

Module - I Introduction to Compiler

What is a compiler?

A compiler is a program that reads a program written in one language and translates it into an equivalent program in another language.

A compiler also reports errors present in the source program as a part of its translation process.

Why should we study compiler design?

- Compilers are everywhere!
- Many applications for compiler technology
 - Parsers for HTML in web browser.

• Interpreter for javascript / flash.

• Machine code generation for high level languages.

• Software testing.

• Program optimization.

• Design of new computer architecture.

• Compiler in-the-loop hardware development.

• Compiled simulation

• Used to simulate designs written in VHDL.

• No interpretation of design, hence faster.

Evaluation of compiler

Computers are electronic machines or in other words they understand only zeros and ones. Therefore in early days of computation instructions to a computer were given using perforated paper tapes and punch cards. In these

the presence of a ~~hole~~ hole is representation of zero and absence of them means one. Now writing down instructions in the form of ones and zeros is pretty hectic.

Punched card if converted to binary using ASCII standard

P - 1010000
Q - 1000011

R - 111010
S - 1100011

T - 1101110
U - 1110010

V - 1100011

W - 1101000

X - 1100101

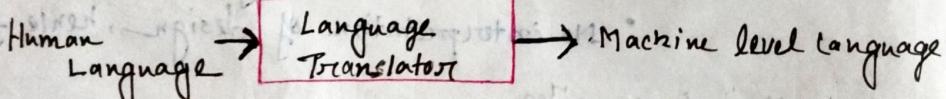
Y - 1100100

This is the format which computers can process directly.

But as it is almost impossible. Even if it becomes

possible for any human to encode instructions like this,

it is undoubtedly time consuming. For this we need language translators, in simpler terms we will say punch cards, and the middleman that is the language translator will convert that into the strings of ones and zeros.



(High level language)

1. Assembler :

Mov R1, 02H

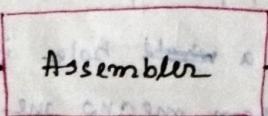
Mov R2, 03H

ADD R1, R2

STORE X, R1

Assembly

Language



01000100

100100110

010001100

101011100

Machine Code

None assembly language has a very strong correspondence b/w the instructions of it and the architectural specification of the computer. That means in order to write instructions for the computer in assembly language the coder needs to have a proper knowledge about the computer's architecture first.

