

# COMPILER CONSTRUCTION

Compiler

It is a program that translates one <sup>programming</sup> language into another <sup>programming</sup> language.

Programming language code

↓ → (source code) → input of compiler

Compiler

↓

Programming language code (target code) → output

H.L.L → High Level Language

↓

Compiler

↓

L.L.L → Low level language (typically assembly)

Assembler

Assembly code ⇒ Assembler ⇒ Machine code

Types of Compilers

1) L.L.L → Decompiler → H.L.L

2) H.L.L → Transpiler → H.L.L

↳ also called source to source compiler

Interpreter

→ line by line translation (statement by statement)

Compiler vs Interpreter

- If you want fast compiling, a compiler is more suitable
- If you want more portability, an interpreter is more suitable

## Compiler

## Interpreter

→ Whole code translation → line by line (on the fly)

→ Compile once & run always → translate on every run  
(Overhead bcz you are translating each line & executing it)

→ Faster → Slower

↳ bcz already compiled

↳ lots of optimizations

→ Shows all errors but → stops on the first error  
↳ is vague in its telling (precisely tells the error)  
↳ as it tries to fix errors (as it takes assumptions)

→ less portable (compiled → more portable (can run  
differently on different OS) (hardware) on different OS) (translating code  
→ Less RAM → More RAM)

