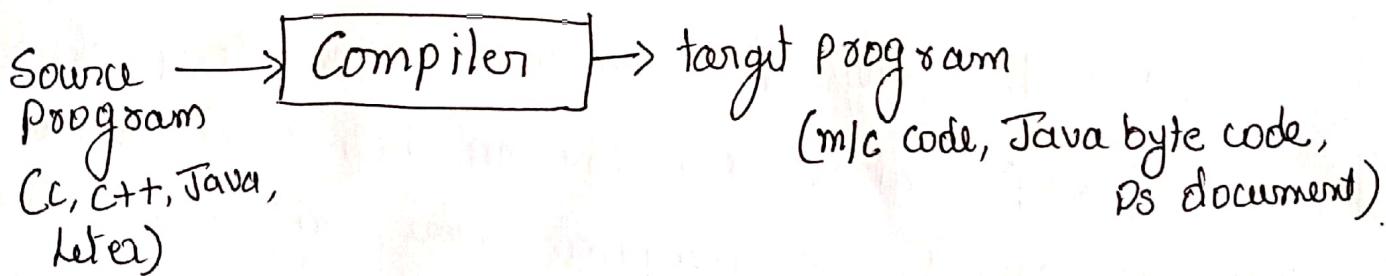


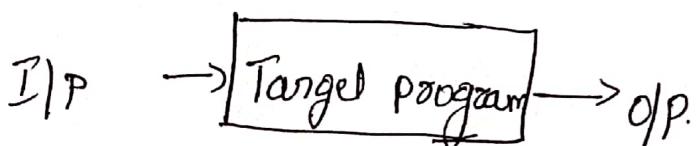
Translator: Is a program (slow) Converts a high level language program in to a Target code. Compiler & Interpreters are the translators.

Compiler: is a program that reads source program and translates into equivalent program in another language i.e target language.

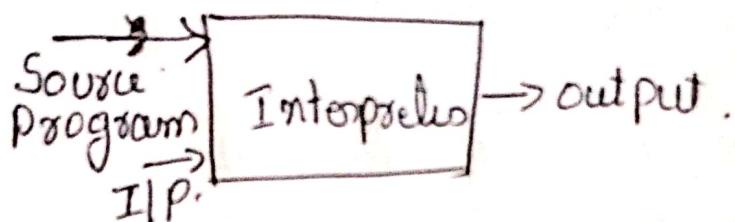
Target code can be Assembly code or machine code. An important role of the Compiler is to report any errors in the Source program that it detects during the translation process.



If the target program is an executable m/c code, it can then be called by the user to process inputs and produce outputs.



Interpreter: Is also a language processor. Instead of producing a Target Program as a translation, An interpreter appears to directly execute the operation specified in source program on inputs supplied by the user.

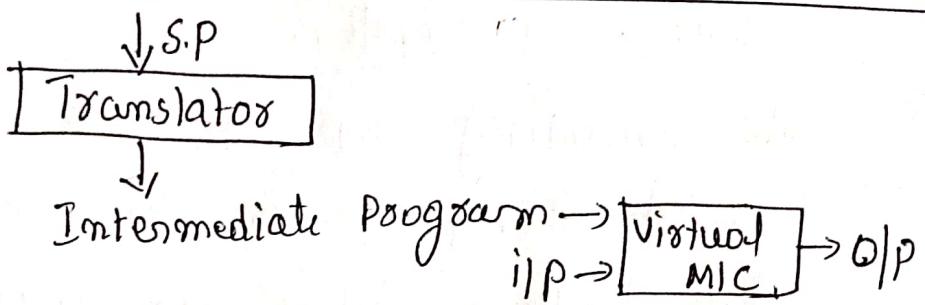


The m/c lang Target Program produced by a compiler is usually much faster than an interpreter at mapping input to output. An interpreter, however, can usually give better error diagnostics than a compiler.
It executes source program statement by statement.

Hybrid Compiler

Java language processors combine compilation and Interpretation as shown below. a Java S.P may first be compiled into an intermediate form called byte codes. The byte codes are then interpreted by a virtual m/c.

Hybrid compiler



Interpreter Vs Compiler.

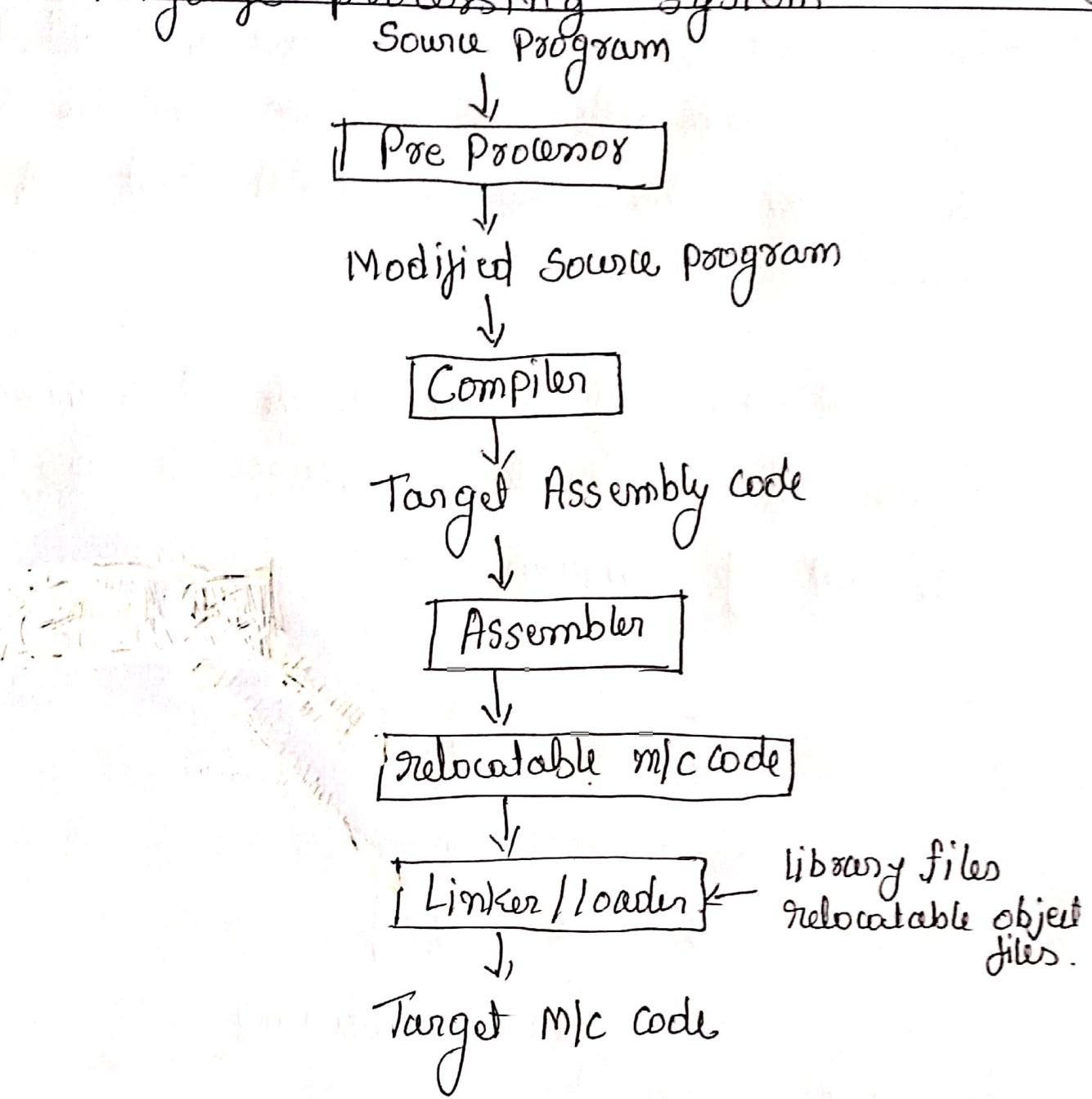
- * Compiler translates the entire program into target program as then it executes it. whereas interpreter executes instruction of sp one by one.
- * Compilers are 5 times faster than the interpreters.
- * Interpreters requires much more memory than the m/c code generated by the compiler.
- * Interpreters are easy to write and can provide better error message.
- * Compiler produces efficient object code whereas no intermediate code is generated by interpreter.
- * Difficult to use compiler for beginners compare to interpreters.

Components of Compilers

While converting S.P into T.P in addition to compiler we need other programs like preprocessor, assemblers, loaders, linkers, editors as shown below.

Preprocessor expands macros, brings different files which has different app modules of same application together. The modified S.P is then fed to a compiler. The compiler may produce an assembly language program as its output. The assemble will process this assembly code & produces relocatable m/c code.

Large programs are often compiled in pieces, so the relocatable m/c code may have to be linked together with other relocatable object files and library files into the code that actually runs on the m/c. The linker resolves external memory address, where where the code in one file may refer to a location in another file. The loader then puts together all of the executable object files into memory for execution.



