

Text Books

- ✓ Compilers "Principles, techniques and tools" - Aho, Lam, Sethi and Ullman - Pearson Publication
- ✓ Compiler Design - Dr D. G. Kaele (University Science Press)
- ✓ Compiler Design - Rajesh K. Maurya - Wiley Press, New Delhi
- ✓ Compiler Design - Jay Kant P. S. Yadav, Rajeev Garg (Vayu Education of India)

Unit 1 - Introduction to compiler

Introduction :-

High level language play a vital role in design of software. But these languages in which we want to express our computational needs are significantly different from the language recognized by real machine. It is here that language translator plays an important role. Let us discuss what is translator?

Language translator or language processor :-

A language translator or language processor (or simply translator) is a program that translates an input program written in a programming language (called source program) into a functionally equivalent program in another language (called target program). Apart from translation, a translator should have error detection capability. Any violation of source language specification would be detected and reported to the programmer. The working of translator shown in Fig 1.1.

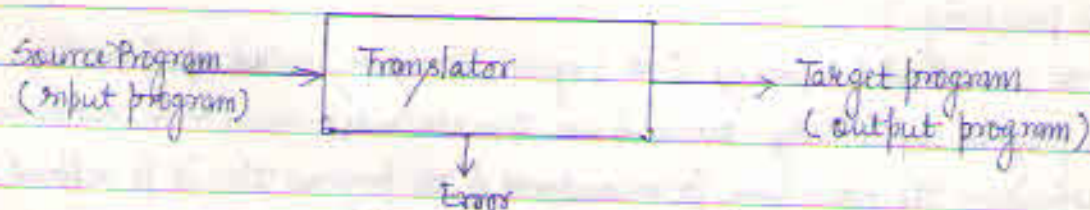


Fig 1.1: Language translator or processor

Language processing system (causing of compiler)

As we know that any computer system is made of hardware and software. The hardware understands a language of 0,1 which human cannot understand so we write program in high level language, which is easier for human to understand and remember. These prog

ram are then fed into a series of utilities (programs) and operating system components to get the desired code that can be used by the machine. This is known as language processing system. Fig 1.2 illustrates language processing system.

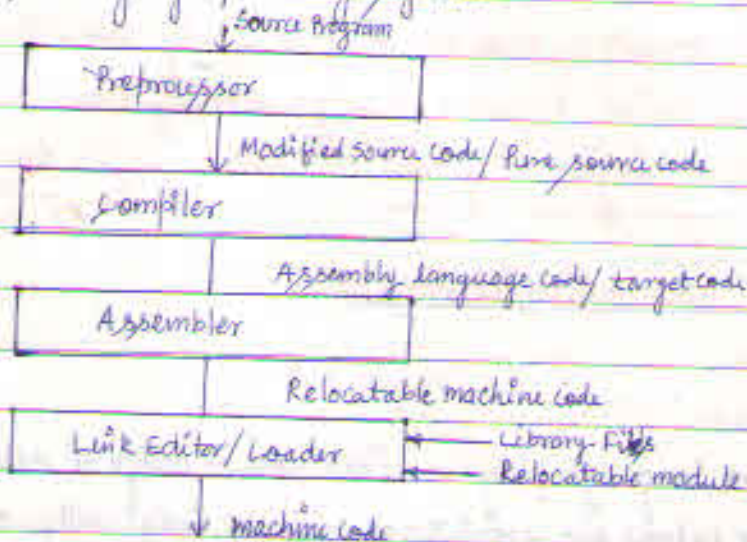


Fig 1.2 : Language Processing System

(i) **Preprocessor** :- It takes the source code as input and produces modified source code as output.

Before compilation, the source program is preprocessed by preprocessor to prepare it for the compilation. The preprocessor program creates modified source program from the original source program by replacing preprocessor directives with the suitable content. The new source program acts as an input to the compiler. The preprocessor performs various tasks as given here :-

(a) **Macro processing** :-

It permits the user to include macros in the program. Macros are smallest of instructions that are used in a program repetitively. Macro have two attributes macro name and macro definition. Whenever the macro name is encountered in the program then it is replaced by macro definition (set of statements corresponding to macros).

