

## Unit - 1

#) **Translator**

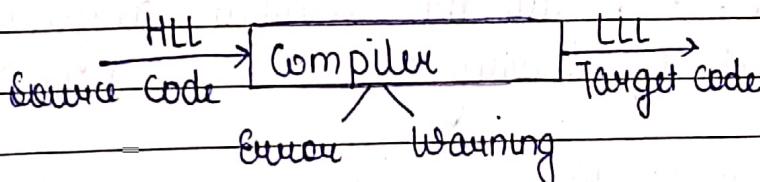
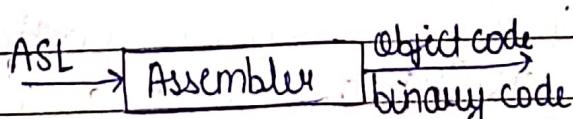
A software system which convert the source code from one form of language to another form of language is called translator.

Types of translator :-

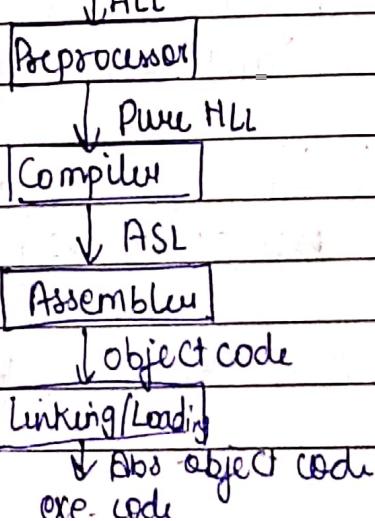
There are two types of translator :-

1. Compiler

2. Assembler

#) **Compiler** - A translation which convert the source code from high level language to low level language is called compiler.#) **Assembler** - A assembler is a software which converts assembly code into object code or binary code is called assembler.

Functions of language processing system (LPS) :-



Object code

Relocatable      Absolute object code.

- i) Preprocessor:- It includes all header files and also evaluates if any MACRO is included. The preprocessor is also called as MACRO evaluator. Preprocessing is optional i.e. if any language which does not support #include and MACRO's preprocessor is not required.
- ii) Compiler:- It takes the preprocessor as input and converts into assembly code (low level lang.)
- iii) Assembler:- It converts the assembly language into object code or binary code or machine code
- iv) Linking and Loading provides your functions
  - a) Allocation
  - b) Relocation
  - c) Linker
  - d) Loader

Allocation - Getting the memory positions from operating system and storing the object data on object code.

Relocation - Mapping the relative address to the physical address and relocating the object code.

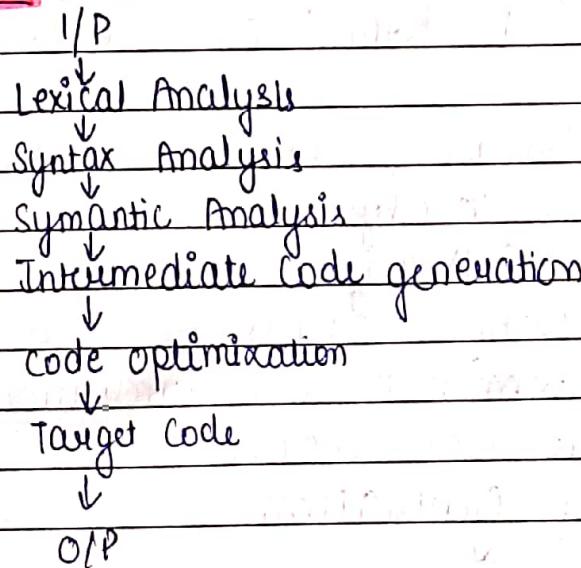
Linker - Combines all the executable object modules to one single executable file.

