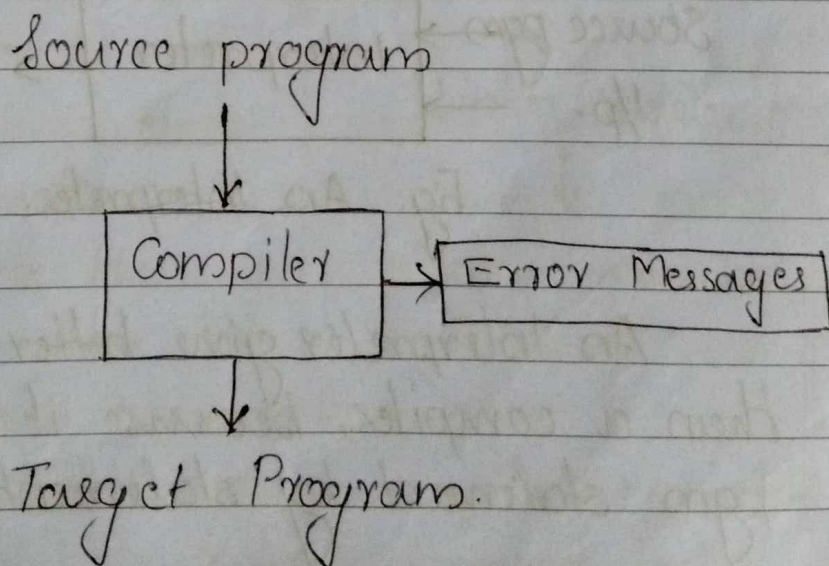


MODULE - 1

Introduction to Compilers:- Phases of a Compiler - Analysis and synthesis phases - Lexical analysis and its role - Review of finite automation, and Regular Expressions - Specification of tokens using regular expressions - Implementing lexical analyser using finite automation - Design of lexical analyser using LEX.

1.1 COMPILER

- Is a language processor
- A compiler is a program that can read a program in one language (i.e. the source language) and translate it into an equivalent program in another language (i.e. the target language)



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 machine language program, it can then be called by the user to process i/p's and produce o/p's.

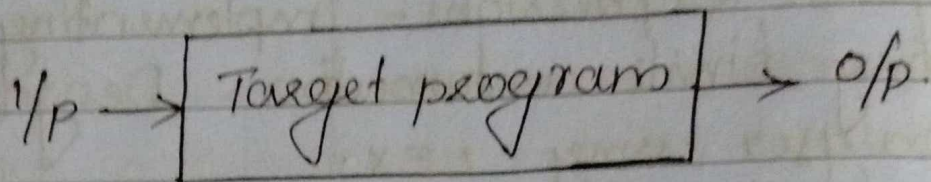


Fig: Running the target program.

1.6.1 INTERPRETER:

- It is a language processor.
- An Interpreter directly executes the ops specified in the source program on inputs supplied by the user.