Structure of a Compiler:

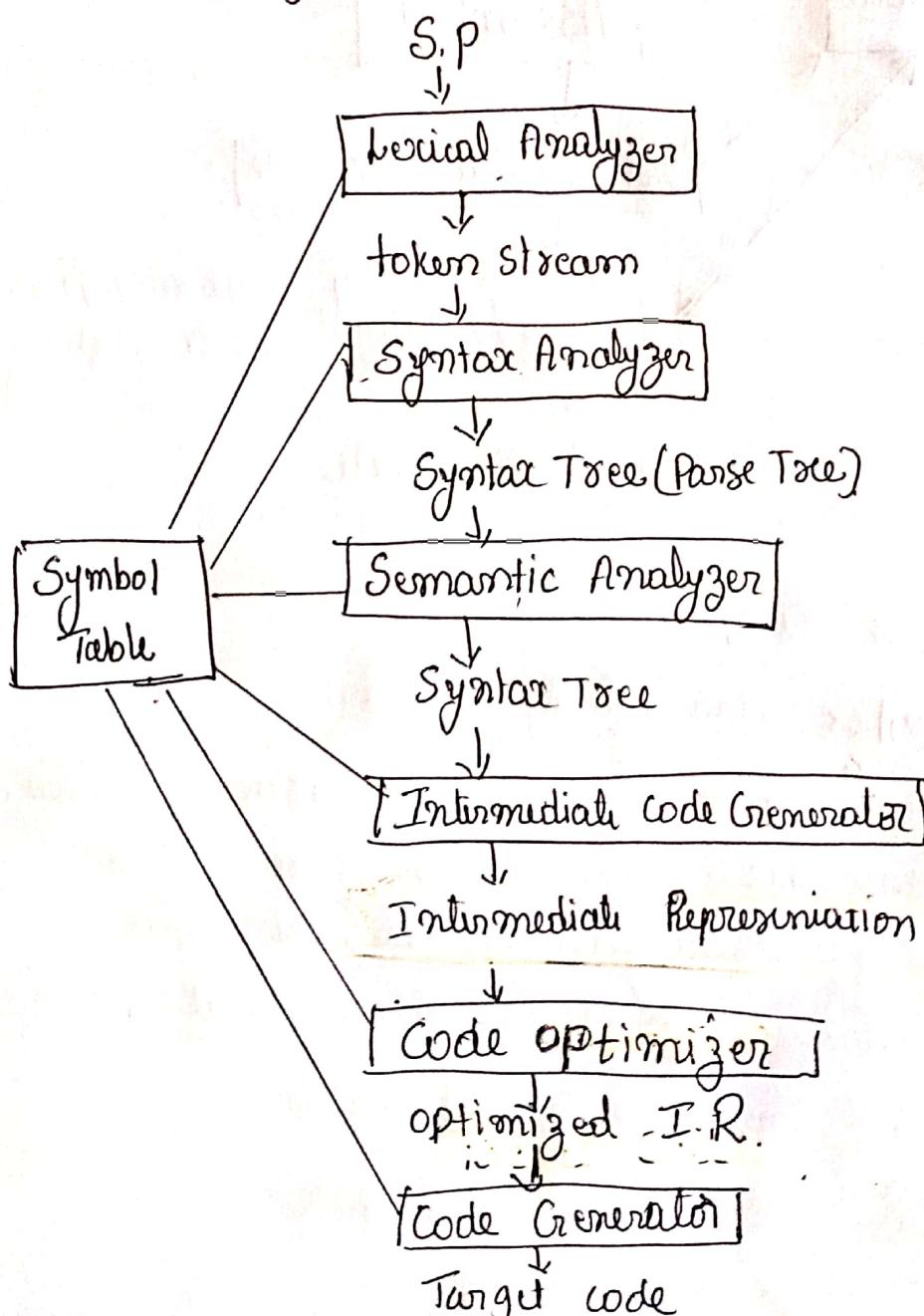
The Compiler has mainly 2 parts.

Analysis: Breaks up the S.P into constituent pieces and imposes a grammatical structure on them. Then it uses this structure to create an intermediate representation of S.P. It also collects information about program and stores it in the "Symbol table".

It is called front end of compiler.

② Synthesis Part: It constructs the desired target program from the intermediate representation and the information in the symbol table. This is called Back end of the Compiler.

The different phases of a compiler. The 1st 3 phases forms analysis part and last 3 phases forms synthesis phase. Part of Compiler.



Lexical Analysis Phase

The Lexical Analyzer reads the stream of characters from the source program and groups the characters in to meaningful sequence called lexemes (units).

For each lexeme, the Lexical Analyzer produces a token of the form <token name, attribute value> & passes it into 2nd phase of compiler i.e Syntax Analysis.

The first component of token is token name which is an abstract symbol used during parsing. The second component is attribute value that points to an entry in the symbol table for this token.

e.g: $a = b + c * 12$

<id₁, 1>, <=,>, <id₂, 2>, <+>, <id₃, 3>, <*> <60>

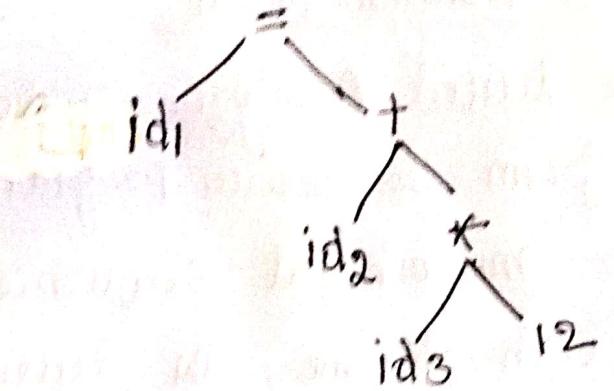
Q) Syntax Analysis (Parsing) : It uses 1st component of the token produced by lexical analyzer. to create a tree like intermediate representation that depicts the grammatical structure of the token stream. It is called Syntax tree, where internal nodes represents operation and the children nodes represents arguments of the operation. This phase is used to find syntax errors in S.P. It uses Context Free Grammar for generating parse tree.

Ex:

$$E \rightarrow E | id | E + E | EXE | num$$

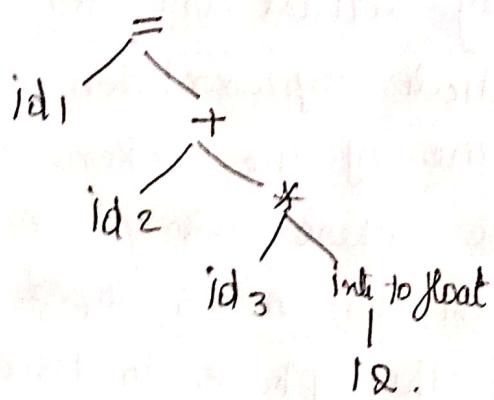
T

Grammar for Expression.



3) Semantic Analysis: This phase uses syntax tree and the information in the symbol table to check the spj. It also gathers type information semantic errors. It also gathers type information and saves it in either the syntax tree or the symbol table, for subsequent use during intermediate code generation. The important part of semantic analysis is type checking, where the compiler checks that each operator has matching operands. If the operator is applied to floating pt and integer, the compiler may convert or coerce int to floating pt numbers.

Ex:



Intermediate code generation:

After syntax & semantic analysis, some compilers generate an explicit intermediate representation

(10)

representation of the source program, which is a program for an abstract machine. This intermediate representation must have 2 properties; it should be easy to produce, & easy to translate into the target program.

The intermediate representation can have a variety of forms. one of which is an "Three address code" which is like an assembly language for a m/c in which every memory location can act like a register. Three address code consists of a sequence of instructions, each of which has at most 3 operands.
 ex. for the statement $(id_1) := (id_2) + (id_3) \times id_4$

$$\left. \begin{array}{l} \text{temp1} := \text{id}_2 \text{ to real } (\#0) \\ \text{temp2} := \text{id}_3 * \text{temp1} \\ \text{temp3} := \text{id}_2 + \text{temp2} \\ \text{id}_1 := \text{temp3.} \end{array} \right\} \rightarrow ①$$

This intermediate form has several properties.

- ① each instruction has at most one operator in addition to the assignment.
- ② compiler must generate a temporary name to

hold the value computed by each instruction.

- ③ Some "three address" instructions have fewer than 3 operands ~~or~~ 1st to last instruction

In general, these representations must do more than compute expressions; they must also handle flow-of-control constructs & procedure calls.