Loader - Loading the executable file into permanent

## #) Difference Between Compiler & Interpreter

<u>Compiler</u>	<u>Interpreter</u>
1) A compiler translates the entire source code in one run	1) An interpreter translates the entire source code one line by line.
2) It consumes less time.	2) It consumes more time.
3) It is more efficient	3) It is less efficient
4) CPU utilization is more	4) CPU utilization is less as compared to compiler
5) The Compiler is larger	5) The Interpreter are smaller than compiler.
6) It is not flexible	6) It is flexible
7) The localization of errors is difficult	7) The localization of error is easier then the compiler
8) Both syntactic & semantic errors can be checked simultaneously	8) Only syntactic errors are checked.

## #) Design of Compiler :-



## Phases of compiler :-

Phases - Converting the source code from one form of representation to another form of representation is called phase.

The phases can be divided into two parts:-

- i) Analysis phase
- ii) Synthesis phase

Analysis phase - It breaks the source program into constituent pieces and imposes a grammatical structure on them which further uses the structure to create an intermediate representation of the source program.

- It is also termed as front end of compiler.
- Information about the source program is collected and stored in a data structure called symbol table.

## Synthesis Phase -

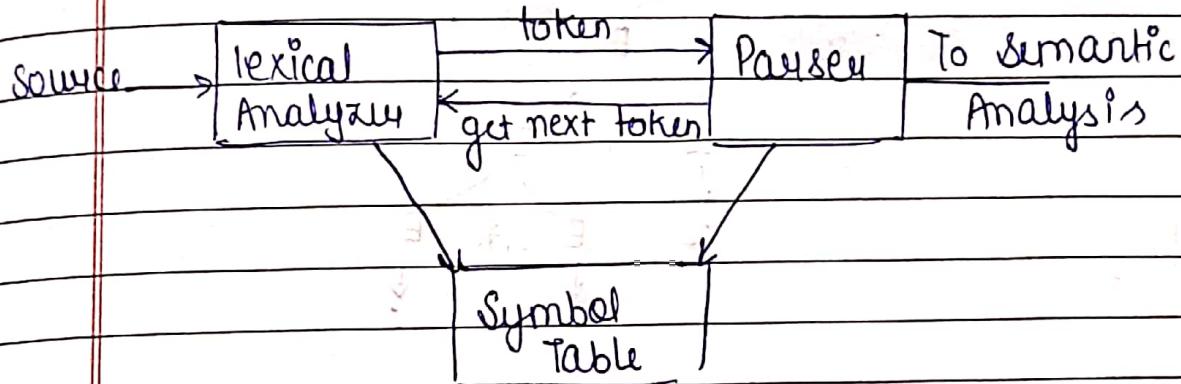
- Synthesis part takes the intermediate representation as input and transforms it to the target program.
- It is also termed as back end of compiler.

The design of compiler can be decomposed into :-

- 1) Lexical Analysis
- 2) Syntax Analysis
- 3) Semantic Analysis
- 4) Intermediate Code generation
- 5) Code optimization
- 6) Code Generation

## i) Lexical Analysis :-

- Lexical Analysis is the first phase of compiler which is also termed as scanning.
- Reads the input characters and those characters are grouped to form a sequence called lexem, which produces token as output uses for syntax analysis.
- Strips out source program comments & white space in form of blank, tab.



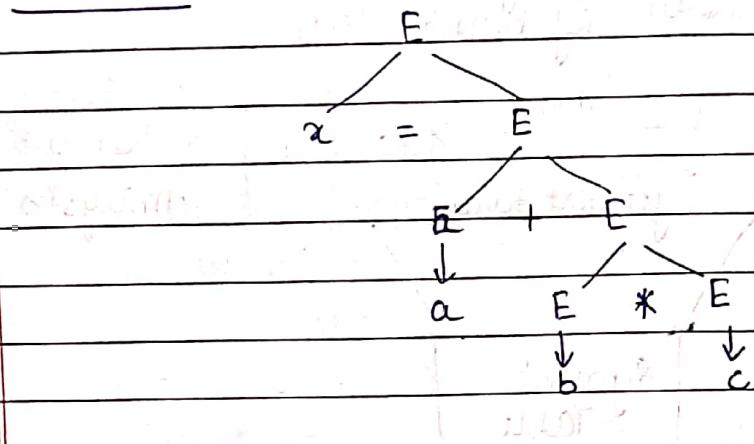
- Token: Token is a sequence of characters that represent lexical unit which matches with the pattern, such as keywords, operators, identifiers etc.
- Lexeme: It is a sequence of characters in the source program that matches the pattern for a token.
- Pattern: Pattern describes the rule that the lexemes of a token take. It is structure that must be matched by strings.
- Once a token is generated the corresponding entry is made in the symbol table.

