# Compiler. Construction tools

Compiler Construction tools can use Morden s/w dvmt., environments Containing tools such as

- language editors
- debuggers
- Version Managers.
- profilers
- text harnesses.

and so on.

} general s/w dvmt., tools.

* Along with the above tools, other more specialized tools have been created to help implement various phases of a Compiler.

* These tools use ~~Specifili~~ Specialized PLs fer implementing Specific Components uring many Sophisticated and Complex algorithms.
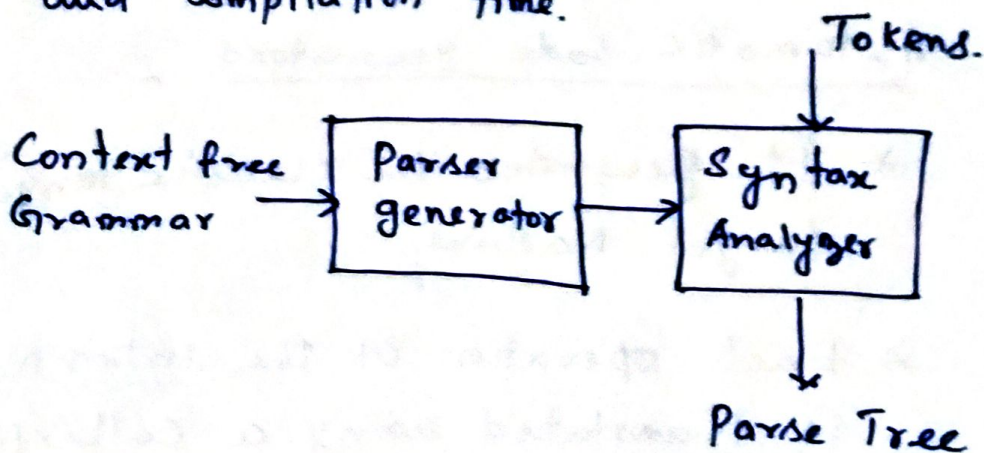
(X·) * Most of the tools provide high degree of abstraction Means the details of implementation is hedden from outside world.

# Compiler Construction tools (or) Writing tools

* The Compiler writer can use some specialized tools that help in implementing various phases of a Compiler.

* These tools assist in the creation of an entire Compiler or its parts.

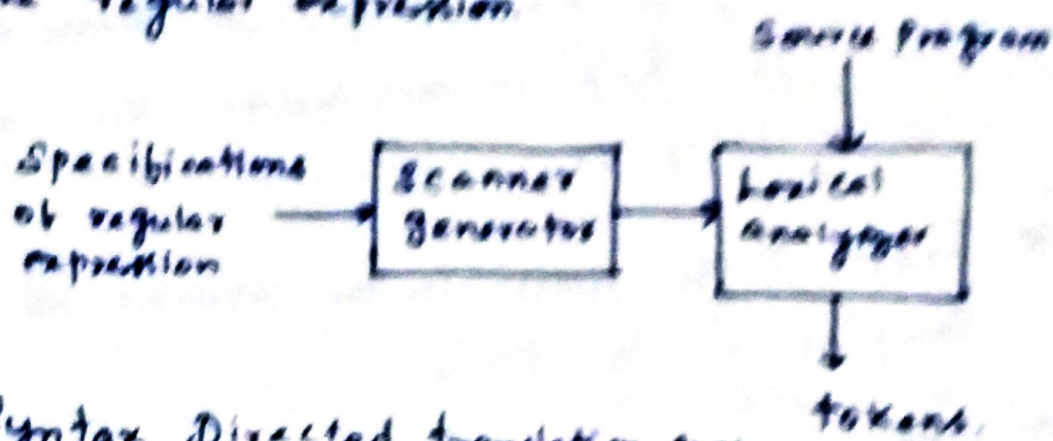* Some Commonly used Compiler Construction tools include.:

## 1. Parser Generator

It produces syntax analyzers (Parsers) from the input that is based on a grammatical description of PL or on a Context free grammar. It is useful as the Syntax analysis phase is highly Complex and Consumes more Manual and and Compilation time.

Context free Grammar → Parser generator → Syntax Analyzer

Tokens → Syntax Analyzer → Parse Tree

## 2. Scanner Generator

It generates lexical analyzers from the input that Consists of regular expression

description based on tokens of a language.
It generates a finite automaton to recognize
the regular expression

Specifications of regular expression → [Scanner generator] → [Lexical analyzer] ← Source Program → tokens.

3. **Syntax Directed translation engines**

* It generates intermediate code with three
  address format from the input that consists
  of a parse tree.

* These engines have routines to traverse the
  parse tree and then Produces the intermediate
  code.

* In this each node of the parse tree is associated
  with one or more translations.

4. **Automatic Code generators**

* It generates the Machine language for a
  target Machine.

