

# Compiler Design Theory

# Compiler Design Theory

PHILIP M. LEWIS II  
DANIEL J. ROSENKRANTZ  
RICHARD E. STEARNS  
General Electric Company



**ADDISON-WESLEY PUBLISHING COMPANY**

Reading, Massachusetts • Menlo Park, California

London • Amsterdam • Don Mills, Ontario • Sydney

Copyright © 1976 by Addison-Wesley Publishing Company, Inc. Philippines copyright 1976 by Addison-Wesley Publishing Company, Inc.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior written permission of the publisher. Printed in the United States of America. Published simultaneously in Canada. Library of Congress Catalog Card No. 75-9012.

ISBN 0-201-14455-7  
ABCDEFGHIJ-HA-798765

### 编译程序设计理论

该书由IBM公司发起,组织了14名专业人员编纂了12本一套的丛书,取名“系统程序设计丛书”。这里介绍的“编译程序设计理论”是该丛书之一,它取材于大学使用过多年的教材,内容包括设计编译程序或其它语言加工程序所必须的某些基本数学理论以及如何实际运用这些实际理论。作者以“自动机”和形式语言理论作为数学概念的基础,并着重介绍了自动翻译程序设计的综合过程。作者认为:只要掌握了上述两种概念即可设计出教学用或应用的编译程序。并以此概念设计了两个商用编译程序为范例。所以此书内容是理论和实践并重,有较好的参考价值,可供那些对于程序设计语言以及数学方法较为熟悉的计算机软件工作者研读。全书共15章:①引论,②有限状态计算机,③有限状态计算机的实现,④ MINI-BASIC 词法框,⑤下推计算机,⑥与上下文无关的语法,⑦句法制导处理,⑧由顶向下的处理,⑨表征文法的由顶向下的处理,⑩ MINI-BASIC句法框,⑪由底向上的处理,⑫替换-识别处理,⑬替换-归约处理,⑭ MINI-BASIC编译程序的代码产生程序,⑮目标代码优化概论。附录A: MINI-BASIC 语言手册;附录B: 关系式;附录C: 语法的翻译。

## THE SYSTEMS PROGRAMMING SERIES



- \*The Program Development Process  
Part I—The Individual Programmer Joel D. Aron
- The Program Development Process  
Part II—The Programming Team Joel D. Aron
- \*The Design and Structure of  
Programming Languages John E. Nicholls
- Mathematical Background of  
Programming Frank Beckman
- Structured Programming. Harlan D. Mills  
Richard C. Linger



- \*An Introduction to Database Systems C. J. Date
- Compiler Engineering Patricia Goldberg
- Interactive Computer Graphics Andries Van Dam
- \*Sorting and Sort Systems Harold Lorin
- \*Compiler Design Theory Philip M. Lewis II  
Daniel J. Rosenkrantz  
Richard E. Stearns



- Recursive Programming Techniques\* William Burge
- Conceptual Structures: Information  
Processing in Mind and Machines John F. Sowa

\*Published

## IBM EDITORIAL BOARD

Joel D. Aron, Chairman  
Edgar F. Codd  
Robert H. Glaser\*  
Charles L. Gold  
James Griesmer\*  
Paul S. Herwitz

James P. Morrissey  
Ascher Opler\*  
George Radin  
David Sayre  
Norman A. Stanton (Addison-Wesley)  
Heinz Zemanek

\*Past Chairman

# Foreword

The field of systems programming primarily grew out of the efforts of many programmers and managers whose creative energy went into producing practical, utilitarian systems programs needed by the rapidly growing computer industry. Programming was practiced as an art where each programmer invented his own solutions to problems with little guidance beyond that provided by his immediate associates. In 1968, the late Ascher Opler, then at IBM, recognized that it was necessary to bring programming knowledge together in a form that would be accessible to all systems programmers. Surveying the state of the art, he decided that enough useful material existed to justify a significant codification effort. On his recommendation, IBM decided to sponsor The Systems Programming Series as a long term project to collect, organize, and publish those principles and techniques that would have lasting value throughout the industry.