## 2) Syntax Analysis :-

Syntax Analyzer verify the grammatical mistake of the source code. To verify the syntax of the code the language must be defined by CFCs. Syntax Analyzer take the stream of tokens as I/P & generates the parse tree.

$$eg \rightarrow x = a + b * c$$

Parse Tree :-



## 3) Syntactic Analysis :-

Syntactic Analysis verify the meaning of each and every sentence by performing type checking. Syntax Analyzer just verifies whether the operator is operating on required no. of operands or not and does not looking into working.

## 4) Intermediate Code Generator :-

The source code is converted into intermediate representation to make code generation process simple & easy and to achieve the platform independently.

Ex →

$$\begin{aligned}
 &x = a + b * c \\
 &t_1 = b * c \\
 &t_2 = a + t_1 \\
 &x = t_2
 \end{aligned}
 \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Intermediate code}$$

### 5) Code optimization

Reducing the number of instructions without affecting the outcome of the source program is called optimization.

It is of two types :-

- (i) Machine independent optimization
- (ii) Machine dependent optimization

### 6) Target Code

The optimized source code should be converted into assembly code.

Ex:-

$x = a + b * c$		
MUL R <sub>1</sub> , R <sub>2</sub> // $b * c$	R <sub>0</sub> → a	
ADD R <sub>0</sub> , R <sub>1</sub> // $a + b * c$	R <sub>1</sub> → b	
MOVE R <sub>2</sub> , x // $x = a + b * c$	R <sub>2</sub> → c	

Ques

why is Intermediate Code Generation required?

