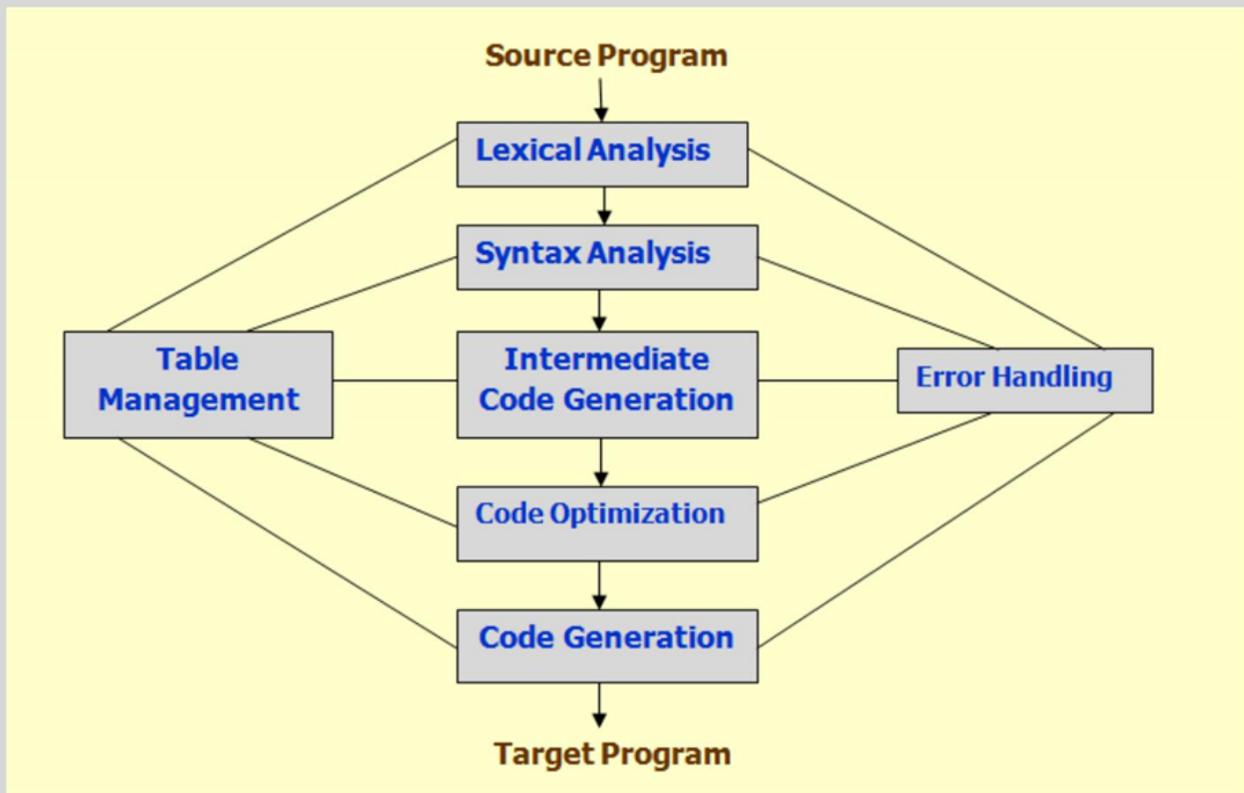


THE DUMMIES' GUIDE TO COMPILER DESIGN



Rosina S Khan

Dedicated to You,

The Valued Reader

Copyright © 2018 by Rosina S Khan. All rights reserved. No part(s) of this eBook may be used or reproduced in any form whatsoever, without written permission by the author.

