

大学本科计算机专业应用型规划教材
丛书主编：高林

编译原理及实现

孙悦红 编著

清华大学出版社



大学本科计算机专业应用型规划教材

丛书主编：高林

编译原理及实现

孙悦红 编著

清华大学出版社
北京

内 容 简 介

本书以通俗易懂的语言介绍编译原理的理论和常用的方法与技术，并着重介绍各种编译方法的实现途径。全书共分 10 章，包括形式语言基础、词法分析、语法分析、语义分析及代码生成、符号表管理、运行时的存储分配，以及代码优化等。考虑目前学计算机专业的学生对 C 语言比较了解，本书中以 C 语言为雏形设计了一种 TEST 语言，并在介绍全书内容时，都用 TEST 语言进行分析与实现，使编译原理的抽象性通过 TEST 语言编译器的实现而具体化，从而使读者轻松掌握编译原理。

本书理论与实践并重，内容深入浅出，便于自学。每章后都提供了适量的习题。

本书可作为高等学校计算机专业的教材，也可供从事计算机应用和开发的人员使用。本书还配有教学辅助课件及书中所有程序示例，需要者可与作者(sun_yh@tom.com)联系。

版权所有，翻印必究。举报电话：010-62782989 13501256678 13801310933

本书封面贴有清华大学出版社防伪标签，无标签者不得销售。

本书防伪标签采用清华大学核研院专有核径迹膜防伪技术，用户可通过在图案表面涂抹清水，图案消失，水干后图案复现；或将表面膜揭下，放在白纸上用彩笔涂抹，图案在白纸上再现的方法识别真伪。

图书在版编目(CIP)数据

编译原理及实现/孙悦红编著. —北京：清华大学出版社，2005. 4

(大学本科计算机专业应用型规划教材/高林主编)

ISBN 7-302-10307-0

I. 编… II. 孙… III. 编译程序—程序设计—高等学校—教材 IV. TP314

中国版本图书馆 CIP 数据核字(2005)第 001292 号

出版者：清华大学出版社

地 址：北京清华大学学研大厦

<http://www.tup.com.cn>

邮 编：100084

社 总 机：010-62770175

客户服务：010-62776969

组稿编辑：谢 琛

文稿编辑：汪汉友

印 刷 者：北京市世界知识印刷厂

装 订 者：三河市化甲屯小学装订二厂

发 行 者：新华书店总店北京发行所

开 本：185×260 印张：16.75 字数：392 千字

版 次：2005 年 4 月第 1 版 2005 年 4 月第 1 次印刷

书 号：ISBN 7-302-10307-0/TP·7020

印 数：1~5000

定 价：23.00 元

大学本科 计算机专业应用型规划教材

编 委 会

主 编：高 林

副 主 编：王 利 鲍 洁

委 员：(按姓氏笔画为序)

王宝智	古 辉	孙悦红	安淑芝
肖 刚	陈 明	张 玲	张建忠
周海燕	赵乃真	娄不夜	顾巧论
崔武子	鲍有文		

策划编辑：谢 琛 汪汉友

丛书序

大学本科计算机专业应用型规划教材



为适应我国“以信息业带动工业化,发挥后发优势,实现社会生产力的跨越式发展”以及大力发展制造业和优化产业结构的要求,应用型人才培养已成为高等学校人才培养的重要任务。

以微电子技术为基础、计算机技术为主体的信息技术,是当前人类社会中发展最快、渗透性最强、应用面最广的先导技术。信息技术的广泛应用推动着以信息产品制造业、软件业、信息系统集成业和信息咨询服务业为主体的信息产业的发展。在新的世纪里,信息已成为重要的生产要素和战略资源,信息技术成为先进生产力的代表,信息产业将发展成为现代产业的带头产业,人类即将跨越工业时代进入信息时代。因此,信息化成为当今世界经济和社会的发展趋势,大力推进社会和国民经济信息化是推进我国社会主义现代化建设的重要任务。计算机和信息技术的发展不仅需要大批专业技术人才,而且还产生了一批新的职业岗位,毋庸置疑,信息及其相关职业将成为未来最紧缺的职业。

