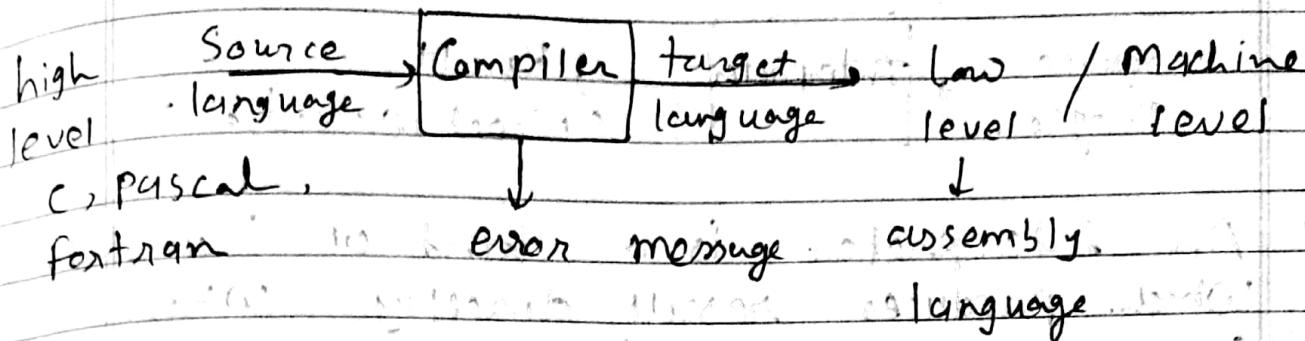


Chapter: 1 Overview of transition process

Compiler : A program which takes a program written in one language and converts it into another language.



→ Compiler is a program that can read a program in one language - source language - and translate it into an equivalent program in another language.

→ Compiler Converts higher language into lower level language or machine language.

for example :

There is a C program : Called HelloWorld.c

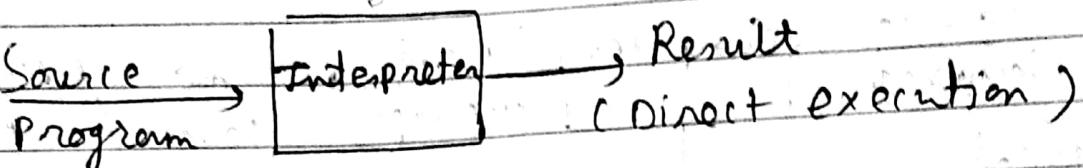
```
#include <stdio.h>
{ printf("HelloWorld"); }
```

So, Computer can not understand higher level languages. So, Compiler Converts it into binary language in the form of 0 and 1. It means Compiler creates object file. Linker creates the executable file from object file.

Example of Compiler:

gcc, javac, Microsoft visual Studio.

line by line execution: Interpreter : Convert one line at a time.



- An Interpreter is a kind of translator which produces result directly when Source language and data is given to it.

Example of Interpreter:

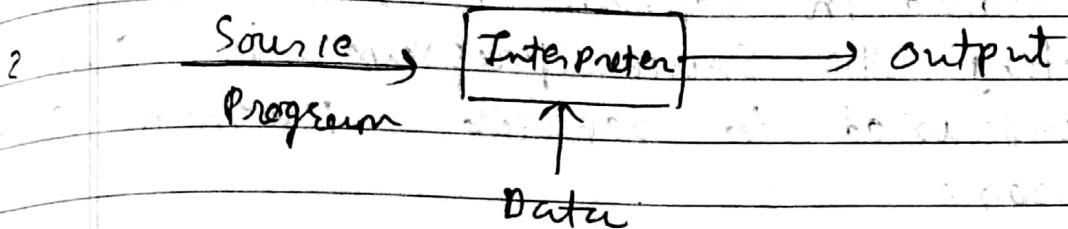
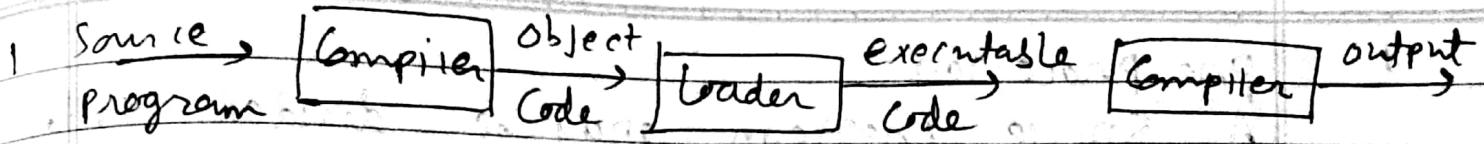
Ruby, Python.

difference between Compiler and Interpreter

- 1 In the process of compilation the program is analyzed only once and then code is generated.
- 2 The source program gets interpreted every time it is to be executed.
- 1 Compiler is efficient than interpreter.
- 2 Interpretation is less efficient than compiler.
- 1 Compilers produce object code.
- 2 Interpreters do not produce object code.