Ans

If we generate machine code directly from source code, we may face 2 problems

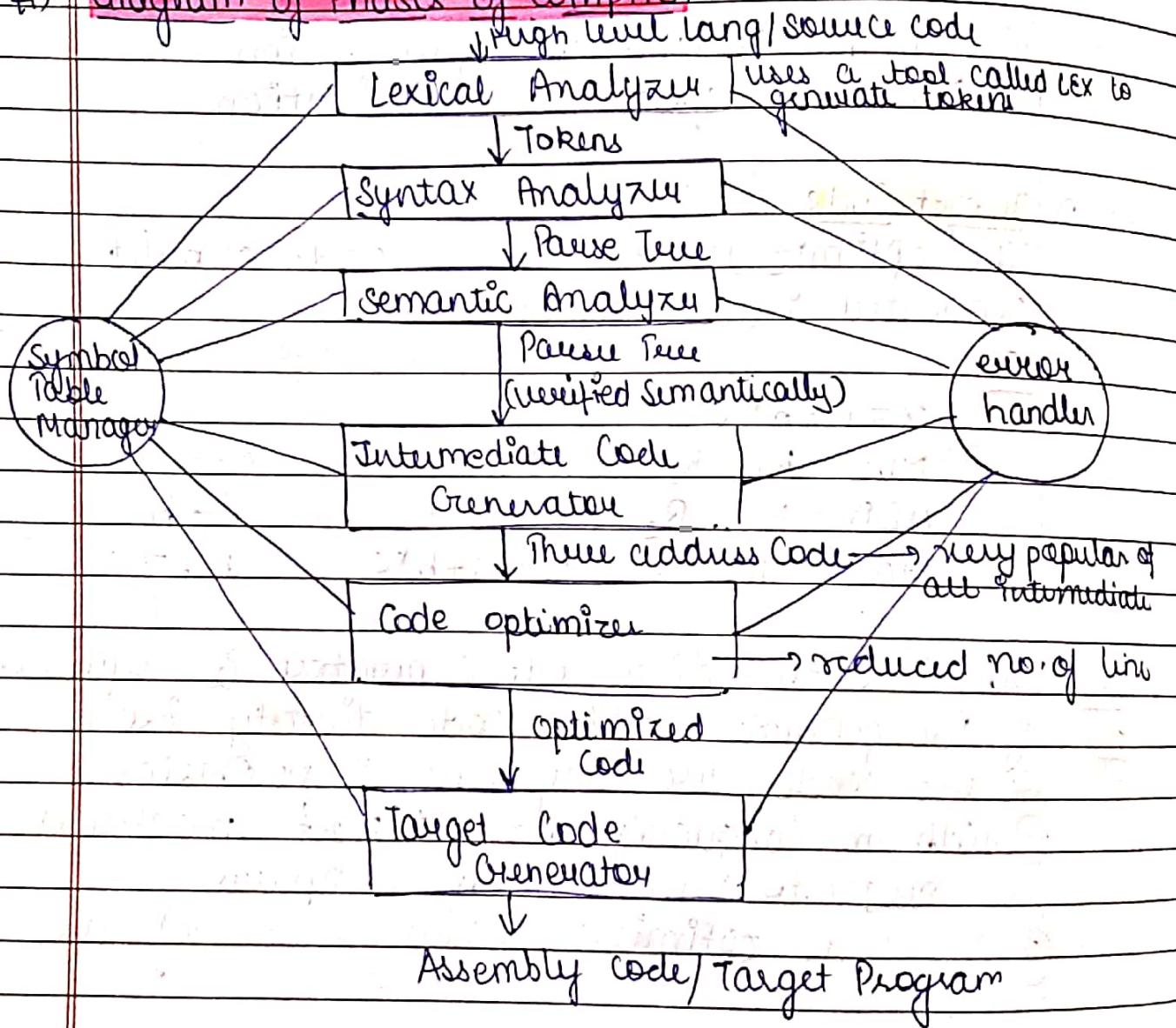
(1) With m languages and n target machines, we need to write  $m \times n$  compilers.

(2) The Code optimizer which is one of the

longest and very difficult to write components of the compiler can't be reused.

By converting source code to an intermediate code, a machine independent optimization can be done. An intermediate code must be easy to produce and easy to translate to machine code, shouldn't contain any machine specific codes in a sort of universal assembly language.

## #) Diagram of Phases of Compiler



#) Lexical Analysis

A Lexical Analyzer takes a high level programming language as input in the form of sequence of characters and produce a sequence of tokens by removing any white space chars and comment line.

So, the LA produce errors for illegal symbols. In such cases, it skips character in the input until a wellformed token is found.

Q. what is Token?

A. Token is a sequence of characters that represent lexical unit which matches with the pattern. Ex- float, identifier, =, +, ;

Q. what is pattern?

A. The set of strings for which some token is produced ex- float, =, +, ;

Q. what is Lexeme?

A. The sequence of characters match by pattern to form the corresponding token.  
Ex = "float", "silicon", "=", "+", ";"

\* Each token consisting of a token name & an optional attribute value.

\* Regular definition (modified version of regular expressions) are used for specification of

tokens & from these definitions, IA can be generated automatically using LEX.

### functions of lexical analyzer

- 1) Read input characters from source program
- 2) Dividing the programs into valid tokens (tokenization)
- 3) Remove white spaces characters
- 4) Remove comments (/\*, \*/) from the program
- 5) Helps to identify tokens into the symbol table
- 6) It generates lexical errors

### Lexical Errors:

A character sequence which is not possible to scan any valid token is called a lexical error.

Lexical phase error can be:-

- Spelling error
- Exceeding length of identifier
- Appearance of illegal character

e.g. -  $x = 1x ab$

\* Lexical phase error is found during the execution of the program

### Advantages

- Lexical Analyzer method is used by programs like compiler which can use the parsed data from
- It is used by web browser to format and display a web page with the help of parsed

- data from javascript, HTML, CSS.
- It helps to construct a specialized and potentially more efficient processor for a task

### Disadvantages of Lexical Analyzer

- You need to spend significant time reading the source program and partitioning it in the form of tokens.
- Some Regular Expression are quite difficult to understand compared to PEG or EBNF rules.
- More effort is needed to develop & debug the lexer and its token descriptions.

