



普通高等教育规划教材

编译原理及实现技术

刘磊 等编著

BIANYIYUANLIJISHIXIAN JISHU

 机械工业出版社
CHINA MACHINE PRESS



普通高等教育规划教材

编译原理及实现技术

刘 磊 金 英 郭德贵
张 晶 杨 冬 赵东范 编著



机械工业出版社

编译原理是计算机学科的一门重要专业基础课。本书旨在介绍编译程序设计的基本原理、实现技术、方法和工具。本书充分考虑了教师便于教学, 学生便于自学的问题, 在介绍基本原理和实现技术中, 注重循序渐进、深入浅出, 每一章节都提供了编译程序实现的具体实例, 每章末尾都给出了丰富的习题以辅助学生更好地掌握编译过程。

本书包含了编译程序设计的基础理论和具体实现技术, 主要内容有: 形式语言与自动机理论、词法分析、语法分析、语义分析、中间代码生成、中间代码优化和目标代码生成等编译过程。

本书设有支持网站: <http://softlab.jlu.edu.cn>, 欢迎读者访问。

本书可作为大专院校计算机专业的教材, 也可作为计算机工程技术人员参考书。

图书在版编目 (CIP) 数据

编译原理及实现技术/刘磊等编著. —北京: 机械工业出版社, 2005.7
普通高等教育规划教材
ISBN 7-111-17116-0

I. 编... II. 刘... 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 089427 号

机械工业出版社 (北京市百万庄大街 22 号 邮政编码 100037)

责任编辑: 王保家 版式设计: 张世琴 责任校对: 张 媛

封面设计: 张 静 责任印制: 杨 曦

北京机工印刷厂印刷·新华书店北京发行所发行

2005 年 8 月第 1 版第 1 次印刷

787mm × 1092mm $1/16$ · 11.5 印张 · 278 千字

定价: 17.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

本社购书热线电话 (010) 68326294

封面无防伪标均为盗版

前 言

编译原理是计算机学科的一门重要专业基础课，学习编译课程，不仅可以掌握编译程序本身的实现技术，而且能够提高对程序设计语言的理解，提高开发大型软件的能力，提高元级程序的设计能力，提高抽象思维能力。

编译程序是计算机系统软件的重要组成部分，其基本原理和实现技术也适用于一般软件的设计和实现，而且在软件工程、软件自动化、程序分析等领域有着广泛的应用。通常把编译程序视为高级语言到机器语言的转换程序，而这种转换不是结构上的变换，而是基于语言语义的等价变换，因此编译程序设计的难度和复杂性是很高的。同时，编译原理也是一门对实践性要求较高的课程。本书充分考虑了教师便于教学，学生便于自学的问题，循序渐进地介绍了编译程序设计的基本原理、主要实现技术、基本设计方法和一些自动构造工具，深入浅出地介绍了完整的编译程序构造和实现过程，使学生能够掌握编译的整体结构。

本书共分 10 章。第 1 章介绍了与编译程序相关的基础知识。第 2 章作为编译程序的理论基础，简单介绍了形式语言、有限自动机理论和正则表达式等基础知识。第 3 章以正则表达式、有限自动机为工具，讨论了词法分析程序的设计与实现，并简要介绍了词法分析器生成器 LEX 的基本原理和使用方法。第 4 章介绍了自顶向下的语法分析方法的基本思想，并讨论了递归下降语法分析方法和 LL(1) 语法分析方法的实现技术。第 5 章介绍了自底向上语法分析方法的基本思想，并详细讨论了 LR(R) 语法分析的基本原理和实现方法，同时简单介绍了流行的语法分析器生成器 YACC、Bison 等工具。第 6 章专门介绍语义分析，包括标识符、类型、值的内部表示及其构造，符号表的构造及其管理。第 7 章介绍了中间代码生成，包括常用中间代码结构、表达式的中间代码、下标变量的中间代码以及语句的中间代码。第 8 章介绍了中间代码优化的基本方法，重点讨论了常量表达式优化、公共表达式优化和循环不变式优化。第 9 章介绍编译程序运行时的存储空间组织与存储分配技术，重点讨论了运行时的存储结构、存储分配、过程活动记录以及变量访问环境等。第 10 章介绍了目标代码生成的基本技术，重点讨论了中间代码到目标代码的翻译。

