中难免有错误或不妥之

处，敬请读者批评指正。

译 者
于1985年8月

• •

前　　言

这本书是为高年级大学生或研究生的编译理论课程写的，可以作为一学期或两学期的教材。书中着重对编译这个实际课题进行理论分析，这样做的出发点有三个方面：

(1) 在一个象计算机科学这样迅速发展的领域，正确有效的教学法是在课程中突出基本思想，而不是强调细节。我们希望在本书中介绍的算法和概念，对下一代计算机和程序设计语言还适用，至少是其中的一部分对编译程序设计以外的领域也能用。

(2) 编译程序的设计已经进展到能把一个编译程序的许多部分分割开来，分别进行设计优化的阶段。这样，为从事优化的人员提供适当的数学工具就很重要。

(3) 一些最有用和效率最高的编译算法，如 LR(k) 剖析法，需要很多的数学基础才能完全理解。因此，对于编译程序的设计者来说，一个基本要求就是要有较好的理论基础。

在不省略与编译有关的较难定理的同时，我们尽量使这本书通俗易懂。为此，列举了大量的例子。每个例子都采用较简单的文法，而不是实际中遇到的复杂文法。希望即使在纯粹理论推导很难懂的情况下，这些例子也足以解释清楚基本概念。

本书的用法

这本书是在给普林斯顿大学和史蒂文斯理工学院的高年级大学生和研究生授课的讲义基础上写成的。以这本书为教材，既讲授过一学期的课，也讲授过两学期的课。作为一学期的课程，不包括本书的第 0, 第二, 第八章，其余各章的内容大部分要详细讲述，但必须先修包括有限自动机和上下文无关语言的课程。

作为连续两学期的课程，第 I 卷的大部分内容都在第一学期讲授，第 II 卷中除第八章以外，大部分内容在第二学期讲授。两

学期的课程比一学期的课程更注重证明和证明的技巧。

显然，书中有些章节是更为重要的，我们想就第 I 卷中各部分内容的重要程度作一简单评论。一般来说，将大部分证明省略掉大致是可取的。显然把所有主要结果的证明都包括到书中来对于深刻理解论题无疑是必要的，但我们估计编译方面的许多课程很难达到这样的深度，而且只要肤浅地了解一些证明，对于编译也能得到一定程度的理解。

第 0 章(数学基础)和第一章(编译概貌)几乎都是基本的材料，也许 1.3 节是例外，这一节包括了句法分析在编译以外的其它应用。

我们相信在第二章(语言理论)中介绍的每个概念和定理都可以在其余九章的一些地方找到用处。但有些材料在编译课中可以省略，如 2.2.1 节中关于正规表达式方程的较难部分。其次是关于右线性文法的某些材料，尽管有办法得到它们与有限自动机的等价性，但也可以省略掉。此外还可以省略 2.4.5 节中达到格雷巴赫范式的 Rosenkrantz 方法。

第三章(翻译)的概念对于本书的其余部分是非常基本的，但 3.2.3 节关于句法制导翻译的层次这部分内容比较难，可以省略。

我们认为 4.1 节关于句法分析的回溯方法与 4.2 节的列表方法相比，是较为次要的。

第五章(单路剖析)的大部分是基本的。我们建议 LL 文法(5.1 节)、LR 文法(5.2 节)、优先文法(5.3.2 节和 5.3.4 节)和算子优先文法(5.4.3 节)应优先考虑。如有必要，其它几节可以省略。

第六章(回溯算法)与第五章的大部分内容相比或与 4.2 节相比，是较为次要的。如果要给予选择，我们宁愿要 6.1 节而不要 6.2 节。

本书的组织

整个句法分析、翻译和编译的理论以两卷的形式出现，即句法分析(第 0 至第六章)和编译(第七至第十一章)。第二卷的主题是

句法分析程序的优化、确定性句法分析理论、翻译、簿记和代码优化。这两卷形成一个整体，页数是连贯编号的，在第 II 卷中有两卷的文献目录和索引。

习题和文献注释在每一节（编号为 *i.j*）的末尾。除了公开问题和待研究的问题以外，我们用*号表示习题的难度。单*号的题目要求深入思考来求解，双*号的题目需要更加深入的思考。

以这本书为教材的课程，最好配备一个程序设计实验室，在那里设计和实现几个编译程序的部件。本书一些章节的末尾附有程序设计练习题，可以作为这种程序设计实验室的实验项目。