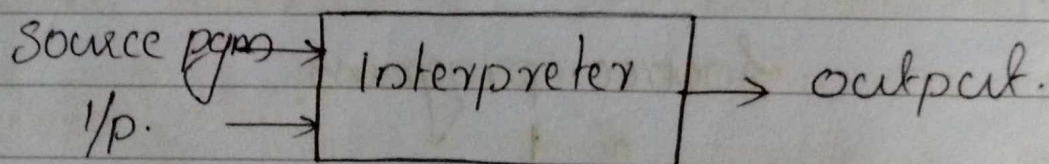
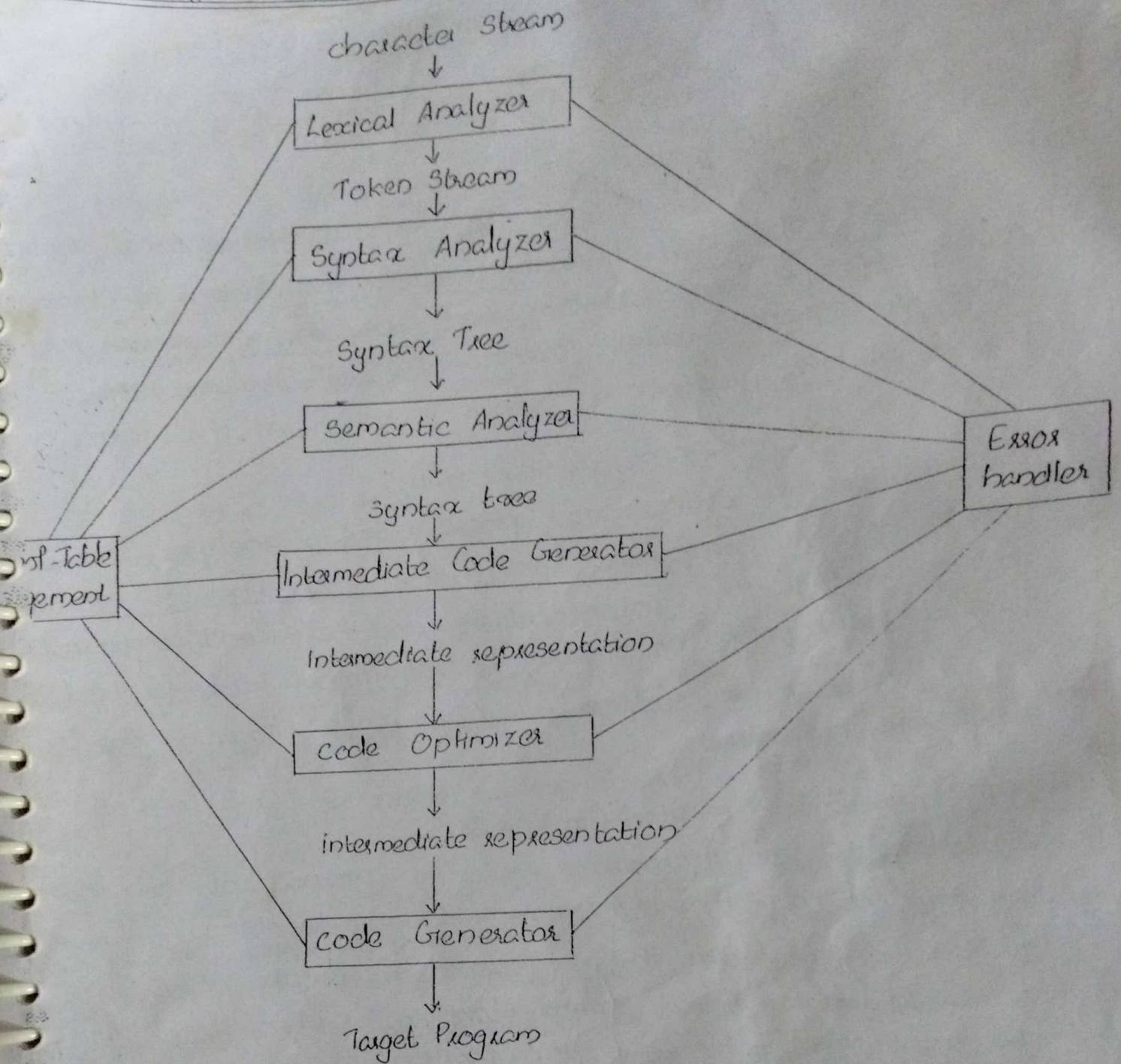


Fig: An Interpreter.

An Interpreter give better error diagnostics than a compiler, because it executes the source program statement by statement.

Phases of a Compiler.



Compiler can be broadly divided into two

phases

1. Analysis phase [often called front end of the compiler]
2. Synthesis phase [often called back end of the compiler].

The analysis phase breaks up the source program into pieces (tokens) and imposes a grammatical structure on them. It then uses this structure to create an intermediate representation. The analysis phase also collects information about the source program and store it in a data structure called a symbol table. If the analysis part detects that the source program is either syntactically ill formed or semantically unsound, then it produces error message.

1.2.1 Lexical Analysis / Scanning

- Performed by lexical analyser or lexer.
- The first phase of a compiler is called lexical analysis or scanning.
- The lexical analyser reads the stream of characters making up the source program and groups the

characters into meaningful sequences called lexemes.

- For each lexeme, lexical analyzer produces as output a token of the form

$\langle \text{token-name, attribute-value} \rangle$

which it passes to the next phase which is syntax analysis.

Here,

token-name is an abstract symbol that is used during syntax analysis.

attribute-value \Rightarrow points to an entry in the symbol table for this token.

eg: If the source program contains the following assignment

$\text{position} = \text{initial} + \text{rate} * 60.$

(The characters in this assignment could be grouped into the following lexemes and mapped into the following tokens passed on to the syntax analyzer.)

After lexical analysis, the sequence of tokens are

$\langle \text{id}, 1 \rangle \quad \langle = \rangle \quad \langle \text{id}, 2 \rangle \quad \langle + \rangle \quad \langle \text{id}, 3 \rangle \quad \langle * \rangle \quad \langle 60 \rangle$

- id is an abstract symbol standing for identifier.

