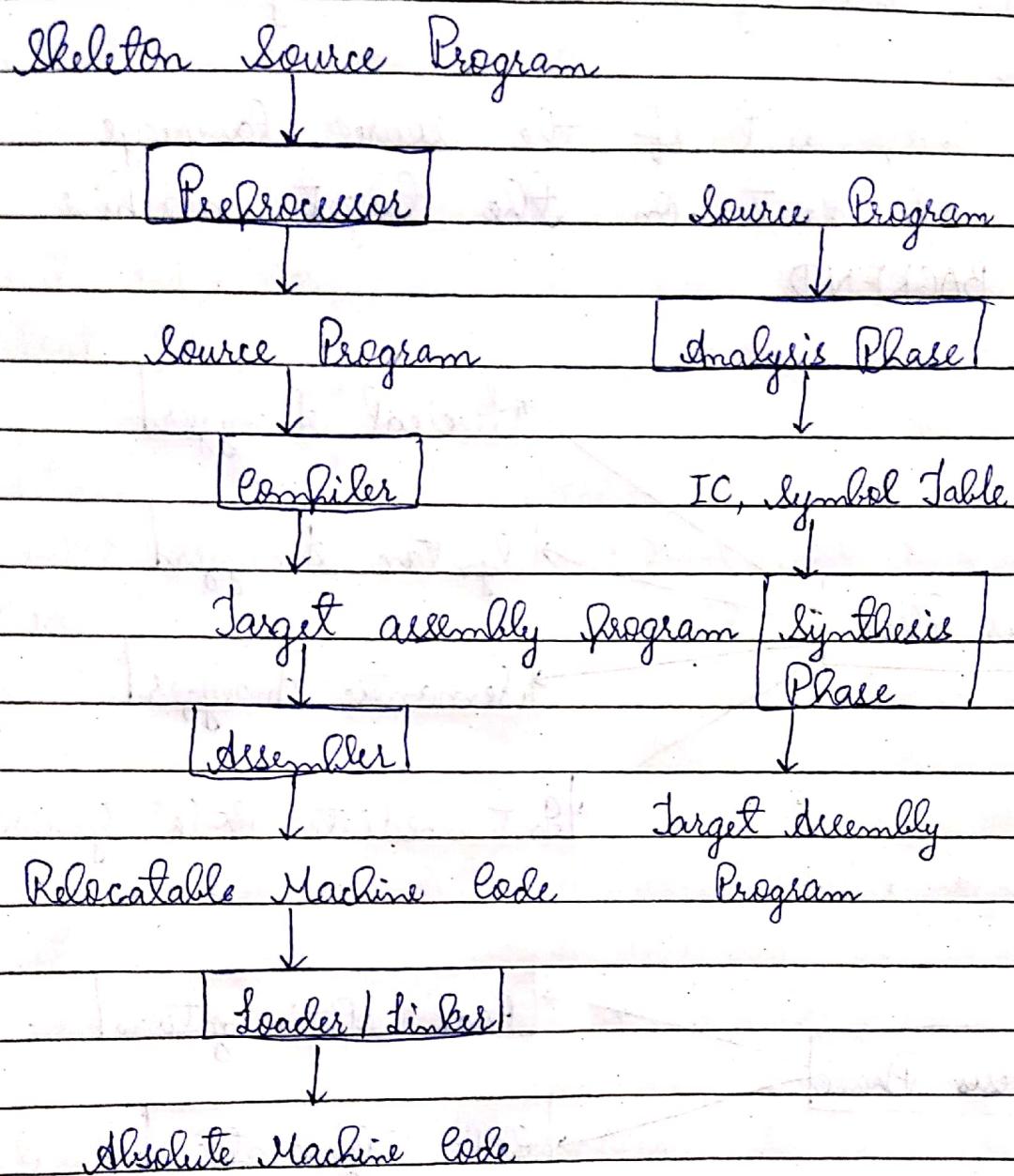


## \* Phases of Compiler

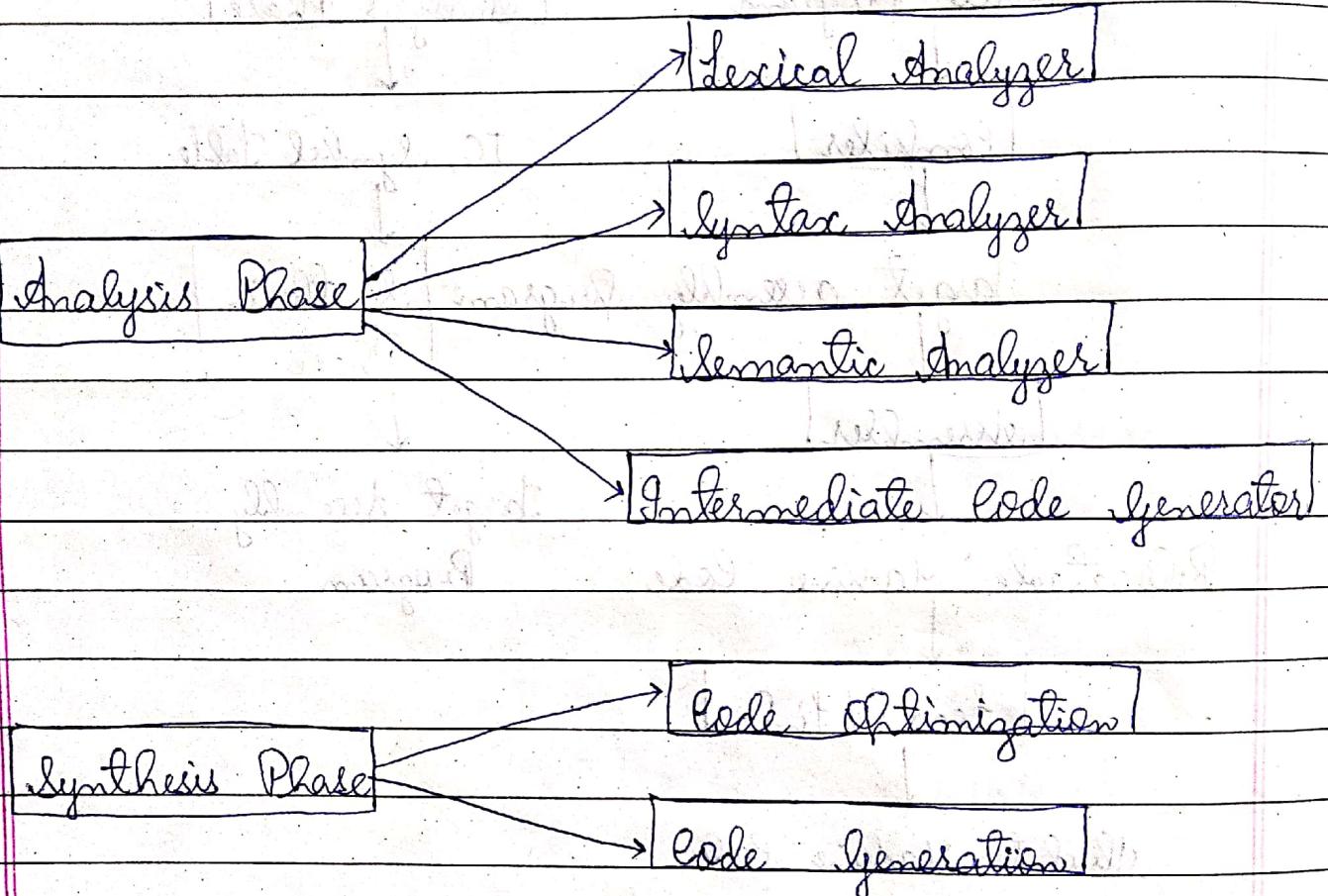


### Analysis Phase:

1. Perform analysis of input code
2. Generate errors and warnings
3. Largely dependent on the source language
4. Largely independent of the target machine
5. Creates FRONTEND

## Synthesis Phase:

1. Takes IC, uses symbol table and generates target program
2. Largely independent of the source language
3. Largely dependent on the target machine
4. creates BACKEND



## 1. Lexical Analysis:

- Grouping into tokens (sequence of characters with meaning)
- Identifies valid words or lexemes in source program
- Removes white spaces, tabs, newlines
- Passes sequence of tokens to syntactic analysis phase
- Implemented as a finite automata
- Error generated: invalid identifier

→ Example :

(i)  $\text{position} = \text{initial} + \text{rate} * 60$

Output of lexical phase:

$\langle \text{id}, 1 \rangle \Rightarrow \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle * \rangle \langle \text{number}, 60 \rangle$

(ii)  $\text{int } n1, n2;$

$\text{float } f;$

$n1 = n1 + f^* 5 - n2;$

Output of lexical phase:

$\langle \text{int} \rangle \langle \text{id}, 1 \rangle \langle \text{id}, 2 \rangle \langle ; \rangle \langle \text{float} \rangle \langle \text{id}, 3 \rangle \langle ; \rangle$   
 $\langle \text{id}, 1 \rangle \Rightarrow \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle * \rangle \langle \text{number}, 5 \rangle$   
 $\langle - \rangle \langle \text{id}, 2 \rangle \langle ; \rangle$

Limitations:

- Can identify validity of lexemes (individual words) only
- Not powerful to analyze expression or statement  
(Can't match a pair of parenthesis)

## 2. Syntax Analysis (Hierarchical analysis, Parsing)

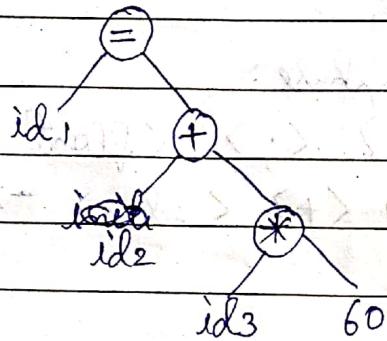
- Takes sequence of tokens from lexical analyzer
- Checks syntactic correctness of expression / blocks / statements / functions / program
- Successful if it can identify grammar rules (CFG, CSG) for a sequence of tokens.
- Generates a parse tree from a sequence of tokens
- If parse tree is complete  $\Rightarrow$  construct is correct
- If parse tree is incomplete  $\Rightarrow$  syntax error in the construct

- uses recursive grammar rules
- Errors generated: missing parenthesis, missing semicolon, missing operand

Example:

- Position = initial + rate \* 60

Tokens: Position (id1), '='; initial (id2), '+', rate (id3)  
 '\*', 60



### 3. Semantic Analysis

- If program is syntactically correct, then check for semantic (meaning) correctness of the program
- Semantics depends on the programming language
- Checks for semantic errors
- Gathers type information
- Uses parse tree generated by the syntax phase
- Common checks are:

- Type of variables and type-casting
- Applicability of operators on operands - each operator has operands permitted by the source language  
 (String operators, only int index in array)
- Scope of variables and functions
- Determine definition of variables and functions  
 (function, overloading)

Example:

```
int a=5, b=2;  
float f;  
f = a/b;  
f = 5/2.0; //inbuilt-type casting
```

→ Errors:

- (i) float array index
- (ii) cannot convert 'int' to 'string'
- (iii) undeclared variable

#### 4. Intermediate Code Generator

→ After semantic analysis, Intermediate Code (IC) of source program will be created (low-level or machine-like IR)

→ Properties:

- (i) Easy to produce
- (ii) Easy to translate into target code

→ Variety of forms:

- (i) Three Address code (3A code)
- (ii) Prefix representation

(i) Three address code (3A code)

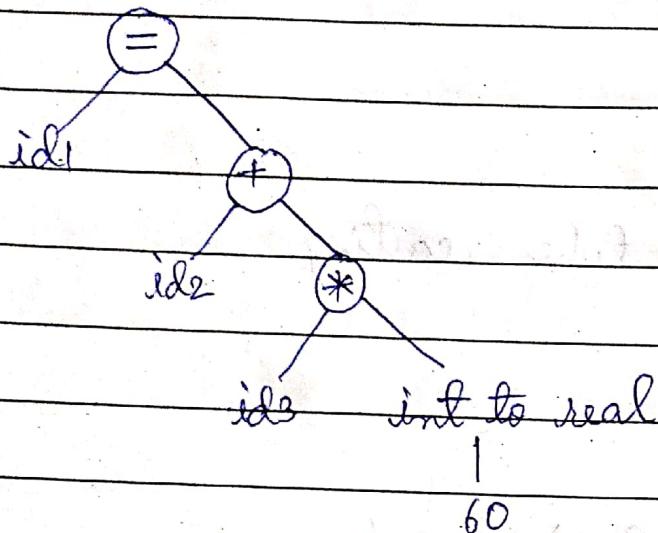
→ Each instruction has

- (a) Maximum three operands
- (b) Maximum one operator in addition to assignment

→ Compiler will decide order of evaluation

→ Compiler will generate a temp name to hold computed results

→ Example:  $\text{id}_1 = \text{id}_2 + \text{id}_3 * 60$



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

$t_2 = \text{id}_3 * t_1$

$t_3 = \text{id}_2 + t_2$

$\text{id}_1 = t_3$

## 5. Code Optimization

→ Improves the Intermediate code for

(i) Faster execution of machine code.