The Series consists of an open-ended collection of text-reference books. The contents of each book represent the individual author's view of the subject area and do not necessarily reflect the views of the IBM Corporation. Each is organized for course use but is detailed enough for reference. Further, the Series is organized in three levels: broad introductory material in the foundation volumes, more specialized material in the software volumes, and very specialized theory in the computer science volumes. As such, the Series meets the needs of the novice, the experienced programmer, and the computer scientist.

Taken together, the Series is a record of the state of the art in systems programming that can form the technological base for the systems programming discipline.

*The Editorial Board*

# Preface

This book is intended to be a text for a one- or two-semester course in compiler design at the senior or first-year-graduate level. It covers the basic mathematical theory underlying the design of compilers and other language processors and shows how to use that theory in practical design situations.

The applicable mathematical concepts come from automata and formal language theory. We have developed these concepts in a rigorous but non-formal style to make them understandable to a wide range of readers, including those who are not mathematically oriented. We believe that automata and formal language concepts constitute an excellent basis both for teaching compiler design and for designing real compilers. We ourselves have designed two commercial compilers based on this theory.

In our selection and presentation of material we emphasize “translation,” in contrast to just “parsing.” The formal concept of a syntax-directed attributed translation is used to specify the input-output performance of various language processors.

Another concept we emphasize is that of an “automaton.” We use such automata as finite-state machines and pushdown machines as basic building blocks for compilers. We emphasize synthesis procedures for designing an automaton to perform a specified translation.

The material in this book constitutes an essentially complete design theory for the lexical and syntax portions of a compiler. The use of attributed translations allows us to include in the syntax-box design a good deal of what is often characterized as “code generation” or “semantics.” The book also includes additional material on code generation and a brief survey of code optimization.

The subject of run-time implementation is not discussed. Although this topic is of considerable importance in deciding what code a compiler should generate, we believe it is not a part of "compiler design theory" but rather a separate subject that should be treated in a course dealing specifically with programming language structures.

Because the automata theory in this text has been selected for its relevance to compiler design, certain basic automata-theory concepts have been bypassed; consequently, this book cannot be used as the exclusive text in an automata theory course. However, students who have taken a course using this book should find a subsequent course in automata-theory greatly simplified. Conversely, some training in basic automata theory permits students to cover the material in this book more rapidly; thus the text can be used either prior to or after taking an introductory course in automata theory.

*Compiler Design Theory* is completely self-contained and assumes only the familiarity with programming languages and the mathematical sophistication commonly found in juniors or seniors.

After an introductory chapter, Chapters 2, 3, and 4 cover finite-state machines and other topics relevant to lexical processing. Chapters 5 and 6 introduce pushdown machines and context-free grammars.

If students have had an introductory course in automata theory, much of the material in Chapters 2 through 6 can be omitted and the remainder covered very rapidly emphasizing the applications to compilers. In any case, Sections 2.7 through 2.11 can be omitted.

Chapter 7 introduces the ideas of translations and attributed translations; this material should be mastered before progressing.

Chapters 8, 9, and 10 cover top-down processing while Chapters 11, 12, and 13 deal with bottom-up processing. These portions of the book are independent and the instructor can elect to cover either or both of them. In any case, Sections 8.7, 10.5, 12.6, and 13.6 can be omitted.

To demonstrate the theoretical concepts in a "real" design situation, the book includes the design of a compiler for a subset of BASIC. The language was selected to be sufficiently complex to illustrate the concepts presented in the book and to have a trivial run-time implementation (since, as we have said, we believe the implementation of language features to be a separate subject). However the language has a variety of syntactic and semantic features including a syntactically recursive control structure (a FOR loop). In Chapter 4, we design a lexical box for the compiler and in Chapters 10 and 12 we design syntax boxes that operate top-down and bottom-up, respectively. In Chapter 14 we design a code generator. Either this design or some extension of it can be implemented in a laboratory portion of the course.



Chapter 15 contains a brief survey of code optimization.

The book also contains three appendices: Appendix A is a language manual for MINI-BASIC; Appendix B discusses those aspects of mathematical relations needed for various test and design procedures; Appendix C presents several methods for transforming a given programming-language grammar into one of the special forms presented in the text.

The material in this book has been taught for several years in one-semester first-year graduate courses at Rensselaer Polytechnic Institute in Troy, N.Y. and at the State University of New York at Albany. It has also been taught as a one-semester undergraduate elective at Union College in Schenectady, N.Y. We would like to thank the students at these institutions whose occasional bewildered looks have motivated us to do several rewrites of the material.