计算机和信息技术与应用的人才需求将呈多元化、多层次趋势,表现在科学、技术、产业、应用、服务诸多方面。不仅需要从事科学、技术研发的人才,而且更需要把研发成果转变为现实产品的技术和管理人才;不仅要有能从事计算机和信息科学、技术工作的人才,而且更需要能从事计算机和信息产业、应用、服务工作的人才;以及在各类人才中的精英人才、领军人物。这实际是对我国计算机和信息类高等教育改革提出了新的要求和新的课题,要求我国高等教育进行结构调整,满足人才培养的多元化,大力培养具有计算机和信息技术专长的应用型人才——他们是这些领域的技术专家和管理专家,可以在相应的行业、企业担任各种技术工作。

目前,我国高等教育中应用型人才培养模式相对落后,如何发展应用型教育已成为课程改革的主要任务。本套教材是以培养计算机和信息类专业本科应用型人才为目的进行的课程与教材改革尝试。在本套教材的策划过程中,清华大学出版社多次组织了由行业企业专家和有丰富教学经验的一线教师参加的研讨会,对应用型高等教育的规律和在计算机教学中的体现进行了深入的研讨。在此基础上我们力求能从整体上把握计算机和信息类应用型人才培养的特征,并体现在这套教材的编写过程中:在教材编写的指导思想上,力求在保持学科科学性的同时,体现工程和技术学科的系统性;在教材



的内容组织上,尽量采用以问题为中心的写作方法,加强案例性教学;在理论联系实际和加强能力培养方面,增加方案性设计习题和实际训练性题目,以培养学生的专门技术能力和完成实际工作任务的能力。

计算机和信息类应用型教材编写还处于改革的初步尝试阶段,希望使用这套教材的教师也能够参与到教材建设工作中来,并提出宝贵意见,以便推动课程改革并提高教材质量。

高 林

2004 年 5 月

前 言

编译原理及实现

编译原理是高等学校计算机专业的必修专业课之一,是一门理论与实践并重的课程。编译原理介绍程序设计语言翻译的原理、技术及实现,对引导学生进行科学思维,提高学生解决实际问题的能力有重要的作用。

在我国高等教育逐步实现大众化后,越来越多的高等学校将会面向国民经济发展的第一线,为行业、企业培养各级各类高级应用型专门人才。而受我国传统历史文化思想的影响,重理论、轻实践的观念在高教界仍较普遍,使我们培养的很多人才不适应社会需求,造成毕业生的结构性就业困难,这也迫使很多高等学校走向应用型教育,培养应用型人才。目前国内大多数的编译原理教材偏重于理论,对实现技术介绍得较少,使学习者感到抽象、难以理解;而且教材篇幅厚重,由于授课时数的限制,以及学生接受能力的差异,教科书的内容往往不能充分利用。根据这种现状,我们编写了本书,目的在于加强对学生应用能力的培养,使学生不仅具备理论知识,更要具备应用能力,使所能为所用,以适应新经济时代对人才的需要,满足就业要求。本书以通俗易懂的语言介绍编译原理,包括词法分析、语法分析、语义分析及代码生成、符号表管理、运行时的存储分配、代码优化等,并着重介绍各种编译方法的实现途径。考虑目前学计算机专业的学生对 C 语言比较了解,书中以 C 语言为雏形设计了一种 TEST 语言,建立该语言的词法、语法、语义文法规则,系统介绍编译过程的各个部分。包括词法分析、语法分析、语义分析及代码生成、符号表的建立及存储分配、错误处理都用具体的实例进行分析与实现。并针对 TEST 语言中的典型语句,深入讲解如何具体用 C 语言编程实现词法分析、语法分析以及语义分析和代码生成,摆脱以往编译教材的抽象性以及理论与实际的脱节,使编译原理的抽象性通过 TEST 语言的编译器实现而具体化,从而使学习者轻松掌握编译原理。

全书共分 10 章,大约需要 70 课时,其中包括 20 课时的上机。第 1 章对编译过程、编译程序的逻辑结构以及编译程序各组成部分的功能进行了概述;第 2 章介绍了文法和语言,它为后面各章的学习奠定了理论基础;第 3 章介绍了词法分析程序的设计原理,包括适合手工设计和自动生成词法分析程序的方法,以及 TEST 语言的词法分析程序的具体编程实现;第 4 章、第 5 章分别介绍了自顶向下和自底向上的语法分析方法,主要介绍了递归下降分析法、LL(1) 分析法以及 LR 分析法,同时介绍了 TEST 语言的递归下降分析