### 1) Specifications of Tokens

These are the types of specifications of tokens :-

i) Alphabets - Any finite set of symbols  $\{0, 1\}$  is a set of binary alphabets  $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F\}$  is a set of hexa decimal alphabets  $\{a-z, A-Z\}$  is a set of English language alphabets.

ii) Strings - Any finite sequence of alphabets is called a strings. Length of the string is the total no. of occurrence of alphabets e.g. the length of the string is the total number of occurrence of alphabet.

3) Language - A language is considered as a finite set of alphabets. Computer languages are considered as finite sets and mathematically set operations can be performed on them. Finite language can be described by means of regular expressions.

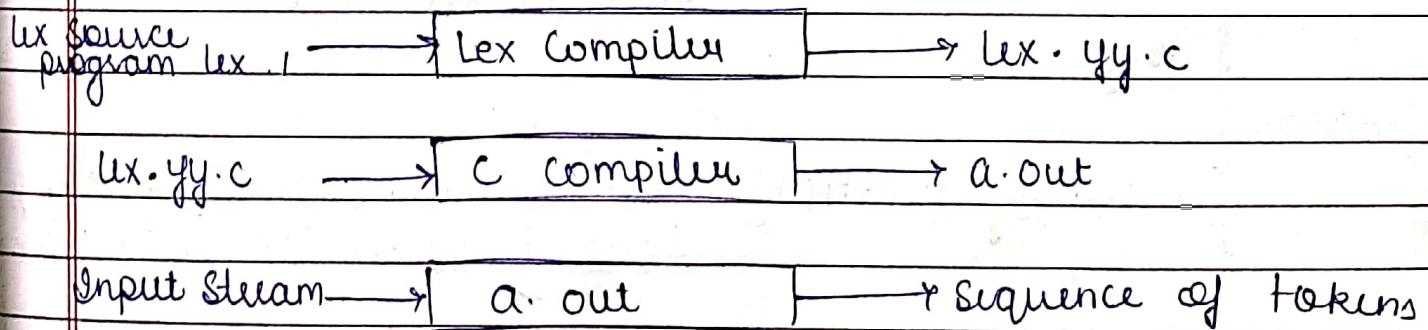
Operations on languages :-

- Union of two languages  $L \cup M$  is written as  $L \cup M = \{s | s \text{ is in } L \text{ or } s \text{ is in } M\}$
- Concatenation of two languages  $L \& M$  is written as  $L \cdot M = \{st | s \text{ is in } L \text{ and } t \text{ is in } M\}$
- The Kleene Closure of a language  $L$  is written as  $L^* = \text{zero or more occurrence of language } L$

4) Regular Expression - Regular expression have the capability to express finite languages by defining a pattern for finite strings of symbols. The grammar defined by regular expression is known as regular grammar. The languages defined by regular grammar is known as Regular language. Regular languages are easy to understand and have efficient implementation.

## #) Lexical Analysis Generator (LEX)

- LEX is a tool / computer program which generates lexical Analyzers. It is used by YACC parser generator.
- LEX written by Mike Lesk and Eric Schmidt and describe in 1975
- Its main purpose is to break up an input stream into tokens (sequence of tokens)
- It reads the input stream and produces the source code as output through implementing the lexical analyzer
- It is a tool for automatically generating a lexer.



- Lex can be used with a parser generator to perform lexical analysis.

Function of Lex through diagram :-

- I) Firstly lexical analyzer creates a program lex.l in the lex language. Then lex compiler runs the lex.l program and produces a c program lex.yy.c

- 2) finally C compiler runs the lex.yy.c program and produces an object program a.out  
3) a.out is a lexical analyzer that ~~program~~ transforms an input stream into a sequence of tokens.

### Structure of LEX :-

A lex program is separated into three sections by % % delimiters.