本书是根据作者多年的教学经验和科研经历编写而成的，由刘磊教授确定内容的取舍和组织，并最后定稿。杨冬编写第 1、2、10 章；郭德贵编写第 3~5

章以及第 8 章的部分内容；张晶编写第 6、9 章；金英编写第 7 章；赵东范编写第 8 章的部分内容；刘磊参与了各章的编写工作。

本书设有支持网站：<http://softlab.jlu.edu.cn>，欢迎读者访问。

由于编者水平有限，书中难免还存在一些缺点和不妥之处，殷切希望广大读者批评指正。

编 著 者

目 录

前言

第 1 章 编译引论	1
1.1 程序设计语言和编译程序	1
1.2 编译程序的结构	2
1.2.1 编译程序的构成	2
1.2.2 遍	4
1.2.3 编译程序的前端和后端	4
1.3 编译程序和程序设计环境	5
1.4 编译程序的实现	6
习题	7
第 2 章 形式语言与自动机理论	
基础	8
2.1 语言和文法	8
2.1.1 基本概念	8
2.1.2 文法分类	11
2.1.3 推导和归约	13
2.1.4 语法树与文法二义性	14
2.1.5 文法等价变换	16
2.2 有限自动机 (FA)	20
2.2.1 确定有限自动机 (DFA)	21
2.2.2 非确定有限自动机 (NFA)	23
2.2.3 DFA 与 NFA 的等价	24
2.2.4 DFA 的化简	26
2.3 正则表达式	28
习题	32
第 3 章 词法分析	33
3.1 词法分析介绍	33
3.1.1 词法分析程序的功能	33
3.1.2 词法分析程序的接口	33
3.2 词法分析程序的设计	34
3.2.1 单词分类	34
3.2.2 单词的内部表示	34
3.2.3 单词的形式描述	34
3.2.4 自动机的实现	35
3.3 词法分析程序的实现	37
3.3.1 实现词法分析程序应注意的问题	37

3.3.2 标识符表和常量表	39
3.3.3 单词结构	39
3.3.4 实现算法	39
3.4 词法分析程序自动生成	41
3.4.1 LEX 简介	41
3.4.2 LEX 工作原理	41
3.4.3 LEX 源文件结构	42
3.4.4 LEX 系统中的正则式	43
3.4.5 LEX 的使用方式	46
3.4.6 应用实例	46
习题	47
第 4 章 语法分析——自顶向下	
分析方法	48
4.1 语法分析程序介绍	48
4.1.1 语法分析程序的功能	48
4.1.2 语法错误类别及错误处理	48
4.1.3 自顶向下语法分析的基本思想	49
4.1.4 三个重要的集合	51
4.1.5 自顶向下语法分析条件	52
4.2 递归下降法	53
4.2.1 递归下降法语法分析原理	53
4.2.2 递归下降法语法分析程序的构造	54
4.3 LL (1) 分析法	56
4.3.1 LL (1) 分析法原理	56
4.3.2 LL (1) 分析表的构造	57
4.3.3 LL (1) 驱动程序构造	58
4.4 自顶向下分析程序的自动生成	59
习题	60
第 5 章 语法分析——自底向上分析	
方法	61
5.1 自底向上语法分析方法介绍	61
5.2 简单优先分析	62
5.2.1 简单优先文法及其优先关系矩阵的构造	62
5.2.2 简单优先分析算法	63
5.3 LR (k) 分析法	64