Code optimization:

This phase attempts to improve the intermediate code, so that fast-running machine code will result. Some optimizations are trivial ~~e.g.~~, a natural algorithm generates the intermediate code ^{as above}, using one instruction for each operator in the tree representation after semantic analysis, even though there is a better way to perform the same calculation; using the 2 instructions.

$$\begin{aligned} \text{temp1} &:= \text{id3} * \text{id2} \\ \text{id1} &:= \text{id2} + \text{temp1} \end{aligned} \quad \left. \right\} \textcircled{2}$$

There is nothing wrong with this algorithm, since the problem can be fixed during

(11)

during the code - optimization phase. i.e the compiler can deduce that the conversion of 60 from int to real representation can be done once & for all at compile time, so the int to float opt can be eliminated.

temp3 is used only once, to transmit its value to id1, it then becomes safe to substitute id1 for temp3, whereupon the last statement of above is not needed & the below code is results in (2)

[Note: It optimizes the code in terms of no. of instructions]

Code generation:

The final phase of the compiler is the generation of target code, consisting of relocatable m/c code or assembly code. Memory locations are selected for each of the variables used by the program. Then, intermediate instructions are each translated into a sequence of m/c instructions that perform the same task.

Ex. above code (2) can be translated as shown below using 2 registers R1 & R2.

LDF R₂, id₃

MOLF R₂, R₂, #12.0

LDF R₁, id₂

ADDF R₁, R₁, R₂

STF id₁, R₁

1st operand is destination & F specifies floating point number
specifies that 12.0 is a constant.

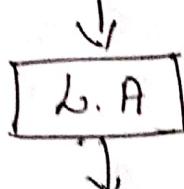
Symbol Table:

The symbol table is a data structure containing a record for each variable name, with fields for the attributes of the name such as storage allocated, its type, its scope and in case of procedure names, numbers and types of its arguments, the method of passing each argument (By value or By reference) and the type returned.

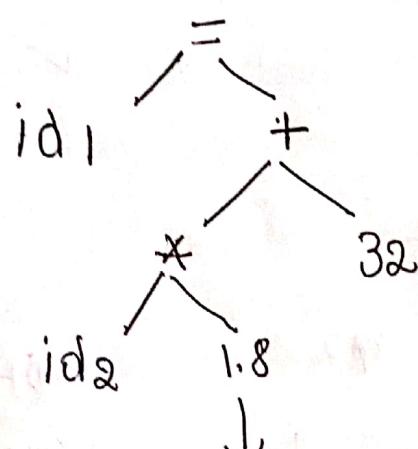
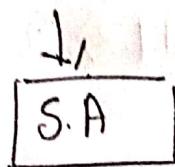
ex:-

1	a	Real	...
2	b	"	...
3	c	"	...

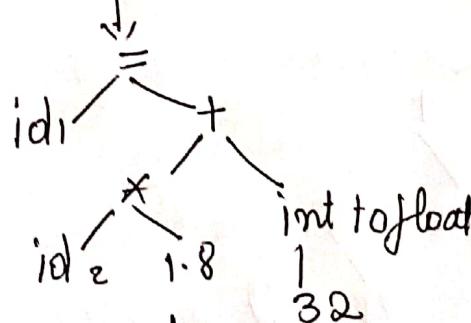
$$F = C * 1.8 + 32$$



$\langle \text{id}_1, 1 \rangle \Leftrightarrow \langle \text{id}_2, 2 \rangle \langle \times \rangle \langle 1.8 \rangle \langle + \rangle \langle 32 \rangle$



Semantic Analysis



$$\begin{aligned} t_1 &= \text{id}_2 * 1.8 \\ t_2 &= \text{int to float}(32) \end{aligned}$$

$$t_3 = t_1 + t_2$$

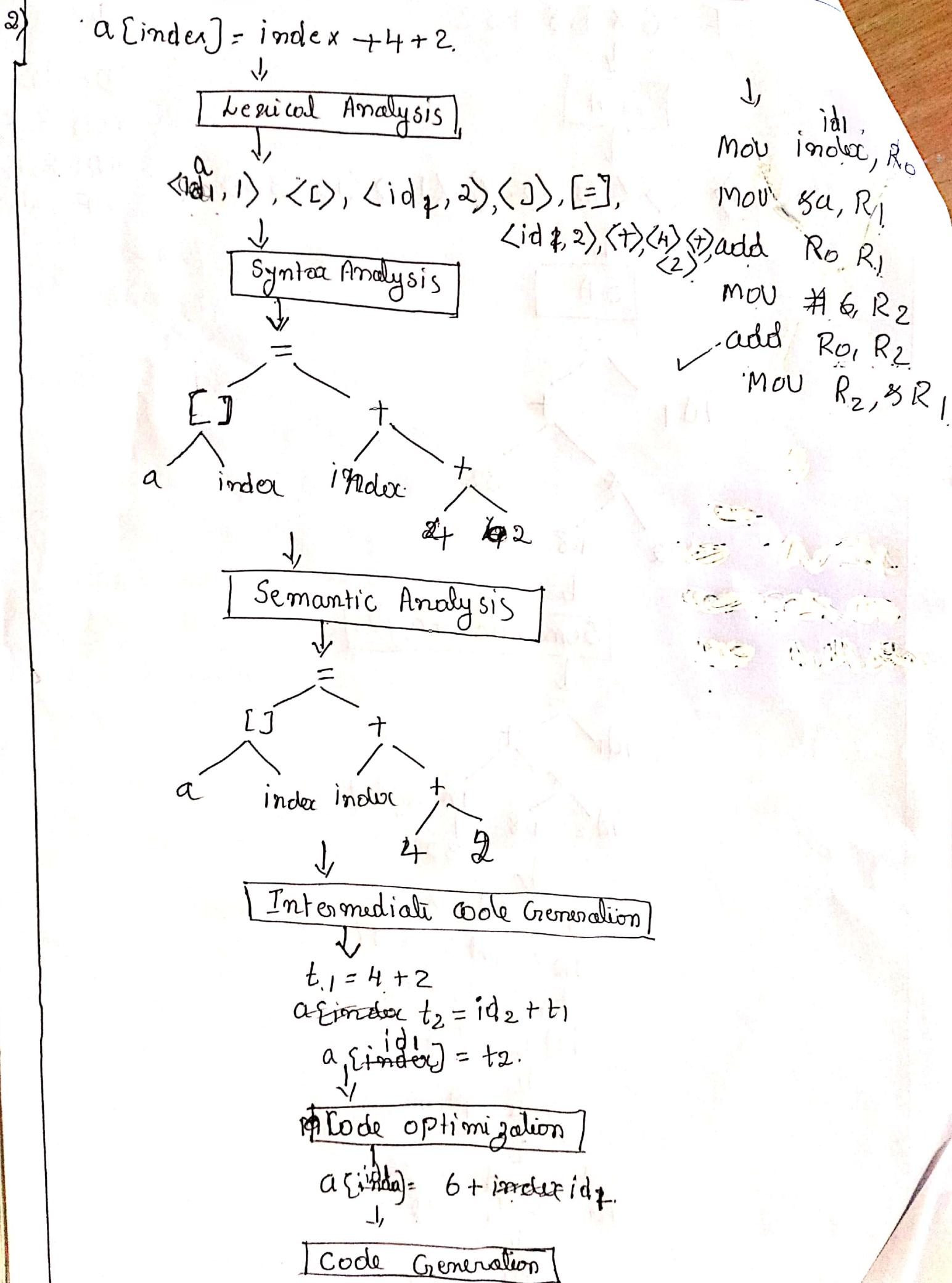
$$\text{id}_1 = t_3$$

↓
Code Optimizer

$$\begin{aligned} t_1 &= \text{id}_2 * 1.8 \\ \text{id}_1 &= t_1 + 32.0 \end{aligned}$$

↓
Code generator

↓
 LDF R₂, id₂
 MULF R₂, R₂, 1.8
 ADDF R₂, R₂, 32.0
 STF id1, R₂



Grouping of Phases into passes

The discussion of phases deals with the logical organization of a compiler. In an implementation, activities from diff phases are grouped together into pass that reads an i/p file & writes an o/p file. ex 1st 4 phases are grouped into one pass & rest as another pass.

Compiler Construction Tools

The writer (s/w developer) of Compiler as to use diff s/w tools such as lang editors, debuggers, version managers, profilers, test harnesses & so on. In addition to these general s/w development tools, often more specialized tools have been created to help implement various phases of compiler.

1. Parser generators: that automatically produce syntax analyzers from a grammatical description of a P.L.
2. Scanner generators: produces lexical analyzers from a R.E. description of the tokens of a lang
3. Syntax-directed translation engines: that produce collection of routines for walking a P.T & generating intermediate code
4. Code-generator: produces a code generator from a collection of rules for translating each opt of the intermediate lang into M/C lang
5. Data-flow analysis engines: gathers info about how values are transmitted from one part of a Prog to

each other part. It is a key part of code optimization.

6. Compiler-Construction toolkits: that provide an integrated set of routines for constructing various phases of a compiler.

The Evolution of P.L:

- * During 1950's 1st people friendly P.L were developed in our mnemonic assembly lang.
- * Then the high level lang like Fortran for scientific computation, cobol for business data processing & Lisp for symbolic computation. as been developed during half of 1950's. The philosophy behind these lang was to create higher level notations with which programmers could more easily write numerical computations, business applications & symbolic programs.
↳ symbolic programs.
- * Then usually these diff lang are classified among diff generations.
↳ 1st generation languages are m/c lang, 2nd gen the assembly lang, 3rd generation the higher level lang like Fortran, cobol, lisp, p, c, C++, Java
↳ 4th gener lang are lang designed for specific applications like NOMAD for report generating, SQL for db queries & postscript for text formatting. The 5th gen lang are logic constraint based lang like prolog & OPS5.