1 Complex program

2 Simple program



1 example : UPS debugger

2 example : Turbo C Compiler

Cousins of Compiler : / Translators :

Language processing System : typical Compilation

1 Context of Compiler :

Source Program



Preprocessor

↓ modified source program

Compiler

↓ target assembly program

Assembler



↓ relocatable machine code

Loader/linker



Absolute machine Code

1 Pre processor :

A Source program may be divided into modules stored in separate files. The task of collecting the source program is sometimes entrusted to separate program called pre processor.

Preprocessors produce input to compilers.

Example :

C preprocessor : is not a part of the Compiler, but is a separate step to in the compilation process.

All Preprocessor Commands begin with # symbol.

#include : insert a particular header from another file.

#include <stdio.h>

These directives tell the CPP to get stdio.h from System libraries and add text to the current source file.

→ Preprocessor also allows a user to define macros that are shorthand for longer constructs.

2 Compiler : The modified source program is then fed to the Compiler. The Compiler may produce an assembly language program as its output because assembly language is easier to produce as output and easier to debug.

Object Code : is portion of machine code that has not been linked yet into target program.

PAGE NO.	
DATE	

3 Assembler :

- Some Compilers produce assembly code as in figure (a). that is parsed to the assemble for further processing.
- other Compilers performs the job of the assembler producing relocatable machine code that can be directly parsed to the Loader / link - editor.
- Assembly Code is mnemonic version of the machine code in which names are used instead of binary codes for operations and names are also given to memory addresses.

Mov a, R₁.

And #2, R₁.

Mov R₁, b.

- Assembler produces relocatable machine code as output.

relocatable machine code : object code

4 Loader and link editors

- Loader performs the two functions: loading and link editing
- The process of loading consists of taking relocatable machine code, altering the relocatable addresses and placing the altered instructions and data in memory at the proper locations
- The link editor allows us to make a single program from several files of relocatable machine code
- Linker / Loader then generates absolute / target machine code

absolute machine code : executable code or target program.

Hierarchy of language :

There are three types of language:

- 1 Machine language } low level language
- 2 Assembly language }
- 3 High-level language

low level languages are closer to the languages used by Computer, while high-level languages are closer to human languages.

language : specifies a set of instructions that can be used to produce various kinds of output.

PAGE NO.	
DATE	

1 Machine language :

- Machine language is a collection of binary digits or bits that the computer reads and interprets.
- Machine languages are the only languages understood by computers.
- But Machine languages are almost impossible for humans to use because they consist entirely of numbers.

Example :

Let us say that an electric toothbrush has a processor and main memory. Processor can rotate the bristles left and right, and can check the on/off switch. The machine instructions are one byte long and correspond to the following machine operations:

Machine instruction	Machine operation
0000 0000	Stop
0000 0001	Rotate bristles left
0000 0010	Rotate bristles right
0000 0100	go back to start of program
0000 1000	Skip next instruction if switch is off

2 Assembly language :

- is a low level programming language for a computer or other programmable device or microprocessors.

Assembly language consist of series of mnemonic instructions followed by a list of data, arguments or parameters.

- These are translated by an assembler into machine language instructions that can be loaded into memory and executed.

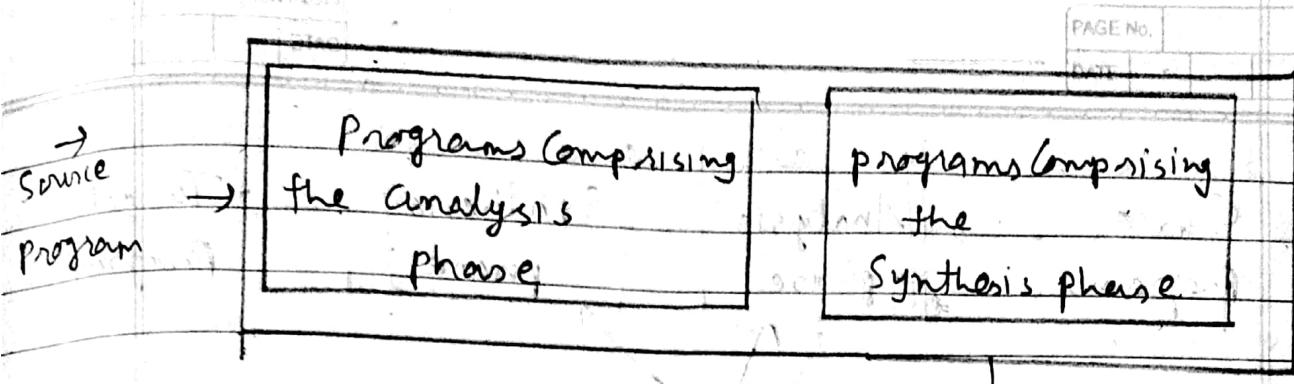
3 High level language :