In the 50's two new language translators came into picture

II. Interpreter

~~same~~ PHP

RUBY

PYTHON

Javascript

* one line at a time

1. Translates the program one statement at a time.

2. Considering it scans code one line at a time, errors are shown line by line.

3. Due to interpreters being slow in executing the object code, it is preferred less.

4. It does not convert source code into object code instead it scans it line by line.

5. The machine code is nowhere stored.

6. Less efficient.

7. CPU utilization is less.

Example - PHP, Python, Javascript, Ruby

III. Compiler

C

C++

C#

ERLANG

* whole program at once

1. The compiler scans the whole program in one go.

2. As it scans the code in one go, the errors (if any) are shown at the end together.

3. The main advantage of compilers is its execution time.

4. It converts source code into object code.

5. The machine code is stored in disk storage.

6. More efficient.

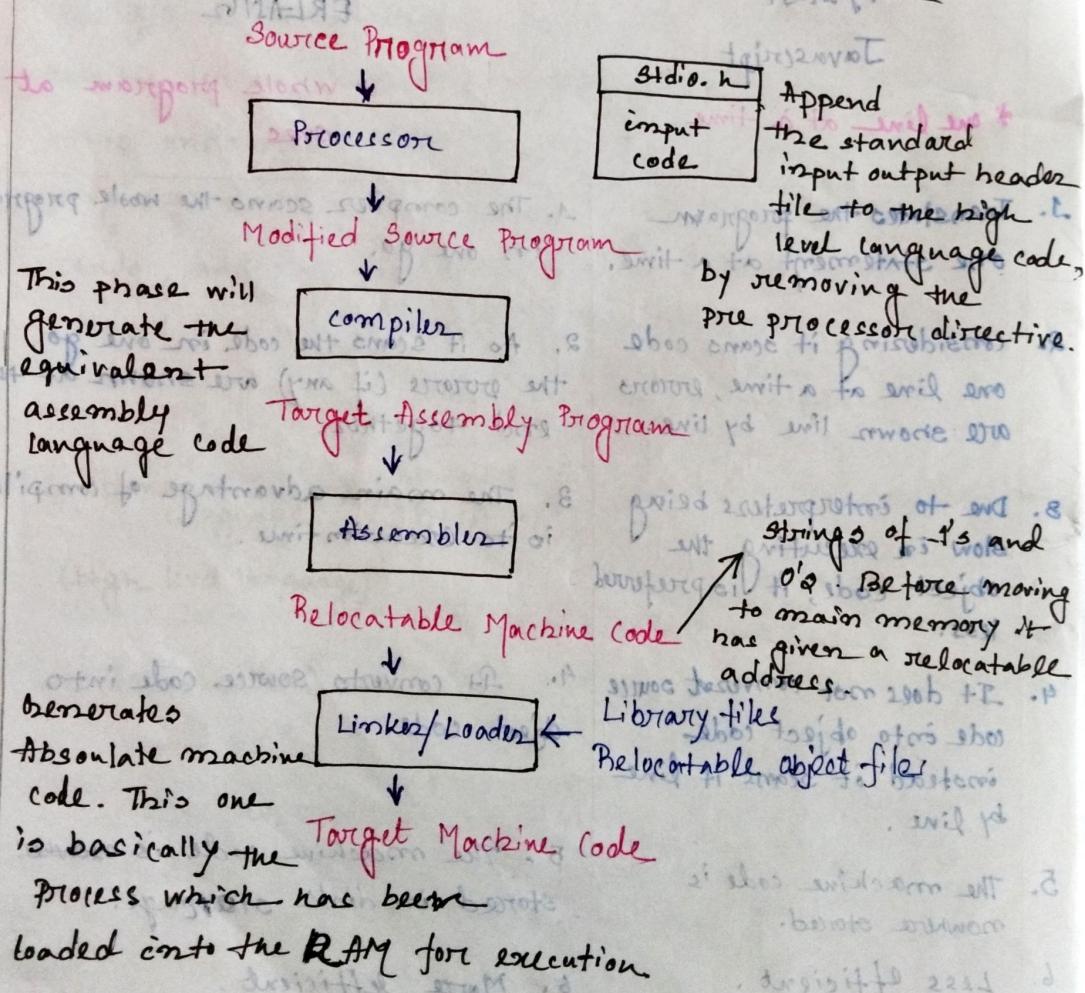
7. CPU utilization is more.