<http://rosinaskhan.weebly.com>

Contents

Preface	8
CHAPTER 1	11
Introduction to Compilers	11
1.1 What actually is a Compiler?	11
1.2 Other Translators	11
1.3 What is the Significance of Translators?	12
1.4 Macros	13
1.5 High-Level Languages	14
1.6 The Structure of a Compiler	15
CHAPTER 2	18
Lexical Analysis	18
2.1 How a Lexical Analyzer Works	18
2.2 Input Buffering	19
2.3 Preliminary Scanning	20
2.4 A Simple Design of Lexical Analyzers	21
2.5 Regular Expressions	26
2.6 Finite Automata	29
2.7 Nondeterministic Automata	29
2.8 Converting an NFA to a DFA	31
2.9 Minimizing the Number of States of a DFA	41
2.10 A Language for Specifying Lexical Analyzers	43
2.11 Implementation of a Lexical Analyzer	48
CHAPTER 3	52
Syntax Analysis - Part 1	52
3.1 Context-Free Grammars	52
3.2 Derivations and Parse Trees	55
3.3 Regular expressions vs Context-Free Grammar	58

3.4 Further example of Context-Free Grammar	59
CHAPTER 4	62
Syntax Analysis - Part 2	62
4.1 Shift-Reduce Parsing	62
4.2 Operator-Precedence Parsing	64
4.3 Top-Down Parsing	67
4.4 Recursive-Descent Parsing	70
4.5 Predictive Parsers	72
CHAPTER 5	79
Syntax Analysis - Part 3	79
5.1 LR Parsers	80
5.2 CLOSURE	81
5.3 Parsing Table	87
5.4 Moves of LR parser on an Input String	90
CHAPTER 6	92
Syntax-Directed Translation	92
6.1 Syntax-Directed Translation Scheme	92
Semantic Actions	92
6.2 Implementation of Syntax-Directed Translators	95
6.3 Intermediate Code	100
6.4 Postfix Notation	101
6.5 Parse Trees and Syntax Trees	105
6.6 Three-Address Code	106
CHAPTER 7	114
Run-Time Storage Organization	114
7.1 Memory Layout of an Executable Program	114
7.2 Run-Time Stack	115
7.3 Storage Allocation	118

7.4 Accessing a Variable	121
7.5 The Code for the Call of $Q(x,c)$	122
7.6 The Code for a Function Body	123
CHAPTER 8	124
Intermediate Representation (IR) Based on Frames	124
8.1 Example of IR Tree	125
8.2 Expression IRs	125
8.3 Statement IRs	126
8.4 Local Variables	127
8.5 L-values	128
8.6 Data Layout : Vectors	128
CHAPTER 9	130
Type Checking	130
9.1 Type Checking	130
9.2 Type Systems	131
9.3 Type Expressions.....	131
9.4 Type Expressions Grammar	132
9.5 A Simple Typed Language	132
9.6 Type Checking Expressions	133
9.7 Type Checking Statements	133
CHAPTER 10	135
Code Optimization	135
10.1 Introduction to Code Optimization	135
10.2 Loop Optimization	135
CHAPTER 11	143
Code Generation	143
11.1 Problems in Code Generation	143
11.2 A Machine Model	145

11.3 The Function GETREG.....149

Appendix: A Miscellaneous Exercise on Compiler Design

154

About the Author..... 160

Further Free Resources 161

Preface

While students in Computer Science and Engineering (CSE) field or any other equivalent field program in high-level languages and run their programs in editors using a compiler, they do need to understand the mysteries and complexities about how a compiler functions. And the best way to do that is to grasp the underlying principles and actually design a compiler in lab. This is where this book comes into the picture. In a simple, lucid way, the content of this book is made available to the students of CSE or any other equivalent program so that they can understand and grab all the concepts behind Compiler Design conveniently and thoroughly.

Now the principles and theory behind designing a compiler presented in this book are nothing new and they are presented as they have always been known but the real difference lies in the fact that they have been outlined in a really simple and easy-to-understand way. Now I have collected some of the resources from varying sources and assembled with mine to make the flow of reading logical, comprehensible and easy to grasp.

Some of these resources are:

[1] Aho A. V., Lam M. S., Sethi R., Ullman J.D., *Compilers: Principles, Techniques & Tools*, Pearson Education, 2nd Ed., 2007

[2] Aho A. V., Ullman J. D., *Principles of Compiler Design*, Addison-Wesley/Narosa, Twenty-third Reprint, 2004

[3] Dr Fegarsas's Lecture Notes on *Compilers*, CSE, UTA

Organization

Let me now explain the organization of this book.