{declaration}  $\Rightarrow$  declaration of variable, constant  
% %

{translation rules}  $\Rightarrow$  have the pattern function  
contains rules in the form of  
regular expressions  
% %

{auxiliary functions}  $\Rightarrow$  can be compiled separately &  
loaded with lex program

## #) Symbol Table Management.

Symbol Table is an important data structure created and maintained by compiler in order to store information about the occurrence of various entities such as variable names, function names, objects, classes,

Symbol table is used by both analysis & synthesis parts of a compiler. It is used for scope management.

### Purpose :-

- To store the names of all entities in a structured form at one place.
- To verify if a variable has been declared.
- To implement type checking, by verifying assignment and expression in the source code are semantically correct.
- To determine the scope of a name (scope resolution).

### #) Finite Automata

- Finite Automata are used to recognize patterns.
- It takes string of symbol as input and changes its state accordingly. When the desired symbol is found, then the transition occurs.
- At the time of transition, the automata either move to the next state or stay in the same state.
- Finite Automata has two states, Accept State or Reject State.
- Finite State Automata are used to design L.A.
- F.S. Automata are machine that accepts regular languages.

A Finite Automata is a collection of five-tuple  $(Q, \Sigma, \delta, q_0, F)$  where :-

$Q$  is finite set of states

$\Sigma$  is finite set of the input symbol

$q_0$  is initial state

$F$  is final state

$\delta$  is transition function

It is of two types :-

- 1) DFA
- 2) NFA

### 1) DFA (Deterministic Finite Automata)

- DFA refers to when the machine is reading an input string one symbol at a time.
- In DFA, there is only one path for specific input from the current state to next state.
- DFA does not accept the null move.
- DFA can contain multiple final states. It is used in lexical Analysis in Computer.

It consists of 5 tuples  $(Q, \Sigma, q_0, F, \delta)$ .

Transition function can be defined as :-

$$\delta: Q \times \Sigma \rightarrow Q$$

### 2) NFA (Non-Deterministic Finite Automata)

- It is easy to construct an NFA than DFA for a given Regular language.
- The finite automata are called NFA when there exist many paths for specific input.

from current state to next state

- Every NFA is not DFA but each NFA can be translated into DFA.

- It accepts Null move

- It contains multiple next states.

Transition Function  $\rightarrow S : Q \times \Sigma \rightarrow 2^Q$

& it contains five tuples  $(Q, \Sigma, q_0, F, \delta)$ .

## #) Three Address Code

Three Address Code is a simple, intermediate code to construct and convert to machine code. It employs just three places and one operator to describe an expression and the value computed at each instruction is saved in a compiler-generated temporary variable.

- It is used form by the optimizing compilers.
- The given expression is broken down into several separate instructions. These instructions can easily translate into assembly language.
- It is a combination of assignment and a binary operator.
- The three address code has -three operands

General Representation :-

$$a = b \text{ op } c$$

operator

a, b, op c are operands

## ① Quadruple -

It is structure with consist of 4 fields namely  
op, aug1, aug2, and result  
op denotes operator  
aug1, aug2 denotes operands  
result is to store the result of expression

	operator
	aug1
	aug2
	result

### Advantages

- 1) Easy to rearrange code for global optimization
- 2) One can quickly access value of temporary variables using symbol table

### Disadvantages

- 1) Contain lot of temporaries
- 2) Temporary variable creation increases time and space complexity

Ex →

$$a = b^* - c + b^* - e$$

The three address code :-

$$t_1 = \text{uminus } c$$

$$t_2 = b^* t_1$$

$$t_3 = \text{uminus } c$$

$$t_4 = b^* t_3$$

$$t_5 = t_2 + t_4$$

$$a = t_5$$

#	op	Aug 1	Aug 2	Result
0	unimux	c		$t_1$
1	*	$t_1$	b	$t_2$
2	unimux	c		$t_3$
3	*	$t_3$	b	$t_5$
4	+	$t_2$	$t_4$	$t_5$
5	=	$t_5$		a

### Quadtuple Representation

#### ② Tuples -

This representation doesn't make use of extra temporary variable to represent a single operation instead when a reference to another tuple's value is needed, a pointer to that tuple is used. So, it consists of three fields namely op, aug1, aug2.