* Each operation of the intermediate language
  is translated using a collection of rules
  and then is taken as an input by the code
  generator.

* A template matching process is used.

* An intermediate language statement is replaced by its equivalent machine language statement using templates.

5. **Data flow analysis engines.**

    * It is used in code optimization.

    * Data flow analysis is a key part of the code optimization that gathers information that is the values that flow from one part of a program to another.

6. **Compiler Construction toolkits.**

    * It provides an integrated set of routines that aids in building compiler components in the construction of various phases of compiler.

⏺

## Ex - 2

$$x = a + b * c ;$$

$$\Downarrow$$

$l(l+d)^* \rightarrow$
**Pattern.**

```
Lexical Analyser.
```

→ once the LA receives this i/p, the i/p is Converted into a stream of tokens. Not only this, it removes the white spaces and skipping the comments.
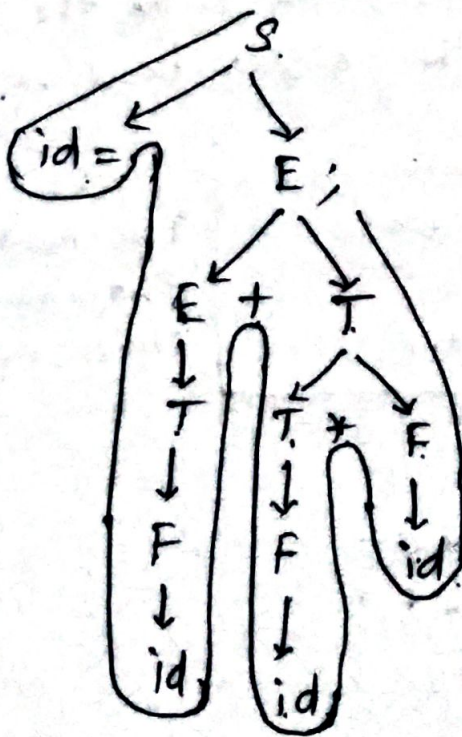
$$\Downarrow$$

$$id . = id + id * id$$

$$\Downarrow$$

```
Syntax Analyser
```

← Contex free Grammar.

$$\boxed{\begin{array}{l} S \rightarrow id = E ; \\ E \rightarrow E + T / T \\ T \rightarrow T * F / F \\ F \rightarrow id . \end{array}}$$ rules.

**Parse Tree.**

S → id = , E ;
E → E + T
E → T, T → T * F
T → F, T → F, F → id
F → id, F → id
id, id

The Rules the PL can be entirely represented in some **Productions** rules.

S → A statement can be identifier = expression followed by ;

E ⟹ An expression can be Expression + Term or Term.

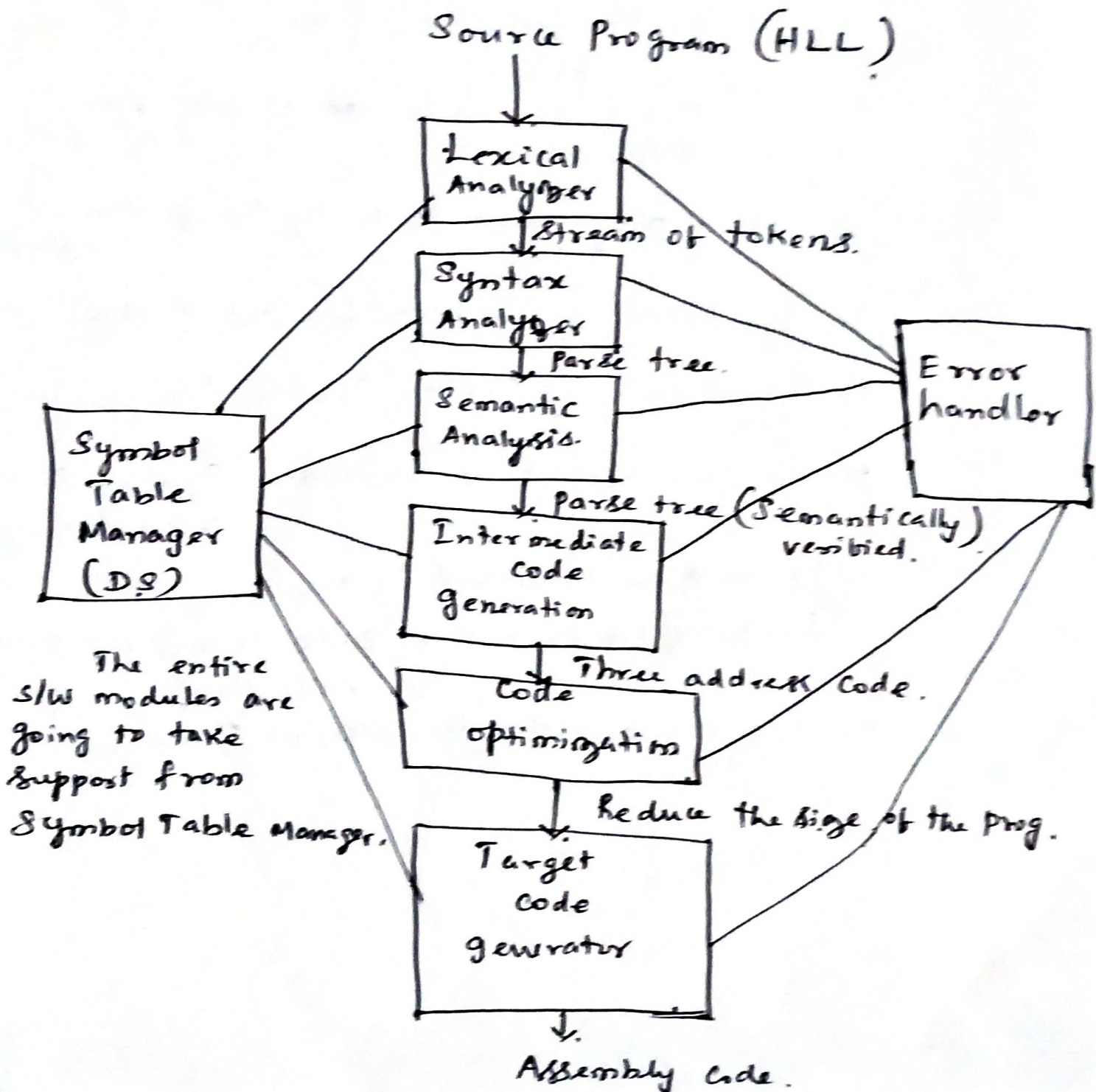T → A Term can be Term * Factor or Factor