Chapter 1 is an introductory chapter explaining compilers, translators, their significances and structure of a compiler.

Chapter 2 illustrates lexical analyzers which take input from source programs and produce group of characters called tokens, how they work and function and finally their implementation introducing such concepts as regular expressions, nondeterministic finite automata and deterministic finite automata.

Chapter 3 covers syntactic analysis which groups tokens into syntactic structures such as expressions and statements. For this, we use concepts such as context free grammar, derivations and parse trees.

Chapter 4 continues with syntactic analysis further covering shift-reduce parsing, operator-precedence parsing, top-down parsing, recursive-descent parsing and predictive parsers.

Chapter 5 further continues with syntactic analysis, portraying a special kind of bottom-up parser, the LR parser which scans input from left to right and how they help with syntactic analysis.

Chapter 6 outlines syntax directed translation that introduces intermediate code generation, which is actually an extension of context-free grammars.

Chapter 7 mainly covers storage of variables within program code in a run-time stack.

Chapter 8 explains intermediate representation (IR) specification in areas of frame layout.

Chapter 9 portrays the role of a type checker in the design of a compiler.

Chapter 10 includes code optimization in order to improve the code space and time-wise before the final code generation.

Chapter 11 introduces code generation in machine language format, the final phase of a compiler.

There is also an Appendix at the very end outlining a miscellaneous exercise on compiler design on which students can work out throughout the whole semester in parallel with theory lectures.

Acknowledgments

I deeply acknowledge my heart-felt thanks to the authors and publisher of Compilers: Principles, Techniques & Tools and Principles of Compiler Design as well as Dr Fegaras, for using and collecting their resources and combining them with mine and hence the birth of this very book.

Last but not the least I am thankful to my family for all their support while writing out this book.

CHAPTER 1

Introduction to Compilers

1.1 What actually is a Compiler?

A translator converts one program in a specific programming language as input to a program in another programming language as output. If the source language is a high-level language such as C++ or Java and the object or target language is assembly language or machine language, then such a translator is known as a compiler.

The function of the compiler takes place in two steps:

- 1) First, the source program is compiled or translated into the object program.
- 2) Second, the resulting object program is stored in memory and executed.

1.2 Other Translators

Certain translators transform a programming language into a less complex language called intermediate code, which can be executed directly with a program called interpreter. Here we can interpret the intermediate code acting as some sort of machine language.

Interpreters are smaller in size than compilers and help in the implementation of complex programming language structures. However, the main downside of interpreters is that they take more execution time than compilers.

Besides compilers, there are other translators as well. For instance, if the source language is assembly language and the target language is machine language, the translator is known as an assembler. As another instance, if a translator converts a high-level language to another high-level language, it is termed as a preprocessor.

1.3 What is the Significance of Translators?

We know machine languages are only sequences of 0's and 1's. If we program an algorithm in machine language, it not only becomes tedious but also becomes prone to making errors. All operations and operands must be numeric and therefore it becomes difficult to distinguish them, which is a serious downside. Another problem that arises is that it becomes inconvenient to modify them. So under these circumstances, machine languages are not reliable to start coding and this is exactly where the picture of high-level languages comes in.

A family of high-level languages has been invented so that the programmer can code in a way nearer to his thought processes, ideas and concepts rather than always think at the machine language level and code, which is almost always impossible. A step away from the machine language is the assembly language which uses mnemonic codes for both operations and data addresses. Thus, a programmer could write ADD X, Y in assembly language rather than use sequences of 0's and 1's for the same operation using machine language. However, the computer only understands machine language and so the assembly language needs to be translated to machine language and the translator which carries out this function is known as an assembler.

