

# Compiler Design

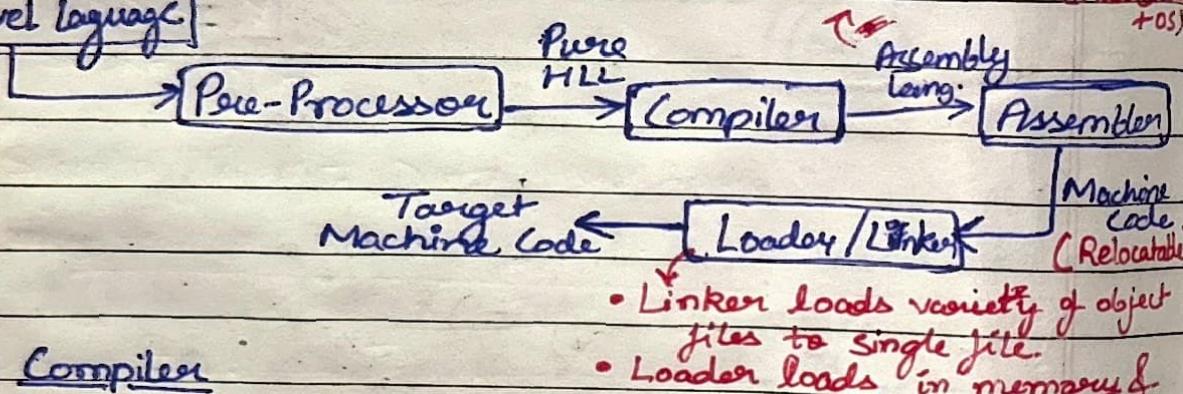
## Introduction to Compilers

### Language Processing System

Programs written in high-level languages are converted to machine usable in a series of steps.

- Neither a binary file nor High level
- Unique for every program (Hardware + OS)

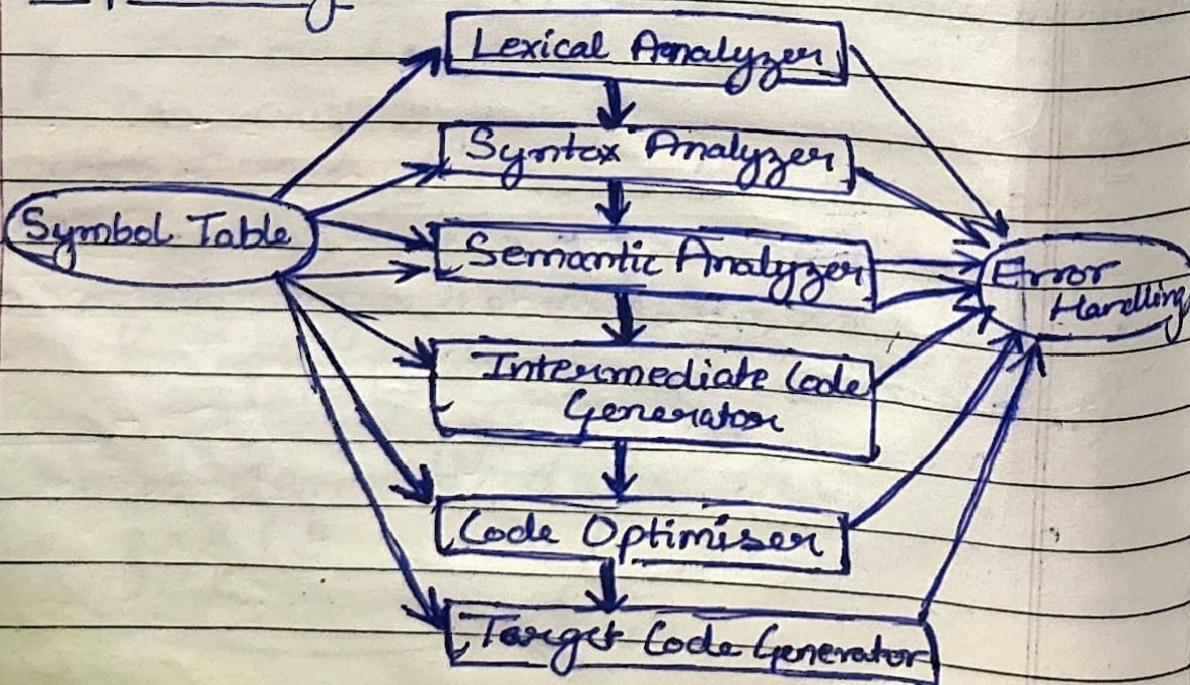
[High Level Language]