致谢

在准备这本书的过程中，许多人曾仔细阅读手稿的各个部分并给予很大帮助。我们特别感谢 David Benson, John Bruno, Stephen Chen, Matthew Geller, James Gimpel, Michael Harrison, Ned Horvath, Jean Ichbiah, Brian Kernighan, Douglas McIlroy, Robert Martin, Robert Morris, Howard Siegel, Leah Siegel, Harold Stone 和 Thomas Szymanski，以及审稿人 Thomas Cheatham, Michael Fischer 和 William McKeeman。我们也从用过有关讲义的许多学生那里得到重要的意见，如他们中间的 Alan Demers, Nahed El Djabri, Matthew Hecht, Peter Henderson, Peter Maika, Thomas Peterson, Ravi Sethi, Kenneth Sills 和 Steven Squires。

我们还要感谢 Hannah Kresse 和 Dorothy Luciani 对手稿的出色打字，以及在手稿准备过程中贝尔电话实验室所提供的支持帮助。Dennis Ritchie 和 Kenneth Thompson 设计的 PDP-11 计算机的操作系统 UNIX 的运用，加速了手稿中一些部分的准备。

A. V. 阿 霍
J. D. 厄尔曼

目 录

第 0 章 数学预备知识	1
0.1 集合论的一些概念.....	1
0.1.1 集合	1
0.1.2 集合的运算	3
0.1.3 关系	5
0.1.4 关系的闭包	7
0.1.5 次序关系	9
0.1.6 映射.....	11
习题.....	12
0.2 符号串的集合	16
0.2.1 符号串.....	16
0.2.2 语言	17
0.2.3 语言的运算.....	18
习题.....	20
0.3 逻辑的一些概念	21
0.3.1 证明.....	21
0.3.2 归纳证明.....	22
0.3.3 逻辑联结词.....	23
习题.....	24
文献注释.....	27
0.4 过程和算法	27
0.4.1 过程.....	28
0.4.2 算法.....	29
0.4.3 递归函数.....	30
0.4.4 过程的阐明.....	31
0.4.5 问题.....	32

J.4.6 波斯特对应问题.....	35
习题.....	36
文献注释.....	39
0.5 图论的一些概念	40
0.5.1 有向图.....	40
0.5.2 有向无圈图	43
0.5.3 树	43
0.5.4 有序图	45
0.5.5 涉及有向无圈图的归纳证明	47
0.5.6 来自偏序的线性次序.....	47
0.5.7 树的表示	49
0.5.8 图上的路径	51
习题.....	54
文献注释.....	56
第一章 编译导引.....	57
1.1 程序设计语言	57
1.1.1 程序设计语言的阐明	57
1.1.2 句法和语义	59
文献注释.....	61
1.2 编译概貌	62
1.2.1 编译程序的部件	62
1.2.2 词法分析	63
1.2.3 簿记	66
1.2.4 句法分析	68
1.2.5 代码产生	69
1.2.6 代码优化	75
1.2.7 误差的分析和挽回	77
1.2.8 小结	79
习题.....	80
文献注释.....	82
1.3 句法分析和翻译算法的其它应用	83

1.3.1 自然语言	83
1.3.2 模式的结构描述	84
文献注释.....	88
第二章 语言理论基础.....	89
2.1 语言的表示法	89
2.1.1 出发点	90
2.1.2 文法	90
2.1.3 有限制的文法	97
2.1.4 识别程序	100
习题	103
文献注释	109
2.2 正规集及其产生程序和识别程序.....	110
2.2.1 正规集和正规表达式	110
2.2.2 正规集与右线性文法	117
2.2.3 有限自动机	120
2.2.4 有限自动机与正规集	126
2.2.5 小结	130
习题	130
文献注释	134
2.3 正规集的性质.....	134
2.3.1 有限自动机的极小化	134
2.3.2 正规集的抽吸引理	138
2.3.3 正规集类的闭包性质	129
2.3.4 关于正规集的可判定性问题	141
习题	143
文献注释	149
2.4 上下文无关语言.....	149
2.4.1 派生树	150
2.4.2 上下文无关文法的变换	155
2.4.3 乔姆斯基范式	164
2.4.4 格雷巴赫范式	165
2.4.5 达到格雷巴赫范式的另一方法	172