1.4 Macros

Macros are statements nearer to assembly language statements but different from them in that they use a single memory address along with the operation code. For example, our previous assembly code, ADD X, Y could be broken down into three macro operations, LOAD, ADD and STORE – all using single memory addresses as shown below:

```
MACRO      ADD2  X, Y

           LOAD  Y

           ADD   X

           STORE Y

ENDMACRO
```

The first statement gives the name ADD2 to the macro along with its dummy arguments X, Y. The next three statements define the macro, assuming the machine has only one register, the Accumulator, the other name for Register A.

LOAD Y is equivalent to $Y \rightarrow \text{Acc}$

which means content of memory address Y is transferred to the Accumulator.

The next statement ADD X is equivalent to $\text{Acc} + X \rightarrow \text{Acc}$

which means content of the Accumulator is added to the content of memory address X and the result is stored in the Accumulator.

The third statement STORE Y is equivalent to Acc -> Y

which means the content of the Accumulator is stored in memory address Y.

In this way the assembly code ADD X, Y is broken into macro statements and the same Add operation happens in the latter case. This is useful when the number of registers in the machine is limited.

In the above code, we defined a macro. Now I shall explain how we can use a macro.

After definition of ADD2 X, Y, a macro use happens when we come across the statement ADD2 A, B. In ADD2 A, B statement, A and B replace X and Y respectively and is translated to:

```
LOAD B
```

```
ADD A
```

```
STORE B
```

Thus macro use is like a function call to the function definition as in a high-level language such as C, C# or Java.

1.5 High-Level Languages

A high-level programming language makes the programming task simpler, but it also introduces a problem. We now need a program to convert to a language the machine understands - in other words, the machine language. In that case, this program becomes a compiler similar to the assembler for an assembly language.

A compiler is more complex to write than an assembler. Sometimes compilers have appended with them an assembler so that the compiler produces assembly code initially, which is then assembled, loaded and converted into machine language.