We would like to express our appreciation to the following people who read early versions of the manuscript and made helpful comments: John Hutchison, Michael Hammer, Stephen Morse, John Johnston, Donna Phillips, Daniel Berry, Alyce Orne, Gary Fisher, Walter Stone, James Roberts, and Robert Blean.

We are also grateful to the management of the General Electric Research and Development Center and especially to Richard L. Shuey and James L. Lawson who established the free and stimulating environment in which our work was done and provided the time and opportunity to write this book.

*Schenectady, N. Y.*  
*December 1975*

P.M.L.  
D.J.R.  
R.E.S.

# Contents

## CHAPTER 1 INTRODUCTION

1.1	Language Processors . . . . .	1
1.2	A Naive Compiler Model . . . . .	2
1.3	Passes and Boxes . . . . .	6
1.4	The Run-Time Implementation . . . . .	7
1.5	Mathematical Translation Models . . . . .	7
1.6	The MINI-BASIC Compiler . . . . .	8

## CHAPTER 2 FINITE STATE MACHINES

2.1	Introduction . . . . .	11
2.2	Finite-State Recognizers . . . . .	12
2.3	The Transition Table . . . . .	14
2.4	Exits and Endmarkers . . . . .	16
2.5	Design Example . . . . .	19
2.6	The Null Sequence . . . . .	22
2.7	State Equivalence . . . . .	24
2.8	Testing Two States for Equivalence . . . . .	27
2.9	Extraneous States . . . . .	31
2.10	Reduced Machines . . . . .	33
2.11	Obtaining the Minimal Machine . . . . .	34
2.12	Nondeterministic Machines . . . . .	38
2.13	Equivalence of Nondeterministic and Deterministic Finite-State Recognizers . . . . .	42
2.14	Example: MINI-BASIC Constants . . . . .	45
2.15	References . . . . .	51

### **CHAPTER 3**

#### **IMPLEMENTING FINITE STATE MACHINES**

3.1	Introduction .....	61
3.2	Representing the Input Set .....	62
3.3	Representing the State .....	64
3.4	Selecting the Transitions .....	64
3.5	Word Identification—Machine Approach .....	67
3.6	Word Identification—Index Approach .....	72
3.7	Word Identification—Linear-List Approach .....	74
3.8	Word Identification—Ordered-List Approach .....	74
3.9	Word Identification—Hash-Coding Approach .....	77
3.10	Prefix Detection .....	81
3.11	References .....	84

### **CHAPTER 4**

#### **MINI-BASIC LEXICAL BOX**

4.1	The Token Set .....	87
4.2	The Identification Problems .....	90
4.3	The Transliterator .....	94
4.4	The Lexical Box .....	97

### **CHAPTER 5**

#### **PUSHDOWN MACHINES**

5.1	Definition of a Pushdown Machine .....	109
5.2	Some Notation for Sets of Sequences .....	116
5.3	An Example of Pushdown Recognition .....	119
5.4	Extended Stack Operations .....	121
5.5	Translations with Pushdown Machines .....	125
5.6	Cycling .....	128
5.7	References .....	129

### **CHAPTER 6**

#### **CONTEXT-FREE GRAMMARS**

6.1	Introduction .....	135
6.2	Formal Languages and Formal Grammars .....	135
6.3	Formal Grammars—An Example .....	136
6.4	Context-Free Grammars .....	139
6.5	Derivations .....	141
6.6	Trees .....	143
6.7	A Grammar for MINI-BASIC Constants .....	148
6.8	A Grammar for S-Expressions in LISP .....	150
6.9	A Grammar for Arithmetic Expressions .....	151

6.10	Different Grammars for the Same Language .....	151
6.11	Regular Sets as Context-Free Languages .....	153
6.12	Right-Linear Grammars .....	155
6.13	Another Grammar for MINI-BASIC Constants .....	161
6.14	Extraneous Nonterminals .....	163
6.15	A MINI-BASIC Grammar for the MINI-BASIC Language Manual. . .	168
6.16	References .....	172