For example - #define identifier replacement strings
(macro name) (macro definition)

eg. - #define Array_Size 50 is the macro, and the program statements
 int array1 [Array_Size];
 int array2 [Array_Size];

The preprocessor would replace Array_Size in body of source code and result of macro processing is
 int array1 [50];
 int array2 [50];

(b) File inclusion :- It permits the user to include the header file in the program and user can make use of function defined in these header file.

when preprocessor encounters # include directive in the program text, it would immediately replace it by the entire content of specified file.

(c) Rational preprocessor :- These preprocessor enhances the capabilities of older languages by providing the programmer with built in macros for constructs that support flow of control and data structuring facilities. Some older languages did not have constructs such as while statements and if statement in them; these can be easily augmented through rational preprocessors.

(d) Language extension :-

Language extensions are preprocessing facilities that enhance the capabilities of the language through concepts such as built in macros. The language extension statements are translated into procedure calls or routines to perform desired operations, such as data base access or similar activities.

For example :- the language SQL is database query language embedded in C. Statements beginning ^{with} ## are taken by preprocessor to be database statements and statement related to C are translated into procedure calls or routines that perform database access.

Note :-

- (i) If program is too large, it may be stored in different files then the preprocessor will take of all these files.
- (ii) preprocessor stripped out all comments in a program.
- (iii) preprocessor is optional eg the language that donot support # symbol are not used preprocessor.

(ii) Compiler :- It takes modified source code as input and produces target code as output. Target code may or may not be assembly language code.

Note :- compilation is also optional eg HTML, JavaScript, SASS

(iii) Assembler :- Some compiler produces assembly language code that is passed to the assembler for further processing. Other compiler directly produces relocatable machine code that can directly passed to the linker. Assembly language code is mnemonics version of machine code in which names are used instead of binary codes, for

operations and names are also given for memory addresses.

(IV) Link editor and loader:-

The large source programs are compiled in small pieces by compiler. To run the target machine code of any source program successfully, there is a need to link the relocated machine language code with library files and other relocatable object files. So loader and linker program programs are used for link editing and loading of relocatable codes. Link editor creates a single program from several files of relocatable machine code. Loader reads the relocatable code and alters the relocatable addresses. To run the machine language program, the code with altered data and comments is placed at correct location in memory.

Compiler

A compiler is a computer program that translates source code written in high level language like C, C++, etc. into functionally equivalent target language code. The target language is usually assembly language or machine language. During the translation process, compiler also informs the programmer about the error present in the source program (see Fig. 1.3).

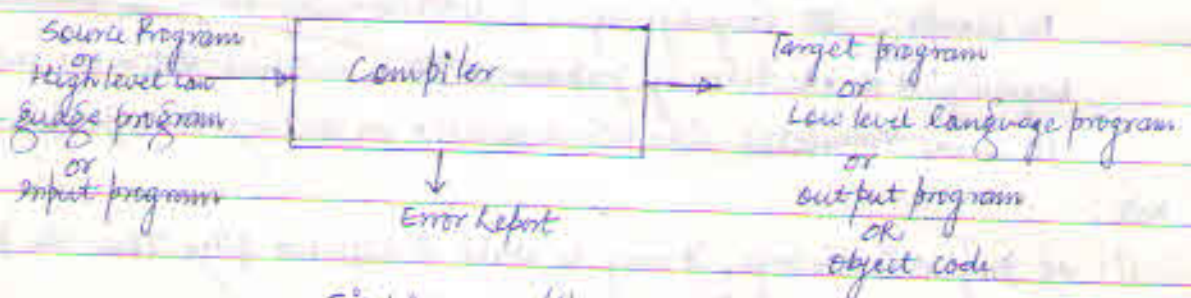


Fig 1.3 - Compiler

Ex: What is compiler? What is its primary function? What is its secondary function?

Compiler converts source program written in high level language into functionally equivalent machine or object or target code. This target code will be called by the user to process the input to produce the output.

The compiler executes the entire code at a time, if any error occurs in the code, all of them will be reported at a time.



Fig 1.4: compilation model

Interpreter:- In interpreter, data and source program are input to interpreter. Instead of producing any object code or executable code as in compiler, the interpreter produces the result by performing operations of the source program on its data. This model is represented as shown in Fig. 1.5



Fig. 1.5 Interpretation model

Interpreter directly executes the source program line by line according to the given inputs. The translation and execution of each statement are carried out side by side. There is no separate execution of the source program. The line by line execution of the source program provides better debugging environment (error finding) than a compiler because it can check for errors like out of bound array indexing at run time.