- High level languages allow us to write computer code using instructions resembling everyday spoken language (for ex: print, if, while) which then translated into machine language to be executed.
- Programs written in high level language need to translated into machine language before they can be executed.
- Some programming language use Compiler to perform this translation and others use an interpreter.

Components of Compilers :

- 1 Single pass Compiler:
- 2 Multi pass Compiler:

Pass: one complete scan of source program is called pass.



- In a Single pass Compiler analysis of Source program is immediately followed by Synthesis of equivalent target statement/program.
- It is difficult to Compile the Source program into Single pass due to :

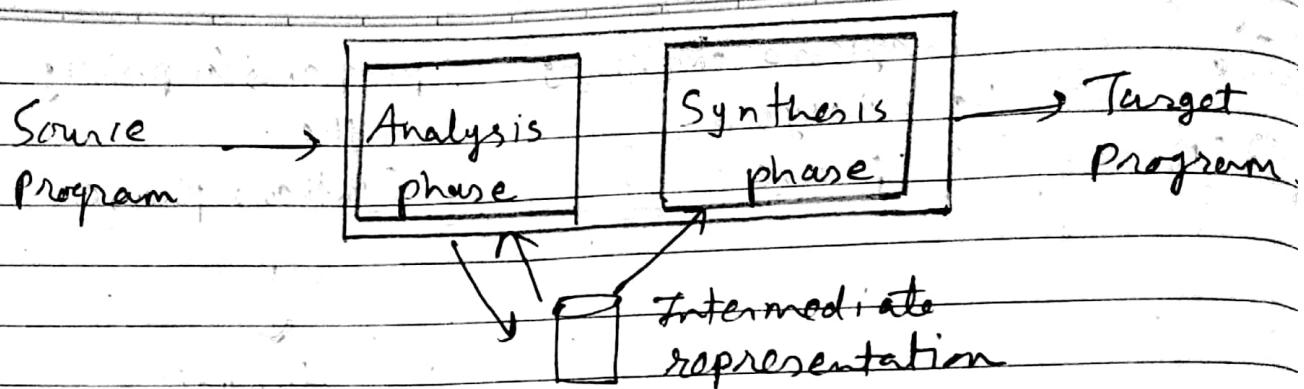
forward Reference: a forward-reference entity is a reference to the entity which precedes its definition in the program.

→ using transient identifiers before declaration.

Multipass Compiler:

- In pass I : perform analysis of the Source program and note relevant information.

In pass II : generate target code using information noted in pass I / perform Synthesis of the target Program.



Intermediate representation:

Compiler generate an easy to represent form of source language which is called IR.

→ difference between single pass and multi pass

Compiler:

- 1 One pass Compiler is a Compiler that passes through the Source Code of each Compilation unit only once.
- 2 A multi pass Compiler is a type of Compiler that processes the Source Code or AST of program several times.

1 faster

2 Slower

- 1 Narrow Compiler
- 2 Wide Compiler

- 1 Language like pascal can be implemented with single pass compiler
- 2 languages like Java requires a multi-pass compiler.

Difference between Compiler and Assembler

- 1 high level language to machine code
- 2 mnemonic of code to machine code

1 Types of Compiler : Single pass compiler,
Multi pass compiler

2 Types of assembler : single pass assembler
Multi pass assembler

1 C Compiler

2 8085, 8086 Instruction Set

Cross Compiler :

A cross compiler is a compiler capable of creating executable code for a platform other than the one on which the compiler is running.