(2)

one more way we can classify the lang as
Imperative lang ^{like C, C++, C#} or declarative lang, object oriented
lang, ML, Haskell, prolog, C++, C#, Java,
scripting lang or JavaScript, Perl, PHP, Python,
Ruby or Tel. Programs written in scripting lang
are much shorter than equivalent prog written in
lang like C.

* The Science of Building A Compiler.

Modelling in CD & Implementation

Some of fundamental models used in CD are

- * R.E which describes lexical units of prog (k.w, identifiers etc) & also describes algorithms used by the compiler to recognize those units.
- * CFG: used to describe ^{synthetic} structure of P.L. such as the nesting of parenthesis or control constructs.
- * Trees are also imp models for representing the structure of prog & their translation into object code.

The Science of Code optimization:

The 'optimization' in CD refers to the attempt that a compiler makes to produce code that is more efficient than the obvious code.

The compiler optimization must meet the following design objectives.

- The optimisation must be correct, i.e. ^{preserve} the meaning of the compiler prog.

- The optimizer must improve the performance of many programs.
- The compilation time must be kept reasonable.
- The engineering effort required must be manageable.

The use of a rigorous mathematical foundation allows us to show that an optimization is correct so that it produces the desirable effect for all possible inputs.

The diff. models used [such as graphs, matrices, & linear programming] are necessary if the compiler is to produce well optimized code.]

Applications of Compiler Technology:

1. Implementation of High-level P.L : (C, C++, Java etc)
2. Optimizations for Computer Architectures.

④ Parallelism: It can be found at several levels at the Instruction level where multiple ops are executed simultaneously and at the processor level, where diff threads of same application are run on diff processors. All modern mp have instruction level parallelism, whereas this parallelism is hidden from the programmer. Progs are written as if all inst were executed in sequence; the h/w checks for dependencies & issues them in parallel when possible. In some cases there will be ~~the h/w~~ H/w scheduler that can change the inst ordering to increase the parallelism in the program. Whether the h/w reorders the inst or not, compiler can rearrange the inst to make instruction level parallelism more effective.

Memory Hierarchies:

b) M.H consists of several levels of storage with diff speeds & sizes. with the level closest to the processor being the fastest but smallest.

A processor usually have small no. of registers consisting of hundreds of bytes, several levels of caches containing kilobytes to megabytes, physical memory containing mega bytes to gigabytes & finally secondary storage that contains gigabytes & beyond. Correspondingly the speed of access b/w adjacent levels of the hierarchy also differ. While compilers traditionally focus on optimizing the processor execution, more emphasis is now placed on making the memory hierarchy more effective.

③ Design of New Computer Architecture

In the early days of computer architecture design, compilers were developed after the m/c were built.

Now, that has been changed because the performance of a computer system is determined not by its raw speed but also by how well compilers can exploit its features. Thus in modern computer architecture development, compilers are developed in the processor design stage, before: of how compiler influenced the " " of computer Architecture was the invention of the RISC (Reduced Instruction Set Computer) architecture over last 3 decades, many architectural concepts have

been proposed like d/a flow m/c's, vector m/c, Every long Inst word) m/c, SIMD (single inst, multi processor with shared or distributed memory d/a multiprocessor with shared or distributed memory

The development of each of these architectural concepts was accompanied by the research & development of corresponding compiler tech.

(A) Program Translation:

The compiler not only used to convert H.L.L to LLL it is also used to translate b/w diff kinds of lang. ex.

(B) Binary Translation:

Compiler Tech can be used to translate the binary code for one m/c to that of another, allowing a m/c to run prog originally compiled for another inst set. Binary Translation Tech has been used by various computer companies to increase the availability of SW for their m/c. In particular, because of the domination of the x86 PC market, most SW titles are available as X86 code. Binary Translators have been developed to convert X86 code into both Alpha & Sparc code.

(C) H/w Synthesis:

How the ~~diff~~ most SW's written in HLL, even H/w design are mostly described in HLL like Verilog & VHDL techniques to translate design at higher levels, such as its behavior or J/n but also exist.

(D) Database Query Interpreters:

or SQL are used to search db. DB queries consist of predicates containing relational & boolean operators.

they can be interpreted or compiled into commands to search a db for records satisfying that predicate.

(4)

compiled simulation:

Simulation is a general tech used in many scientific & engineering disciplines to understand a phenomenon or to validate a design. Inputs to a simulator usually include the description of the design or specific ip parameters for that particular simulation run. Simulations can be very expensive & takes more time. Instead of writing a simulator that interprets the design, it is faster to compile the design to produce m/c code that simulates that particular design natively. Compiled simulation can run orders of magnitude faster than an interpreter-based approach. Compiled simulation is used in many state-of-the art tools that simulate designs written in verilog or VHDL.

(5) SW productivity Tools:

The diff prog we write may contain diff errors. Testing is the primary kic for locating errors in programs. one more approach is to use data flow analysis to locate errors statically. It finds errors in all possible ^{executing} paths. Many of the data flow analysis techniques, originally developed for compiler optimization, can be used to find errors like type checking, bounds checking, memory management errors.

Programming Lang Basics

① Parameter Passing Mechanisms:

Call-by-value: Here, the actual parameter is evaluated (if copied) or copied (if variable). The value is placed in the location belonging to the corresponding formal parameters of the called procedure. This method is used in C, C++ & Java.

Here changes made to the formal parameters does not effect the actual parameters. But if we pass name of array like a then if we do $a[i]=2$ it also changes $a[i]=2$ in actual parameter. Array name is p'th same case in Java when we pass object name.

Call-by-reference: Here the address of the actual parameter is passed to the callee as the value of the corresponding formal parameter. The changes made to the formal parameters also effects the actual parameters. here the address of actual parameters are sent.

If the actual parameter is an expression, however then the exp is evaluated before the call, so its value stored in a location of its own. changes to the formal parameters change the value in its location, but can have no effect on the d/a of the caller.

call by name: A 3rd mechanism - call-by-name was used in the early programming lang Algol 60. It requires that the callee execute as if the actual parameters were substituted literally for the formal parameters in the code of the callee, as if the formal parameters were a macro standing for the a.p (with renaming of local names in the called procedure, to keep them distinct) when the A.p is an expression rather than a variable, some unintuitive behaviour occurs which is one reason this mechanism is not used.

Programming ~~Lang~~ Basics

* Static/ Dynamic Distinction.

→ with in class.

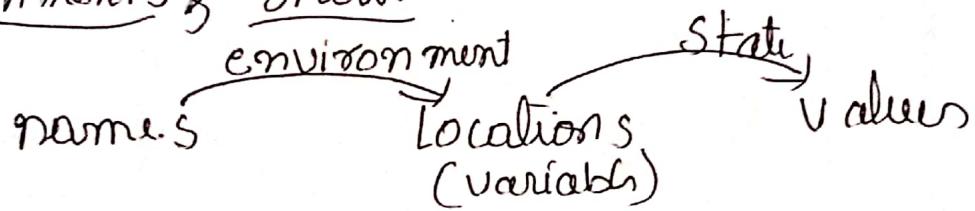
Public static int x;

Static ^{scope} means only once the memory is allocated & it is fixed any updatations we are doing will affect only that location.

Dynamic ~~scope~~ scope means several declaration of x.