F → can be an identifier.

The Syntax error is detected by Syntax Analyser if the input is not according to the Grammar given.

$$\Downarrow$$

```
Semantic Analyser
```

$$\Downarrow$$

# Phases of ~~Computer~~. Compiler.

## Phase ?

is a **logically interrelated operation** that takes source prog. in one representation and produce of output in another representation.
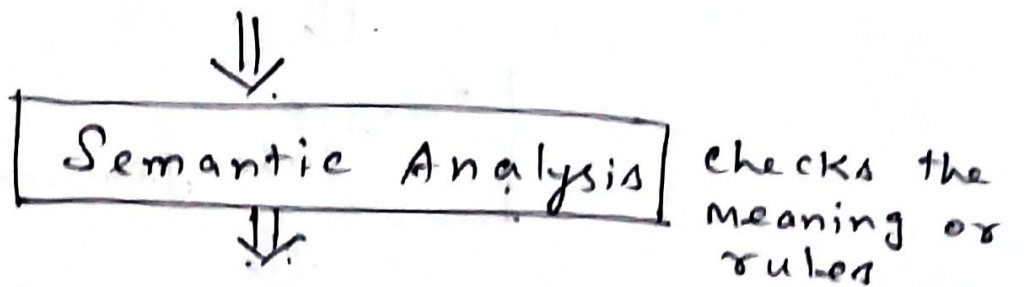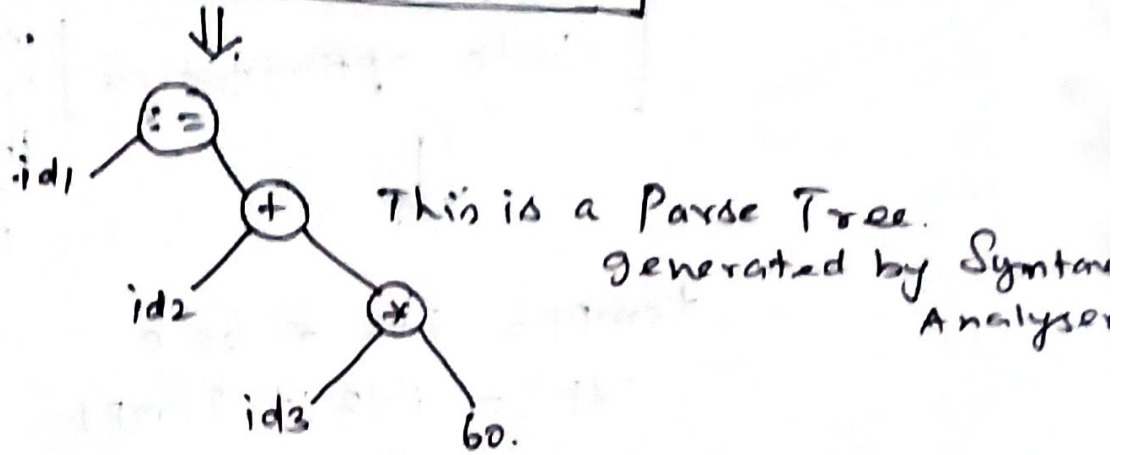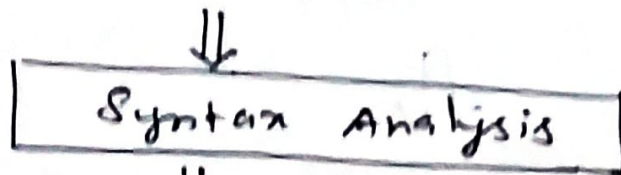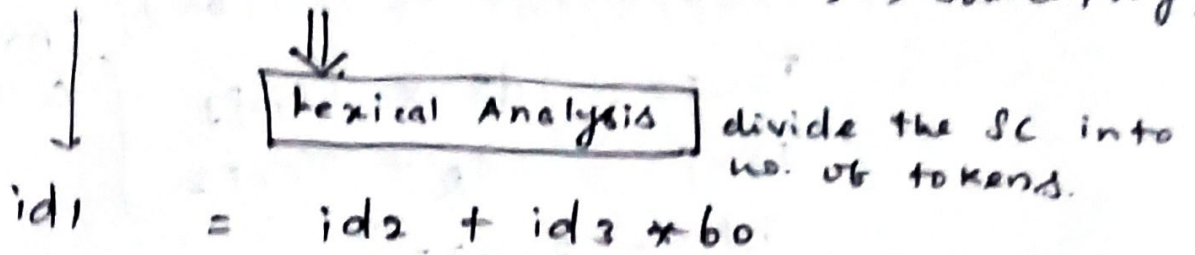
Source Program (HLL)