## Assembler vs Compiler

→ Assembler is a compiler BUT it translates assembly to machine code

→ Compiler can compile into other languages as well

## Just-in-time Compiler

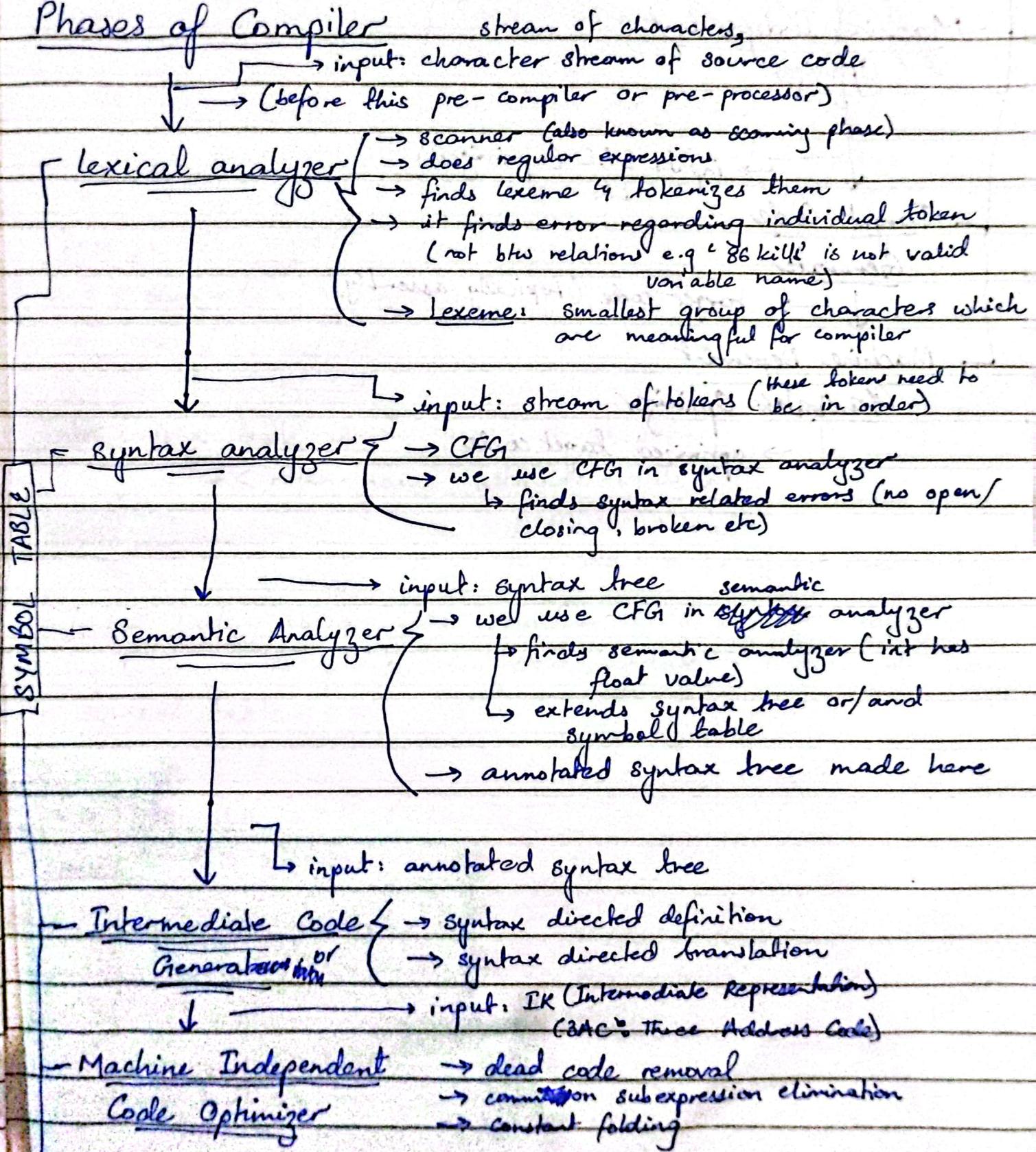
→ on the fly

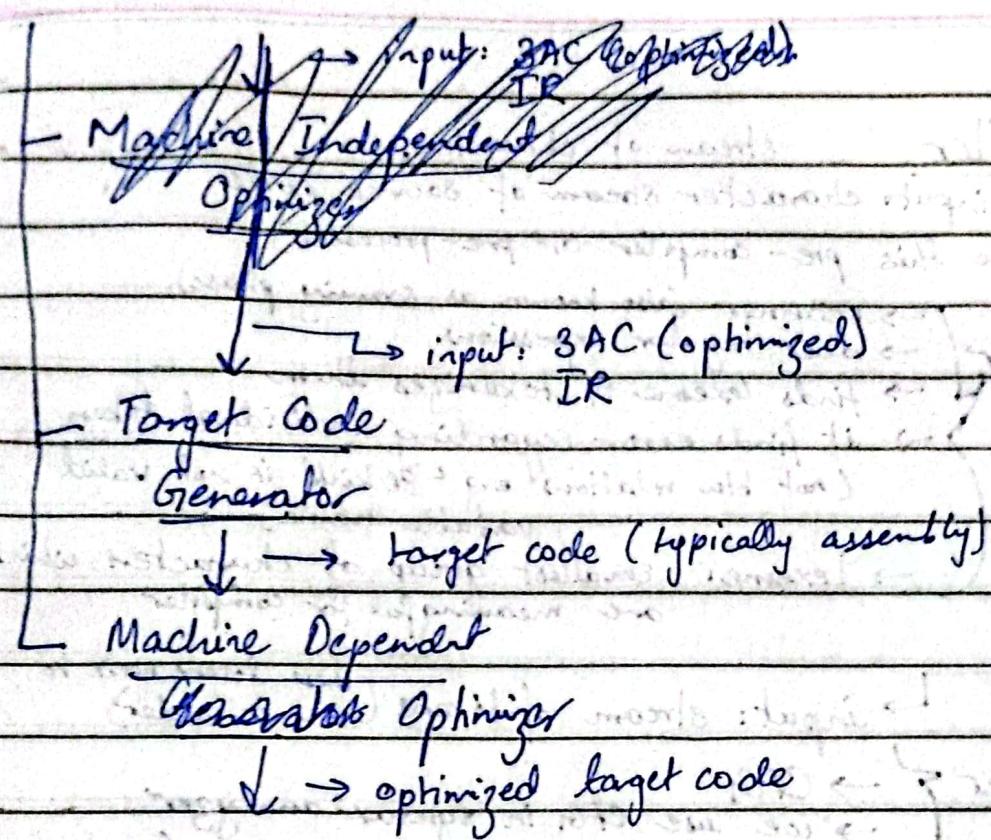
→ stores blocks of code in cache (functions, loops)

→ interpreter + compiler

→ interprets some lines & compiles loops, functions etc.