The objects of class will have diff ^{copies of} x. ~~locations~~

Environments & States



When program runs we need to check whether changes it affects the value of data elements or affects names of d/a. ex $x = y + 1$ changes the value denoted by x.

environments change according to the scope rules of a lang.

```

int i; //global
void f()
{
    int i; //local
    i=3; //use of local i
}
x = i+1; //use of global i
  
```

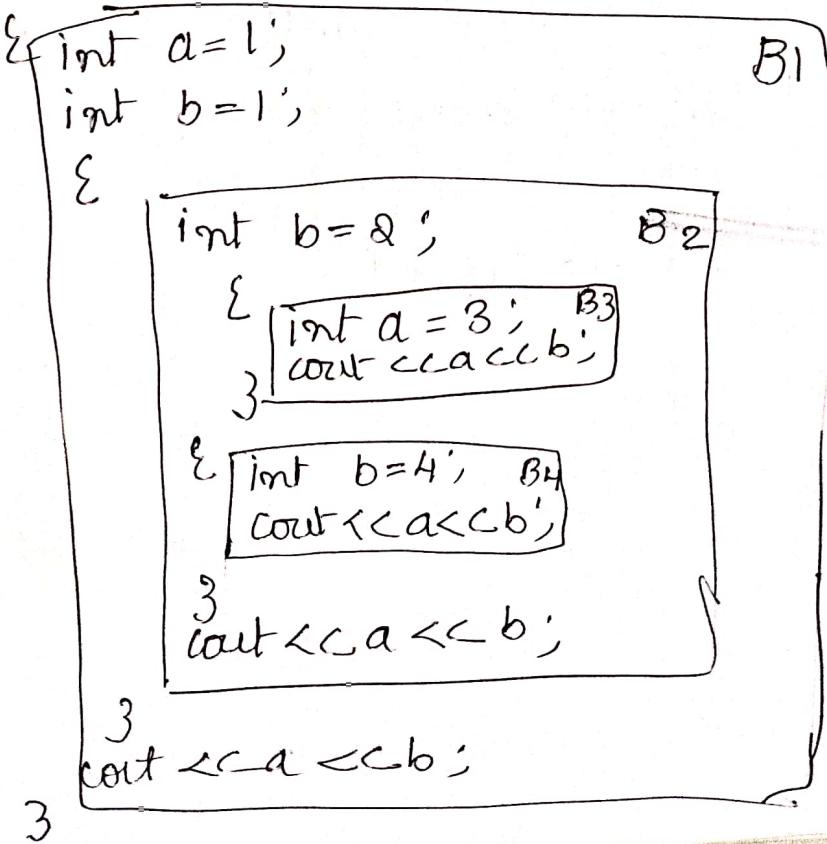
environment & static mapping are usually dynamic. Only if we used static keyword or global variable then environment mapping is static. & the static mapping is dynamic 'cause we can't tell the value in a location until we run the program. But if we have declaration like #define Arraysize 1000

This is static static mapping.

Static scope & Block structure:

blocks in C++ are represented by { & }. The block is the grouping of declarations & statements.

main()



Declaration Scope

int a=1;	B ₁ -B ₃
int b=1;	B ₁ -B ₂
int b=2;	B ₂ -B ₄
int a=3;	B ₃
int b=4;	B ₄

Exercise : 1.6.1 For block structured code below indicate the values assigned to ~~w, x, y, z~~ w, x, y, z.

```

int w, x, y, z;
int i=4; int j=5; □
{
    int j=7;
    i=6;
    w = i+j; ⑬
}
x = i+j; ⑭
{
    int i=8;
    y = i+j; ⑮
}
z = i+j; ⑯⑰

```

```

int w, x, y, z;
int i=3; int j=4;
{
    int i=5; ⑯
    w = i+j;
}
x = i+j; ⑰
{
    int j=6;
    i=7;
    y = i+j; ⑮
}
z = i+j; ⑯

```

```

#define a(α+1)
int x=2;
void b() { int x=1; printf("y.d\n", a); }
void c() { printf("y.d\n", a); }
void main() { b(); c(); } *

```

O/P 2, 3.

Lexical Analysis

(21)

This chapter deals with techniques for specifying & implementing Lexical Analyzers. A simple way to build L.A. is to construct a diagram that illustrates the structure of the tokens of the source lang. & then to hand-translate the diagram into a program for finding tokens.

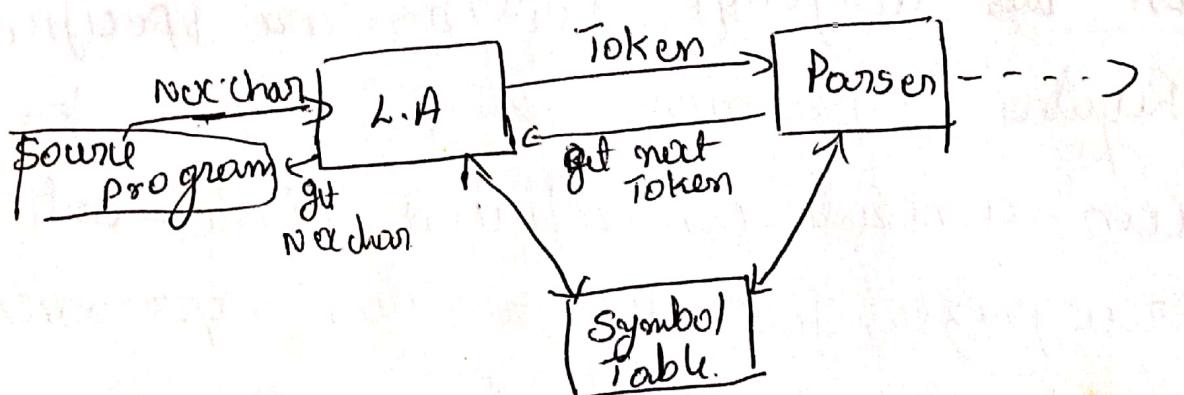
Since pattern directed programming is widely useful, we introduce a pattern-action language called lex for specifying lexical analyzers. In this language, patterns are specified by regular expressions, & a compiler for lex can generate an efficient Finite-automaton recognizer for the regular expressions.

A software tool lexical-analyzer generator can be used for constructing the lex L.A. which uses best known pattern-matching algorithms & thereby creates efficient L.A.'s for people who are not experts in pattern-matching techniques.

The Role of the Lexical Analyzer:

It is the 1st phase of a compiler, which reads the I/O characters & produce as output a sequence of tokens that the parser uses for syntax analysis. usually it is implemented as a subroutine or co-routine of a Compiler Parser. Upon receiving a "get next token" command from the parser, the L.A reads the I/O characters until it can identify the next token.

Interaction b/w L.A & Parser



Since L.A is a part of Compiler that reads the source text, it may also perform certain secondary tasks at the user interface like skipping out from the S.P comments, white space in the form of blanks, tab to newline characters, keeps track of .number of new

The programming language is defined by

1. Syntax

Decides whether sentence in a lang is well formed.

2. Semantics

Determines the meaning, if any, of a syntactically well formed sentence.

3. Grammer:

A formal sm that provides a generative finite description of the lang

Lexical Analyzers or Parsers of a compiler handle the syntax of the programming language.

Token examples

Let us consider the program segment.

```
void main()
{
    printf ("Hello world\n");
}
```

The tokens of this program segment are:

1. void,
2. main,

Token Type

ID

NUM

REAL

~~KEYWORDS~~

~~SYMBOLS~~

Example

foo, ~~114~~, a, temp...

73, 0, 00, 515, +2...

66.1, .5, 10, 1e67 5.5e-10...

IF DO WHILE INT...

, (comma), !=(not equal), ((paran)...

How are tokens formally defined & Recognized?

By using R-E to define a token as a formal Regular language.

A language in real life is made of

1. words made up of alphabets &

2. sentences made up of words arranged according to the grammar of that language.

~~App~~ Formality of languages.

(E) Alphabets \rightarrow A Finite (non-empty) set of symbols

(A) String \rightarrow A finite sequence of symbols from an alphabet which includes even the empty sequence. (E)

language \rightarrow A set of finite strings. The set of all possible finite strings of elements of alphabet Σ (including) is denoted by Σ^* .

Finite specification of language is possible with
be

1. Automaton: a Recognizer; a m/c that accepts all strings in a language (or rejects all other strings)

2. Grammar: a generator; a spm for producing all strings in the lang (or no other strings)

Note: A lang may be specified by many diff grammar
to automato But. A grammar or automato
specifies only one lang.

line characters seen so that a line number can be associated with an error message.

In some compilers L.A is in charge of

~~making a copy of S.P with the errors msg's marked in it. If the S.P supports some macro processor jmp, then these predecessors of jⁿ may also be implemented as L.A analysis takes place~~

L.A are divided into a cascade of 2 phases, "Scanning" & "Lexical Analysis": where Scanner performs simple jⁿ like eliminate blanks from I/P & L.Analysis generates tokens.

Issues in Lexical Analysis:

There are several reasons for separating the analysis phase of compiling into L.A & P Lexical Analysis & Parsing.

1. Simplify the Design of Phases: This is a most imp consideration of separation which makes the phases simple ex: a parser embodying the conventions for comments &

while space is significantly more complex one that can assume comments & white space have already been removed by a L.A.

2. Compiler efficiency is improved:

A separate L.A. allows us to construct a specialized & potentially more efficient process for the task. A large amount of time is spent reading the S.P. & partitioning it into tokens. Specialized buffering tech. for reading I/P characters & processing tokens can significantly speed up the performance of a compiler.

3. Compiler portability is enhanced:

I/P alphabet peculiarities & other device-specific anomalies can be restricted to the lexical analyzer. The representation of symbols ↑ in pascal can be isolated in the L.A.

Tokens, Patterns, lexemes:

In general, there is a set of strings in the IIP for which the same token is produced as O/P. This set of strings is described by a rule called a pattern associated with the token. The pattern is said to match each ~~each~~ string in the set. A lexeme is a sequence of characters in the SP i.e. matched by the pattern for a token. ex. in Pascal

Const pi = 3.1416; [^{printf ("Total=%d\n",} _{ID} ^{score.} _{ID})]

The subString pi is a lexeme for the token "identifier".