5.3.1 LR (k) 分析法的工作过程	65	7.9 控制语句的中间代码生成	119
5.3.2 LR (0) 分析法	66	7.9.1 goto 语句和标号定位的中间代码	120
5.3.3 SLR (1) 分析法	72	7.9.2 条件语句的中间代码	120
5.3.4 LR (1) 分析法	75	7.9.3 while 语句的中间代码	121
5.3.5 LALR (1) 分析法	78	7.10 过程/函数声明的中间代码生成	122
5.3.6 小结	80	习题	123
5.4 自底向上分析程序的自动生成	82	第 8 章 中间代码优化	124
习题	83	8.1 优化方法概述	124
第 6 章 语义分析和符号表	85	8.2 基本块划分	126
6.1 语义分析概述	85	8.3 常量表达式局部优化	127
6.2 标识符的内部表示	86	8.4 公共表达式局部优化	128
6.3 类型的内部表示	90	8.5 循环不变式外提	130
6.4 值的内部表示	93	8.5.1 循环不变式外提概述	130
6.5 符号表的组织和管理	94	8.5.2 循环不变式外提原理	133
6.5.1 符号表的总体组织	94	8.6 其它各类优化介绍	134
6.5.2 符号表项的排列	94	习题	136
6.5.3 符号表的建立与查找	96	第 9 章 运行时存储空间的组织与管理	137
6.6 嵌套式符号表和分程序结构的管理	98	9.1 目标程序运行时的活动	137
6.6.1 简单 C 语言的符号表	98	9.1.1 过程的活动	137
6.6.2 Pascal 语言的符号表	99	9.1.2 名字的作用域和绑定	138
6.6.3 带有分程序结构的符号表	101	9.1.3 过程活动记录	138
6.7 标号的语义分析	103	9.1.4 抽象地址分配	140
习题	104	9.1.5 参数传递	142
第 7 章 中间代码生成	106	9.2 运行时存储器的划分	144
7.1 常用的中间代码结构	106	9.3 静态存储分配	145
7.1.1 后缀式	106	9.4 简单的栈式存储分配	146
7.1.2 抽象语法树和有向无环图	106	9.5 嵌套式语言的栈式存储分配	147
7.1.3 三地址中间代码	107	9.5.1 局部 Display 表方法	149
7.2 语法制导方法概论	108	9.5.2 静态链方法	151
7.3 类型检查和类型转换	110	9.5.3 全局 Display 表方法	152
7.4 中间代码生成中的几个问题	110	9.6 堆式动态存储分配	154
7.4.1 语义信息的获取和保存	110	9.6.1 堆区的分配策略	154
7.4.2 语义栈 Sem 及其操作	111	9.6.2 堆区的回收策略	156
7.4.3 常用的语义子程序	111	9.7 过程调用中几种特殊情况的处理	157
7.5 表达式的中间代码生成	112	9.7.1 非正常出口语句	158
7.6 下标变量的中间代码生成	114	9.7.2 形式过程语句	158
7.6.1 下标变量的地址	114	9.7.3 过程的奇特型调用	158
7.6.2 下标变量的四元式结构	115	习题	159
7.6.3 下标变量的中间代码生成	115	第 10 章 目标代码生成	161
7.6.4 下标变量中间代码生成举例	116	10.1 目标代码生成介绍	161
7.7 赋值语句的中间代码	117		
7.8 过程调用和函数调用的中间代码	118		

10.1.1 代码生成器的输入和输出	161	10.4.4 条件语句四元式的翻译	167
10.1.2 指令选择	161	10.4.5 循环语句四元式的翻译	168
10.2 虚拟机	162	10.4.6 标号语句四元式和 goto 语句四元式 的翻译	169
10.3 寄存器的分配	164	10.4.7 过程、函数说明语句四元式的 翻译	170
10.3.1 单寄存器机器的寄存器分配	164	10.4.8 过程、函数调用语句四元式的 翻译	170
10.3.2 多寄存器机器的寄存器分配	164	习题	171
10.4 四元式到目标代码的翻译	165	参考文献	173
10.4.1 表达式四元式的翻译	165		
10.4.2 赋值语句四元式的翻译	166		
10.4.3 输入输出语句四元式的翻译	166		

第 1 章 编译引论

1.1 程序设计语言和编译程序

众所周知，计算机系统的工作是由事先设计好的程序来控制的。人们首先按自己的需要把让计算机做的工作编写成计算机程序，并把程序送入计算机，然后启动计算机执行程序，从而完成人们设想的要计算机完成的工作。

在计算机发展的初期，人们直接使用机器语言编写程序。机器语言由能被计算机直接执行的机器指令组成，每条机器指令是一串二进制代码，用机器语言编写出来的程序是一串二进制代码序列。这种机器语言程序能被计算机直接识别和执行，但是它很不直观，可读性很差，编写过程费时费力，非常容易出错，而且出错后难于调试和修改。基于上述原因，人们引入了汇编语言。

汇编语言是符号化了的机器语言，也就是引入一些助记符表示机器指令中的操作码和存储地址。汇编语言程序的语句基本上与机器指令对应，却不能被机器硬件直接识别和执行。汇编语言大大提高了编程的速度和准确度，在代码需要极快的运行速度和极高的简洁度时尤为如此，人们至今仍在使用着它。但是，汇编语言也有很多缺点：编写起来仍然很不容易；阅读和理解很难；程序严格依赖于特定的计算机，不便于移植。因此，人们又进一步引进了高级语言。

高级语言不依赖于具体的计算机，它以较接近于自然语言的方式来描述程序的操作。显然，这种程序十分好读，编写起来也很容易，而且还具有通用性，可以在不同机器上运行，便于移植。

