

16/1/18

MODULE - 1



KTU STUDENTS

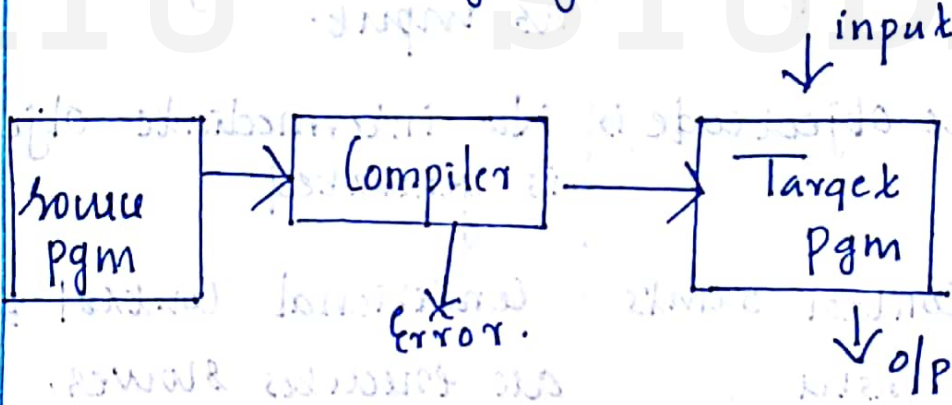
Compiler: It is a pgm that reads a pgm written in ^{high level} ~~the~~ language (source language) and translates it into an equivalent pgm in another language [target language].

An imp. role of compiler is to report any errors in a source pgm that it detects during the translation process.

NOTE

Compiler is a system software

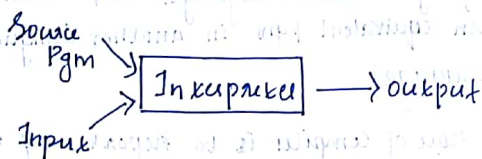
Compiler is a language translator



eg: C compiler, Java c, Dempo C, ...

Interpreter: An Interpreter another common kind of language processor. Instead of producing a target pgm an interpreter directly execute the operatⁿ specified by in the ~~at~~ source program.

The m/c language target program produced by a compiler is usually much faster than an interpreter, as mapping inputs to outputs.



Difference b/w Compiler & Interpreter

Compiler	Interpreter
<ul style="list-style-type: none"> It takes entire program as input. Intermediate object code is generated Conditional control stmts are executed faster M/c requirement is more (time object code is generated) Pgm need not be compiled everytime 	<ul style="list-style-type: none"> It takes single instruction as input. No intermediate object code is generated Conditional control stmts are executed slower. M/c requirement is less Everytime higher level pgms is converted into lower level Pgms

Errors are displayed after entire pgm is checked

eg: C, C++ languages uses Compilers.

Errors are displayed for every instruction interpreter

Pgming languages like Python, Ruby uses interpreter

Concept of a Compiler:

Source Program

Preprocessor

Source Program

Compiler

Assembly Code

Assembler

Relocatable Machine Code/object code

Loader/Linker

Executable machine code/Absolute code.

Assembler: is a language translate assembly code into machine language.

2 marks
The programs which assist the compiler to convert a skeleton source code into executable form make the context of a compiler.

① PREPROCESSOR: scans the source code and includes the header files which contain various relevant information for various functions.

* preprocessor performs macroexpansion also.

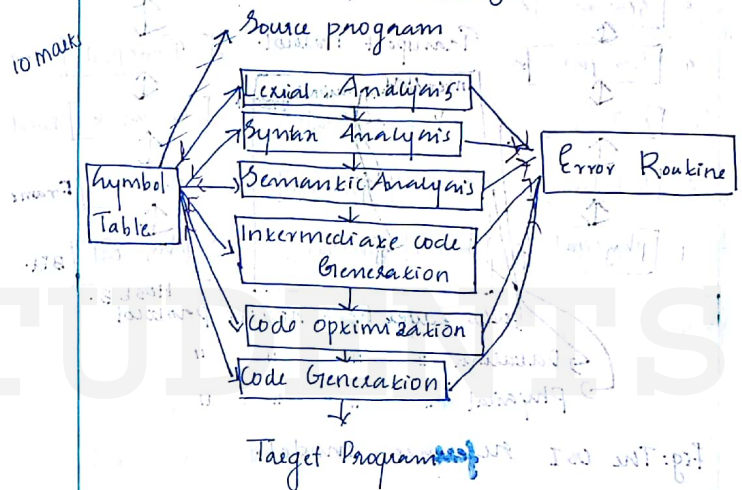
② COMPILER: passes the source code through various phases and generates the target assembly code.