### Compiler

A program that translates code written in high level language to low level language (e.g. assembly lang.) to create an executable program.

### Compiler Design



Date / /

2 Phases → Analysis Phase  
→ Synthesis Phase

Code check

Lexical Analyzer: Clears input text by removing comments, whitespaces etc. & breaks # input text to tokens.

int a = 10; → int, a, =, 10, ;  
Tokens

Syntax Analyzer: Also called Parser, it creates syntax tree & uses production of context free grammar to construct parse tree.

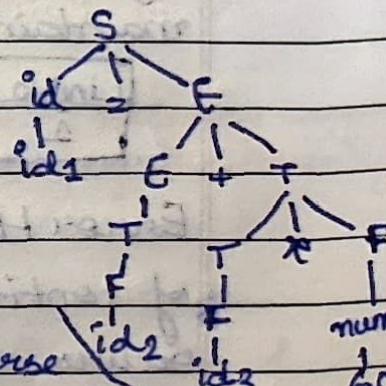
$S \rightarrow id = E$

$E \rightarrow E + T / T \quad id_1 = id_2 + id_3 * 60$

$T \rightarrow T * F / F$

$F \rightarrow id / num$

Reduced Production Rules



Semantic Analyzer: Verifies if parse tree is meaningful or not. Also produces verified / annotated parse tree.

Intermediate Code Generator:

Generates intermediate code which is machine independent. To build new compiler steps upto 3 here are some & last 2 part can be coded machine specific.

$$\begin{aligned} x &= y + z * 60 \\ t_1 &= z * 60 \\ t_2 &= y + t_1 \\ x &= t_2 \end{aligned}$$

Code Optimizer :- Transform the code so that it consumes fewer resources & produces more speed.

Target Code Generator :- Machine understandable code.

MOV R<sub>1</sub>, Z

- MUL R<sub>1</sub>, 60 ←  $(0 = 0 \text{ tri})$

ADD R<sub>1</sub>, Y

STORE X, R<sub>1</sub>

Symbol Table :- Data structure being used & maintained by compiler.

Line No	Keyword	Identifier	Constant	Operator
1	int	x	10	;

Error Handler :- Sub-routine to take care of continuation of compilation even if error occurs. After phase 3 if error handler object is empty then source code is error free.

3 types of errors → Lexical, Syntax, Semantic

int a=10  
int

a+b=c;  
c=a+b

a = 'C' + 1.2

### Phases & Passes

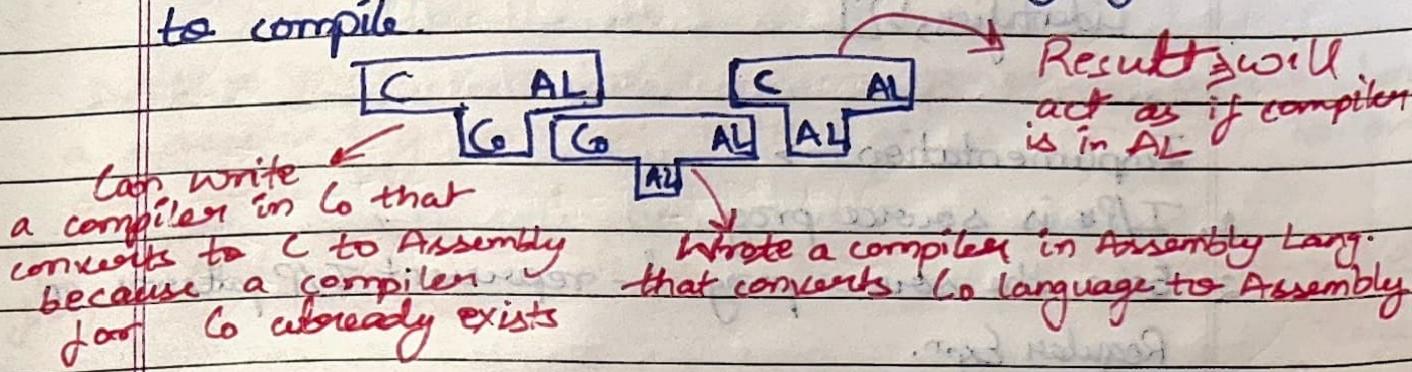
Phase - Distinguishable stage, which takes I/P from previous stage, processes & gives O/P that can be used by next stage.