尽管用汇编语言和高级语言编写程序比用机器语言方便得多，但是机器硬件只能识别机器语言，因此必须有这样的一个程序，它把用汇编语言或高级语言编写的程序转换成等价的机器语言程序。我们把这种执行转换功能的程序统称为翻译程序 (Translator)，其中汇编语言的翻译程序称为汇编程序 (Assembler)，高级语言的翻译程序称为编译程序 (Compiler)，也称为编译器。

编译程序的输入对象称为源程序 (Source Program)，它是用高级语言编写的程序；输出对象称为目标程序 (Object Program)，目标程序可以是机器语言程序、汇编语言程序或用户自定义的某种中间语言程序。

编译程序支持的源程序的执行分为两个阶段：编译阶段和运行阶段。编译阶段对整个源程序进行分析，翻译成等价的目标程序，翻译的同时做语法检查和语义检查，凡是有错误的源程序均指出其错误；运行阶段在运行子程序的支持下执行目标程序。运行子程序是为了支持目标程序的运行而开发的程序，例如系统提供的标准函数和其它目标程序所调用的程序等。其过程如图 1-1a 所示。

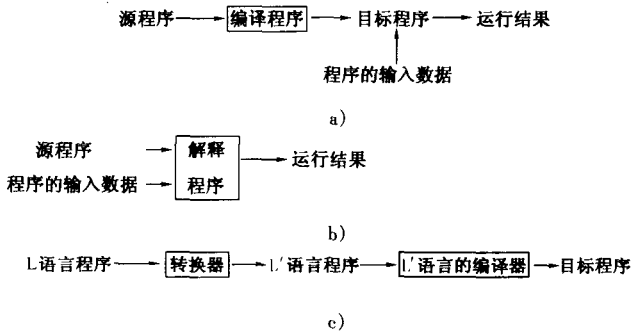


图 1-1 高级程序设计语言的实现方式

a) 编译实现方式 b) 解释实现方式 c) 转换实现方式

高级语言程序的执行也可以采用另一种方式，即并不把源程序先翻译成目标程序然后再执行目标程序，而是一个语句一个语句的读入源程序，边翻译边执行，在翻译过程中不产生目标程序。这种处理程序称为解释程序（Interpreter）。解释程序的工作结果是源程序的执行结果而不是目标程序。其过程如图 1-1b 所示。

从理论上说，任何程序设计语言都可以被解释或被编译，在实际应用中，解释方式特别适合于交互式语言；而且解释方式允许程序执行时改变自身，比如调试程序。这种情形编译程序不易胜任，因为它需要动态编译，而解释程序可以毫无困难地胜任；此外，解释程序不依赖于目标机，因为它不生成目标代码，可移植性优于编译程序。但是和编译程序相比，解释程序开销大，运行速度慢得多，因此，大多数高级语言主要采用编译方式。解释方式和编译方式的最根本区别在于：在解释方式下，并不生成目标代码程序，而是直接执行源程序本身。

实现高级语言的方式除了上文提到的编译方式和解释方式，还有一种转换方式。假如要实现 L 语言，现在已有 L' 语言的编译程序，就可以先把用 L 语言编写的程序转换成等价的 L' 语言的程序，再利用 L' 语言的编译程序实现 L 语言，其过程如图 1-1c 所示。在实际应用中，这种转换方式使用得较少。

1.2 编译程序的结构

1.2.1 编译程序的构成

编译程序是比较庞大的系统软件，它的处理过程复杂，不同的编译程序实现的方法千差万别，构造原理各异，但编译程序的基本逻辑功能及必须完成的处理任务是非常相似的。一般编译程序要完成以下工作：词法分析、语法分析、语义分析、中间代码生成、中间代码优化、目标代码生成、表格管理和出错处理。它们之间的关系如图 1-2 所示。其中除了中间代码部分外，其它所有工作都是任何编译程序都得完成的工作，但完成的方式可能不同，有的可能把几个工作合起来交错完成。不过从原理的角度说，图 1-2 所示的结构是比较理想的逻辑结构，它具有一般性。