实现;第 6 章介绍了语法制导翻译的概念以及属性翻译文法;第 7 章介绍符号表的组织与管理;第 8 章介绍存储组织与分配技术;第 9 章介绍了语义分析及代码生成的概念和技术,并以 TEST 语言为范例,实现语义分析并同时生成抽象机汇编目标代码;第 10 章主要介绍了局部优化和循环优化常采用的方法。另外,附录中列出了 TEST 语言的文法规则、词法分析程序、语法分析程序、语义及代码生成程序以及 TEST 抽象机模拟器完整程序。每章后都提供了适量的习题,使学习者通过适量的练习掌握书中内容。

本书是作者多年教学实践的汇集、提炼,同时也参考了许多国内外的参考书。本书还配有相应的教学辅助课件,以及词法分析、语法分析和语义分析方法的演示程序(包括递归下降、LL(1)、LR 分析法和可在 DOS 环境下运行的 LEX 与 YACC),有需要者可与作者联系,地址为:sun_yh@tom. com。

鉴于作者水平有限,书中难免有错误和不妥之处,恳请读者批评指正。

编 者

目 录

编译原理及实现

第 1 章 编译概述	1
1. 1 程序设计语言	1
1. 2 翻译程序	2
1. 3 编译程序的组成	3
1. 3. 1 词法分析	4
1. 3. 2 语法分析	4
1. 3. 3 语义分析及中间代码生成	5
1. 3. 4 代码优化	5
1. 3. 5 目标代码生成	6
1. 3. 6 符号表管理	6
1. 3. 7 错误处理	7
1. 4 编译程序的结构	7
1. 4. 1 单遍编译程序	7
1. 4. 2 多遍编译程序	7
1. 4. 3 编译程序分遍的优缺点	8
1. 4. 4 “端”的概念	8
1. 5 编译程序的前后处理器	9
1. 5. 1 预处理器	9
1. 5. 2 汇编程序	9
1. 5. 3 连接加载程序	10
1. 6 TEST 语言与编译器	10
1. 6. 1 TEST 语言	10
1. 6. 2 TEST 编译器	11
1. 6. 3 TEST 机	11
习题	11
第 2 章 文法和语言	12
2. 1 字母表和符号串	12
2. 1. 1 字母表	12
2. 1. 2 符号串	13



2.1.3 符号串及其集合的运算	13
2.2 文法	14
2.2.1 文法形式定义	14
2.2.2 文法的 EBNF 表示	16
2.3 推导	17
2.3.1 直接推导定义	17
2.3.2 推导定义	17
2.3.3 规范推导	18
2.4 句型和句子	18
2.5 语言	19
2.6 递归规则与递归文法	20
2.6.1 递归规则	20
2.6.2 递归文法	20
2.7 短语、简单短语和句柄	21
2.8 语法树	21
2.9 子树与短语	22
2.10 由树构造推导过程	23
2.11 文法的二义性	23
2.12 有关文法的实用限制	25
2.13 文法和语言分类	26
习题	27
 第 3 章 词法分析	28
3.1 词法分析的功能	28
3.2 程序语言的单词符号种类及词法分析输出	29
3.3 正则文法及状态图	30
3.3.1 状态图	30
3.3.2 状态图的用法	31
3.4 词法分析程序的设计与实现	32
3.4.1 TEST 语言的词法规则及状态图	32
3.4.2 TEST 语言词法分析程序的构造	34
3.4.3 TEST 语言的词法分析程序实现	35
3.5 正则表达式	37
3.5.1 正则表达式定义	37
3.5.2 正则文法到正则表达式的转换	38
3.6 有穷自动机	39
3.6.1 确定的有穷自动机	39
3.6.2 不确定的有穷自动机	42

3.6.3 NFA 到 DFA 的转化	43
3.6.4 正则表达式与有穷自动机的等价性	46
3.6.5 确定的有穷自动机的化简	48
3.6.6 根据 DFA 构造词法分析程序	51
3.7 词法分析程序的自动生成器 LEX	52
3.7.1 用 LEX 语言表达正则表达式	52
3.7.2 LEX 源程序结构	54
3.7.3 使用 LEX 生成 TEST 语言的词法分析程序	57
习题	60
第 4 章 语法分析——自顶向下分析	61
4.1 自顶向下分析方法	61
4.2 FIRST 集合和 FOLLOW 集合	62
4.2.1 FIRST 集合定义及构造方法	62
4.2.2 FOLLOW 集合定义及构造方法	63
4.3 递归下降分析	64
4.3.1 递归下降分析的基本方法	64
4.3.2 递归下降分析中存在的问题及解决方法	64
4.3.3 TEST 语言的递归下降分析实现	69
4.4 LL(1) 分析方法	70
4.4.1 LL(1) 分析的基本方法	71
4.4.2 LL(1) 分析表的构造方法	74
4.4.3 LL(1) 分析的主要问题及解决方法	75
习题	76
第 5 章 语法分析——自底向上分析	78
5.1 规范推导、规范句型和规范归约	78
5.2 自底向上分析方法的一般过程	79
5.3 LR 分析方法	80
5.3.1 LR 分析器逻辑结构	80
5.3.2 LR 分析表构成	80
5.3.3 LR 分析过程	82
5.4 LR(0) 分析器	83
5.4.1 活前缀和可归前缀	83
5.4.2 LR(0) 项目	84
5.4.3 构造识别活前缀的有穷自动机	86
5.4.4 LR(0) 分析表的构造	90
5.4.5 LR(0) 分析器的工作过程	92