↓

```
┌─────────────┐
│ Lexical     │
│ Analyzer    │
└─────────────┘
```
↓ Stream of tokens.

```
┌─────────────┐
│ Syntax      │
│ Analyzer    │
└─────────────┘
```
↓ Parse tree.

```
┌─────────────┐
│ Semantic    │
│ Analysis.   │
└─────────────┘
```
↓ Parse tree (Semantically) verified.

```
┌─────────────┐
│ Intermediate│
│ Code        │
│ generation  │
└─────────────┘
```
↓ Three address Code.

```
┌─────────────┐
│ Code        │
│ optimization│
└─────────────┘
```
↓ Reduce the size of the prog.

```
┌─────────────┐
│ Target      │
│ Code        │
│ generator   │
└─────────────┘
```
↓

Assembly code.

Symbol Table Manager (DS)

Error handler

The entire s/w modules are going to take support from Symbol Table Manager.

# Explanation with example

**Ex.**

Position := initial + rate * 60 ; → Source Prog.

↓

| Lexical Analysis | divide the SC into no. of tokens.

id1 = id2 + id3 * 60

⇓

| Syntax Analysis |

⇓

id1 ⟨:=⟩
⟨+⟩
id2 ⟨*⟩
id3 60.

This is a Parse Tree.
generated by Syntax
Analyser

⇓

| Semantic Analysis | checks the meaning or rules

⇓

Suppose rate is float
type but 60 is int.
So, convert 60 into float.

id1 ⟨:=⟩
id2 ⟨+⟩
id3 ⟨*⟩
int to real
↓
60.0

⇓

```
        ⇓
┌──────────────────┐  Generates
│  Intermediate    │  a three address code.
│     Code         │
│  Genatation.     │
└──────────────────┘
        ⇓
```

$Temp_1 = int$ to $real\,(60)$  ⎫ Three address
$Temp_2 = id_3 * Temp_1$         ⎬ Code
$Temp_3 = id_2 + Temp_2$         ⎭
$id_1 = temp_3$

⇓

```
┌──────────────────────┐  removes unnecessary
│  Code optimization   │  information to
└──────────────────────┘  increase the
                          speed of the
                          program.
```

⇓

$temp_1 = id_3 * 60.0.$
$id_1 = id_2 + Temp_1$

⇓

```
┌──────────────────────┐
│  Target              │
│      code generator  │
└──────────────────────┘
```

⇓

MOV F    60.0, R1   ⎫
MUL F    id3, R1    ⎬ Assembly
MOV F    id2, R2    │ Language
ADD F    R1, R2     │ Prog.
MOV F    R2, id1.   ⎭

## Ex:-

$$\# = A + b * c ;$$

⇓

$$l(l+d)^* \rightarrow \boxed{\text{Lexical Analyser}}$$
Pattern.

Once the LA receives the i/p, the i/p is converted into a stream of tokens. Not only this, it removes the white spaces and skipping the comments.

⇓

$$id = id + id * id$$

⇓

$$\boxed{\text{Syntax Analyser}}$$

⇓

Context free Grammar

$$S \rightarrow id = E ;$$
$$E \rightarrow E + T / T$$
$$T \rightarrow T * F / F \quad \text{rules.}$$
$$F \rightarrow id .$$

**Parse Tree**



The above the PL can be entirely represented in some productions.