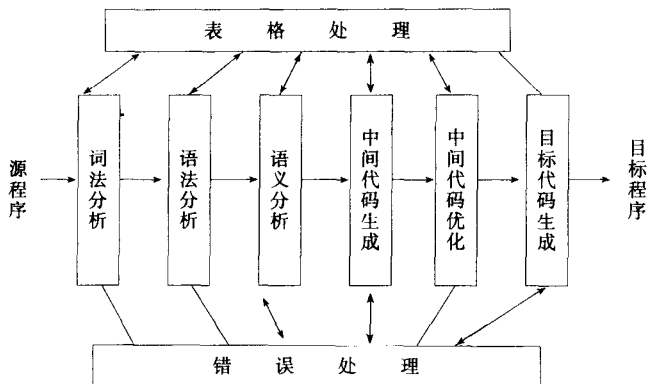


图 1-2 编译器的功能结构图

1. 词法分析阶段

词法分析的任务是扫描源程序的 ASCII 码序列，识别出一个个具有独立意义的最小语法单位，即单词，如关键字（if、then、for、while 等）、标识符、常数、运算符和界符（标点符号、左右括号等等），并把每个单词的 ASCII 码序列替换为所谓的 TOKEN 形式。其中 TOKEN 是单词的内部表示，它的结构没有统一的规定，依赖于具体语言和具体编译器，也依赖于词法分析器的具体功能，但一般包括单词属性标识和单词内码两个部分。在词法分析阶段还要检查括号类配对等词法错误并去掉源程序中注释。词法分析阶段不依赖于语言的语法定义。

2. 语法分析阶段

语法分析的任务是根据程序设计语言的语法规则，把词法分析的结果分解成各种语法单位，同时检查程序中的语法错误。语法分析的扫描对象有两种可能：一种是将词法分析程序作为独立的一遍运行，扫描整个源程序的 ASCII 码序列，将之转换为 TOKEN 序列，输出到一个中间文件，该文件作为语法分析程序的扫描对象继续编译的过程；更一般的情况是将词法分析程序设计成一个子程序，每当语法分析程序需要读取单词时，则调用该子程序。这种设计方案中，词法分析程序和语法分析程序处于同一遍，可以省去中间文件。

3. 语义分析阶段

这一阶段的任务是对语法分析所识别出的各类语法范畴，分析其含义，并进行静态语义检查，如变量是否定义、类型是否匹配等。这一阶段依循的是语言的语义规则，通常使用属性文法描述语义规则。

4. 中间代码生成

在进行了上述的语法分析和语义分析阶段的工作后，有些编译程序将源程序变成一种内部表示形式，这种内部表示形式叫做中间代码。使用中间代码的主要好处是：便于移植、便于修改、便于优化。这种中间代码的形式有很多种，常见的有后缀式（栈式）中间代码、三地址中间代码（三元式和四元式）、图结构中间代码（树，DAG）。其中后缀式中间代码是最早使用的一种中间代码，现在很少使用，目前使用的主要是后两种。

5. 中间代码优化

此阶段的任务是对前阶段产生的中间代码在不改变源程序语义的前提下进行加工变换，使生成的代码更为高效，缩短运行时间或节省存储空间。主要的优化方式有常量表达式优化、公共子表达式优化、不变表达式的循环外提和削减运算强度等。

6. 目标代码生成

这一阶段的任务是把中间代码变换成特定机器上的机器指令代码或汇编指令代码。这是编译的最后阶段，因为目标语言的关系而十分依赖于硬件系统。如何充分利用寄存器、合理选择指令、生成尽可能短而有效的目标代码，都与目标机的结构有关。

生成的目标代码如果是汇编指令代码，则需经由汇编程序处理后才能执行；生成的目标代码如果是绝对指令代码，则可直接投入运行；生成的目标代码如果是可重定位的指令代码，那么目标代码只是一个代码模块，必须由连接装配程序将输入/输出模块、标准函数等系统模块与目标代码模块连接在一起，才能形成一个绝对指令代码程序以供执行。大多数现代实用的编译程序生成的目标代码都是这种可重定位的指令代码。

7. 表格管理

编译程序在对源程序的分析过程中，需要创建和管理一系列的表格，以登记源程序的各类信息和编译各阶段的进展情况。例如，编译程序需要知道变量的类型、数组的大小、函数的参数个数和类型等，这些信息一般可以从源程序中得到。随着编译过程的进行需要不断的建表、查表和填表，或修改表中的某些数据，或从表中查找相关信息，这些活动贯穿整个编译过程的始终。因此，合理的设计和使用表格是编译程序构造中的一个非常重要的问题。不少编译程序都设立一些专门子程序（称为表格管理程序），它们专门负责管理表格。

8. 错误处理

一个编译程序不仅能对书写正确的程序进行翻译，而且能对出现在源程序中的错误进行处理。错误包括词法错误、语法错误、静态语义错误、动态语义错误，其中动态语义错误只能在运行目标程序时才能发现。在编译程序的各个阶段都要有错误处理部分，词法错误和语法错误都集中一次完成检查，而语义检查则分散在以后的各个阶段在完成别的工作时顺便完成。