The main drawback of interpreter is that the execution time of interpreter program is generally slower than that of a compiled program because the program needs to be translated every time it is executed.

The interpreter has certain advantages over compiler^{as} it can handle certain language features which can not be compiled e.g. language like APL are normally interpreted as it involves features about the data such as size and shape of arrays which cannot be deduced at compile time. Another advantage is that interpreters can be made portable as they do not produce machine code programs. They also save the time required for assembling and linking a real program.

Difference between compiler and interpreter

S.No	Compiler	Interpreter
1.	It is translator.	It is not a translator. It is a simulator.
2.	In compilation process, run time and compile time are different.	In interpretation process, run time and compile time are same.
3.	The compiler reads input program and translates it only once.	The interpreter reads the input program over and over to compute results.

4- It produce object code (or .exe file)	4- It donot produce object ^{code} file. It simply runs the program.
5- The execution time of the compiler is less than that of interpreter because compiler reads the input program only once and occupies less space in memory.	5- The execution time of interpreter is more than that of the corresponding object program because in the interpretation process each line must be scanned and parsed prior to the execution and it occupies more memory.
6- Debugging is very slow	6- Debugging is very fast.
7- Size and complexity are more	7- Size and complexity are less.
8- It is used to designing stand alone applications	8- It is not used to design stand alone applications.
9- The compiler sees the sequence of source lines as printed not as executed	9- The interpreter sees that the statements in order of execution rather than as printed.
10- Compiler reads the whole source program code at once and report all error at last	10- Interpreter reads a statement from input execute it on data, if any error occurs an interpreter stop execution and report it.
Example - C/C++	Example: Command language, script language, batch language (LISP, PYTHON)

Part of compiler

As we know and illustrated in Fig 1.2, the compiler takes the preprocessed file as the input and translates it into an equivalent assembly language file or machine code file. here we will get an overview of how a compiler translates a preprocessed input file into a assembly code file or machine language file.

The translation of input source program file (preprocessed) into target assembly language file can be divided into two stages called as:

- ✓ Analysis of a source program (front end)
- ✓ synthesis of a source program (back end)

The analysis part or front end of compiler transforms the input source program into intermediate code. The intermediate code (sometimes called intermediate representation (IR)) is a machine independent representation of input source program.

The synthesis part of the compiler takes machine independent intermediate code and generates the target assembly language code. The synthesis part or back end deals with machine specific details like registers, number of allowable operators and so on. Fig 1.6 illustrates the two stage design approach of a compiler.



Fig 1.6 - The analysis - synthesis model model of compiler

The Main advantage of having two stage approach are:

- Take the front end of a compiler and attach different back end to produce a compiler for the same language on different machine (see Fig 1.7 (a))
- Take the common back end for different front ends to compile several different language on the same machine (see Fig 1.7 (b))