$S \rightarrow$ A statement can be. rules.
identifier = expression followed by ;

$E \rightarrow$ An expression can be expression + Term or Term.

$T \rightarrow$ A Term can be Term * Factor or Factor

$F \rightarrow$ can be an identifier.

The Syntax error is detected by Syntax Analyser if the input is not according to the Grammar given.

⇓

$$\boxed{\text{Semantic Analyser}}$$

⇓

Semantic Analyser is going to verify the parse
Tree whether the Parse tree is semantically
correct or not.

eg

$$id = id + id * id.$$

$\downarrow$

left hand side should be a variable
can't be a constant or can't be a fun name
or can't be an array name.
So, ~~left~~ left hand sing has to be a variable
which is compatible with the type of
variable in the right hand side.

$\Downarrow$

| Intermediate Code Generator | The most popular Intermediate code is three address code. |

$\Downarrow$

$$t_1 = b * c;$$
$$t_2 = a + t_1;$$
$$x = t_2;$$

$\Downarrow$

| Code optimization | $\rightarrow$ reduce the ~~the~~ no. of li |

$\Downarrow$

$$t_1 = b * c;$$
$$x = a + t_1;$$

$\Downarrow$

| T C G |

$\Downarrow$

MUL   R1, R2    $\Big| a \rightarrow R_0$
add    R0, R2    $\Big| b \rightarrow R_1$
Mov   R2, X     $\Big| c \rightarrow R_2$

①

# Compiler Design

## Importance of Subject

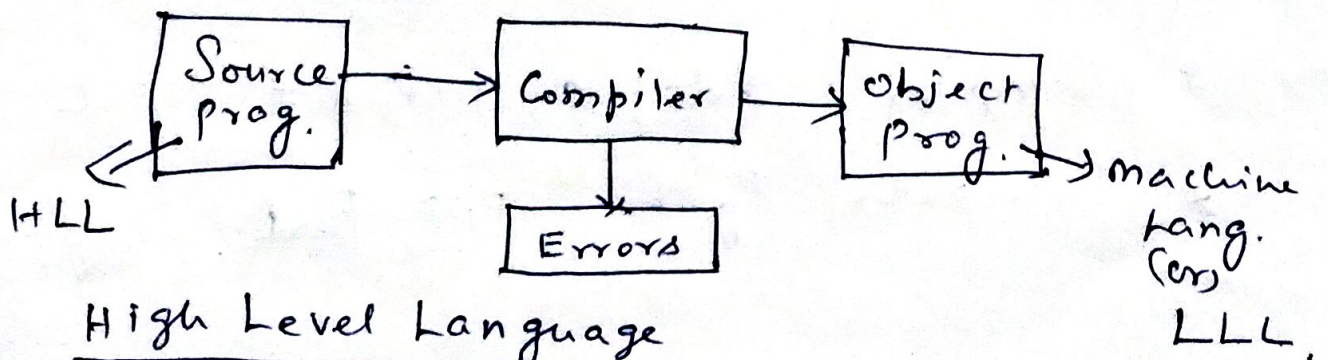* Compilers provide you with the _theoretical_ and _practical_ knowledge that is needed to implement a Prog. Language.

## Contents

* Structure of Compiler
* Interpreter
* Compiler Vs Interpreter
* Language Processing System.

## Introduction to Compiler.

A Compiler is a translator/s/w that _Converts_ the program written in high-level language (Source Language) to Low Level Language (object / Target / Machine Language).

Block Diagram.



```
   ┌─────────┐      ┌──────────┐      ┌─────────┐
   │ Source  │─────▶│ Compiler │─────▶│ object  │
   │ prog.   │      │          │      │ Prog.   │──▶ machine
   └─────────┘      └────┬─────┘      └─────────┘    Lang.
     HLL                 ▼                            (or)
                    ┌─────────┐                        LLL.
                    │ Errors  │
                    └─────────┘
```

## High Level Language

A Program written in English Language English Language is familer with us to understand easily.

We can easily & understand the Syntax, Simantic, and Structure.

* Machines do not understand HLL. They can understand LLL.

So Compiler Converts HLL to LLL.

Compiler and Interpreter:

* To Convert Source Code (SC) into machine code, we use either a Compiler or an interpreter.

* Compiler transforms Code written in a high-level Prog. Lang. into the Machine Code at once before the Prog. runs.

* where-as an interpreter Converts each high level Program statement, one by one into the Machine Code, during execution of prog.

* Compiler : Ada, c, C++, C#, COBOL.
* Interpreter : Python, PHP, Perl, Ruby.