习题	176
文献注释	180
2.5 下推自动机.....	180
2.5.1 基本定义	180
2.5.2 下推自动机的变形	186
2.5.3 PDA 语言和 CFL 的等价性	191
2.5.4 确定性下推自动机	200
习题	207
文献注释	209
2.6 上下文无关语言的性质.....	209
2.6.1 奥登引理	210
2.6.2 上下文无关语言类的闭包性质	214
2.6.3 一些可判定性结果	217
2.6.4 确定性 CFL 的一些性质	220
2.6.5 多义性	221
习题	226
文献注释	230
第三章 翻译理论	232
3.1 翻译的形式方法.....	232
3.1.1 翻译和语义	233
3.1.2 句法制导的翻译模式	235
3.1.3 有限变换器	244
3.1.4 下推变换器	249
习题	255
文献注释	259
3.2 句法制导翻译的性质.....	259
3.2.1 特征化语言	259
3.2.2 简单 SDT 的性质	264
3.2.3 SDT 的层次.....	265
习题	273
文献注释	274

3.3 词法分析	275
3.3.1 正规表达式的一种扩展语言	276
3.3.2 间接词法分析	278
3.3.3 直接词法分析	282
3.3.4 有限变换器的软件模拟	285
习题	286
文献注释	287
3.4 句法分析	287
3.4.1 句法分析的定义	287
3.4.2 自上而下的剖析	289
3.4.3 自下而上的剖析	293
3.4.4 自上而下和自下而上剖析的比较	296
3.4.5 文法的覆盖	301
习题	303
文献注释	306
第四章 一般句法分析方法	307
4.1 回溯剖析法	307
4.1.1 下推变换器的模拟	308
4.1.2 自上而下剖析法的非正式描述	311
4.1.3 自上而下的剖析算法	316
4.1.4 自上而下剖析程序的时、空复杂度	324
4.1.5 自下而上剖析	329
习题	335
文献注释	341
4.2 列表的剖析方法	342
4.2.1 科克-杨格-卡萨米算法	342
4.2.2 厄利剖析方法	349
习题	362
文献注释	364
第五章 单路无回溯剖析法	365
5.1 LL(k)文法	366

5.1.1	$LL(k)$ 文法的定义	366
5.1.2	预测剖析算法	370
5.1.3	$LL(k)$ 定义的含义	375
5.1.4	剖析 $LL(1)$ 文法	379
5.1.5	剖析 $LL(k)$ 文法	382
5.1.6	对 $LL(k)$ 条件的检验	392
	习题	397
	文献注释	404
5.2	确定性自下而上剖析法	405
5.2.1	确定性移动-缩减剖析法	405
5.2.2	$LR(k)$ 文法	408
5.2.3	$LR(k)$ 定义的含义	418
5.2.4	对 $LR(k)$ 条件的检验	430
5.2.5	$LR(k)$ 文法的确定性右剖析程序	432
5.2.6	$LL(k)$ 和 $LR(k)$ 剖析程序的实现	437
	习题	437
	文献注释	440
5.3	优先文法	440
5.3.1	形式的移动-缩减剖析算法	441
5.3.2	简单优先文法	444
5.3.3	扩展优先文法	452
5.3.4	弱优先文法	458
	习题	467
	文献注释	469
5.4	移动-缩减可剖析文法的其它类型	470
5.4.1	有界右文关联文法	470
5.4.2	混合策略优先文法	480
5.4.3	算子优先文法	483
5.4.4	弗洛伊德-伊文思产生式语言	488
5.4.5	本章小结	493
	习题	495
	文献注释	500

第六章 回溯量有限制的剖析算法	501
6.1 回溯量有限制的自上而下剖析法	501
6.1.1 TDPL	502
6.1.2 TDPL 与确定性上下文无关语言	513
6.1.3 TDPL 的推广	517
6.1.4 识别 GTDPL 语言的时间复杂度	523
6.1.5 GTDPL 程序的实现	526
习题	533
文献注释	536
6.2 回溯量有限制的自下而上剖析法	536
6.2.1 非规范剖析法	536
6.2.2 双堆栈的剖析程序	538
6.2.3 科默劳尔优先关系	542
6.2.4 科默劳尔优先性的检验	544
习题	552
文献注释	553
参考文献	554
索引	566