如果源程序有错误，编译程序应设法发现错误，并把相关错误信息报告给用户。这部分工作是由专门的错误处理程序来完成的。一个好的编译程序应能最大限度地发现源程序中的各种错误，准确地指出错误的性质和发生错误的位置，并且能将错误所造成的影响限制在尽可能小的范围内，使得源程序的其余部分能继续被编译下去，以便进一步发现其它可能的错误。

1.2.2 遍

前面介绍的编译过程的各个阶段仅仅是逻辑功能上的一种划分。具体实现上，受不同源语言、设计要求和计算机硬件条件的限制，往往将编译程序组织成若干遍（Pass）。所谓“遍”就是对源程序或源程序的中间表示形式从头到尾扫描一次，并做加工处理，生成新的中间结果或目标程序。既可以将编译过程中的几个不同阶段合为一遍，也可以把一个阶段的工作分为若干遍。例如，词法分析这一阶段可以作为单独的一遍，但更多的时候是把词法分析程序作为语法分析程序的子程序来加以调用，将词法分析阶段和语法分析阶段合并为一遍。

1.2.3 编译程序的前端和后端

我们有时从概念上把编译程序划分为前端和后端。编译程序的前端主要由与源语言有关但与目标机无关的那些部分组成，通常包括词法分析、语法分析、语义分析、中间代码生

成，与目标机无关的中间代码优化部分也可包含在前端，当然前端也包括相应部分的错误处理。

编译程序的后端包括与目标机有关的中间代码优化部分和目标代码生成等。一般来说，这些部分与源语言无关而仅仅依赖于中间语言。很明显后端是面向目标语言的，而前端则不是，它几乎独立于目标语言。

一种典型的做法是取已有的编译程序的前端，改写其后端以生成不同目标机上相同源语言的编译程序。如果编译程序的后端是经过精心设计的，这样的改写无需太大的工作量，可以很方便地将程序设计语言移植到不同体系结构的计算机上。

1.3 编译程序和程序设计环境

除了编译程序之外，程序设计开发人员通常还需要一些其它的工具，如编辑器、连接程序、装入程序、调试工具和项目管理程序等。编译程序同这些程序一起，构成了一个完整的程序设计环境。

1. 编辑器 (Editor)

编译器通常接受由任何生成标准文件（例如 ASCII 文件）的编辑器所编写的源程序，而且编辑器和编译器及其它程序通常捆绑在一起，构成交互开发环境。有些编辑器还根据源语言的特点，具有语法制导的结构化编辑功能。除了一般文本编辑器的功能外，还可对正在编辑的文本进行分析、提示，能自动地提供关键字和与其匹配的关键字，如 if 和 then、begin 和 end 的配对。这样可以减少程序中的语法错误，提高效率和质量。

2. 预处理器 (Preprocessor)

预处理器是翻译工作开始之前由编译器调用的独立程序，它所做的工作包括删除源程序中的注释、执行宏替换以及包含文件的嵌入等。

3. 连接程序 (Linker)

连接程序负责将分别在不同的目标文件中编译或汇编的代码集中到一个可执行文件中，并将目标和标注库函数的代码以及计算机的操作系统提供的资源连接在一起。连接程序严重依赖于操作系统和处理器。

4. 装入程序 (Loader)

装入程序用来把程序加载到内存储器中，以便执行。由于用户的程序经汇编或编译后生成的目标代码通常采用相对地址的形式，它的起始地址是不确定的。这样的代码被称为可重定位的。装入程序可处理所有的与指定的基地址或起始地址有关的可重定位的地址，它使得可执行代码更加灵活。

5. 调试程序 (Debugger)

调试程序是在被编译了的程序中判定执行错误的程序，它也经常与编译器一起放在集成开发环境 (IDE) 中。运行一个带有调试程序的程序与直接执行程序不同。这是因为调试程序保存着大量的源代码信息（例如行数和变量名等）。它还可以在预先指定的位置（称为断点）暂停执行，并提供有关已调用的函数和变量的当前值等信息。

在高级语言发展的早期，这些工具都是独立的，缺乏整体性。随着程序设计语言的发展，这些工具往往被集成在一起，构成了所谓的集成开发环境 (IDE)，它集编辑、编译、调

试、连接、运行等功能于一体。在这种集成开发环境中，编译程序起核心作用。

1.4 编译程序的实现

编译程序是一个相当复杂的系统程序，通常有上万甚至几万条指令。随着编译技术的发展，编译程序的开发周期也在逐渐缩短，但仍然需要很多人年，而且工作很艰巨，正确性也不易保证。