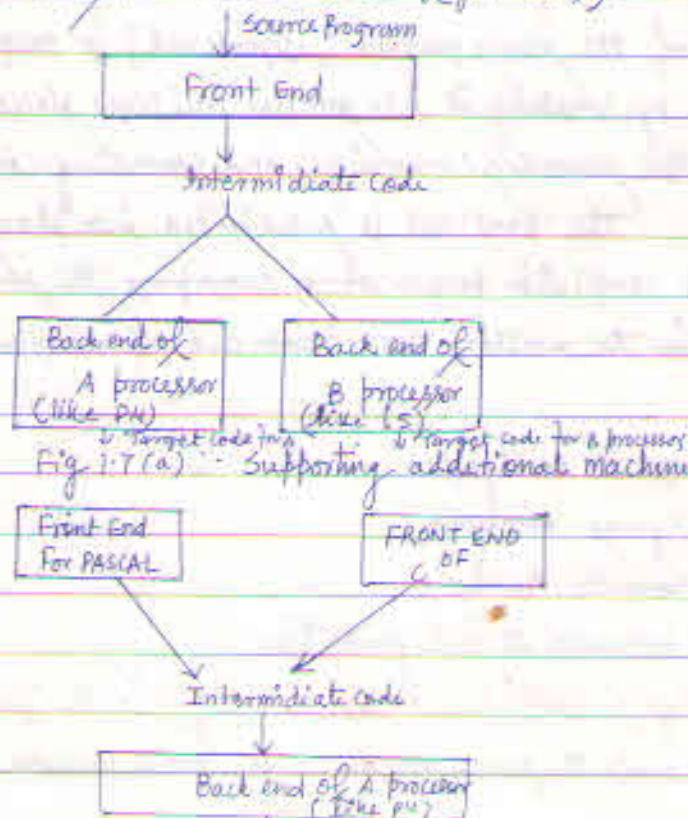


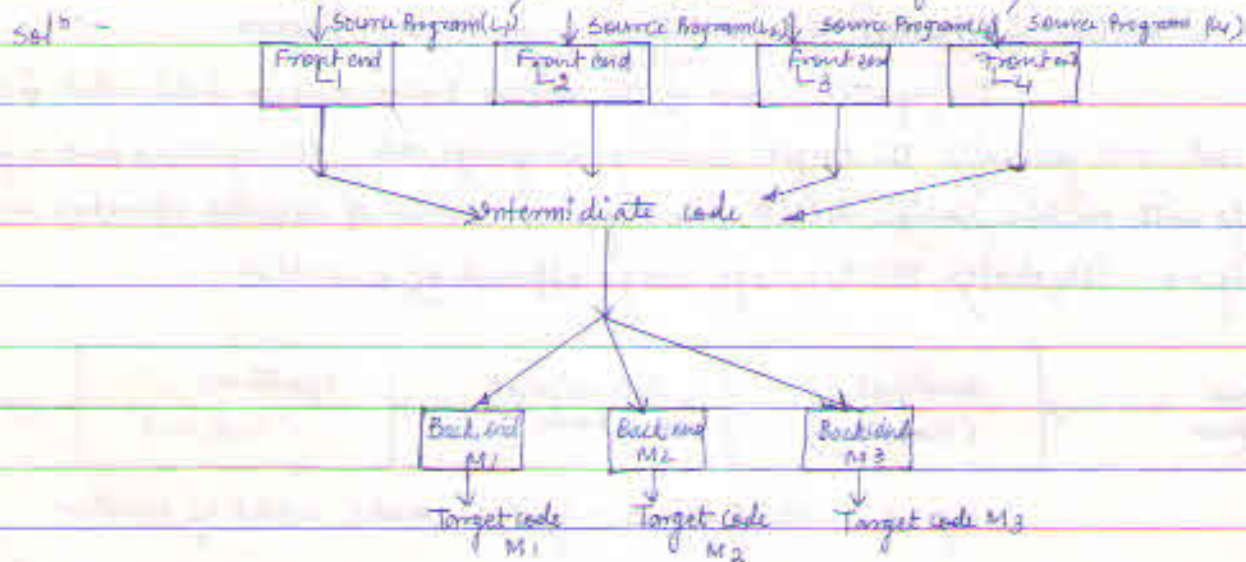
Fig 1.7 (a) - Supporting additional machine by adding back end.

Fig 1.7 (b) - Supporting an additional language by adding front end.

Date _____

This organization of compiler is called contemporary compiler.

Example :- Suppose we want to create compilers for four languages L_1, L_2, L_3, L_4 over three different machines M_1, M_2, M_3 . How many codes we have to write down?



So instead of producing $4 \times 3 = 12$ compiler code we will simply write down four code for front end of every language and three back end for each of machine. So total code required to construct the compiler for each language over every machine will be $4 + 3 = 7$.

(1) Analysis of source program :-

The analysis of the source program (Front end) is responsible for analysing the input source by breaking it into smaller entities (pieces) and checking syntax and verifying the semantics (meaning) and generating intermediate code.

The front end is subdivided into phases. A phase is an independent task in the compilation process which transforms the source program from one representation to another. The front end of compiler consists of the following phases

- ✓ Lexical Analysis
- ✓ Syntax Analysis
- ✓ Semantic Analysis
- ✓ Intermediate code generation

(i) Lexical analysis (also known as scanning) is the first phase of compiler. Lexical analyzer or scanner reads the source program in the form of character stream and groups the logically

related characters together that are known as lexemes. For each lexeme, a token is generated by the lexical analyzer.