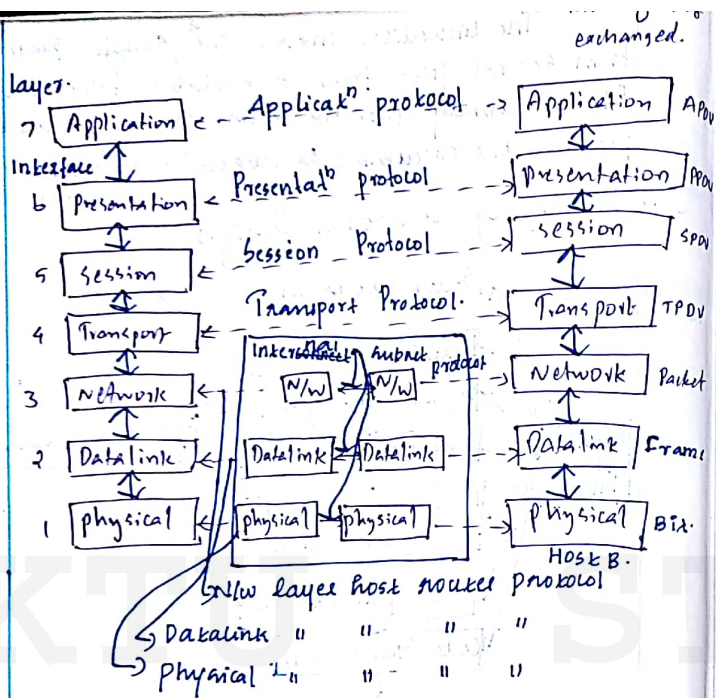
③ ASSEMBLER: converts the assembly code into relocatable m/c code or object code. Although, this code is in zero and one form but it can't be executed Bcz this code has not been assigned the actual m/c address.

④ LOADER OR LINKER (LINK EDITOR): it performs two functions. The process of loading consist of taking m/c code, altering relocatable address & placing the altered instructions in memory at proper location.

The linker makes the single program from several files from relocatable codes. These files are library files the program needs. The loader produces the executable or absolute machine code.

→ Phases of Compiler Design:-





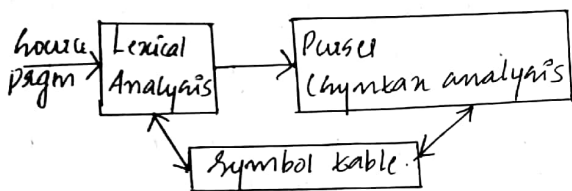
A compiler operates in phases. A Phase is a logically interrelated operation, that takes source program in one representation and produces output in another representation. The phases of compiler are shown in above diagram. There are two major phases of compilation.

1. **Analysis phase:** [Machine independent but language dependent]
 2. **Synthesis Phase:** [Machine dependent and language independent]
- Analysis phase consists of lexical analysis, syntax analysis, semantic analysis, intermediate code generation.
- Synthesis phase consists of code optimization and code generation.

Phase-1: Lexical Analysis

A lexical analyzer needs the stream of character making up the source program and groups the character into meaningful sequences of lexemes.

For each lexeme lexical analyzer produces a token of the form [token name, attribute value]. Tokens are pass to subsequence phase for syntax analysis.



Example:

Sample statement newval = oldval + 12.

TOKENS:	lexeme	token
	newval	id 1
	=	Assignment operator
	oldval	id 2
	+	Add operator
	12	Number

lexical analyzer translates while space also removes error.