## Compiler Vs Interpreter.

| Compiler | Interpreter. |
|---|---|
| * Scans the entire Prog. and translates it as a whole into Machine code. | * Translate prog., one st., at a time. |
| * Take a large amt., of time to analyze the SC. | * Take less amt., of time to analyze the SC. |
| * The overall exe., time is Comparatively faster than interpreters. | * The overall exe., time is Comparatively slower than Compilers. |
| * Requires More Memory (generates intermediate code, so it requires more amts of memory) | * Memory efficient. (No need to generate intermediate code.) |

Have you heard any PL name which uses both Compiler and Interpreter.?

## Language Processing System.

The task of Compilation will not be done by Compiler along. There are other s/w modules which will accompany Compiler in this process.

So, in order to know what are the s/w modules, we need need to know what is Lang. Processing System. is.

Language Processing System Consists of 4 s/w modules.

1. Pre-Processor
2. Compiler
3. Assembler
4. Loader / Linker.

HLL
↓↓
| Pre-Processor |
⇓ Pure HLL.
| Compiler |
⇓ Assembly Lang.
| Assembler |
⇓ Relocatable Machine Lang.
| Loader / Linker. |
⇓

Absolute Machine Code.

Let us have a detailed descriptions about each s/w Modules.

1. <u>Pre-Processor</u>.

The first s/w module is Preprocessor whose input is HLL and output is Pure HLL.

Now what is this pure HLL?

Consider a C program. Basically your C prog should start with a line # include < stdio.h > and # include < math.h >.

These header files are also called as Pre processor Directives

```
Removes : -
Directives like
    # include
    # define
Performs : -
Macro
expansion
File  Inclusion
```

Function calling is over head but Macro calling is simple

The Pre Processor is going to remove the above # include, # define by including files related to them. This is called as <u>File inclusion</u>

ie whatever the files you want to include Pre processor is going to <u>Substitute</u> that the entire files in your source prg.

Pre processer will also do Macro expansion

2. **Compiler**

The input to the Compiler is Pure HLL which means the prog. will not contain any # lines (or) tags.

Compiler is going to Convert Pure HLL into assembly Language.

This Assembly Lang. is not entirely $0^s$ and $1^s$ and not entirely HLL.

It is somewhat intermediate.

3. **Assembler**

Assemblers of one platform. will not work with any other Platform

The o/p of the Assembler will be machine codes.

Basically Machine codes of two types.
1. Relocatable Machine Code
2. Absolute Machine code.

Relocatable Machine code is. loaded
at any point of your computer. and
you can run it.

This ~~to~~ Relocatable Machine
code is given as input to the next phase
Loader/linker.

## 4. Loader/Linker.

The linker will link variet of
object files into one file and the
loader will load this ~~into~~ files into
the memory.

By this the execution of prog
will be completed.

Now we will focus on Compiler
and not on all s/w modules.

Before discussing the phases of
a Compiler, Let us know some thing
about
* word
* Sentences
* Languages

\* <u>Word</u> can be defined as a set of chars. which gives a meaning.

eg "boy" → word.
↳ defined from the set of chars from the alphabets available in English Lang.
↳ It denotes Masculine gender child.

\* <u>Sentences</u>

is a set of words which gives a meaning.

eg "The boy is running"

- The sentence should follow the grammar
- Here the Grammar defines the way in which the sentence can be formed.

- If a sentence does not follow any Grammar, then it can be called as a grammatically incorrect sentence.

\* <u>Languages</u>

- The Lang. is a tool used to Communicate with others
- It is defined over the sentences in turn words, which in turn defined over the chars.

The following example will help us in understanding Concepts of Compiler Construction.

Consider the following Sentence in English.

I AM GOING TO MARKET.

I → Subject

AM → Auxiliary verb.

Go ting → Verb + ing

To MARKET → object.

This simple sentence is <u>Syntactically</u>, <u>Grammatically</u> Correct.

e.g

The boy is going to hastel

↳ This Sentence is not Syntactically correct.

⇓ Syntax error.

e.g

The boy is go to hostel.

⇓

Grammatically incorrect.

we have a Sentence which is Syntactically and Grammatically Correct.

"The boy is going to hostel"

Now what is to be done?

The Answer is, The boy must do some action.

ie The boy must go to the hostel.

The <u>action</u> is attached to the <u>Sentence</u>

-x- It is clear that an action is associated with a Sentence which is Syntactically, Grammatically Correct.

The above simple analogies will explain how the ~~computer~~ Compiler works in very simpler terms.

| | | |
|---|---|---|
| Language | → | Any HLL. |
| word | → | String |
| Sentence | → | Statement in a HLL. |
| Set of Sentences | → | Set of Statements ie a Prog. in a HLL. |

To carryout certain task a <u>sentence</u> must be written, which should satisfy the foll:

    1. Syntactically Correct

    2. Grammatically Correct

    3. An Action must be associated

    4. It must be understood by the Executor for execution.

               (or)
    " must be made ~~by~~ ready for execution. ".