A stream of tokens is generated as the output of the lexical analysis phase which acts as input for the syntax analysis phase. Tokens can be of different types, namely keywords, identifiers, constants, punctuation symbols, operator etc. The syntax of any token is $\langle \text{token_name}, \text{value} \rangle$

where token name is the name or symbol which is used during syntax analysis phase and a value is the location of that token in the symbol table.

4. Syntax analysis phase (parsing.)

Syntax analysis phase also known as parsing. Parser uses token name taken from from the token stream to generate the output in the form of tree like structure known as syntax or parse tree.

Syntax analysis checked, if they form a valid sequence as defined in programming language. So parse tree illustrate the grammatical structure of the token stream.

Note: - Due to parse tree arrangement of stream of tokens. This phase is also known as hierarchical analysis.

(3) **Semantic analysis** :- In semantic analysis, we check if the syntactically correct statements make a meaningful reading. For example, a statement in the input source program " $x = y + 2$ " would not make a meaningful reading if say x is the name of function or array and y is a float type of identifier. This statement might be syntactically acceptable by the productions of the context free grammar in syntax analysis but would not hold out during semantic analysis because the data types of x and y are not compatible.

In natural language parlance, this is very similar to having a grammatically correct sentence, but devoid of meaning. For example, the syntax rule would also accept a sentence, "The bicycle rides the boy". This sentence does not make sense so we would reject it during semantic analysis.

In similar way a c statement " $\text{myfun} = y$ " where " myfun " is the name of function and y , a float type identifier might be acceptable in syntax analysis but would be rejected in semantic analysis, since the data types of left hand

side and right hand side do not match.

Most of semantic analysis revolves around such type checking. Other tasks of semantic analysis involves detection of undeclared variable, access violations (classes accessed correctly or not) and so on.

The result of semantic analysis is annotating the parse tree with more information on the data types.

4- Intermediate code generation:-

In intermediate code generation phase, the parse tree representation of the source code is converted into low level machine like representation (easy to generate from source program and easy to convert (translate) into machine language) ✓

In intermediate code generation phase, we walk through the annotated parse tree and generate intermediate code. The three address (or postfix) is one of the common forms of intermediate code. Three address code is a sequence of instructions each of which can have at most three operands.

(ii) Synthesis of source program (Back end): +

The synthesis of source program (Back end) of compiler is responsible for translating the machine independent intermediate code into target code (either machine code or assembly code). The back end of the compiler depends on target processor, where final binary would be executed.

The back end or synthesis of part of the compiler consists of the following phases.

- ✓ code optimization
- × code generation

(5) Code optimization:-

Code optimization phase, which is an optional phase, performs the optimization of the intermediate code; optimization means making the code shorter and less complex, so that it can execute faster and takes lesser space, so that it can execute faster and takes less space. The output of code optimization phase is also an intermediate code, which performs the same tasks as the input code but requires lesser space and time.

(6) Code Generation:- Code generation phase translates the intermediate code representation of the source program into the target program (machine or assembly code). If the target

program is machine language code, the code generator produces the target code by assigning registers, and memory locations to store variables defined in the program and to hold the intermediate computation results. The machine code produced by code generator can be directly executed on the machine.

In case of assembly code, mnemonics code is generated which is translated by assembler and other utility in the machine language code further.

Besides these, there are two other phases which interact with all the other phases of compiler. These are:

- ✓ Symbol table Management
- ✓ Error detection and handling

(a) Symbol table Management :-

A symbol table is a data structure that is used by compiler to record and collect information about source program constructs like identifiers (variables) and subroutines (procedure) in the source program.

For identifiers (variables), in addition to its name, table contains information about the type, scope and storage space of identifiers.

For subroutines (procedures) it records the number and types of its parameters, the way each parameter passed and return type of the subroutines (if there is any).

Symbol table management is an important function of each phase of the compiler. During lexical analysis phase, lexical analyzer identifies tokens and put them into symbol table. About type and other information or any other information is entered by syntax and semantic analysis phase.

For ex:-

For following code segment in C, the symbol table will be given below
`int a;` // now it will read each identifier & entry will be
`float b, c;` done in symbol Table.