Phase-2: Syntax Analysis

Syntax analysis is also called parsing the parser uses token produced by the lexical analyzer to create a tree in intermediate representation that depicts the grammatical structure of the token string.

A typical represent is the syntax stream in which each internal node represent operation and children of the node represent argument of the operation.

eg: syntax correspond to
newval = oldval + 12

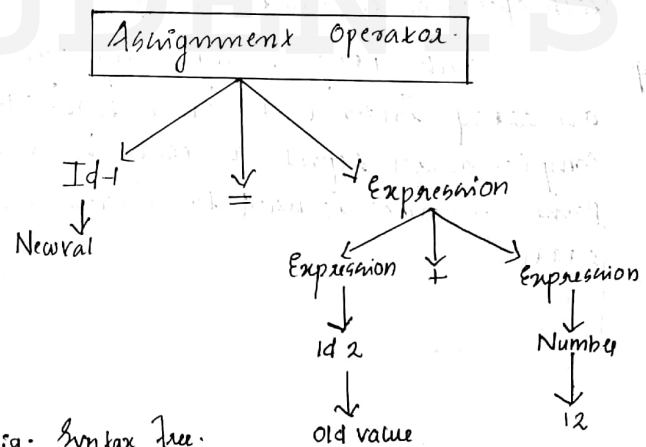


Fig: Syntax Tree.

Phase-3 Semantic Analysis

The semantic analyzer uses the syntax tree and the information in the symbol table to check the source program for semantic consistency with the language definitions.

It gathers information from either in the syntax tree or symbol table and saves it in symbol table for subsequent use during intermediate code generation. An important part the semantic analysis is type checking, where the compiler checks that each operation on matching operations.

eg:

Each programming language requires an array index to be an integer. The compiler must report an error if a floating point number is used to index an array.

Phase-4: Intermediate Code Generation

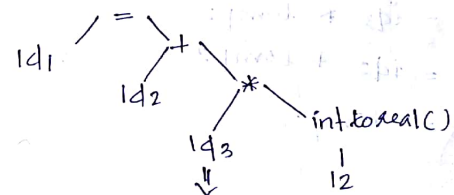
After the semantic analysis of source program many compilers an explicit low level or machine layer intermediate representation. This intermediate representation should have two important properties

1. It should be easy to produce
2. It should be easy to translate into the target machine.

The considered intermediate code for three address code, which consists of three operands per instruction each operand can act like a register. This phase bridges analysis & synthesis phase of translator.

eg:

$$\text{Newval} = \text{oldval} + \text{fact} * 12$$
$$ld\ 1 \quad ld\ 2 \quad + \quad ld\ 3 \quad * \quad 12$$



kemp-1 = int to real (12)

kemp-2 = id3 * kemp1

kemp-3 = id2 + kemp2

id1 = kemp3

intermediate code.

Three address code
operand=3.

Phase-5: Code Optimization

The compiler looks at the large segment of the formula to decide how to improve performance. The machine independent on optimization attempt to include the intermediate code so that better target code will be generated.

Usually better means: faster, shorter, target code that consumes less power.

eg: The above intermediate code will optimize

for kemp = id3 * 12.0

id1 = id2 + kemp2.

kemp1 = int to real (12)

kemp2 = id3 * kemp1

kemp3 = id2 + kemp2

id1 = kemp3

Phase 6: Code Generation

The last phase of translation is code generation. It takes intermediate representation of source program and maps it into the target language.

eg: id1 = id2 + id3 * 12

mov R1, id3

mul R2, #12

mov R2, id2

ADD R1, R2

mov id2, R1

25/1/17

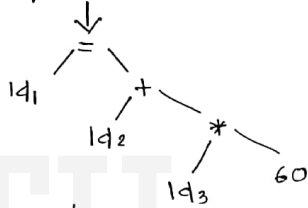
Example:

position = initial + rate * 60

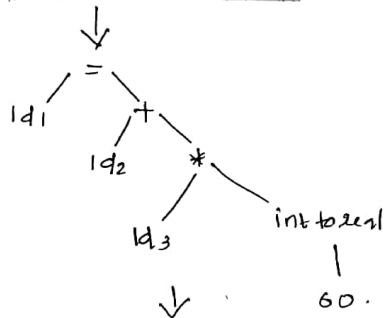
↓
Lexical analyzer

↓
 $id_1 = id_2 + id_3 * 60$

↓
Syntax analyzer



↓
Semantic Analyzer



Symbol table

1	Position	id1
2	Initial	id2,...
3	rate	id3,...

Intermediate code generator

temp1 = integer(60)

temp2 = id3 * temp1

temp3 = id2 + temp2

id1 = temp3

↓

Code generator optimizer

↓

temp1 = id3 * 60

id1 = id2 + temp1

↓

Code Generator

Mov id3, R2

Mov 60, R1

Mov id2, R1

Add R2, R1

Mov R1, id1

Symbol Table Management:

Symbol table is a data structure containing a record for each variable name with fields for attributes of the variable name.

The data structure should be designed to allow the compiler to find the record for each variable name quickly or to store/retrieve data from it quickly. These attributes may provide information about storage allocated for a name, its datatype, its scope (where it in the pgm is value may be used and in the case of procedure names the things such as name of arguments, its type, the method of passing each argument).

Position	Id ₁ , ...
Initial	Id ₂ , ...
Rate	Id ₃ , ...
⋮	⋮

Error handling routines:

One of the most imp. functions of a compiler is the detection & reporting of errors in the source pgm. The error msg should allow the programmer to determine exactly where the errors have occurred. Errors may occur in all

or the phases of a compiler. Whenever a phase of the compiler discovers an error, it must report the error to the error handler, which issues an appropriate diagnostic message. Both of the symbol table management & error handler routines interact with all phases of the compiler.

Lexical Errors:

e.g. SA (Invalid variable name)

First character should be alphabet followed by alphabet/digit.

Syntax Error:

missing of braces, ;

Semantic Error:

Invalid array declaration a[10.5]

Division by zero.

char a;

int b, c

c = a + b

One pass

It passes through the source code of each compilation phase only once. Thus efficiency

is limited because they don't produce intermediate. OPC is very common because of its simplicity they are faster than multipass compilers. Also known as narrow compilers.

eg Pascal MC compiler.

→ Multipass compilers:

The i/p is passed through certain phases in one pass. Then the o/p of previous phases is passed through another phases in second pass until the desired o/p is generated. It requires less mly because each pass takes o/p of previous phase as i/p. It may create one/more intermediate code. Also known as wide compiler.

eg Modula 2.

Front end & back end of a compiler.

The phases of compiler are collected into front end and back end (analysis phase).

The front end consists of those phases that depend primarily on the source program. These normally include lexical

analysis, syntactic analysis, semantic intermediate code generator.

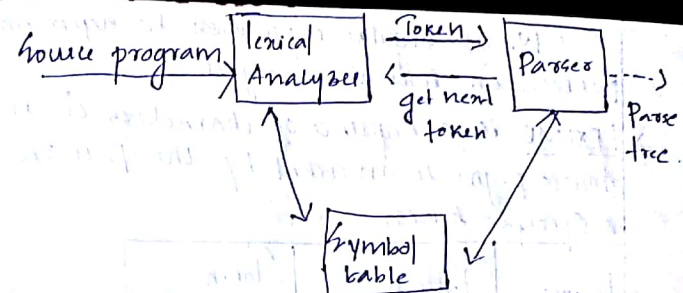
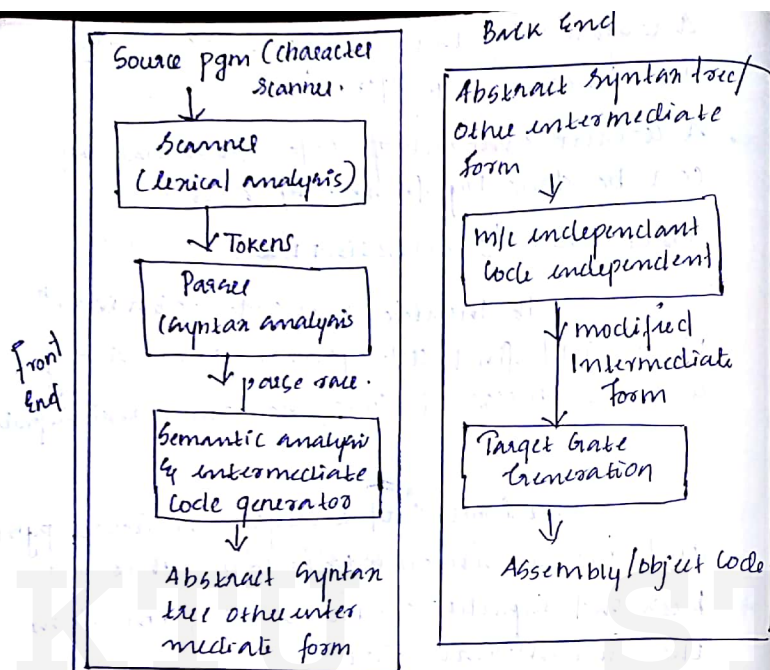
A certain amount of code optimisation can be done by front end as well.