The st, of a prog. must satisfy the foll: in a Compiler.

    1. Syntactically Correct (Lexical Analyser)

    2. Grammatically Correct (Syntax Analyser)

    3. An action must be associated (Syntax Direct Translation)

    4. It must be understood by the Executor for execution (Code generation and execution)

  * Code optimization is optional phase.

# Symbolic Assembly Language

* The most immediate step away from machine language is Symbolic Assembly Language.

* In this Language, a programmer uses <u>Mnemonic names</u> for both <u>operation Code</u> and <u>data addresses</u>.

* Thus a programmer. Could write

$$ADD \quad X, \quad Y. \longrightarrow \text{Assembly Language,}$$

instead of $\underset{\downarrow}{0110} \quad \underset{\downarrow}{001110} \quad \underset{\downarrow}{010101} \longrightarrow$ Machine Language.

* A Computer can't execute a program written in assembly Language.

* That program has to be first translated to Machine Language, which the Computer Can understand. The program that performs this translation is the ~~assembly~~. <u>Assembler.</u>

# Macros

* A Macro facility is a text replacement capability.

* There are two aspects to Macros.
  - definition
  - use.

To illustrate the utility of Macros, consider the foll. situation.

## Macro definition with two-address ADD inst.

```
MACRO      ADD2  x, y
           LOAD  y
           ADD   x
           STORE y
ENDMACRO.
```

These three Sts define the Macro. ie they give its translation.

ADD2 ⟶ name of the MACRO.

x, y ⟶ dummy arguments (formal parameters)

Having Defined ADD2 in this way, we can then use it as an ordinary assembly Lang., op-code.

## For example

If the St, ADD2 A, B is encountered somewhere after the definition of ADD2, we have a Macro Use.

Here the Macro processor substitutes for ADD2 A, B the three St, which form the defition of ADD2, but with the actual parameters A and B replacing the formal parameters x and y respectively.

That is ADD2 A,B is translated to

LOAD   B → Moves a datum from
            memory to a register.
ADD    A → ADDs the contents of
            a memory address. to
STORE  B ⌐ that of a register.
          └→. moves data from a
            register to memory.

## Bootstrapping

- a process by which Simple language is used to translate more complicated Program, which in turn may handle for ~~an even~~ More Complicated program. This Complicated Program can further handle even more Complicated prog. and so on.



fig-1

SL → Compiler → TL

This Compiler is implemented by a Language say I (Intermediate Lang.)

* Writing a Compiler for any high level Language is a Complicated Process.

* It takes lot of time to write a Compiler from Scratch.

* Hence Simple Lang., is used to generate target Code in Some Stages.

* To clearly understand the Bootstrapping tech., Consider the foll., scenario.

* The above fig-1 is represented as T diagram.

fig-2.

$$S \quad T \quad \Rightarrow \quad S_I{}^T$$

$$I$$

* Suppose we want to write a cross Compiler for new language X.

* The implementation Lang., of this Compiler is say Y.

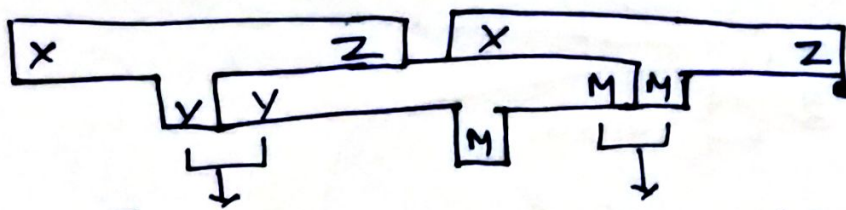* The target code being generated is in Lang., Z.

(ie) We create $\underline{XYZ}$.

* Now if existing Compiler Y runs on Machine M. and generates Code for M.

Then it is denoted as $\underline{YMM}$.

* Now if we run XYZ using YMM then we get a Compiler $XMZ$

* That means a compiler for source Lang., X that generates a target Code in lang, Z. and which Machine M.

The foll., diagram illustrates the above scenario.



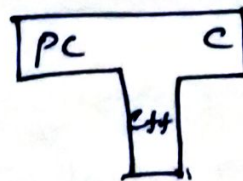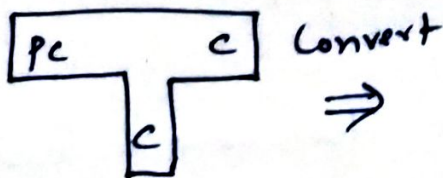These two Languages must be Same

These two Lang., must be Same.