Pass - Traversal of compiler through entire program.

- Single Pass Compiler → Faster, memory efficient, less optimized  
Ex: Tiny C Compiler
- Multi Pass Compiler → Multiple scans, more optimized code, context awareness, more memory usage. Ex → GCC, FORTRAN (only)

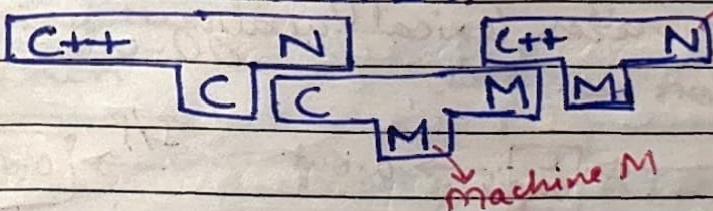
### Bootstrapping

Earlier, compilers were written in assembly language but as languages evolved so did writing compilers. Writing a compiler in some language it intends to compile.



### Cross Compiler

Allows devs to write code on 1 machine & execute it on another.



### Lexical Analyzer

Source Prog → Lexical Analyzer → Token

- Lexeme - Sequence of ~~Prog~~ characters in source code that matches a pattern for token. It is the actual text
- Token - Logical meaning assigned to lexeme

- Lexing can be divided in 2 phases:
- i) Scanning - Breaking input string into lexemes & categorizing into tokens
  - ii) Evaluating - Converting lexemes into processed values.

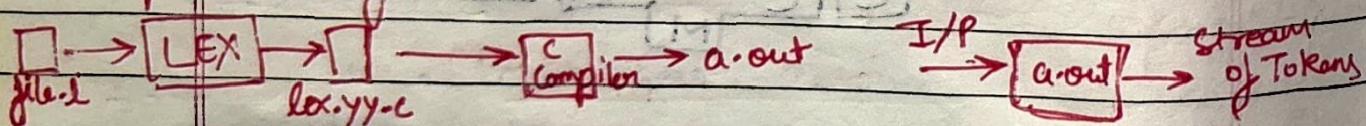
Ex →  $x = a + b;$

[Identifier, x), (operator, =), (Identifier, a), (operator, +),  
 (Identifier, b)]

Implementation →

1. I/P is source prog.
2. Scan the source prog. & represent I/P patterns in Regular Expr.
3.  $R^{\text{Expr}} \rightarrow \text{NFA} \rightarrow \text{DFA} \rightarrow \text{mDFA} \rightarrow \text{lexemes}$

Note → Instead of us writing, a UNIX utility (LEX) exists that generates lexical analyzer using Regular Expression for us.



## Syntax Analysis

### Type 2 Grammar (CFG)

$\alpha \rightarrow \beta$  *Variables*

$\alpha \in V_n$   $| \alpha | = 1$

$\beta \in \{ \Sigma \cup V_n \}^*$

↓ *Terminals*

Formal Grammar represent specification of programming lang. with use of production rule

## BNF Notation

Backus-Naur Form is notation used to express grammar of computer language.