`b = a + b;` // If already entered so it will not be reentered

Entry No	Symbol	Type	Length	...	Address
1	a	Integer	2		
2	b	Float	4		
3	c	Float	4		
4	a*	Float	4		
5	T ₁	Float	4		

Compiler
type cast variable 'a'
to a* in float type
from integer

(b) Error detection & handling

Date _____

Programs are written by programmer (human) and hence cannot be free from errors. Each phase encounters errors. The lexical analyzer detects all errors when the remaining characters in the input do not form any token. The syntax analyzer usually detects a large number of errors where the token stream violates the structural rules of language. The type mismatch type errors are usually detected by semantic analyzer (eg. addition of character variable to an integer variable).

Whenever, a phase encounters an error, it must report it, in addition, it must somehow deal with that error. Error handler is invoked. Error handler generates a suitable message allowing programmer to find out exact location of error.

The analysis and synthesis model of compiler in Fig 1.8, consolidating the phases of front-end and back-end of the compiler.

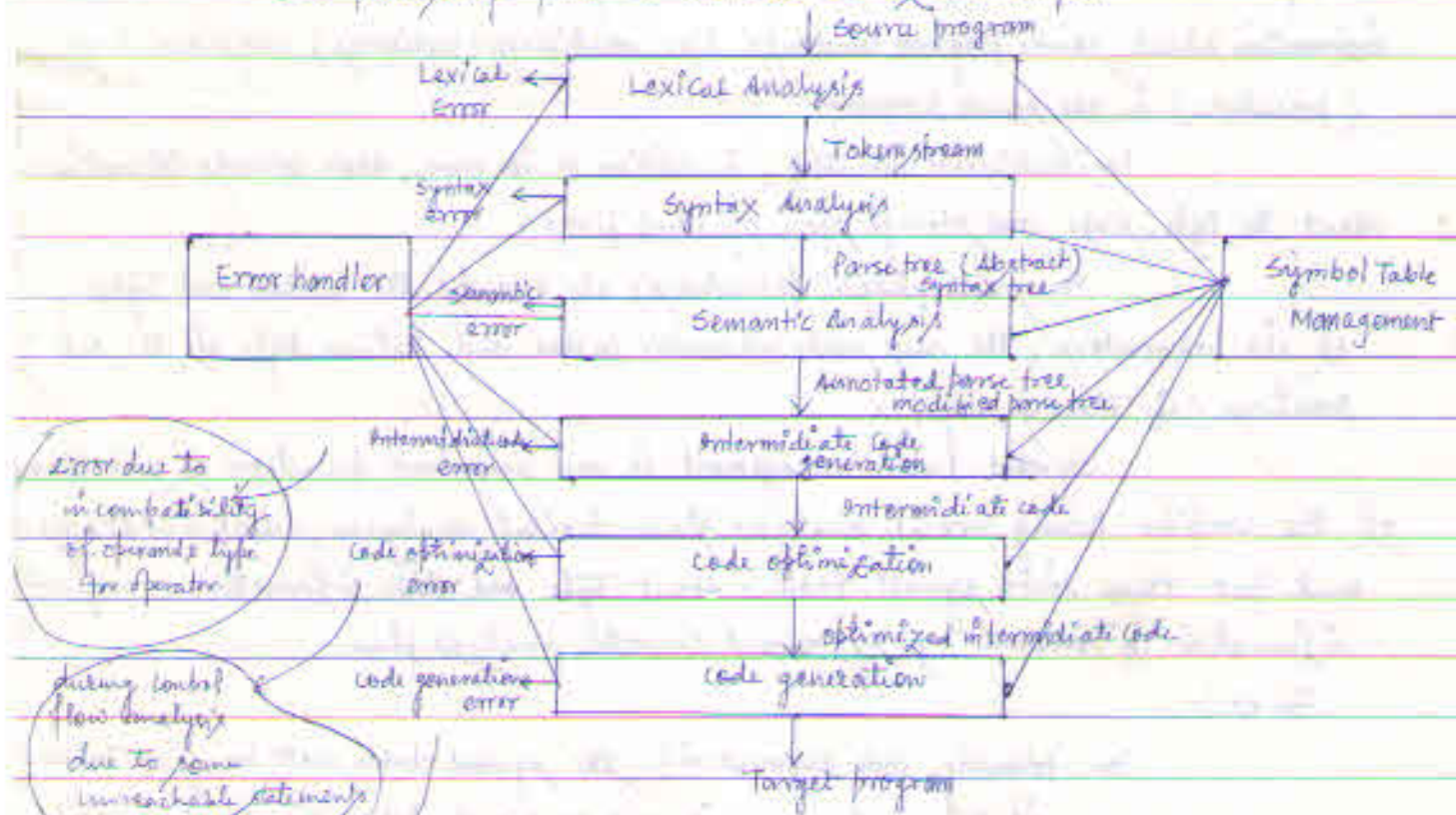


Fig 1.8: Phases of compiler (analysis & synthesis model)