第〇章 数学预备知识

为了使讲解清楚和确切，我们需要精确的且严格定义的语言。这一章就来描述在讨论句法分析、翻译和书中其它论题时要用到的语言，包括集合论初步、图论和逻辑的一些基本概念。已经具有这方面基础知识的读者，可以容易地掠过这一章，而把它当作有关定义和记号的参考资料来对待。

0.1 集合论的一些概念

这一节将扼要地复习集合论中一些最基本的概念：关系、函数、次序和普通集合运算。

0.1.1 集合

以下假设有一些当成原子的确定对象。术语“原子”是一个初级概念，换句话说，我们不再给原子下定义。而且，把什么东西叫做原子，依赖于所论及的范围。在许多情况下，把整数或一个字母表中的字母当成原子是合适的。

我们也给“成员资格”一个抽象的记号。若 a 是 A 的成员(元)，就写成 $a \in A$ 。反之，则写成 $a \notin A$ 。还假设一个原子没有元，就是说，若 a 是一个原子，则对所论及的范围中一切 x 都有 $x \notin a$ 。

我们也考虑一些不是原子的基本对象，称之为集合。若 A 是一个集合，则它的成员或元素是满足 $a \in A$ 的那些对象(不一定是原子)。一个集合的每个元或者是原子，或者是其它集合。我们假定，一个集合的每个元，恰好在该集合中出现一次。若 A 的成员数目是有限的，则称 A 为有限集，通常记为 $A = \{a_1, a_2, \dots, a_n\}$ ，只要 a_1, \dots, a_n 都是 A 的元，且当 $i \neq j$ 时， $a_i \neq a_j$ 。元的次序则是无关

紧要的。例如，我们也可以将上述集合写成 $A = \{a_n, \dots, a_1\}$ 。我们约定用符号 ϕ 表示空集，它是一个没有元的集合。要注意，一个原子也没有元，但空集不是原子，原子也不是空集。

表达式 $\#A = n$ 的意思是：集合 A 有 n 个元。

例 0.1

如果把所有非负整数作为原子，那么 $A = \{1, \{2, 3\}, 4\}$ 是一个集合。 A 的元是 1, $\{2, 3\}$ 和 4。其中元 $\{2, 3\}$ 也是一个集合，它的元是 2 和 3。但原子 2 和 3 本身不是 A 的元。我们可以等价地写成 $A = \{4, 1, \{3, 2\}\}$ 。还注意到 $\#A = 3$ 。 \square

规定集合的一个有效方式是借助于谓词，即一个语句，其中包含一个或多个未知量，每个未知量可取两值之一：真或假。由一个谓词所规定的集合，是使该谓词为真的那些元素组成的。但我们必须仔细选择规定集合的谓词，不然，所规定的集合就可能不存在。

例 0.2

罗素(Russell)悖论就可以说明上述现象。令 $P(X)$ 为谓词：“ X 不是自己的元”，即 $X \notin X$ ，则可设想规定一个集合 Y ，它含有使 $P(X)$ 为真的所有 X ，即 Y 由所有不是自己的元的集合组成。由于最普通的集合似乎都不是自己的元，好象就可以认为集合 Y 存在。

但是，若 Y 存在，则应能回答这样的问题：“ Y 是它自己的元吗？”而这将导致不可能情况。若 $Y \in Y$ ，则 $P(Y)$ 为假，由集合 Y 的定义， Y 不是它自己的元。因此， $Y \in Y$ 是不可能的。反之，设 $Y \notin Y$ ，则由 Y 的定义， $Y \in Y$ 。可见， $Y \notin Y$ 蕴含 $Y \in Y$ 且 $Y \in Y$ 蕴含 $Y \notin Y$ 。而 $Y \in Y$ 与 $Y \notin Y$ 只能有一个为真，两者同时为真是不可能的。因此，只有一条出路，就是承认 Y 不存在。 \square

避免罗素悖论的常规方法是：只用形如“ X 在 A 中且 $P_1(X)$ ”的谓词 $P(X)$ 来规定集合。这里 A 是一已知集合且 P_1 是任一谓词。当对集合 A 已有了解时，只要写 $P_1(X)$ 来代替“ X 在 A 中且 $P_1(X)$ ”就可以了。

若 $P(X)$ 是一个谓词，则用 $\{X | P(X)\}$ 表示使 $P(X)$ 为真的所有