- 1, 2, and 3 are points to the symbol table entries for position, initial and rate respectively.

- +, = and * are abstract symbols for addition, assignment and multiplication operators. They

do not need any attribute value

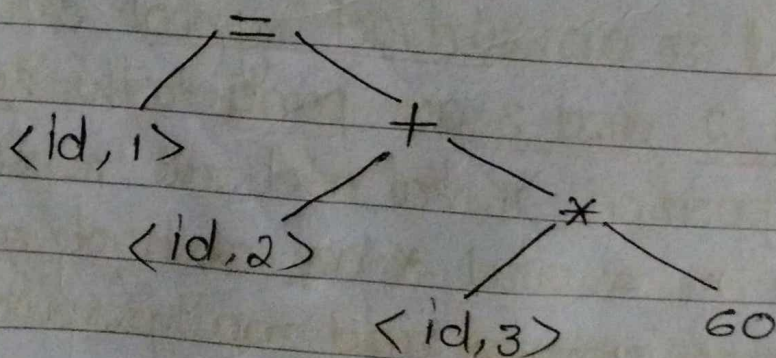
1	Position	---
2	initial	---
3	rule	---

Symbol Table

1.2.2

Syntax Analysis / Parsing

- Performed by syntax Analyzer or Parser
- The second phase of the compiler is called Syntax Analysis or Parsing.
- The parser uses the tokens produced by the lexical analyzer to create a tree like intermediate representation that (called syntax tree) depicts the grammatical structure of the tokens stream.
- In a Syntax tree, each interior node represents an operation and the children of the node represent the arguments of the operation.



1.2.3

Semantic Analysis:

- The Semantic Analyzer uses the syntax tree and the information in the symbol table to check the source program for semantic consistency with the language definition.

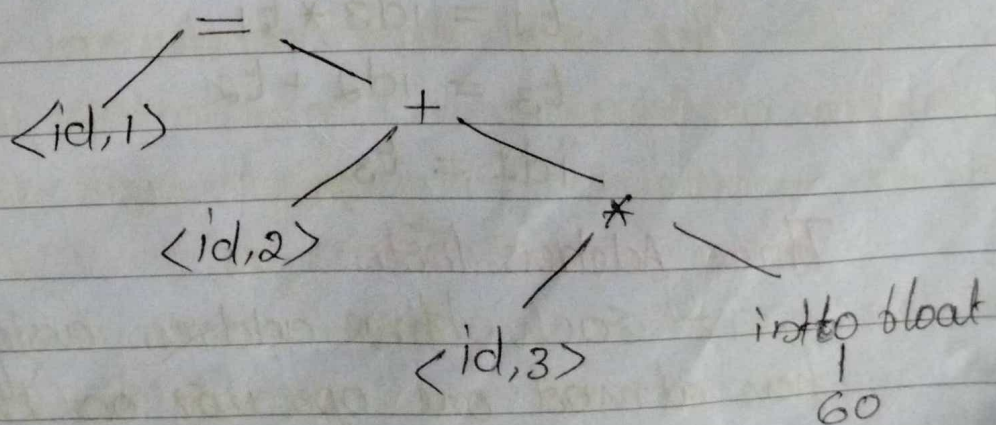
✓ - It also gathers type information and save it in either the syntax tree or symbol table for further use.

- An important part of semantic analysis is type checking, where the compiler checks that each operator has matching operands.

Eg: The compiler must report an error if a floating point no is used to index an array.

- It permits some type conversions called coercions.

Eg: If a binary arithmetic operator is applied to an integer and a floating point number, then the compiler may convert or coerce the integer into a floating point number.



1.2.4

Intermediate Code Generation.

- Syntax trees are a form of intermediate representation; they are commonly used during syntax and semantic analysis.

- Compiler generates an explicit low-level or machine-like intermediate representation.

- This intermediate representation should have two important properties.

* It should be easy to produce

* It should be easy to translate into the target machine language.

- Example of intermediate representation is called three address code, which consists of a sequence of assembly-like instructions with three operands per instruction.

- The o/p of this phase is given below

$$t_1 = \text{inttofloat}(60)$$

$$t_2 = id_3 * t_1$$

$$t_3 = id_2 + t_2$$

$$id_1 = t_3$$

Three Address Instrs:

- Each three address assignment instruction has at most one operator on the right side.

- These instns fix the order in which operations are to be done.

Eg: multiplication precedes the addition in the source prog.

- The compiler must generate a temporary name to hold the value computed by a three-address instruction.