Token	Sample lexeme	Informal description of Pattern
const	const	const
if	if	if
relation	<, <=, >=, =, <?, >	< > <= > = <? > = < ? >
id	pi, count, D2	letter followed by letters & digits
num	3.14, 0, 10, 6.4	any numeric constant
literal	"core dumped"	any character b/w " " except "

Tokens are treated as terminal symbols in the source long, using boldface names to represent tokens. The lexemes i.e strings of characters in the S.P that can be treated together as a lexical unit.

In most P.L the following constructs are treated as tokens: Keywords, operators, identifiers, constants, literal strings, & punctuation symbols such as parentheses, commas, & semicolons. In the example above when the character sequence pi appears in the S.P a token representing an identifier is returned to the parser. The returning of a token is often implemented by passing an integer corresponding to the token. It is this integer i.e referred to in fig 8.2 as boldface id.

A pattern is a rule describing the set of lexemes that can represent a particular token in S.P. To represent ^{Pattern of} some tokens like num, identifiers we use R.E.

management of tokens may be imp in b. Analysis. (94)

Blanks are treated in different ways in diff lang. & it will complicate the task of identifying tokens.

ex: In Fortran

DO 5 I = 1, 25

we can't tell until we have seen the decimal point that DO is not a k.w. but rather part of the identifier DO5I. On the other hand

DO 5 I = 1, 25

we have 7 tokens, corresponding to the k.w. DO, the statement label 5, Identifier I, operator =, Constant 1, the comma & the constant 25. Here we can't be sure until we have seen the comma that DO is a k.w.

Attributes for tokens:

when more than one pattern matches a lexeme, the lexical analyzer must provide additional info. about the particular lexeme that matched to the subsequent phases of the compiler. ex, the pattern num

matches both the strings or !, but it is essential for the code generator to know what string was actually matched.

The lexical analyzer collects info. about tokens into their associated attributes. A token has only a single attribute i.e. pointer to the symbol table entry in which the info. about the token is kept; the pts becomes the attribute for the token. We may be interested in both the lexeme for an identifier or the line number on which it was first seen. Both these items of information can be stored in the symbol table entry for the identifier.

ex. The tokens & associated attributes values for the fortran statement.

$$E = M * C * \underline{2}$$

\langle id, pts to S.T entry for E \rangle

\langle assign-op, \rangle

\langle id, pts to S.T entry for M \rangle

\langle mult-op, \rangle

<id, ptr to S.T entry for c>

<exp- op>

<num, integer value 2>

Note that in certain pairs there is no need for own attribute value; the 1st component is sufficient to identify the token.

Lexical Errors:

Few errors can't be detected by lexical level alone, as L.A has a very localized view of a S.P. If the string fi is encountered in a C program for the first time in the context

fi(a==f(x))

A.L.A can't tell whether fi is a misspelling of the k.w if or an undeclared identifier. Since fi is a valid identifier L.A must return tokens for identifiers so that some other phases of compiler handle any errors.

Suppose L.A is unable to proceed; none of the patterns for tokens matches a prefix of the remaining i/p.

The simplest recovery method is "Panic in recovery, we delete successive characters from the remaining i/p until the LA finds a well formed token." One Progress

Other possible error-recovery actions are

1. deleting an extraneous character
2. Inserting a missing character
3. Replacing an incorrect character by a correct character.
4. Transposing two adjacent characters.

error transformations like these may be tried in our attempt to repair the i/p. The simplest method is to see whether a prefix of the remaining i/p can be transformed into a valid

name by just a single error transformation. This strategy assumes most lexical errors are the result of a single error transformation, an assumption usually but not always, borne out in practice.

The way of handling the errors in a program is to compute the minimum no. of error transformations required to transform the erroneous prog into one i.e syntactically well-formed. (26)

Input Buffering:

There are 3 general approaches to the implementation of a L.A.

1. Use a L.A generator, S.A box Compiler to produce the L.A from a R.E based specifications. Here generator provides routines for reading & buffering the i/p.
2. write the L.A in a conventional sm lang - common language, using the I/O facilities of that lang to read the i/p.
3. write the L.A in assembly lang & explicitly manage the reading of i/p.

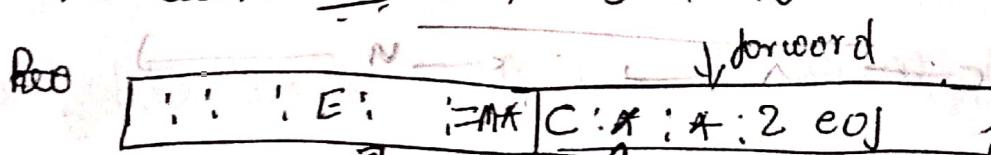
Since L.A is the only phase of the compiler that reads the sp character by character, it

is possible to spend a considerable amount of time in the lexical analysis phase, even though the later phases are conceptually more complex. Thus the speed of L.A. analysis is a concern in compiler design. We use first method i.e. generators to implement L.A.

Buffer Pairs:

For many source languages, there are times when the L.A. needs to look ahead several characters beyond the lexeme for a pattern before a match can be announced. There are diff specialized buffering techniques available but we will see one such technique.

We use a buffer divided into $\frac{N}{2}$ -character halves, where N is the number of characters in one disk block e.g. 1024 or 4096.



Read $\frac{N}{2}$ characters into each half of the buffer with one read command instead of using read command for each I/P character.

If fewer than N characters remain in the i/p then a special character eos is read into the buffer after the i/p characters.

Two pts to the i/p buffer are maintained. The string of characters b/w the 2 pts is the current lexeme. Initially both pts point to the 1st character of the next lexeme to be found. One called forward pt, scan ahead until a match for a pattern is found, once the next lexeme is found the forward pt is set to the character at its right end. After the lexeme is processed, both pts are set to the character immediately past the lexeme, with this scheme the comments & the white space can be treated as patterns that yield no token.

If the forward pt is about to move past the half way mark, the right half is filled with N new i/p characters.

If the forward ptr is about to reach the right end of the buffer, the left half is filled with N new characters as the forward ptr wraps around to the beginning of the buffer.

This buffering scheme works quite well most of the time, but with it the amount of lookahead is limited, as this limited lookahead may make it impossible to recognize tokens in situations where the distance that the forward ptr may travel is more than the length of the buffer. ex - if we see

DECLARE (ARG1, ARG2... ARGn);

in a PL/I program, we can't determine whether Declare is a k.w or an array name until we see the character that follows the right parenthesis. In either case, the name ends at the second E, but the amount of lookahead needed is proportional to the number of arguments,

which in principle is unbounded.

Code to advance forward pt is

if forward at end of first half then begin

reload second half;

forward := forward + 1

end

else if forward at end of 2nd half then begin

reload 1st half

move forward to beginning of
1st half.

end

else forward := forward + 1

Sentinels:

If we use the above scheme, then we must check each time we move the FP that we have not moved off one half of the buffer; if we do, then we must reload the other half, i.e. our code for advancing the forward pts performs tests like above.

Except at the ends of the buffer we
the code above requires 2 tests for each
advance of the forward pts. we can
reduce the 2 sets to one if we extend
each buffer half to hold a sentinel
character at the end. ~~The~~ ^{we} i.e specify
all characters like ej.

came to begin
for was

↓ ↓
`* : : E : : = : M : A : ej) | C : * ; * ; 2 ; ej ;
: : ! : ej)`

with this arrangement of ~~of~~ we can use below
code to advance the F.p.

forward := forward + 1;

if forward \neq ej then begin

if forward at end of 1st half then
begin

reload 2nd half;

forward := forward + 1

end

else if forward at end of 2nd half
then begin

reload 1st half;

move forward to beginning of 1st half

end

else /x ej within a buffer signifying end of if,

terminate L. Analysis

we also need to decide how to process the character scanned by the forward ptr; does it mark the end of a token, does it represent progress in finding a particular keyword or what? One way to structure these tests is to use a case statement, if the implementation ~~language~~ has one. The test

- i) forward \neq eof

can be implemented as one of the diff cases.

Specifications of tokens

Patterns are represented by Regular Expression,

Some Basic formal definitions of language

Alphabet [character] (Σ)

It is a finite (non empty) set of symbols, where symbols can be ~~characters~~ ^{digits or punctuations} or letters.

The Set Σ_0, Σ_1 is the binary alphabet.

ASCII & ~~EBCDIC~~ EBCDIC are 2 examples of computer alphabets

(1) String: It is a finite sequence of symbols from an alphabet. It is also known as sentence or word. $l(S)$ is the length of string i.e. the number of occurrences of symbols in S . ex banana, $l(\text{banana}) = 6$. The empty string, denoted by ϵ , is a special string of length zero.

Language $\{\epsilon^*\}$

A set of finite strings over some fixed alphabets. The set of all possible finite strings of elements of alphabet Σ is denoted by Σ^* . The Abstract lang. like \emptyset , the empty set, or $\{\epsilon\}$, the set containing only the empty string, are called.