Example - C, C++, C#

C - High-level language

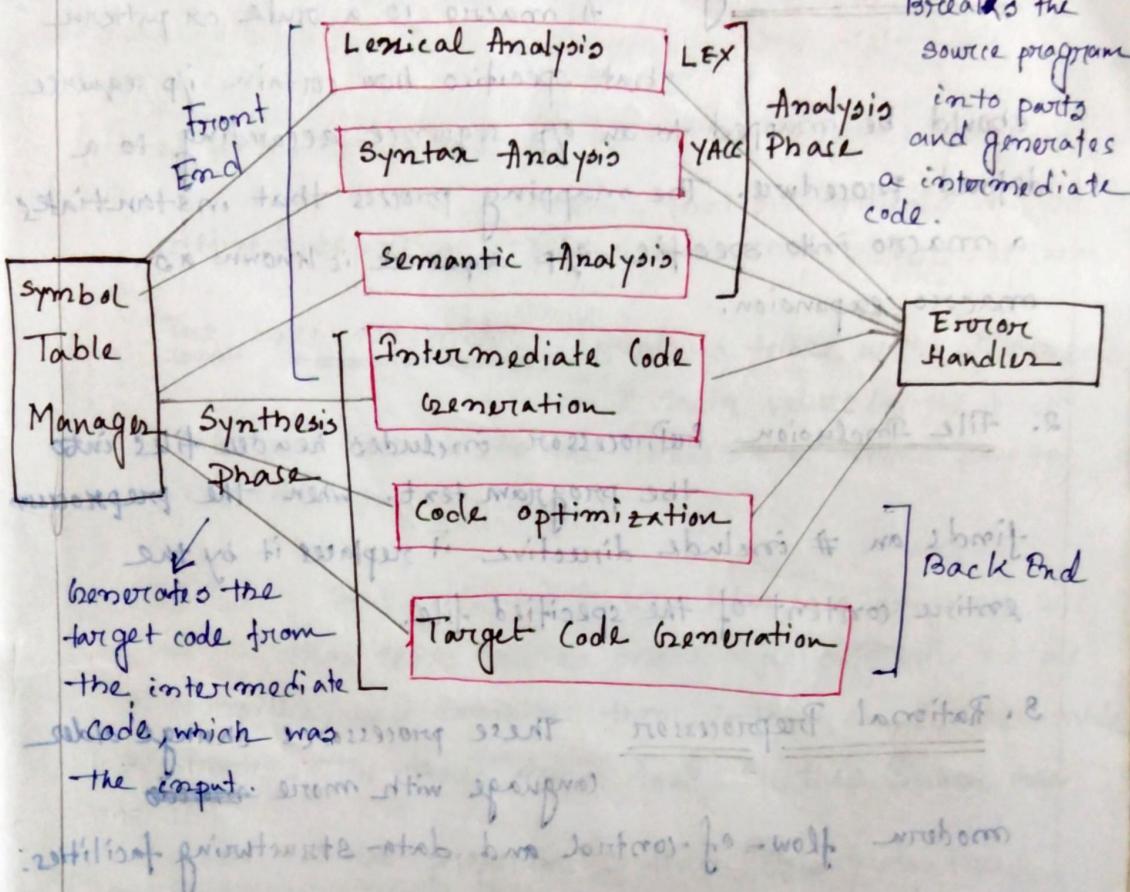
```
#include <stdio.h> → preprocessor directive  
int main() → main function  
{  
    int a, b=2, c=5;  
    a = a+b*c;  
    printf ("The value of a is %d", a);  
    return 0;  
}
```

Language Processing System / Language Translator



Internal Architecture of Compiler

It has 6 different phases.



Cousins of a compiler

1. Preprocessor
 2. Assembler
 3. Loader
 4. Linker

1. Preprocessor A preprocessor is a program that processes its input data to produce output that is used as input to another program. The output is said to be a preprocessed form of the input data, which is often used by some subsequent program like compiler.

They may perform the following functions:

1. Macro processing

A macro is a rule or pattern that specifies how certain i/p sequence should be mapped to an o/p sequence according to a defined procedure. The mapping process that instantiates a macro into specific o/p sequence is known as macro expansion.

2. File Inclusion

PreProcessor includes header files into the program text. When the preprocessor finds an #include directive it replaces it by the entire content of the specified file.

3. Rational Preprocessor

These processors change older language with more modern flow-of-control and data-structuring facilities.

4. Language Extension

These processors attempt to add capabilities to the language by what amounts to built-in macros. For example, the language Equel is a database query language embedded in C.

Assemblers Assembler creates object code by translating assembly instruction mnemonics into machine codes. There are two types of assemblers:

One-pass assembler go through the source code and assume that all symbols will be defined before any instruction that references them.

Two-pass assembler create a table with all symbols and their values in the first pass and then use the table in a second pass to generate code.