## CHAPTER 7 SYNTAX-DIRECTED PROCESSING

7.1	Introduction .....	181
7.2	Polish Notation .....	182
7.3	Translation Grammars .....	183
7.4	Syntax-Directed Translations .....	187
7.5	Example—Synthesized Attributes .....	190
7.6	Example—Inherited Attributes .....	195
7.7	Attributed Translation Grammars .....	197
7.8	A Translation of Arithmetic Expressions .....	201
7.9	Translation of Some MINI-BASIC Statements .....	205
7.10	Another Attributed Translation Grammar for Expressions .....	207
7.11	Ambiguous Grammars and Multiple Translations .....	213
7.12	References .....	216

## CHAPTER 8 TOP-DOWN PROCESSING

8.1	Introduction .....	227
8.2	An Example .....	228
8.3	S-Grammars .....	235
8.4	Top-Down Processing of Translation Grammars .....	239
8.5	$q$ -Grammars .....	245
8.6	LL(1) Grammars .....	252
8.7	Finding Selection Sets .....	262
8.8	Error Processing in Top-Down Parsing .....	276
8.9	Method of Recursive Descent .....	284
8.10	References .....	288

## CHAPTER 9 TOP-DOWN PROCESSING OF ATTRIBUTED GRAMMARS

9.1	Introduction .....	303
9.2	L-Attributed Grammars .....	303
9.3	Simple Assignment Form .....	305
9.4	An Example of an Augmented Machine .....	311
9.5	The Augmented Pushdown Machine .....	320
9.6	Example of a Conditional Statement .....	327

9.7	Example of Arithmetic Expressions .....	332
9.8	Recursive Descent for Attributed Grammars .....	337
9.9	References .....	343

## **CHAPTER 10**

### **MINI-BASIC SYNTAX BOX**

10.1	An LL(1) Grammar for MINI-BASIC .....	355
10.2	The Atom Set and Translation Grammar .....	357
10.3	An L-Attributed Grammar .....	364
10.4	The Syntax Box .....	367
10.5	A Compact MINI-BASIC Expression Processor .....	383

## **CHAPTER 11**

### **BOTTOM-UP PROCESSING**

11.1	Introduction .....	401
11.2	Handles .....	402
11.3	An Example .....	405
11.4	A Second Example .....	411
11.5	Grammatical Principles of Bottom-Up Processing .....	420
11.6	Polish Translations .....	423
11.7	S-Attributed Grammars .....	424

## **CHAPTER 12**

### **SHIFT-IDENTIFY PROCESSING**

12.1	Introduction .....	435
12.2	The SHIFT-IDENTIFY Control .....	436
12.3	Suffix-Free SI Grammars .....	443
12.4	Weak-Precedence Grammars .....	447
12.5	Simple Mixed-Strategy-Precedence Grammars .....	452
12.6	Computing BELOW and REDUCED-BY .....	457
12.7	Error Processing in SHIFT-IDENTIFY Parsing .....	462
12.8	MINI-BASIC Syntax Box .....	469
12.9	References .....	485

## **CHAPTER 13**

### **SHIFT-REDUCE PROCESSING**

13.1	Introduction .....	493
13.2	An Example .....	493
13.3	Another Example .....	504
13.4	LR(0) Grammars .....	513
13.5	SLR(1) Grammars .....	515
13.6	Epsilon Productions .....	520
13.7	Error Processing in SHIFT-REDUCE Parsing .....	526
13.8	References .....	531

## CHAPTER 14 A CODE GENERATOR FOR THE MINI-BASIC COMPILER

14.1	Introduction .....	537
14.2	The Compiler Environment and the Target Machine .....	537
14.3	The Simulation of Run Time .....	538
14.4	Memory Layout .....	539
14.5	Table Entries .....	540
14.6	The GEN Routine .....	542
14.7	The Register Manager .....	544
14.8	Routines for the Atoms .....	545
14.9	Processing Declarations in Block-Structured Languages .....	553
14.10	References .....	555

## CHAPTER 15 A SURVEY OF OBJECT CODE OPTIMIZATION

15.1	Introduction .....	559
15.2	Register Allocation .....	559
15.3	One-Atom Optimizations .....	560
15.4	Optimizations Over a Window of Atoms .....	560
15.5	Optimizations Within a Statement .....	561
15.6	Optimizations Over Several Statements .....	563
15.7	Optimization Over Loops .....	564
15.8	Miscellaneous .....	567
15.9	References .....	568