Example of compilation :-

Let us consider the translation of the following statement

$$x = y + z * 10 \quad \text{--- (1.1)}$$

The internal representation of source program changes with each phase of the compiler.

(1) Lexical Analysis phase

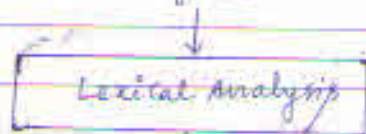
The lexical analyzer builds uniform descriptors for the elementary constituents of the source string. To do this it must identify lexical units in the source string and categorize them into identifiers, constants, keywords, operators etc. This uniform descriptor is called token. A token has the following format

category	lexical value
----------	---------------

For example, when lexical analyzer finds an identifier x , it generates a token, say id and also enters ' x ' into the symbol table, if it is not already there. The token for x will be $\langle id, \#x \rangle$, if x is x^{th} entry in symbol table. The lexical value associated with this occurrence of id contains a pointer (or index) to the symbol table entry for x .

A token is generated / constructed for each identifier / constant as well as each operator in the string. Let it assign id_1 , id_2 , and id_3 to x , y , z respectively, and assign op to '=', $multop$ to '*', $addop$ to '+' and num to 10. New representation of statement (1) after lexical analysis may be given as

$x = y + z * 10$

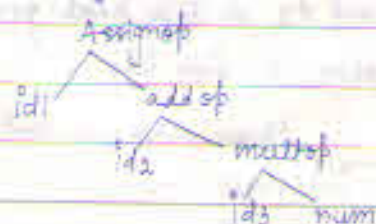
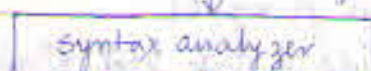


id_1 assignop id_2 addop id_3 multop num (2)

(2) Syntax analysis phase

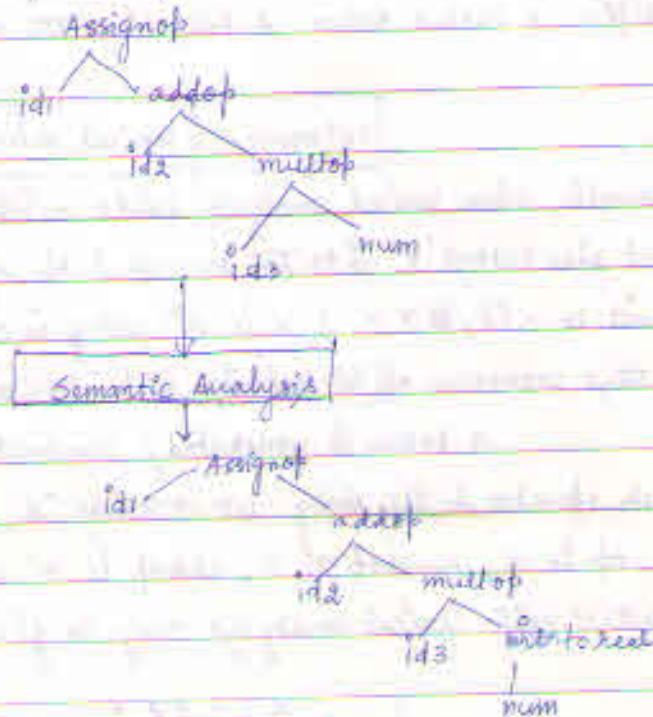
Stream of tokens acts as input for the syntax analyzer. Output of the syntax analyzer is a parse (syntax) tree that acts as input to semantic analyzer. The parse tree for statement (2) is shown in fig.

id_1 assignop id_2 addop id_3 multop num



(iii) Semantic analysis phase.