要实现一个编译程序，通常要做到：

- 对源语言的语法和语义要有准确无误的理解，否则难以保证编译程序的正确性。
- 对目标语言和编译技术也要有很好的了解，否则会生成质量不高的目标代码。
- 确定对编译程序的要求，如搞不搞优化，搞优化搞到哪一级等等。
- 根据编译程序的规模，确定编译程序的扫描次数、每次扫描的具体任务和所要采用的技术。

·设计各遍扫描程序的算法并加以实现。

一般开发编译程序有如下几种可能途径：

1. 转换法（预处理法）

假如我们要实现 L 语言的编译器，现在有 L' 语言的编译器，那么可以把 L 语言程序转换成 L' 语言的程序，再利用 L' 语言的编译器实现 L 语言，这种方法通常用于语言的扩充。如对于 C++ 语言，可以把 C++ 程序转换成 C 程序，再应用 C 语言的编译器进行编译，而不用重新设计和实现 C++ 编译器。常见的宏定义和宏扩展都属于这种情形。

2. 移植法

假设在 A 机器上已有 L 语言的编译程序，想在 B 机器上开发一个 L 语言的编译程序。这里有两种实现方法：

(1) 实现方法一 最直接的解决办法就是将 A 机的代码直接转换成 B 机代码。

(2) 实现方法二 假设 A 机和 B 机上都有高级程序设计语言 W 的编译程序，并且 A 机上的 L 语言编译程序是用 W 语言写的。我们可以修改 L 编译程序的后端，即把从中间代码生成 A 机目标代码部分改为生成 B 机的目标代码。这种在 A 机上产生 B 机目标代码的编译程序称为交叉编译程序（Cross compiler）。

3. 自展法

自展法的实现思想是：先用目标机的汇编语言或机器语言书写源语言的一个子集的编译程序，然后再用这个子集作为书写语言，实现源语言的编译程序。通常这个过程会分成若干步，像滚雪球一样直到生成预计源语言的编译程序为止。我们把这样的实现方式称为自展技术。使用自展技术开发编译器要求这种高级语言必须是能够编译自身的。

4. 工具法

20 世纪 70 年代，随着诸多种类的高级程序设计语言的出现和软件开发自动化技术的提高，编译程序的构造工具陆续诞生，如 Bell 实验室推出的 LEX、YACC 至今还在广泛使用。其中 LEX 是词法分析器的自动生成工具，YACC 是语法分析器的自动生成工具。然而，这些工具大都是用于编译器的前端，即与目标机有关的代码生成和代码优化部分由于对语义和目标机形式化描述方面还存在困难，虽有不少生成工具被研制，但还没有广泛应用。

5. 自动生成法

如果能根据对编译程序的描述，由计算机自动生成编译程序，是最理想的方法，但需要对语言的语法、语义有较好的形式化描述工具，才能自动生成高质量的编译程序。目前，语法分析的自动生成工具比较成熟，如前面提到的 YACC 等，但是整个编译程序的自动生成技术还不是很成熟。虽然有基于属性文法的编译程序自动生成器和基于指称语义的编译程序自动生成器，但产生目标程序的效率很低，离实用尚有一段距离，因此，要想真正的实现自动化，必须建立形式化描述理论。

习 题

1. 典型的编译程序在逻辑功能上由哪几个部分组成？
2. 实现编译程序的主要方法有哪些？
3. 将用户使用高级语言编写的程序翻译为可直接执行的机器语言程序有哪几种主要的方式？
4. 编译方式和解释方式的根本区别是什么？

第 2 章 形式语言与自动机理论基础

2.1 语言和文法

2.1.1 基本概念

人们日常交流所使用的语言，如汉语、英语等，被称为自然语言。这种自然语言过于复杂，难以描述；而人与计算机打交道的程序设计语言是形式化的语言，具有语法严格、结构正规、便于计算机处理等特点。程序设计语言和自然语言也存在共性，即语言的核心均由语法和语义两部分构成。

语法是语言的形式，语义是语言的内容。可以说，语言是由具有独立意义的单词根据一定语法规则组成的句子的集合，用来表达一定的意义。

自从 1956 年语言学家 Noam Chomsky 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。下面我们针对语言和文法进行简单的讨论。

定义 2.1 字母表

字母表 (Alphabet) 是元素的非空有穷集合，字母表中的一个元素称为该字母表的一个字母 (Letter)，也可叫做符号 (Symbol) 或者字符 (Character)。值得注意的是，字母表具有非空性和有穷性。