## Phases of Compiler





## Lexical Analyzer

Phase: stream of characters

Lexical Analyzer

↓  
stream of tokens

Specifications:

Regular Expressions (RE)

What Errors Are Produced?

→ related with individual lexemes

↳ e.g. 87 kills → "illegal identifier name"

Phase 1:

e.g.  $8 \text{sum} = (i + j) * 10$ .

↓↓↓↓↓↓↓↓↓↓  
all are lexemes

Token: is a tuple which consists of two things:

< token-name, <sup>(optional)</sup> attribute-value >

→ this is optional

↳ is index to the symbol table entry

Symbol Table: is a data structure used by all phases of compilation. The phases can store, update & extract info from a symbol table

for this example

e.g.  $\text{sum} = (i + j) * 10$  ←

Tokens for this will be:

<id, 1> <=> <(> <id, 2> <+> <id, 3> <)> <\*> <10> Symbol Table

↳ index to the symbol table entry

↳ identifier

1	Sum	...
2	i	...
3	j	...

→ Lexical analyzer removes spaces

check if type conversion  
into float is done in straight  
symbol analysis

Syntax analyzer's responsibility  
Deals with operators

## Syntax Analyzer

Phase:

↓ stream of tokens

## Syntax Analyzer

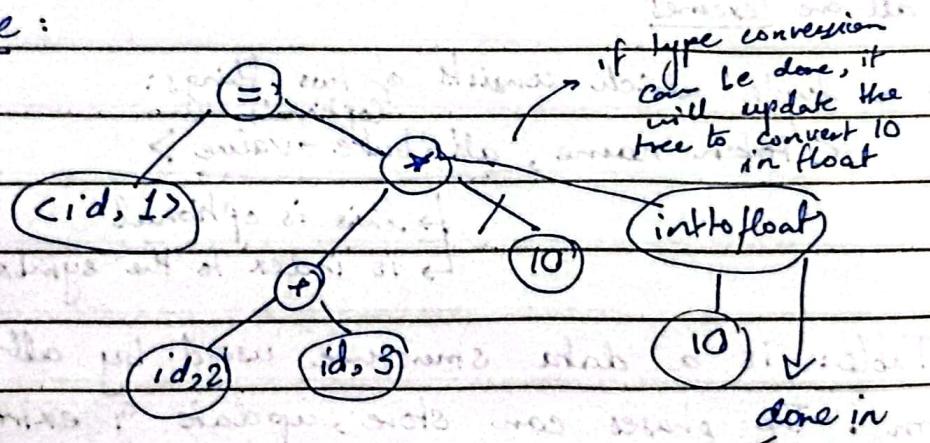
↓  
Syntax tree

Tokens:

e.g.  $\langle \text{id}, 1 \rangle \langle = \rangle \langle ( \rangle \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle * \rangle \langle 10 \rangle$

Phase 2:

## Syntax Tree:



## Semantic Analyzer

↳ updates the symbol table by adding more info in it  
↳ updates type by scope in symbol table by in the tree

Phase:

↓  
Syntax tree

## Semantic Analyzer

↓  
Syntax tree

## Specification:

Context-Free Grammar (CFG)

## Errors:

↳ Type e.g. int i = 3.5

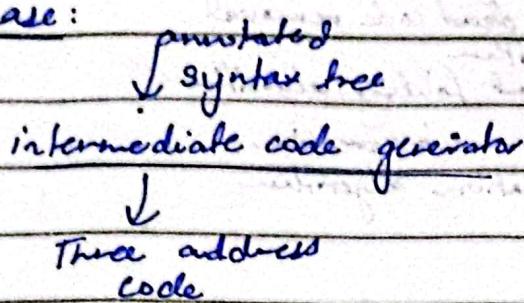
↳ Scope "illegal assignment"

Phase 3:

Type / scope related variables caught here

## Intermediate Code Generation

Phase:



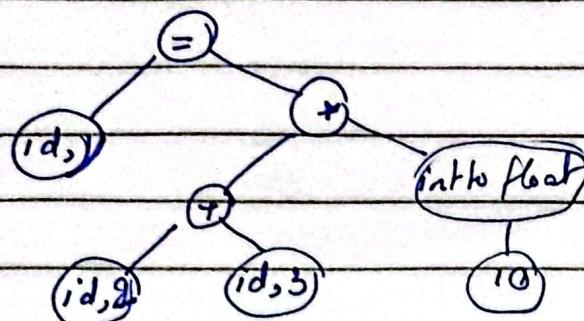
Specifications

- ↳ syntax directed definition
- ↳ syntax directed translation (basically CFG with some code by rules attached with the grammar symbols or productions)

Errors:

- ↳ does not produce many errors

e.g



Phase 4:  
4 Phases 3

$$t_1 = id_2 + id_3$$

$$t_2 = \text{int to float}(10)$$

$$t_3 = t_1 * t_2$$

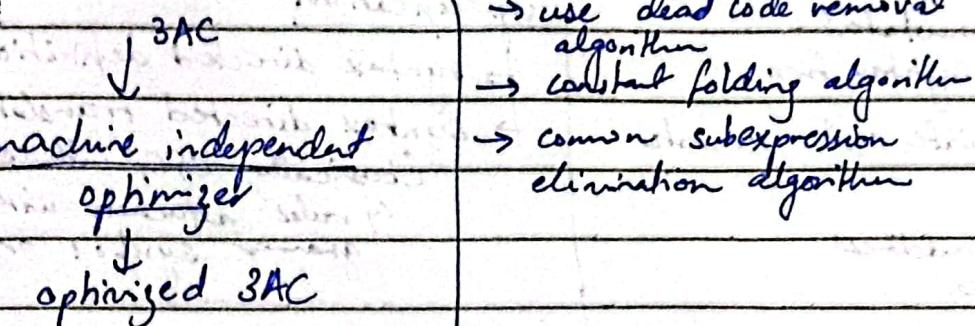
$$id_1 = t_3$$

Advantages:

- ↳ can do optimizations which are machine independent

## Machine Independent Optimizer ( )

Phase:



e.g.  $t_1 = id_1 + id_2$

$t_2 = \text{intToFloat}(10)$

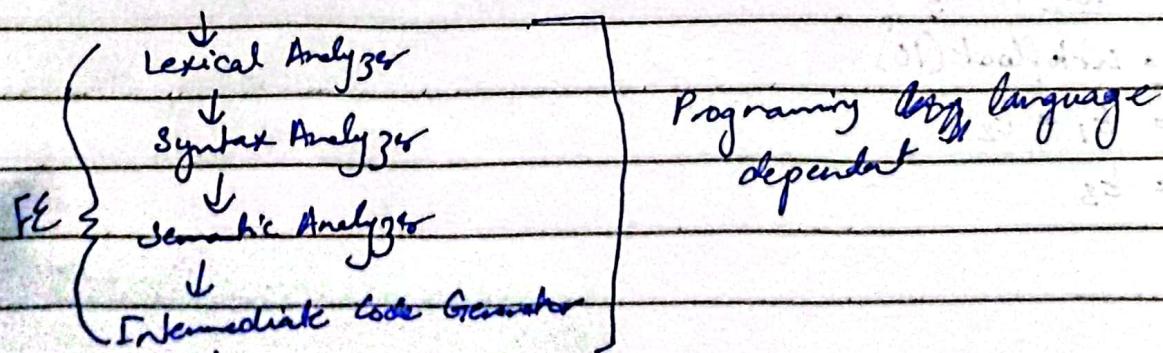
$t_3 = t_1 + t_2$

$id_1 = t_3$

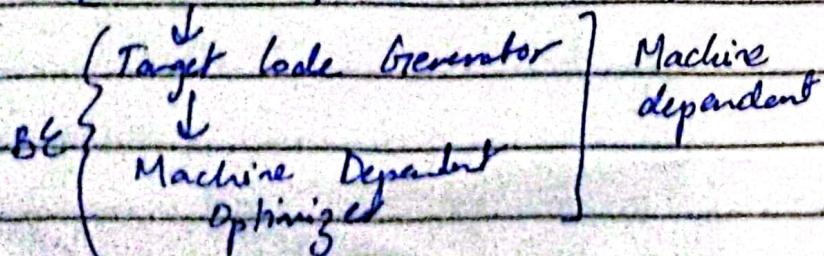
Optimized code:

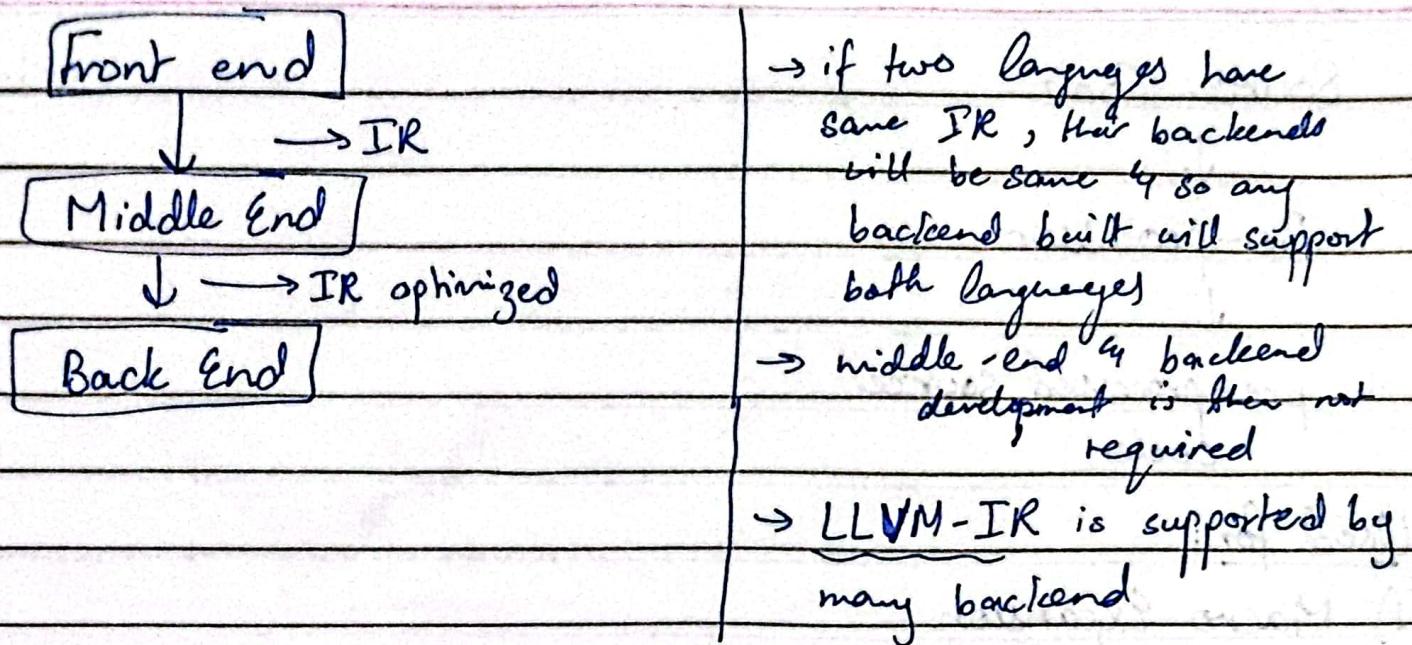
$t_1 = id_1 + id_2$

$id_1 = t_1 * 10.0$



ME [Machine Indep Optimizer] independent





Source Code



Pre-processor



pre-processed source  
code

Used for:

(1) Macro Expansion

→ `#define E 2.71`

→ `#define square(x) (x*x)`

e.g. `i = square(3);`  
`i = 3*3;`

but if `square(3+7)`

→ `i = (3+7 * 3+7)` // incorrect answer

Solution:

`#define square(x) ((x)*(x))`

→ `((3+7) * (3+7))` // correct solution

(2) `#include <stdio.h>` → including libraries

(3) `#ifdef __linux__` → code related to different  
// linux specific code OS's

`#elif _WIN32_`

// windows specific code

`#else`

`#endif`

what is  
unresolved ref,  
relative address ??

Linker/ Loader Object file → used for modular development ??

a.c ⇒ a.o

b.c ⇒ b.o

d.c ⇒ d.o ] → unresolved reference  
relative address

several object files



Linker



single executable