(ii) Space optimization

→ Optimizing compilers

→ Machine independent code optimization

→ Example: (i) 60.0 instead of int to real(60)

(ii) No need of  $t_3$

$t_1 = \text{id}_3 * 60.0$

$\text{id}_1 = \text{id}_2 + t_1$

→ Removes dead code like unreachable code

→ If anything is written after return statement in function

→ Loop optimizations - major source of optimization

→ Saving even a single line in loop improves execution time significantly

→ Example,

$b = 5;$

for ( $i=0; i < 10000; i++$ ) {

$a = b + 100 - i * 2;$

}

Optimized code:

$b = 5;$

$C = b + 100;$

for ( $i=0; i < 10000; i++$ ) {

$a = C - i * 2;$

}

## 6. Code Generation

- Takes an optimized Intermediate code as an input and converts it to the target code
- The target code can be:
  - (i) Relocatable Machine code
  - (ii) Assembly code
- Memory locations are selected for each variables of a program
- Intermediate code is translated to a sequence of machine instructions
- Part of backend of the compiler design
- Target code generation depends on:
  - (i) availability of machine instructions
  - (ii) addressing modes
  - (iii) number of registers (general and special purpose)

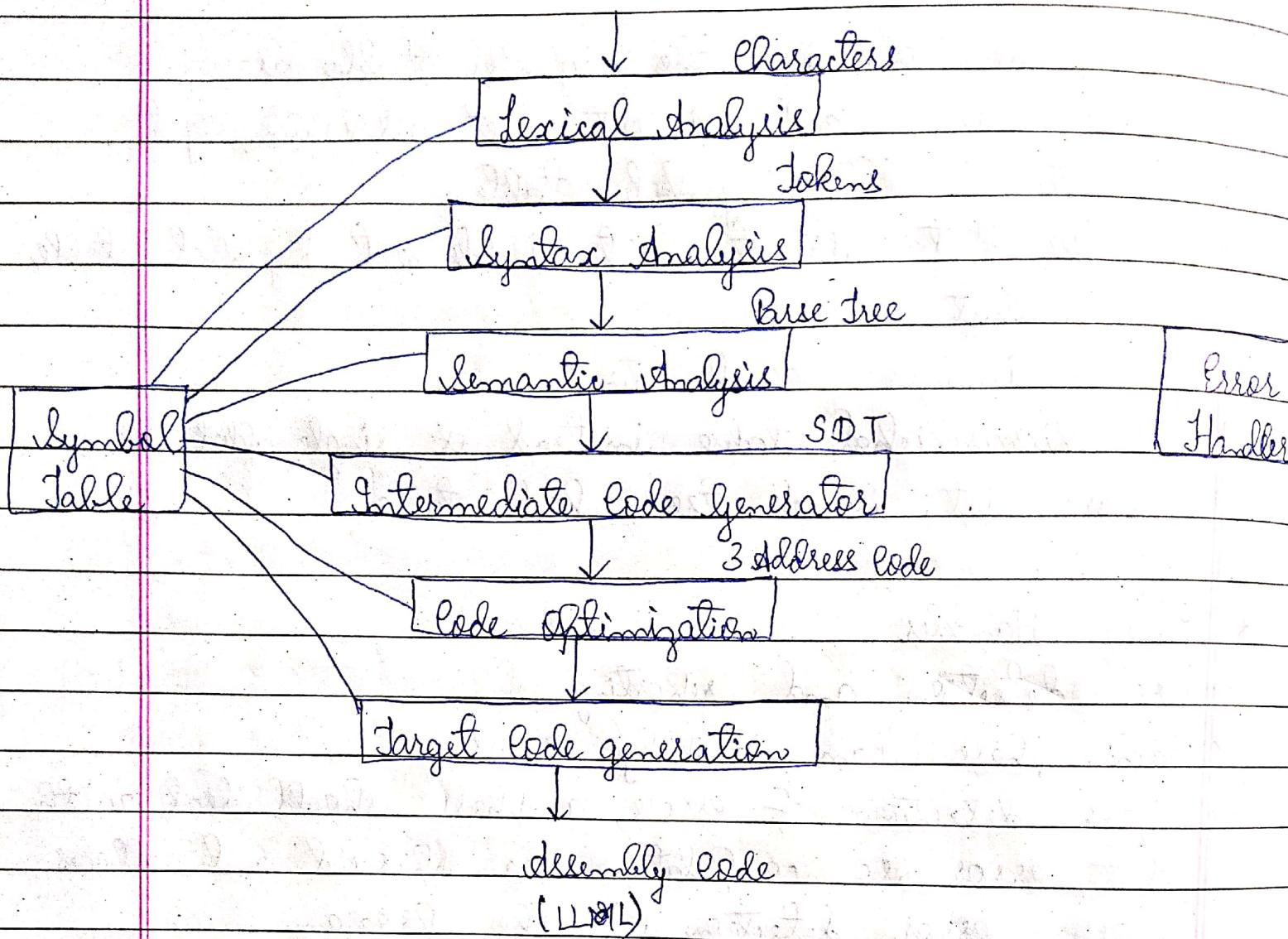
## \* Symbol Table Management:

- Interacts with all phases - used as a reference table by all the phases
- Compiler needs
  - (i) To record the identifiers (variables) used in source program
  - (ii) To collect information about various attributes of each identifier
  - (iii) Attributes for variable names - Name, type, class, size, relative offset within program, scope, value
  - (iv) Attributes for function names - Name, number of arguments, types of arguments, return type, method of passing each parameter
- It is a data structure containing a record for each identifier with fields for the attributes of the identifiers.
- Lexical analyzer: detects identifier and enters in the symbol table, but can't determine all its attributes.
- The remaining phases enter information about identifiers into symbol table and use the entered information in different ways.
- Semantic analyzer and Intermediate code generation needs to know type of variable.
- Code generator: enters and uses details about the storage assigned to identifier.
- During compilation, all phases often look up this table for definition of variable, example variable used is defined before that or not?

- Hence, operations in the symbol table are allocate, free, search, insert, set attribute, get attribute in the symbol table.
- Use the data structure to implement symbol table such that
  - (i) It minimizes search time
  - (ii) the hierarchical table instead of single flat one
  - (iii) linear list, search trees, hash table

### \* Error Handler

- Error detection and reporting
- Each phase can identify errors
- After detection of errors, a phase should deal with that error so compilation can proceed and allows further error detection in the program
- Semantic error - parser can progress
- Syntax error - parser reaches in error state
- If indicate regarding type of error
- Undo some processing already carried out by parser - error recovery
- Error detection and recovery - vital role in overall operations of compiler

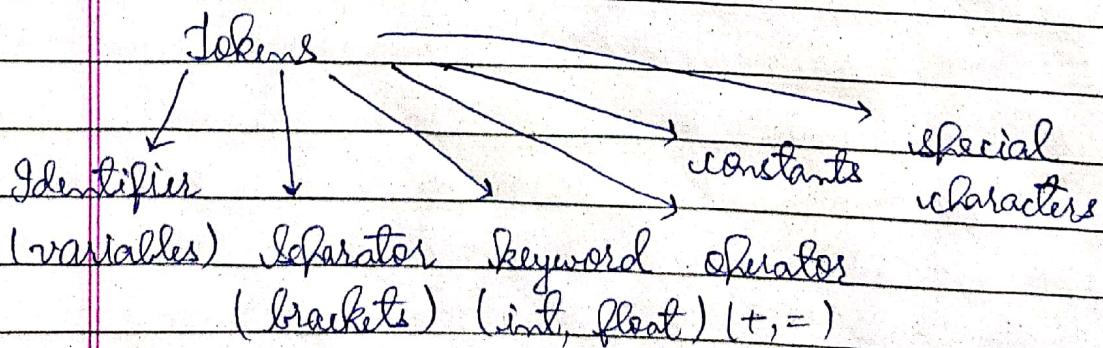


\* Lexical analysis (Lexer, Tokenizer, Scanner)

1. Tokenization

2. Give error message
- Exceeding length
  - unmatched string
  - Illegal character

3. Eliminate comments, white space (tab, blank space, new line)



Example,

1. int max()

{  
5 /\* Find max of a and b \*/ → comment

int a = 30, b = 20; 14.

if (a < b) 20

— return(b); 25

else

return(a); 31

{ (32)

2. printf("i = %.d, &i = %ox") ( i, &i );

3                  4        5      6      7      8      9      10

3. int main()

{ 5

x = y + z; 11

int x, y, z; 18

printf("sum %d %d", x);

{ (26)

4. main() { } 4

a = b ++ + -- - - + ++ = ;

5    6    7    8    9    10    11    12    13    14    15  
15    16    17

printf("%d %d, a, b);

15    16    17