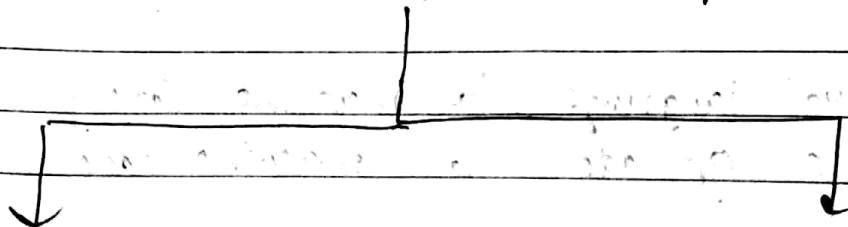
Compiler which run on one machine and produces the target code for another machine

Analysis Synthesis model of Compilation / Structure of Compiler

Phases of Compiler

Write output of phases of a compiler.
for $a = a + b * c / 2$; type of a, b, c are float.

→ two Parts of compilation process.



Analysis phase: The main objective of the analysis phase is to break the source code into parts and then arranging these pieces into a meaningful structure.

Analysis Phase is divided into three Subpart

- 1 Lexical analysis
- 2 Syntax analysis
- 3 Semantic analysis.

Synthesis phase: Synthesis phase is concerned with generation of target language statements which has the same meaning as the source statement.

- 4 Intermediate Code generation
- 5 Code optimization
- 6 Code generation

Source Program



lexical analysis



Syntax analysis



semantic analysis



Intermediate Code



optimization
Code generation



Code generation



target program.

Symbol
table

error detection
and recovery

I Lexical analysis: remove the whitespace or comment in the source code.
 Also called as Scanning or linear Analysis.
 Complete Source Code is scanned from left to right and your source program is broken up into group of strings called token
 for example:

$a = b + c * 60$

In lexical analysis following assignment statement $a := b + c * 60$ would be grouped into the following tokens:

- 1 The Identifier a
- 2 The assignment symbol =
- 3 The Identifier Identifier b
- 4 The plus sign +
- 5 The Identifier c
- 6 The multiplication sign *
- 7 The number 60

token : Sequence of characters that can be treated as single logical entity.

→ Examples : identifiers, Constant, keywords, operators.

Identifiers : The symbol whose meaning is not defined in library.

int a → identifier.

Pattern : Set of rules that describe the token

Lexeme : Sequence of characters in the source program that are matched with pattern

lexical errors:
Keywords: predefined reserved words used in programming that have special meaning to the Compiler.

of token.

lexeme → token

int a

int: keyword

a: identifier

lexeme identifier token pattern: letter followed by letters and digits.

2 Syntax Analysis:

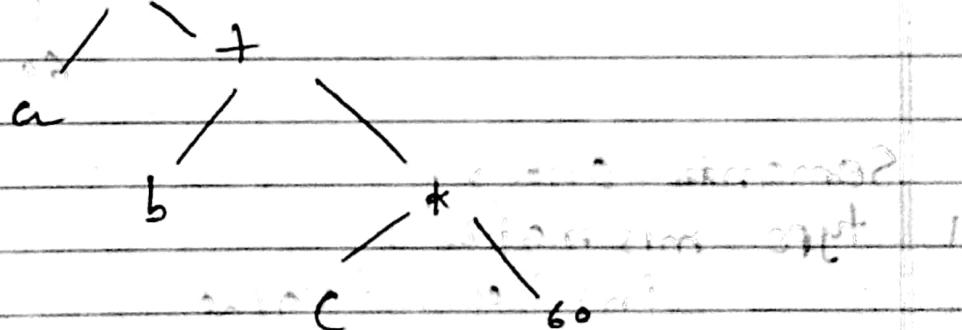
also called as Parsing or Hierarchical analysis

It involves grouping of tokens of the source program into grammatical Phases that are used by the Compiler to synthesize output.

Or generate Syntax errors.

for the expression: $a = b + c * 60$

parse tree: =



→ parsing can be done in two ways: top down and bottom-up.

Here, the expression are:

$b + c$ and $c * 60$

So, we will evaluate the $c * 60$ because '*' has higher priority than '+'.

\uparrow , \star , $/$, $+$, $=$

3 Semantic analysis :

The next phase is Semantic analysis . Which determines the meaning of source string means matching of parenthesis , matching of if ... else statements

→ main ^{use} of Semantic analysis is type checking, type conversion:

~~to~~ → a → +

~~int~~

b

c int to real

| result

60

Semantic errors :

1 type mismatch :

int a = "value"

2 undeclared variable

4 Intermediate Code generation :

It should have two properties :

1 It should be easy to produce

2 easy to translate into the target program

This code is in variety of forms such as three address code , triple , quadruple , posix .