字母表有时也称为符号表，本书用 Σ 表示字母表。

定义 2.2 符号串

由字母表中的符号组成的任何有穷序列称为符号串。符号串还可以称为字符串、字或句子。

例如，0、1、00、11、01、10、100、101、110、... 均是字母表 $\Sigma = \{0, 1\}$ 上的符号串。

在符号串中，符号的顺序是很重要的。在上例中，01 和 10 就不是同一个符号串，101 和 110 也不是同一个符号串。可以使用字母来表示符号串，如 $x = 010$ ，表示字符串 x 是由 0、1、0 三个符号，并按此顺序组成的符号串。

如果某符号串 x 中有 m 个符号，则称其长度为 m ，表示为 $|x| = m$ 。例如符号串 01 的长度为 2，符号串 1011 的长度为 4。

符号串中可以不包含任何符号，此时称其为空符号串或空串，用字母 ϵ 表示，即 $|\epsilon| = 0$ 。 ϵ 是一个特殊的字符串，对任一字母表 Σ ，都有 ϵ 是 Σ 上的符号串。空字符串有时也用 λ 来表示。值得注意的是，集合 $\{\epsilon\}$ 并不等于空集合 \emptyset ，前者是包含一个元素 ϵ 的非空集合，而后者是不包含任何元素的空集合。

只有当两个符号串长度相等且相应位置的符号均相同时,这两个符号串才是相等的。

下面介绍一些关于符号串的运算。

定义 2.3 符号串的连接

设 x 和 y 均是字母表 Σ 上的符号串,它们的连接是把 y 的所有符号顺序接在 x 的符号之后所得到的符号串。例如设 $x = abc$, $y = de$, 则它们的连接 $xy = abcde$ 。容易看出 $|xy| = |x| + |y| = 5$ 。由于 ϵ 是不包含任何符号的字符串,显然有 $\epsilon x = x\epsilon = x$ 。但连接运算一般不满足交换律。

定义 2.4 符号串的方幂

设 x 是字母表 Σ 上的符号串,把 x 自身连接 n 次得到的符号串 z ,即 $z = xx \cdots xx$ (n 个 x),称作符号串 x 的 n 次幂,记作 $z = x^n$ 。根据定义有

$$x^0 = \epsilon$$

$$x^1 = x$$

$$x^2 = xx$$

$$x^3 = x^2x = xx^2 = xxx$$

.....

$$x^n = x^{n-1}x = xx^{n-1} = xx \cdots xx \text{ (} n \text{ 个 } x \text{)}$$

例 2.1 设 $x = 001$, 则有 $x^0 = \epsilon$, $x^2 = 001001$, $x^3 = 001001001$ 。

定义 2.5 前缀和后缀

设 x 是某一字母表上的符号串, $x = yz$, 则 y 是 x 的前缀, z 是 x 的后缀,特别是当 $z \neq \epsilon$ 时, y 是 x 的真前缀; $y \neq \epsilon$ 时, z 是 x 的真后缀。

例 2.2 设 $x = abc$, 则 ϵ , a , ab , abc 都是 x 的前缀,其中 ϵ , a , ab 为 x 的真前缀;而 abc , bc , c , ϵ 都是 x 的后缀,其中 bc , c , ϵ 为 x 的真后缀。

定义 2.6 子字符串

一个非空字符串 x , 删去它的前缀和后缀后所得到的字符串称为 x 的子字符串,简称子串。如果删去的前缀和后缀不同时为 ϵ , 则称该子串为真子串。

例如, 设 $x = abc$, 其子串为 abc , ab , bc , a , b , c , ϵ , 真子串为 ab , bc , a , b , c , ϵ 。值得注意的是 ac 并不是 x 的子串。

下面定义几个符号串的集合运算。

定义 2.7 符号串集合

若集合 A 中的所有元素都是某字母表上的符号串, 则称 A 为该字母表上的符号串集合。

定义 2.8 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, 则定义

$$AB = \{xy \mid (x \in A) \wedge (y \in B)\}$$

运算结果仍然表示符号串的集合。

例 2.3 设 $A = \{a, bc\}$, $B = \{de, f\}$, 则

$$AB = \{ade, af, bcde, bcf\}$$

注意: $\{\epsilon\} A = A \{\epsilon\} = A$, $\emptyset A = a\emptyset = \emptyset$, 其中 \emptyset 为空集。符号串集合的乘积一般不满足交换律。

定义 2.9 符号串集合的方幂