If x & y are strings, their concatenation is given by xy . ex $x=\text{dog}$ $y=\text{house}$ then $xy=\text{doghouse}$. The empty string is the identity element under concatenation.

$$\text{i.e. } S\epsilon = \epsilon S = S.$$

If we think concatenation as a product then string "exponentiation" is defined as follows. Define S^0 to be ϵ & for $i > 0$ define S^i to be $S^{i-1}S$. Since ϵS is S itself, $S^1 = S$. Then $S^2 = SS$ $S^3 = SSS$ & so on.

Example of language

1) language of all strings consisting of n 0's followed by n^{th} for some $n \geq 0$ is

$$\{ \epsilon, 01, 001, \dots \}$$

2) \emptyset is empty language, but $\emptyset \neq \{\epsilon\}$ where former is empty language & later is a language with empty string.

3) set of binary numbers whose value is a prime.

$$\{10, 11, 101, 111, 1011, \dots \}$$

Some common terms associated with parts of a string are given below.

<u>Term</u>	<u>Definitions</u>
-------------	--------------------

Prefix of s A string obtained by removing zero or more leading symbols from the ^{end of} string ex. ban is a prefix of banana.

Suffix of s A string obtained by removing zero or more trailing symbols of s. ex. ana is suffix of banana.

substring of S : A string obtained by deleting some part of S.

Prefix or a suffix from S; ex noun is a prefix of banana. every Prefix or every suffix is a substring, but not every substring of S is a prefix or a suffix of S. for every string S, both ϵ & S & λ are prefix, suffixes & substrings of S.

proper prefix, suffix or substring of S

Any non empty string x i.e. respectively, a prefix, suffix or substring of S. $S \neq x$.

Subsequences of S: Any string formed by deleting zero or more not necessarily contiguous symbols from S. ex abada is a subsequence of banana.

operations on languages:

Operations	Definitions
Union of $L \otimes M$ ($L \cup M$)	$L \cup M = \{S \mid S \text{ is in } L \text{ or } S \text{ is in } M\}$
Concatenation of $L \otimes M$ ($L \cdot M$)	$L \cdot M = \{St \mid S \text{ is in } L \text{ & } t \text{ is in } M\}$
Kleene closure of L (L^*)	$L^* = \bigcup_{i=0}^{\infty} L^i$ [L^* denotes "zero or more concatenations of L "]
+ve closure of L (L^+)	$L^+ = \bigcup_{i=1}^{\infty} L^i$ [L^+ denotes "one or more concatenations of L "]

Example Let $L = \{A, B, \dots, Z, a, b, \dots, z\} \times D = \{0, \dots, 9\}$

Some new languages created from $L \otimes D$ by applying different operators are

1. $L \cup D$ is the set of letters & digits

2. LD is the set of strings consisting of a letter followed by a digit.

3. L^4 is the set of all 4 letter strings

4. L^* is the set of all strings of letters including ϵ , the empty string.

5. $L(L \cup D)^*$ is the set of all strings of letters & digits beginning with a letter.

6. D^+ is the set of all strings of one or more digits.

Regular Expressions: The text strings can be represented by Regular Expressions. R.E can be derived by Finiti Automata but style & notations are different.

The R.E are built another type of language defining notation. R.E n denotes a language $L(n)$.

R.E are built recursively out of smaller R.E. using below rules.

Basis: 1) \emptyset is a R.E denoting an empty language $L(\emptyset) = \emptyset$.

2) G is a R.E denoting $L(G) = \{\epsilon\}$
i.e lang with empty string.

3) a is a R.E denoting lang containing only $\{a\}$ i.e $L(a) = \{a\}$

Induction:

Suppose r & s are R.E denoting $L(r)$ & $L(s)$ respectively then

(32)

- 1) $(r) | (s)$ is a R.E denoting the lang $L(r) \cup L(s)$
- 2) $(r)(s)$ " " " " $L(r)L(s)$
- 3) $(r)^*$ " " $(L(r))^*$

4. (r) is an R.E denoting $L(r)$. it says we can add additional pairs of Parentheses around exp. without changing the lang they denote.

R.E may also contain some unnecessary Pairs of parentheses. we may drop certain Pairs of parenthesis by adopting some conventions like

- 1) $*$ has highest precedence as is left associative
- 2) Concatenation has 2nd highest precedence \rightarrow is left associative
- 3) | has lowest precedence as left associative

$$\text{ex: } (a) | (b)^*(c) \Rightarrow a | b^*c$$

both represent same lang i.e string that are sigl either a single a or 3 or more b's followed by one c.

Let $\Sigma = \{a, b\}$

1. The R.E a/b denotes the lang $\{a, b\}$
2. $(a/b)(a/b)$ denotes $\{aa, ab, ba, bb\}$ lang of all string of length 2 over Σ .
3. a^* denotes the lang consisting of all strings of 0 or more a's i.e $\{\epsilon, a, aa, \dots\}$
4. $(a/b)^*$ denotes the set of all string consisting of zero or more instances of a or b i.e all string of a's & b's $\{\epsilon, a, b, aa, ab, bb, ba, aaa, \dots\}$
OR $(a^*b^*)^*$
5. a/a^*b denotes the lang $\{a, b, ab, aab, aaab, \dots\}$ i.e the string of all strings consisting of zero or more a's & ending in b.

A lang that can be defined by R.E is called Regular Set. Suppose if 2 R.E γ & δ defines same lang then we say both are equivalent $\gamma = \delta$ or $(a/b) = (b/a)$

Regular Definitions:

we can use ~~variables~~ to certain R-E.

$id \rightarrow \text{letter}(\text{letter} | \text{digit})^*$

$\text{letter} \rightarrow A | B | \dots | z | a | b | \dots | z | -$

$\text{digit} \rightarrow 0 | 1 | \dots | 9$

unsigned numbers (integers or Fpt) are strings
such as 5280, 0.01234, 6.336E4 or 1.8E-4

The Regular definitions.

$\text{digits} \rightarrow \text{digit} \text{digit}^*$

~~dig~~ optional fraction $\rightarrow . \text{digits} | \infty$

optional Exponent $\rightarrow (\text{E} (+|-|) \text{digit}) \text{E} \infty$

number $\rightarrow \text{digits} \text{optional fraction} \text{optional}$

1. is not valid.

Exponent

using shorthands we can rewrite the regular def as below

$\text{letter} \rightarrow [A-Za-z-]$

$\text{digit} \rightarrow [0-9]$

$id \rightarrow \text{letter}(\text{letter} | \text{digit})^*$

$\text{digit} \rightarrow [0-9]$

$\text{digits} \rightarrow \text{digit}^*$

$\text{number} \rightarrow \text{digits} (. \text{digits})? (\text{E} (+|-|)? \text{digit})?$

Recognition of tokens:

Take the patterns for all needed tokens & build
piece of code that examines the i/p string & finds
a prefix that is a lexeme matching one of the
patterns.

e.g.: $\text{stmt} \rightarrow \text{if expn then stmt}$

$| \text{if expn then stmt else stmt}$

$\text{expn} \rightarrow \text{term relop term}$

$| \text{term relop term}$

$\text{term} \rightarrow \text{id}$
 number .

The terminals of the grammar, which are if, then, else,
relop, id & number are the names of tokens as far
as the L.R.0 is concerned.

Patterns for

$\text{if} \rightarrow \text{if}$

$\text{else} \rightarrow \text{else}$

$\text{then} \rightarrow \text{then}$

$\text{relop} \rightarrow \{ > | <= | = | < \}$

Lexical Analyzer skips the white spaces by
recognizing the "tokens" ws defined by:

$\text{ws} \rightarrow (\text{blank} | \text{tab} | \text{newline})^+$

Here blank, tab & newline are abstract symbols that we use to express the ASCII characters of the same names. When LA reads the token as it will not return it to the parser, but instead restarts the lexical Analysis from the character that follows the whitespace.

Transition diagram { Finite Automata }

Before constructing an L.A. first we have to convert the patterns (R.E) into transition diagrams which are deterministic i.e. only one edge for one input from each state. (NFA)

The T.D. have a collection of states (nodes) each state represents a condition that could occur during the process of scanning the i/p looking for a lexeme that matches one of the several patterns. We may think of a state as summarizing all we need to know about what characters we have seen b/w token begin pts & forward pts.

Each edges are directed from one state to another, which are labeled by a symbol or set of symbols. If we are in state s & if symbol is a , we look for an edge out of state s labeled by a . If we find such an edge advance the forward pointer to the state of the transition diagram to which that edge leads.