1.6 The Structure of a Compiler

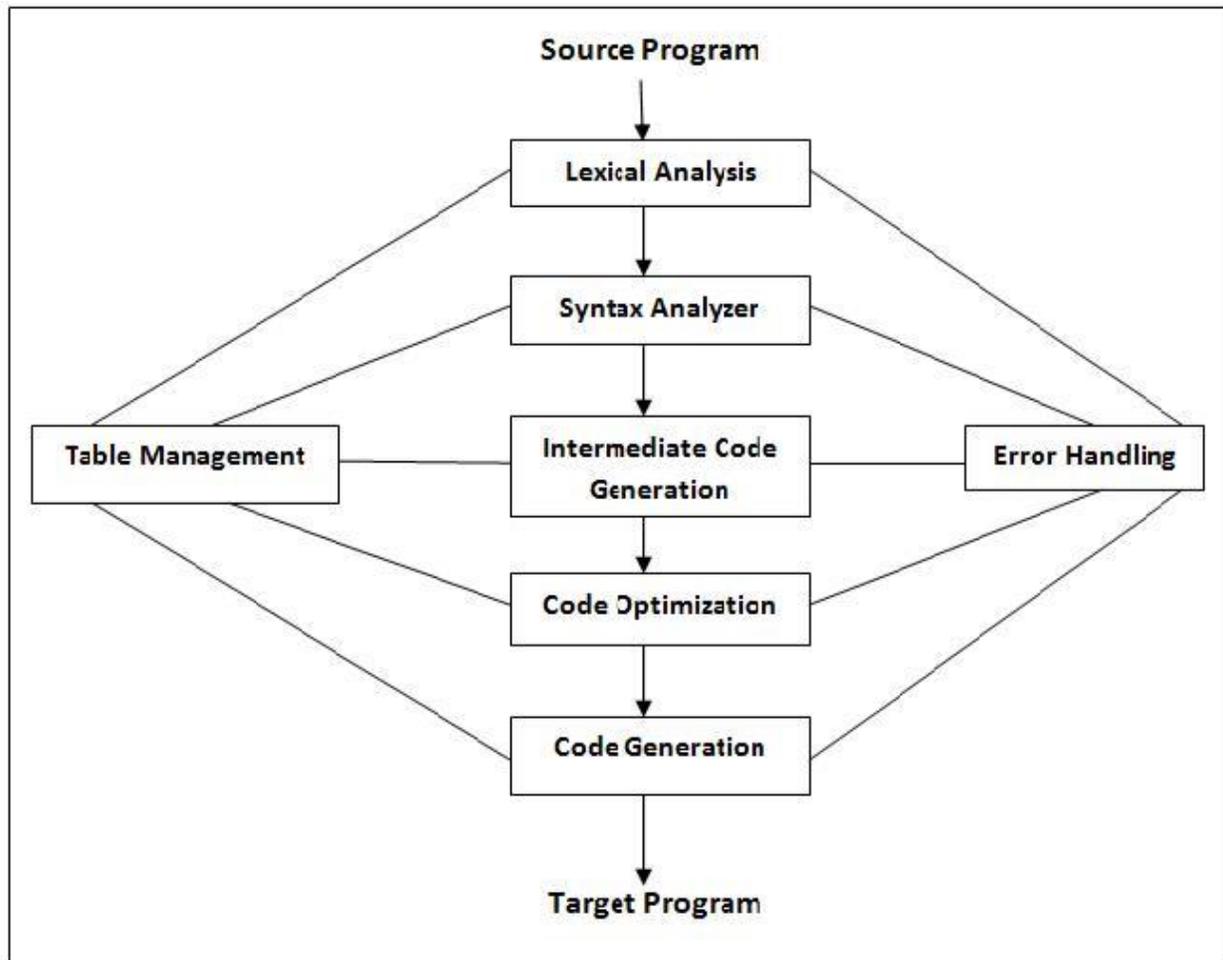


Fig. 1: Phases of a Compiler

As with our definition, a compiler converts a high-level language program into a machine language program. But the whole process does not occur in a single step but in a series of subprocesses called phases as shown in Fig. 1.

Each phase of the diagram becomes a chapter of this book because the phases of the compiler structure leads to the compiler design and that is what this book is about. Now let's go briefly through all the phases of the compiler.

The first phase of the compiler called lexical analyzer or scanner separates the characters of the source program into groups called tokens that logically belong together. Tokens can be keywords such as, DO or IF, identifiers or variables such as, X or NUM, operator symbols such as \leq or + and punctuation symbols such as parenthesis or commas. The tokens can be represented by integer codes for instance, DO can have the integer code 1, '+' can have 2 while 'identifiers' 3.

The syntax analyzer or parser, the next phase of the compiler, groups tokens into syntactic structures. For instance, the three tokens in A+B can be grouped together to form a syntactic structure called an expression. Expressions can further be grouped to form statements. Sometimes the syntactic structure can be represented as a tree whose leaves form the tokens.

The output of syntax analyzer or parser is a stream of simple instructions called intermediate code. The intermediate code generator which produces the intermediate code in the third phase usually use instructions such as simple macros with one operator and a small number of operands. The macro ADD2 statement explained in section 1.4 is

a bright example of this. The main difference between intermediate code and assembly code is that the former does not need to specify registers for each operation unlike the latter.

Code optimization, the next phase after intermediate code generator, is an optional phase that is geared to improving the intermediate code so that the ultimate object program runs faster and/or takes less space.

The final phase, Code Generator, should be designed in such a way that it produces a truly efficient object code, which is challenging, both in practical and theoretical terms.

In table management or bookkeeping portion of the compiler, a data structure called symbol table keeps track of the names (identifiers) used by the program and records important information about each such as, its type whether integer, real etc.

The error handler handles errors as the information passes from the source program through one phase to another.

Both the table management and error handling routines interact with all phases of the compiler.

Thank You for previewing this eBook

You can read the full version of this eBook in different formats:

- HTML (Free /Available to everyone)
- PDF / TXT (Available to V.I.P. members. Free Standard members can access up to 5 PDF/TXT eBooks per month each month)
- Epub & Mobipocket (Exclusive to V.I.P. members)

To download this full book, simply select the format you desire below