eg

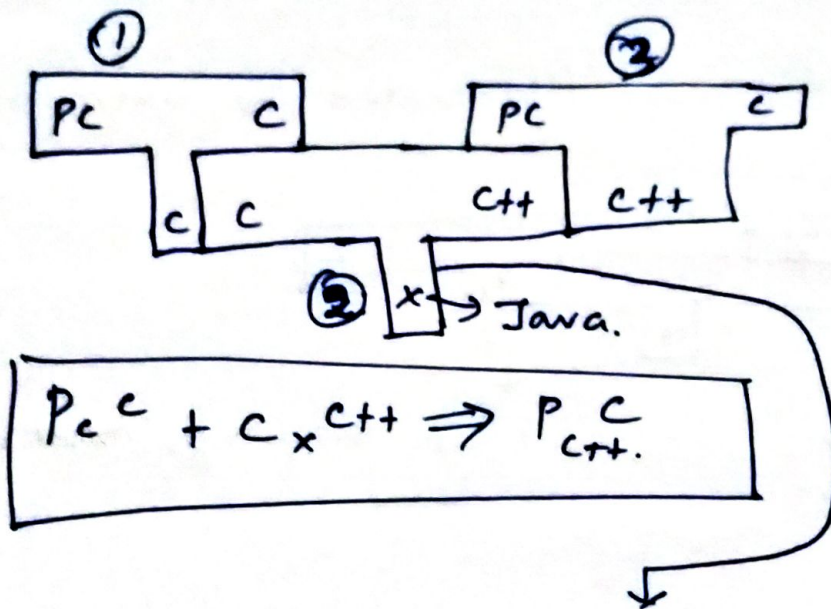We have a Pascal Translator which is written in C Lang.,

Pascal Translator — C Lang.,

I/P. Pascal code → Pascal Translator → C Lang o/P.

Create a Pascal Translator in C++

How will we do It?

Convert ⇒

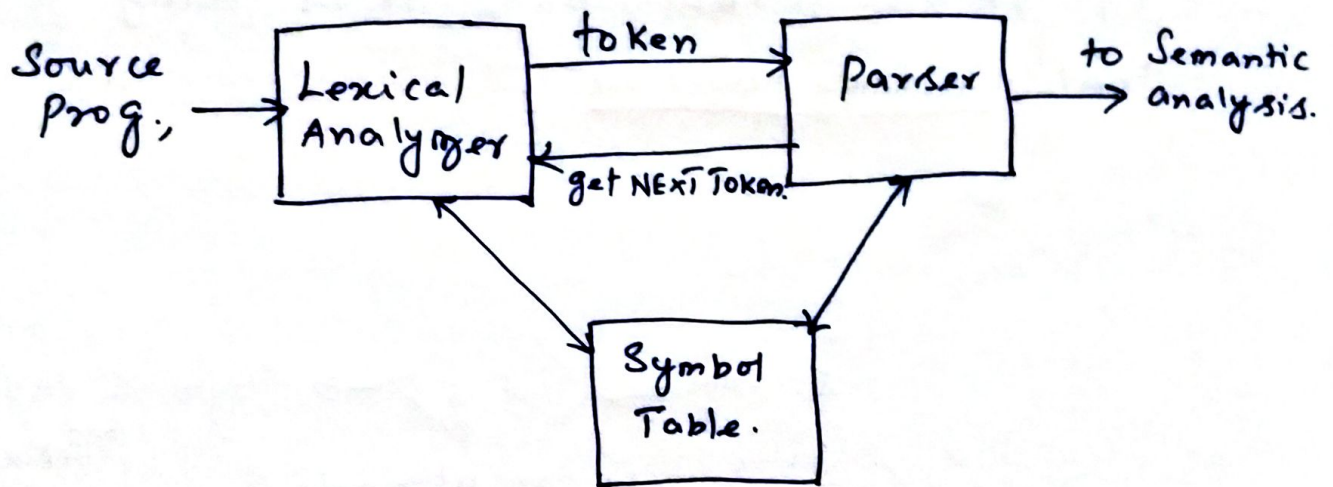Pascal translator written in c Lang and takes Pascal code as input and produce C Lang., as output.

Pascal translator written in C++ and takes Pascal code as input and Produce C Lang., as output.

$$P_c{}^c + c_x{}^{c++} \Rightarrow P_{c++}{}^c$$

To Convert First T diagram to Third T diagram, we should process Second T diagram.

# Role of Lexical Analyzer

* Reads the input characters of the Source Prog., group them into logically cohesive elements called Lexemes.

* Produces a Seq., of tokens for each lexeme in the source Program as output.

* when the lexical analyzer discovers a lexeme constituting an identifier, it enters that lexeme into the Symbol table.

① Source Prog., is the input to Lexical Analyzer.

②. Lexical Analyzer reads Source Prog., char by char until a meaningful token is found.

③ The output of the Lexical Analyser is token.

④. The Lexical Analyzer is accompanied.a Symbol Table.

⑤ whenever the LA identifies an identifier in the Source Prog., it has to check that the identifier is already existed. in the Symbol table or not.

when it reads an identifier SUM, we have to check whether the SUM is already encounted in the Prog.