5.4.6 LR(0)文法	93
5.5 SLR(1)分析器	94
5.5.1 SLR 解决方法的基本思想	96
5.5.2 SLR(1)分析表的构造	97
5.6 LR(1)分析器	100
5.6.1 LR(1)项目	102
5.6.2 LR(1)项目集规范族构造算法	103
5.6.3 LR(1)分析表的构造	107
5.7 LALR(1)分析器	108
5.8 语法分析程序的自动生成工具——YACC	112
5.8.1 YACC 源程序结构	113
5.8.2 YACC 源程序说明部分的组成	113
5.8.3 YACC 源程序的语法规则部分的组成	114
5.8.4 YACC 源程序的程序部分组成	115
5.8.5 二义性文法的处理	117
5.8.6 YACC 示例运行	117
习题	118
第 6 章 语法制导翻译技术	120
6.1 翻译文法	120
6.2 语法制导翻译	122
6.3 自顶向下语法制导翻译	123
6.3.1 递归下降翻译	123
6.3.2 LL(1)翻译器	126
6.4 属性翻译文法	128
6.4.1 综合属性	128
6.4.2 继承属性	130
6.4.3 属性翻译文法定义	131
6.4.4 属性翻译文法举例——算术表达式的翻译	132
6.5 属性文法的自顶向下翻译	133
6.5.1 L-属性翻译文法	134
6.5.2 L-属性翻译文法的翻译实现——递归下降翻译	135
6.5.3 L-属性翻译文法的翻译实现——LL(1)法	139
6.6 自底向上语法制导翻译	142
6.6.1 波兰翻译	142
6.6.2 S-属性文法	144
6.6.3 S-属性波兰翻译文法的翻译实现	145
习题	147

第 7 章 符号表管理技术	149
7.1 何时建立和访问符号表	149
7.2 符号表的组织和内容	150
7.3 符号表上的操作	152
7.4 非块程序结构语言的符号表结构	153
7.5 块程序结构语言的符号表组织	155
7.5.1 块程序结构语言的概念	155
7.5.2 栈式符号表	156
习题	157
第 8 章 程序运行时的存储组织及管理	158
8.1 程序运行时的存储组织	158
8.2 静态存储分配	159
8.3 栈式动态存储分配	160
8.3.1 活动记录	161
8.3.2 运行时的地址计算	163
8.3.3 递归过程的处理	163
8.4 堆式动态存储分配	164
8.4.1 堆分配方式	164
8.4.2 堆式存储管理技术	165
习题	168
第 9 章 语义分析和代码生成	169
9.1 语义分析的概念	169
9.2 中间代码	170
9.2.1 波兰后缀表示	170
9.2.2 N-元表示	171
9.2.3 栈式抽象机及其汇编指令	172
9.3 声明的处理	174
9.3.1 符号常量	174
9.3.2 简单变量	175
9.3.3 数组	177
9.3.4 过程声明	180
9.4 表达式语句	180
9.5 if 语句	185
9.6 while 语句	186
9.7 for 循环语句	188



9.8 write 语句	190
9.9 read 语句	190
9.10 过程调用和返回	191
9.10.1 参数的基本传递形式	191
9.10.2 过程调用	192
9.10.3 过程定义的处理	193
9.10.4 返回语句和过程终止语句	194
9.11 语义分析及代码生成实现	194
9.12 错误处理	194
习题	195
第 10 章 代码优化	196
10.1 局部优化	196
10.1.1 基本块的划分	197
10.1.2 基本块的优化技术	198
10.1.3 基本块的 DAG 表示	199
10.1.4 基本块优化的实现	203
10.2 循环内的优化	205
10.2.1 循环结构的定义	205
10.2.2 循环的查找	206
10.2.3 循环优化的实现	207
习题	213
附录 A TEST 语言文法规则	214
A1 TEST 语言词法规则	214
A2 TEST 的语法规则	214
A3 TEST 的语义和代码生成规则	216
附录 B 词法分析程序	218
附录 C 语法分析程序	221
附录 D 语义及代码生成程序	231
附录 E TEST 抽象机模拟器完整程序	246
参考文献	251