Linker A linker or link editor is a program that takes one or more objects generated by a compiler, and combines them into a single executable program. The three main tasks of the Linker are;

1. Searches the program to find library routines used by program, e.g. printf(), math routines.
2. Determines the memory locations that code from each module will occupy and relocates the instruction by adjusting absolute references.
3. Resolves references among files.

Loader A loader is the part of an operating system that is responsible for loading programs in memory, one of the essential stages in the process of starting a program.

Different Phases of Compiler

$x = a + b * 20;$

Lexical Analysis Phase generates tokens by scanning the whole program.

$x = a + b * 20$

Lexemes

Tokens

x identifier 1

= operator 1

a identifier 2

+ operator 2

b identifier 3

* operator 3

~~20~~ number

20

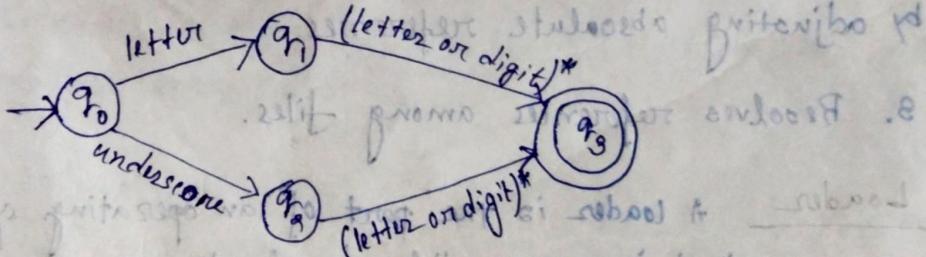
Regular expression for identifier:

$(l+d)^* | - (l+d)^*$

$l : \text{letter}$

$d : \text{digit}$

$- : \text{underscore}$

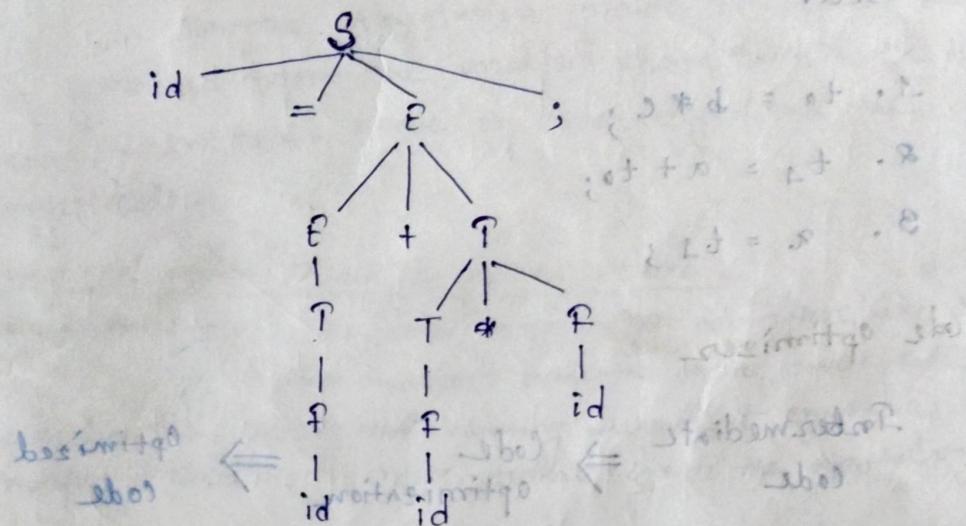


So, for examining the lexemes the lexical analyzer makes use of the type grammar. It tries to

Syntax Analyzer — The stream of tokens is passed to the syntax analysis phase. The syntax analyzer depends on the type 2 or context-free grammar.

$$S \rightarrow id \# E ; \quad E \rightarrow E + T / T$$

CFG
Production Rules of the expression



Parse Tree

Output $\rightarrow id = id + id * id ;$

If the parse tree and the provided expression is not same, then syntax analyzer phase generates error message, or syntax error.