Imp. conventions about the transition diagrams are;

1. Certain states are said to be accepting or final, which indicates that some has been found, although the actual some may not consist of all positions b/w the pt's. represent it by double circle or attach action to be taken to it.
2. If it is necessary to retreat the F.pt's one position (i.e. the some does not include the symbol that got us to the accepting state) then we additnally place a * near the accepting state.
3. one state is considered as start state or initial state it is indicated by an edge labeled

St → ij exp thr stat

main()

{
 do int a, b
 Pointf ("")
 } == 3 a = b + 10 if (a > 10)
 Pointf ("", a)
 else Pointf ("")

Prog → #include main() { declar-list } stat-list }

declar-list → <dec> | <dec-list>; <dec>

dec → <id-list> & <type> <id-list>

type → Integer

id-list → id | (id-list), id

stat-list → stat | (stat-list), stat

stat → <assign> | <pointf> | if-else

assign → id: = <exp>

exp → term | exp + term | exp - term

term → Factor | term * Factor / term di Factor

Factor → id | int | <exp>

Pointf → Pointf ("", id-list)

if → id <exp> then

if exp → Factor term if exp term stat

Factor Use stat | E

exp → + | - | > | >= | < | <= | < | >

Tokens

Pointf, Pointf ("", id-list),
 main(),
 main, g,

3

int

if

else

Pointf

+

÷

l

)

id

in Factor

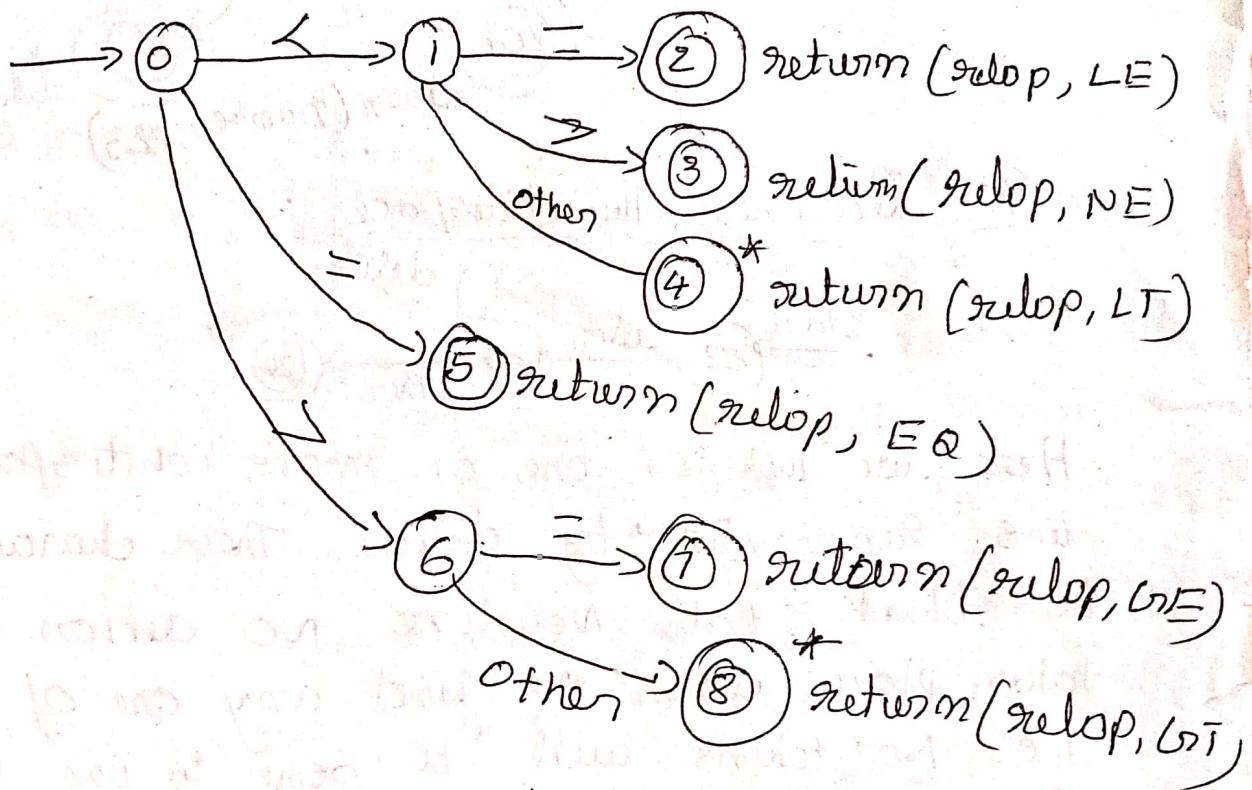
if exp

if exp

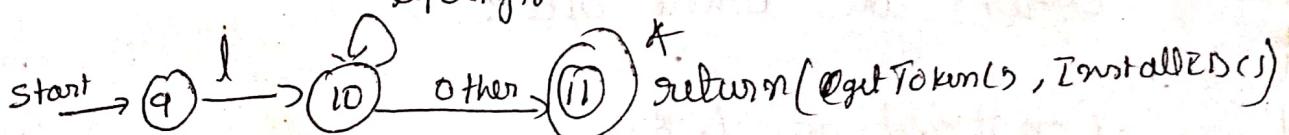
(35)

"start". The T.D always begins in the start state before any i/p symbols have been read.

E8: T.D For relational operator that recognizes the lexemes matching the token relop.



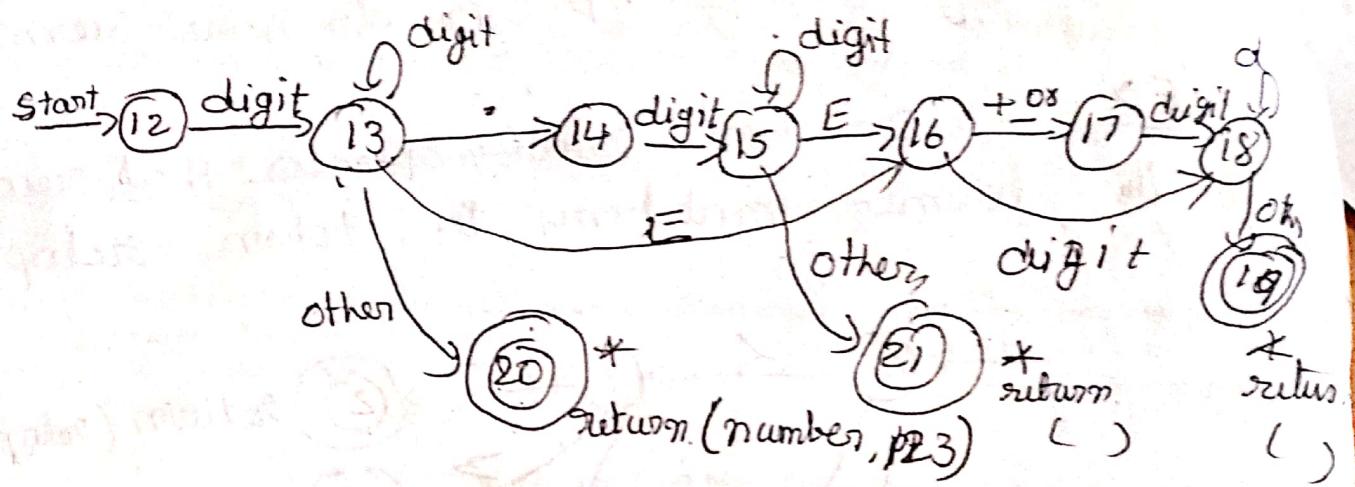
To Recognize an Identifier & K.W.
Iddigit



To Recognize K.W.



To Recognize the Number



To Recognize the whitespace



Here we look for one or more "whitespace" characters represented by `delim`. These characters could be blank, tab, newline. No action will be taken place when we find any one of the above i.e. No tokens will be sent to the parser. Rather we must restart the lexical Analysis after the white space.

Constructing L.A :

To construct an L.A we need the collection of T-Diagrams for different terminals of a grammar. & each state is represented by a piece of code.

ex: Take `getrelop()` which is a C++ fn whose job is to simulate the T. D for relop. and return an object of type TOKEN, i.e pair consisting of token name & attribute value.

`getrelop()` 1st creates a new object `retToken` & initializes the 1st component to `RELOP`.

If `State = 0` then Adj `nextchar()` obtains the next character from the i/p & assigns it to the local variable `c`. we then check `c` for the three characters we expect to find, ex if next character is `=` then we go to state 5.

If the next character is not the one that can begin a comparison operator then the fn `fail()` is called, & it should reset the forward pointer to `locmeBegin`, to allow another T. D to be applied to the true beginning of the unprocessed i/p.

It might then change the value of `State` to be the start state for another T. D, which will search for another token.

Alternatively if there is no other F. D. that remains unused, fail() could initiate an error-correction phase that will try to repair the i/p & find a lexeme.

If we are in state 8 we must retract the i/p ptr one position, which is done by f^n retract(), since state 8 represents the recognition of lexeme \geq we set the second component of the return object, which we names as G₂T

TOKEN getrelop()

```
{ TOKEN getToken = new (RELOP);  
  while (!) /* repeat characters processing until  
           a return or failure occurs */  
    switch (state)  
    { case 0: c = nextchar();  
      if (c == '<') state = 1;  
      else if (c == '=') state = 5;  
      else if (c == '>') state = 6;  
      else fail(); /* lexeme is  
                     not relop */
```