第1章

编译原理及实现

编译概述

程序设计语言之所以能由专用的机器语言发展到现今通用的多种高级语言,就是因为有了编译技术。编译技术是计算机专业人员必须具备的专业基础知识,它涉及程序设计语言、形式语言与自动机、计算机体系结构、数据结构、算法分析与设计、操作系统,以及软件工程等各个方面。现在程序员大多数使用各种高级程序设计语言编写程序,而计算机只能识别用二进制数 0 和 1 表示的指令和数所构成的机器语言程序,用高级语言编写的程序不能直接在机器上运行,要想它运行并得到预期的结果,必须将源程序转换成等价的目标程序,这个转换过程就是所谓的编译。

本章主要介绍程序设计语言编译程序的组成和结构以及编译程序的工作环境,以便读者对编译的基本概念和工作过程有所了解。

1.1 程序设计语言

在计算机发展的初期,人们直接使用机器语言编写程序。机器语言是由二进制数字 0 和 1 表示的机器指令组成,很不直观,而且难写、难读、难记,容易出错,调试极不方便,程序员在程序设计及检查错误方面要花很大的精力;更由于不同类型的计算机使用不同的机器指令,程序员必须针对某种类型的机器编程,编写的程序不适用于移植,因此限制了计算机的推广与使用。

为了便于记忆、阅读和检查,人们用较直观的符号来代替机器指令,进一步发展成为汇编语言。汇编语言采用比较直观且具有含义的指令助记符表示每条机器指令,同时为了方便编程,还提供了若干宏指令对应一组机器指令,从而完成一些特定的功能。但是汇编指令依赖于机器,对问题的描述处于低层次,没有高级语言中的条件、循环等控制结构,编程人员必须考虑寄存器和内存的分配,使用仍不方便,程序设计的效率依然很低。

为了解决这些问题,提高编程效率,人们又发展出了更接近自然语言的高级程序设计语言,如 BASIC、C、Pascal 语言等。这类语言完全摆脱了机器指令的约束,用它编写的程序接近自然语言和习惯上对算法的描述,故称为面向用户的语言或面向过程的语言。后来,又相继出现许多专门用于某个应用领域问题的专业语言。例如用于数据库操作的 SQL(Structured Query Language, 结构化查询语言)语言,这类语言称为面向问题的语言。

汇编语言和机器语言依赖于机器,称为低级语言;而面向用户的语言称为高级语言。高级语言与低级语言相比较具有以下优点:

- (1) 高级程序设计语言不依赖于具体的机器,对计算机了解较少的人也可以学习和使用,有良好的可移植性,在一种类型的机器上编写的程序不做很大改动就能在别的机器上运行;
- (2) 编写高级语言程序时,不用考虑具体的寄存器和内存的分配,不用知道如何实现将数据的外部形式转换成计算机内部形式,也不必了解机器的硬件;
- (3) 每条高级语言语句对应于多条汇编指令或机器指令,编程效率高;
- (4) 高级语言提供了丰富的数据结构和控制结构,提高了问题的表达能力,降低了程序的复杂性;
- (5) 高级语言接近于自然语言,编程更加容易,编写出的程序有良好的可读性,便于交流和维护。

尽管高级语言有这么多优点,但是使用高级语言编写的程序是不能立即在计算机上执行的,它必须经过翻译程序翻译成机器语言程序,计算机才能执行。这种翻译程序就称为编译程序。虽然汇编语言不是高级语言,但是不能在计算机上直接执行,它需要翻译程序将汇编语言编制的程序翻译成机器语言程序,这种翻译程序称为汇编程序。对高级语言来说,其编译程序再加上一些相应的支持用户程序运行的子程序就构成了该语言的编译系统。编译系统是计算机的重要组成部分。本书主要介绍构造编译系统的原理、技术及其实现方法。

1.2 翻译程序

翻译程序扫描所输入的源程序,然后将源程序转换成目标程序。翻译程序的源程序分高级语言源程序和汇编语言源程序两种:

- (1) 如果要翻译的源程序是汇编语言编写的,而目标语言是机器语言,则翻译程序称为汇编程序;
- (2) 如果要翻译的源程序是用高级语言编制的,其翻译后的目标程序是某种具体机器的机器语言或汇编语言,那么这种翻译程序称为编译程序,而实现源程序到目标程序的转换所花费的时间叫做编译时间。

在高级语言程序的编译和运行过程中,源程序和数据是在不同时间处理的。源程序的处理是在编译阶段进行,而数据则是在程序的运行阶段处理。有的编译程序的目标程序是机器语言,则源程序从编译到被执行的过程如图 1-1 所示。

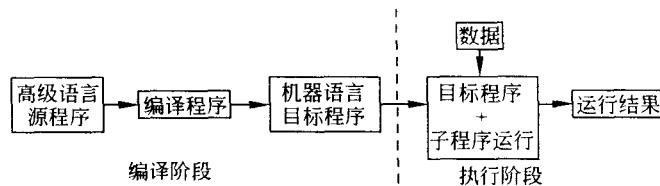


图 1-1 生成机器语言目标程序的编译方式

如果编译程序翻译得到的目标程序是汇编语言程序,则还要经过汇编程序翻译成机器语言程序,这种编译方式的源程序从编译到被执行的过程如图 1-2 所示。

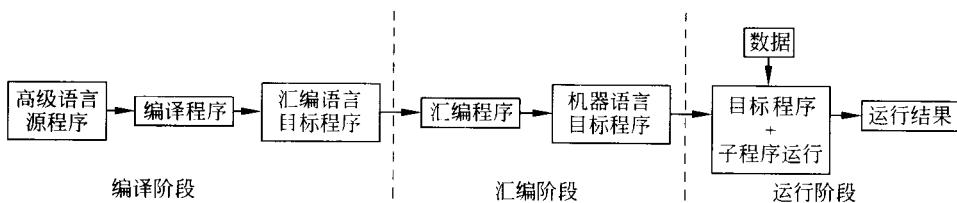


图 1-2 生成汇编语言目标程序的编译方式

还有一种高级语言翻译程序,从源程序的编译到执行只有一个阶段——解释执行阶段,它同时处理源程序和数据,按源程序中语句的动态顺序,逐句地进行分析解释,并立即予以执行,这种翻译程序称为解释程序,运行过程如图 1-3 所示。最常见的高级语言 BASIC 就是在解释环境下运行的。在解释方式下,最终并不生成目标程序,这是编译方式与解释方式的根本区别。解释方式很适合于程序调试,易于查错,在程序执行中可以修改程序,但与编译方式相比,执行效率太低。现在有些语言的集成开发环境提供了两种方式,如 Visual Basic, 调试期间可解释执行源程序,而调好的程序可以编译生成目标程序。

本书虽然主要介绍编译技术及其实现,但同样的技术也适合于解释程序。掌握了编译技术,就不难设计和实现解释程序,而汇编程序的设计与实现和编译程序相比则更为简单。



图 1-3 高级语言的解释方式

1.3 编译程序的组成

为了将高级语言程序翻译成目标程序,编译程序首先必须对高级语言源程序进行分析,然后,产生目标程序。因此,编译程序分成前后两个阶段:分析阶段和综合阶段。分析阶段根据源语言的定义,分析源程序的结构,检查源程序是否符合语言的规定,确定源程序所表示的对象和规定的操作,并以某种中间代码形式表示出来。词法分析、语法分析、语义分析和中间代码生成都属于分析阶段。综合阶段根据分析结果构造所要求的目标代码程序,包括代码优化和目标代码生成。为了记录分析过程中识别出的标识符及有关信息,还需要使用符号表。如果在编译过程中发现源程序有错误,不仅要报告错误,还要进行错误处理,使编译能继续进行。基本的编译程序模型见图 1-4。

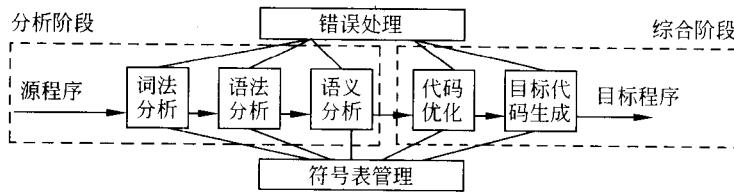


图 1-4 典型的编译程序模型