Symbols  $\rightarrow$  Terminals      Production Rules  $\rightarrow$  non-terminals

Ex: Expression  $\rightarrow$  Expression + Term | Term

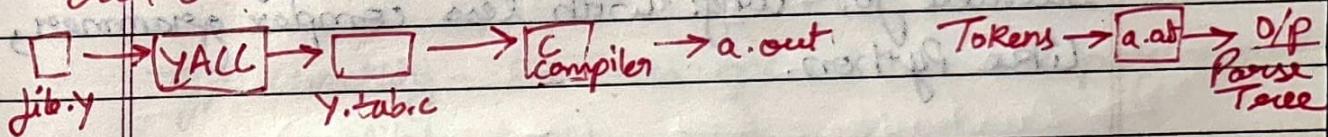
Term  $\rightarrow$  Term \* Factor | Factor

Factor  $\rightarrow$  ( Expression ) | Number

Number  $\rightarrow$  Digit | Number Digit

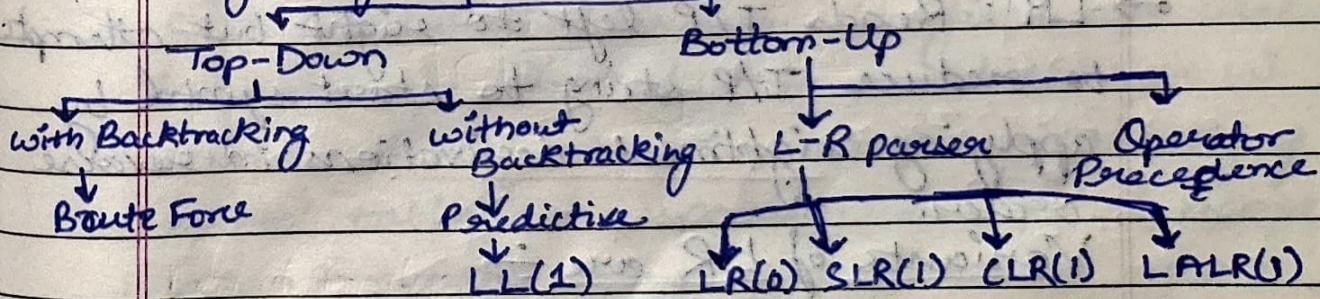
Digit  $\rightarrow$  0|1|2|3|4|5|6|7|8|9

Note  $\rightarrow$  YACC is a tool used to generate parser for a specified grammar instead of us writing it on our own.



## Basic Parsing Techniques

### 2 Types of Parser



Top-Down is method of parsing in which parser starts at root of parse tree & works its way down to leaves.

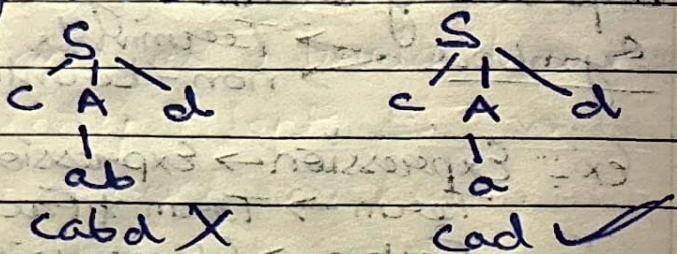
It cannot work with Left Recursive Grammar, Ambiguous grammar.

It can work on Non-Deterministic grammar

using backtracking but that is inefficient.

→ Brute Force : Try all alternatives

$$\text{Ex } S \rightarrow cAd \\ A \rightarrow ab/a \\ w_1 = cad$$



→ LL(0)(1) : Reads the I/P from left to right & constructing a left-most derivation of sentence. Uses a parsing table & stack to decide which production to apply.

Suitable for long. with less complex grammar, like Python.

Bottom-up starts from leaves (I/P tokens) & work their way up to root of tree.

→ LR : Reads I/P left to right but attempt to reduce I/P string to start symbol by applying rightmost derivations in reverse order.