Back end / synthesis phase

It includes the code optimisation phase and final code generation phase along with the necessary error handling and syntax.

The front end analyses the source program and produces intermediate code while the back end synthesises the target program from the intermediate code.

The front end phase consists of those that primarily depend on source program and are independent of target machine. Various phases of a compiler consist of those which depend on target and are independent of source program.



Lexical analyser is also called scanner
 Secondary function performed by lexical analyser are as follows

- 1) Remove comments
- 2) Remove wide spaces in the form of blank
- 3) Tab and new line character
- 3) Report lexical errors.

TOKENS, PATTERNS & LEXEME

- TOKENS are basic building blocks of a source program which are indivisible. Tokens can be classified into keywords, operators, identifiers, separators, strings, numerical constants.
- PATTERN is a set of rule to recognise tokens present in source program.

Lexical Analyser : is the first phase of compiler's main task to read the i/p character (source program) and produce as output a sequence of tokens that the parser uses for syntax analysis. This interaction is summarised in below figure

function of Lexical analyser

We use regular expression to represent pattern in lexical analyser.

→ LEXEME is a sequence of characters in the source prgm is matched by the pattern. For a specific token

lexeme	Pattern	Token
785	$[0-9]^+$	number
Sum	$[a-z]$	id 1
a10	$[a-z][0-9]$	id 2

Specification Of Token.

To specify a token in lexical analyser by using pattern. Regular expression are an important notation for pattern. Specifying pattern.

Definition of Regular Exp

Regular expression over alphabet Σ is defined as follows.

1. ϵ, a, b are primitive regular expression denotes the language $\{\epsilon\}, \{a\}, \{b\}$
Suppose 'R' and 'S' are regular expression denoting the language $L(R)$ and $L(S)$. Then

a) 4.5 is a regular expression

b) 1/5 is a regular expression

c) a^* " " " " " " " " " " " "

d) (1) " " " " " " " " " " " "

eg: Specification of Token, Identifier Using regular expression is as following

id $[A-Za-z][A-Za-z0-9]^+$

Recognise digit and decimal no's using regular expression.

Q Develop pattern to recognise fractional number 10.21

digit $\rightarrow [0-9]$

num $\rightarrow [digit]^+$

frnum $\rightarrow [digit]^+ (\cdot [digit]^+)^?$

? \Rightarrow zero/one time.

Definition of finite automata

Review of finite automata

NFA & DPA

Context free grammar.

Derivation tree or Parse tree

left most derivation & Right Most derivation.

Ambiguous grammar (TOC note)

Q Pattern specification for keywords

keyword \rightarrow "int" / "float" / "while" / "if" / "else" / "for"

Q Pattern specification to identify comment lines.

cmt \rightarrow "/*" . "*" /

[here, . is a regular generator matches
with all characters except newline]

cmt \rightarrow "//" single line.

Q Pattern to identify arithmetic operators.

oprs \rightarrow '+' / '-' / '*' / '/'