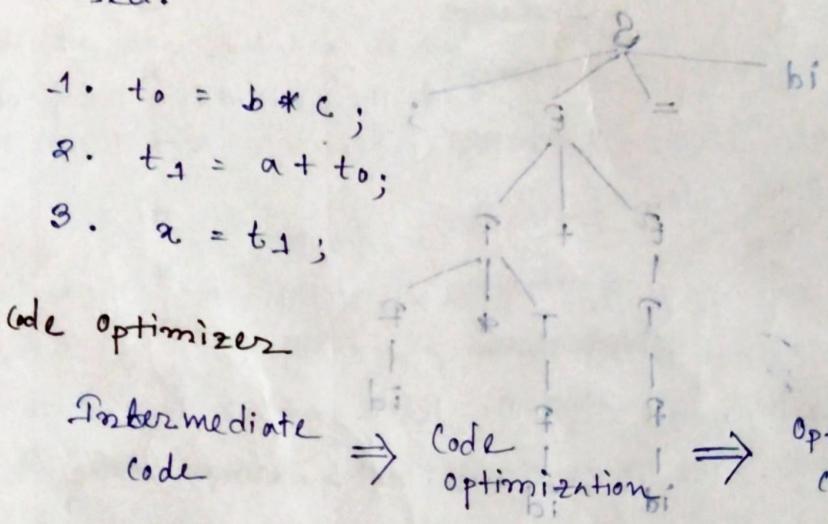
Semantic Analyzer — is responsible for type checking, array bound checking and the correctness of scope resolution. Basically it does the logical analysis of the parse tree. Semantic analyzer detects type mismatch errors, undeclared variables, misuse of reserved words, multiple declaration of a variable, missing out of scope variable, etc.

Semantic analyzer looks for the meaningfulness of the parse tree and verifies that.

Intermediate Code Generators

Semantically
verified Parse Tree \Rightarrow Intermediate
Code Generation \Rightarrow Postintermediate
code

For generating intermediate code, Three Address Code is used.



Code optimization can be machine dependent or machine independent.

1. $t_0 = b + c$: behavior with time and aging effect
 2. $a = a + t_0$: replaces surface next, aging term

Target Code Generator

Optimized code \Rightarrow Target Code generator \Rightarrow Assembly code segment design

mov eax, BWORD PTR [rbp-8]

Symbol Table It gathers information from the analysis phase and those gathered information are used by the synthesis phase. Symbol table is basically a data structure which is created and maintained by compilers in order to store information about the occurrences of various entities such as variable and function name, objects, classes, interfaces etc.

Now the information inside the symbol table is gathered from the analysis phase which in turn is used by synthesis phase in order to generate the target code.

Usage of Symbol Table by various phases

Lexical Analysis - It creates entries for identifiers. Lexical analyzer scans the entire source code line by line and during that scanning whenever it encounters any identifier it creates the entry for that inside the symbol table.

Syntax Analysis - It adds information regarding attributes, like type, scope, dimension, line of reference, line of usage etc.

Semantic Analysis - using available info checks semantics of the identifier created ↓ and updates the symbol-table accordingly during lexical analysis phase

Intermediate code Generation - The available information inside the symbol table, helps the intermediate code generator in adding the temporary variables information.

Code Optimization - The available information stored inside machine dependent the symbol table is used specifically in optimization.

Target code Generator - Generates the target code using address info for identifiers, stored inside the symbol table.

Symbol Table - Entries

int count;

char x[] = "Compiler Design";

Name	Type	Size	Dimension	Line of Declaration	Line of Usage	Address
count	char	1	2	char count;	char x[] = "Compiler Design";	int count;

Symbol Table - operations

- Non Block Structure Language (Fortran)
 - Contains single instance of the variable declaration.
 - Operations:
 - i) insert()
 - ii) lookup()
- Block Structure Language (c, c++)
 - Variable Declaration may happen multiple times.
 - Operations:
 - i) insert()
 - ii) lookup()
 - iii) set()
 - iv) Reset()