Variants of LR are :-

→ SLR : Simple LR

→ LALR : Lookahead LR

→ CLR : Canonical LR (most powerful type of LR)

LR parsers are used in modern prog. language  
Ex - C, C++, Java

## ~~Syntax Directed Translators~~

### Semantic Analysis

Takes Abstract syntax tree (AST) generated by syntax analysis phase as its input & performs type checking, scope resolutions & validate semantic consistency according to language's rules.

AST is traversed & compiler builds & maintains a symbol table that tracks information about variables, functions, classes etc.

Type Checking - Ensures operations are performed on compatible types.

```
int a=5;
float b=3.5;
a=a+b; // Error.
```

Scope Resolution - Ensures identifiers are used in correct context & that variables

→ Scope Hierarchy - Compiler manages hierarchy of symbol tables corresponding to nested scope. When traversal enters new scope a new table is created.

→ Lookup - When identifier is referenced it looks in current scope, if not found moves to next outer scope & so on.

```
int x = 10;
void fool() {
    int y=5;
```

```
    int x=3; // This x shadows outer x
```

```
    y=x+2;
```

## Intermediate Code Generation

- Converts annotated AST into intermediate code, which is abstract representation of source code & acts as a bridge b/w HLL & LL
- It is machine independent & can be translated to different target machine.

## Intermediate Code Forms -

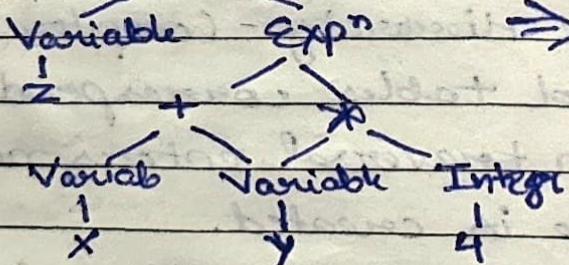
- Three-Address Code : Each instruction has at most 3 operands (2 source, 1 destination)
- Ex →  $t_1 = a + b$
- $t_2 = t_1 * c$
- Quaduples : ~~set~~ (+, a, b, t<sub>1</sub>)
- Directed Acyclic Graph : Nodes represent operations & edges represent data flow

Example →

int main() {

    int x=2, y=3, z;  
    z = x + y \* 4;  
    return z;  
}

Assignment Statement



⇒ Intermediate  
Code gen.

Traverse AST, process each node.

$t_1 := y * 4$

$t_2 := x + t_1$

$z := t_2$

## Target Code Generation

### Code Optimization

Improving Intermediate code so that final machine code runs efficiently. It transforms intermediate code into equivalent but more optimized form.

#### Techniques →

- Constant Folding :  $a = 2 + 3 \rightarrow a = 5$
- Constant Subexpression Elimination :  $b = a * c$        $d = a * c \Rightarrow b = d$
- Dead Code Elimination
- Loop Invariant Code Motion : Code that does not change within loop is moved outside
- Strength Reduction : Replace expensive operations  
 $a * 2 \rightarrow a \ll 1$
- Machine dependent
- Register Allocation : Assign frequently used variable to registers for faster access

### Code Generation

Takes optimized intermediate code & translates to machine code that can be directly executed by target processor.

#### Tasks in code generation -

- Instruction Selection : choose appropriate machine instruction to represent operations in intermediate code.

$$a = b + c$$

addi x1, x2, x3 (RISC-V processor)

- Register Allocation : Assigning ~~registers~~ variables

to registers to minimize memory allocation

$$a = b * c$$

~~D00000~~

a assigned to  $x_1$

b " "  $x_2$

c " "  $x_3$

Machine code

~~mul x1, x2, x3~~

- Instruction Scheduling : Rearranging instruction to improve pipelining & reduce stalls.

$$a = b * c$$

$$d = a + e$$

~~without scheduling~~

~~mul x1, x2, x3~~

~~add x4, x1, x5~~

~~mul x1~~

- Memory Management : Handling memory allocation & deallocation for variables & data structures.