## APPENDIX A MINI-BASIC LANGUAGE MANUAL

A.1	General Form of a MINI-BASIC Program .....	569
A.2	Numbers .....	569
A.3	Variables .....	570
A.4	Arithmetic Expressions .....	570
A.5	Statements .....	571

## APPENDIX B RELATIONS

B.1	Introduction .....	577
B.2	Representing Relations on Finite Sets .....	578
B.3	The Product of Relations .....	580
B.4	Transitive Closure .....	582
B.5	Reflexive Transitive Closure .....	586

**APPENDIX C**  
**GRAMMATICAL TRANSFORMATIONS**

C.1	Introduction .....	591
C.2	Top-Down Processing of Lists .....	591
C.3	Left Factoring .....	594
C.4	Corner Substitution .....	595
C.5	Singleton Substitution .....	598
C.6	Left Recursion .....	599
C.7	Goal-Corner Transformation .....	602
C.8	Eliminating $\epsilon$ -Productions .....	608
C.9	Making Translations Polish .....	611
C.10	Making a Grammar SHIFT-IDENTIFY Consistent .....	612
C.11	References .....	613

<b>BIBLIOGRAPHY .....</b>	<b>619</b>
---------------------------	------------

<b>INDEX .....</b>	<b>637</b>
--------------------	------------

# 1

## Introduction

### 1.1 LANGUAGE PROCESSORS

There is a natural communication gap between man and machine. Computer hardware operates at a very atomic level in terms of bits and registers, whereas people tend to express themselves in terms of natural languages such as English or in mathematical notation. This communication gap is usually bridged by means of an artificial language which allows the human to express himself with a well-defined set of words, sentences, and formulas that can be "understood" by a computer. To achieve this communication, the human is supplied with a user manual which explains the constructs and meanings allowed by the language, and the computer is supplied with software by which it can take a stream of bits representing the commands or programs written in the language by the human and translate this input into the internal bit patterns required to carry out the human's intent.

Existing computer languages vary widely in complexity, including, for example:

- the instruction set of a particular computer, the machine language, which is interpreted by the hardware or micro-programs of the machine itself;

- assembly languages, the "low-level" languages that largely mirror the instruction set of a particular computer;

- control-card and command languages that are used to communicate with an operating system;

- "high-level" languages, such as FORTRAN, PL/I, LISP, etc., which have a complicated structure and do not depend on the instruction set or operating system of any particular machine.



We use the term “language processor” to describe the computer programs that enable the computer to “understand” the commands and inputs supplied by the human. Broadly speaking, there are two types of such language-processing programs: interpreters and translators.

An *interpreter* is a program that accepts as input a program written in a computer language called the *source language* and performs the computations implied by the program.

A *translator* is a program that accepts as input a program written in a *source language* and produces as output another version of that program written in another language called the *object language*. Usually the object language is the machine language of some computer, in which case the program can be immediately executed on that computer. Translators are rather arbitrarily divided into *assemblers* and *compilers* which translate low-level and high-level languages, respectively.

The common mathematical foundation of all language processing is the theory of automata and formal languages. Since the main concern of this book is the design of compilers, we present those portions of this theory that are most relevant to compiler design and show practical methods whereby this mathematics may be applied. Although the theory is presented in the context of compilers, it can be used in the design of any language processor.

## 1.2 A NAIVE COMPILER MODEL

The job of a compiler is to translate the bit patterns that represent a program written in some computer language into a sequence of machine instructions that carry out the programmer's intent. This task is sufficiently complex that understanding or designing a compiler as a single entity is both difficult and cumbersome. Therefore, it is desirable to consider the compilation process as an interconnection of smaller processes whose tasks can be more easily described.

The selection of these subprocesses for any particular compiler may depend on the details of the language being processed and, in any case, can best be done when taking available design theory into account. Hence, we do not want to endorse a specific set of subprocesses. On the other hand, it is impossible to describe a compiler design theory without having some ideas as to a possible internal organization of a compiler. Therefore, in order to establish a frame of reference, we introduce a rather naive but specific model.

In this model, the compiling job is done by a serial connection of three boxes which we call the *lexical box*, the *syntax box*, and the *code generator*. These three boxes have access to a common set of tables where long-term or