#### Disadvantages -

- 1) Temporaries are implicit and difficult to reuse in code.
- 2) It is difficult to optimize because optimization involves moving intermediate code. When a tuple is moved, any other tuple referring to it must be updated also.

Ex→

$$a = b^* - c + b^* - c$$

#	op	Bug 1	Bug 2
0	uminus	c	
1	*	0	b
2	uminus	c	
3	*	2	b
4	+	1	3
5	=	a	4

## #) Difference Between Static and Dynamic Memory Allocation

### Static Allocation

- 1) When the allocation of memory performs at the compile time, then it is known as static memory.
- 2) The memory is allocated at the compile time.
- 3) While execution, the memory cannot be changed.
- 4) It saves running time as it is fast.
- 5) It allocates memory from the stack.

### Dynamic Allocation

- 1) When the memory allocation is done at the execution or run time, then it is called dynamic memory allocation.
- 2) The memory is allocated at the runtime.
- 3) While execution, the memory can be changed.
- 4) It is slower than static memory allocation.
- 5) It allocates memory from the heap.

- |  |  |
|--|--|
| 6) It is less efficient  | 6) It is more efficient                                    |
| 7) The memory allocation is simple   | 7) The memory allocation is complicated.                   |
| 8) It is preferred in array  | 8) It is preferred in linked list.                         |
| 9) Once the memory is allocated, it will remain from the beginning to end of program | 9) The memory can be allocated at any time in the program. |

## ⇒ Optimizing Codes

The code optimization in the synthesis phase is a program transformation technique which tries to improve the intermediate code by making it consume fewer resources so that faster running machine code will result.

### Code optimization objectives :-

- The optimization must be correct, it must not, in any way change the meaning of the program.
- Optimization should increase the speed and performance of the program.
- The compilation time must be kept reasonable.
- The optimization process should not delay the overall compiling process.

## #) Machine Independent Optimization :-

- Machine Independent optimization attempts to improve the intermediate code to get a better target code. The part of the code which is transformed here does not involve any absolute memory location or any CPU register.
- The process of intermediate code generation much inefficiency like : using variable instead of constants, extra copies of variable, repeated evaluation of expression.
- Through code optimization, you can remove such efficiencies and improves code.
- It can change the structure of program sometimes of beyond recognition like : unrolls loops, inline functions, eliminates some variables that are programmer defined.

## #) Machine Dependent Optimization :

Machine dependent optimization is done after the target code has been generated and when the code is transformed according to the target machine architecture. It involves CPU registers and may have absolute memory rather than relative references. Machine-dependent put efforts to maximum advantage of the memory hierarchy.

## #) Peephole optimization

This optimization technique works locally on the source code to transform it into an optimized code. By locally, it means a small portion of the code block at hand. These methods can be applied on intermediate code as well as target code. A bunch of statement is analyzed and are checked for following optimization :-

### 1) Ridundant Instruction elimination -

At the compilation level, the compiler searches for instructions redundant in nature. Multiple Loading and Storing of instructions may carry the same ~~repeated~~ meaning even some of are removed.

### 2) Unreachable Code

Unreachable Code is a part of the program code that is never ~~selected~~ accessed because of programming constructs. Programmers may have accidentally written a piece of code that can never be reached.

### 3) Flow of Control Optimization

There are instances in code where the program control jumps back and forth without performing any significant task. These jumps can be removed.

#### 4) Algebraic Expression Simplification

There are occasions where algebraic expression can be made simple.

#### 5) Strength Reduction

There are operations that consume more space and time. Their strength can be reduced by replacing them with other operations that consume less time and space, but produce same result.

#### 6) Accessing Machine Instructions

The target machine can deploy more sophisticated instructions which can have the capability to perform specific operations much efficiently. If the target code can accommodate these instructions directly, that will only improve the quality of code but also yield more efficient results.