The semantic analyzer determines the meaning of a source string. It finds out that the types of operands in fig are not identical as id3 is real (float) and num is identifier. Therefore, it converts 'num' to type real so that the syntax tree now looks as in figure.



(iv) Intermediate code generation:-

After semantic analysis, the compiler generates an explicit intermediate code representation of source program. This intermediate code could be three address code or postfix (polish) code. But it should be easy to produce and easy to translate into target program.

For syntax tree shown in above figure the following steps must be taken in the intermediate code generation.

- 1- Convert num to real giving t1.
- 2- multiply t1 and id3 in type real giving t2.
- 3- Add id2 and t2 in type real giving t3.
- 4- Finally store t3 into id1.

These steps shown in three address code as follows.

```

t1 = int to real (num)
t2 = id3 * t1
t3 = id2 + t2
id1 = t3
  
```


(v) code optimization phase :-

The code optimizer tries to improve the intermediate code in order to achieve faster running machine code. In our example we can convert num from integer to real once at compile time (no need of conversion at execution time). Also we can substitute $id1$ for $t3$ as $t3$ is used only once to transmit its value to $id1$. Now intermediate code (of previous one) can be rewritten as

$$t_2 = id3 * rnum$$

$$id1 = id2 + t_2$$

(vi) code generation phase

Finally, compiler can generate the target code for the intermediate code. The storage must be allocated and registers must be assigned to the variables and addressing modes to be used for accessing the data must be decided before generating code.

If we use two registers R_1 and R_2 , the code can be written as

MOVE $R_2, id3$

MUL $R_2, rnum$

MOVE $R_1, id2$

ADD R_1, R_2

MOVE $id1, R_1$

The first operand in each instruction specifies destination, 'x' in each instruction indicates destination - real value

Questions :-

Illustrate the compilation process on the following statements

(1) $x = y * z + 10$

(2) $I = (P * X * X * t) / 100$

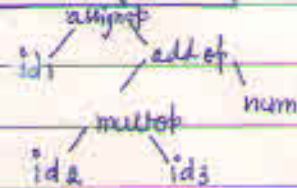
Solⁿ:- Here we illustrate compilation process of each statement

$x = y * z + 10$

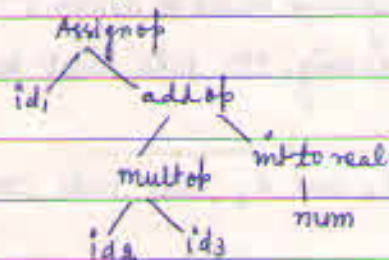
Lexical analysis

id1 assignop id2 multop id3 addop num

Syntax Analysis



Semantic Analysis



Intermediate code generation

 $t_1 = id_2 * id_3$
 $t_2 = \text{int to real}(num)$
 $t_3 = t_1 + t_2$
 $id_1 = t_3$

code optimization

 $t_1 = id_2 * id_3$
 $id_1 = t_1 + rnum$

code generation

```

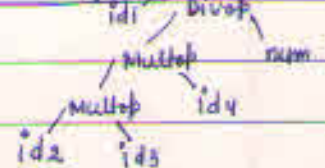
MOVE R2, id3;
MULTF R2, id2;
MOVE R0, rnum;
ADDF R0, R2;
MOVE id1, R0;
  
```

 $I = (p * r * t) / 100$

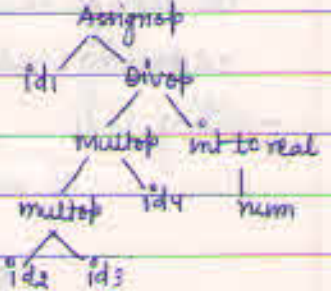
Lexical Analysis

id1 assignop id2 multop id3 multop id4 divop num

Syntax analysis



semantic Analysis



Intermediate code generation

 $t_1 = id_2 * id_3$
 $t_2 = t_1 * id_4$
 $t_3 = \text{int to real}(num)$
 $t_4 = t_2 / t_3$
 $id_1 = t_4$

code optimization

 $t_1 = id_2 * id_3$
 $t_2 = t_1 * id_4$
 $id_1 = t_2 / rnum$

code generation

```

MOVE R1, id2;
MULF R1, id3;
MOVE R2, id4;
MULF R2, R1;
MOVE R2, R2;
DIVF R2, rnum;
MOVE id1, R2;
  
```