Here, we will consider an Intermediate Code in three address code

Three address code consists of instructions each of which has atmost three operands

$t_1 := \text{int to real (60)}$

$t_2 := c + t_1$

$t_3 := b + t_2$

$a := t_3$

- 1) Each instruction has at most one operator in addition to assignment.
- 2) Compiler must generate a temporary name to hold the value computed by each instruction.
- 3) Some instruction may have fewer than three operands.
- 5) Code optimization :
 - used to improve the intermediate Code for faster execution and less consumption of memory.
 - It eliminates the redundant instructions.

$t_1 := c * \text{int to real (60)}$

$a := b + t_1$

6 Code generation:

- In Code generation phase the target code gets generated.
- The Intermediate instruction are translated into sequence of machine instructions.

ex: MOV C, R₁ The # signifies that
MUL #60.0, R₁ 60.0 is to be treated
MOV b, R₂ as a constant.
ADD R₂, R₁
MOV R₁, a

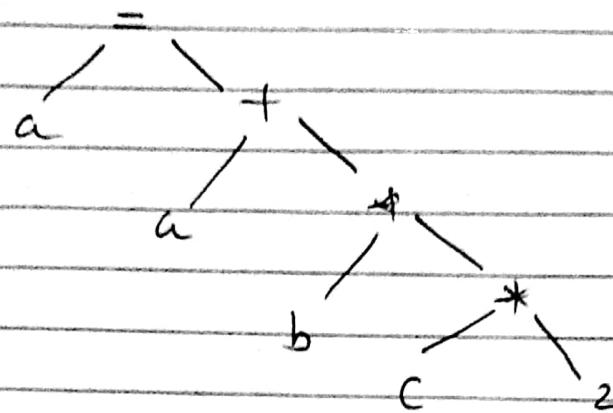
3B-5
for a statement given below, write output of all phases (except of an optimization phase) of a Compiler. Assume a, b, and c of type float

$$a = c + b * (* 2);$$

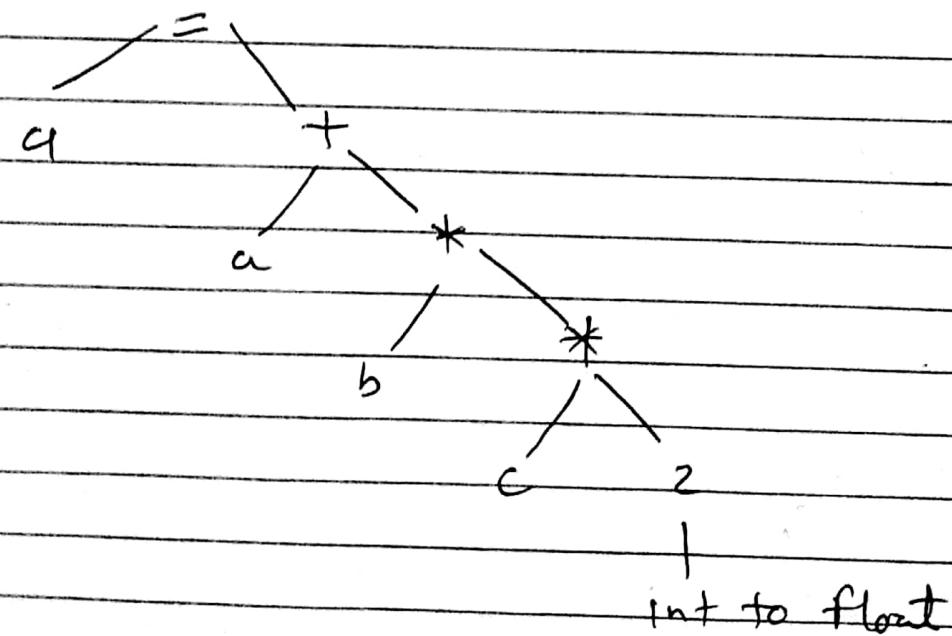
① lexical analysis:

- 1) The identifier a
- 2) The assignment symbol =
- 3) The identifier b
- 4) The multiplication sign *
- 5) The identifier c
- 6) The multiplication sign *
- 7) The number 2.

② Syntax analysis :



③ Semantic analysis :



④ Intermediate code generation :

$t_1 := \text{int to float}(c)$

$t_2 := c * t_1$

$t_3 := b * t_2$

$t_4 := a + t_3$

~~$t_5 := t_4$~~

⑥ Code generation:

```
Mov C, R1  
Mul #2.0, R1  
Mov b, R2  
Mul R2, R1  
Mov a, R2  
Add R2, R1  
Mov R1, a
```

prepared by kinjal patel