- Some three address instns have fewer than three operands.

1.2.5 Code Optimization.

- Code optimization phase attempts to improve the intermediate code so that better target code will result.

- Better means shorter, faster etc.

Eg. inttofloat operation can be eliminated by replacing the integer 60 by the floating point number 60.0

$$t_1 = id3 * 60.0$$

$$id1 = id2 + t_1$$

- Different optimization techniques are :-

Local Optimization, loop optimization, deadcode elimination, strength reduction, frequency reduction etc.

1.2.6

Code Generation:-

- The code generator takes as i/p an intermediate representation of the source prog and maps it to the target language.

- The intermediate instructions are translated into sequence of machine instructions that perform the same task.

- A crucial aspect of code generation is the judicious assignment of registers to hold variables (registers or memory locations are selected for each of the variables used by the prog)

Eg. Using registers R₁ and R₂, the intermediate code might get translated into the machine code

LDF R₂, id₃

MULF R₂, R₂, #60.0 multiplier with floating point constant

LDF R₁, id₂ move id₂ into reg^r R₁

ADDF R₁, R₁, R₂ Adds it with value in reg^r R₂

STF id₁, R₁

- The first operand of each instruction specifies a destination.

- The F in each instr tells that it deals with floating point numbers.

- LDF in first line - loads contents of address id₃ into reg^r R₂.

signifies that 60.0 is to be treated as an immediate constant.

STF - stores the value in reg^y R1 into the address of id1.

1.2.7 Symbol-Table Management

- (An essential feature of a compiler is to record the variable names used in the source program and collect information about various attributes of each name. (such as storage allocated for a name, its type, its scope number and type of arguments etc).)

- A symbol table is a data structure containing a record for each variable name, with fields for the attributes of the name.)

- The data structure should be designed to allow the compiler to find the record for each name quickly and to store or retrieve data from that record quickly.

- These attributes associated with a name provide information about the storage allocated for an identifier, its type, its scope and in case of procedural names such things as the no and types of its arguments, the method of passing each arguments and the type returned.

Error Handler:

- Each phase can encounter errors. Compiler stops when it finds the 1st error.

1. Lexical Analyzer : next token in the source program is misspelled.
2. Syntax Analyzer : syntactic errors such as a missing parenthesis.
3. IC generator : May detect an error operator whose operands have incompatible type.
4. Code Optimizer : May detect that certain blocks can never be reached.
5. Code Generator : May find a compiler created constant that is too large to fit in a word on the target machine.
6. Symbol Table : May discover an identifier that has been multiply declared with contradictory attributes.

position = initial + rate * 60

Lexical Analyzer

$\langle id, 1 \rangle \langle = \rangle \langle id, 2 \rangle \langle + \rangle \langle id, 3 \rangle \langle * \rangle \langle 60 \rangle$

Syntax Analyzer

$\langle id, 1 \rangle = + \langle id, 2 \rangle * \langle id, 3 \rangle 60$

Semantic Analyzer

$\langle id, 1 \rangle = + \langle id, 2 \rangle * \langle id, 3 \rangle \text{inttofloat} 60$

Intermediate Code Generator

```
t1 = inttofloat(60)
t2 = id3 * t1
t3 = id2 + t2
id1 = t3
```

Code Optimizer

```
t1 = id3 * 60.0
id1 = id2 + t1
```

Code Generator

```
LDF R2, id3
MULF R2, R2, #60.0
LDF R1, id2
ADDF R1, R1, R2
STF id1, R1
```

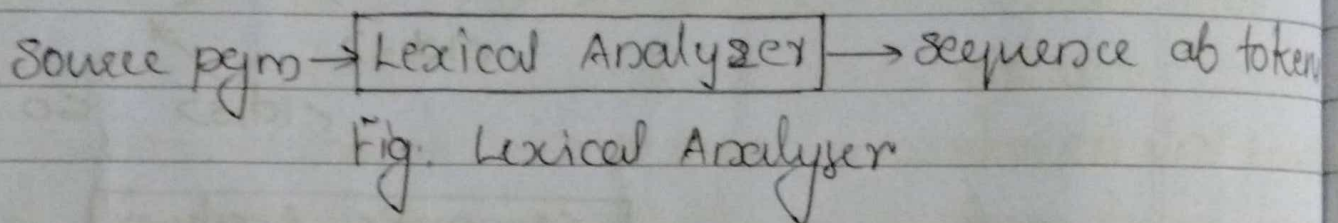
1	position	...
2	initial	...
3	rate	...

SYMBOL TABLE

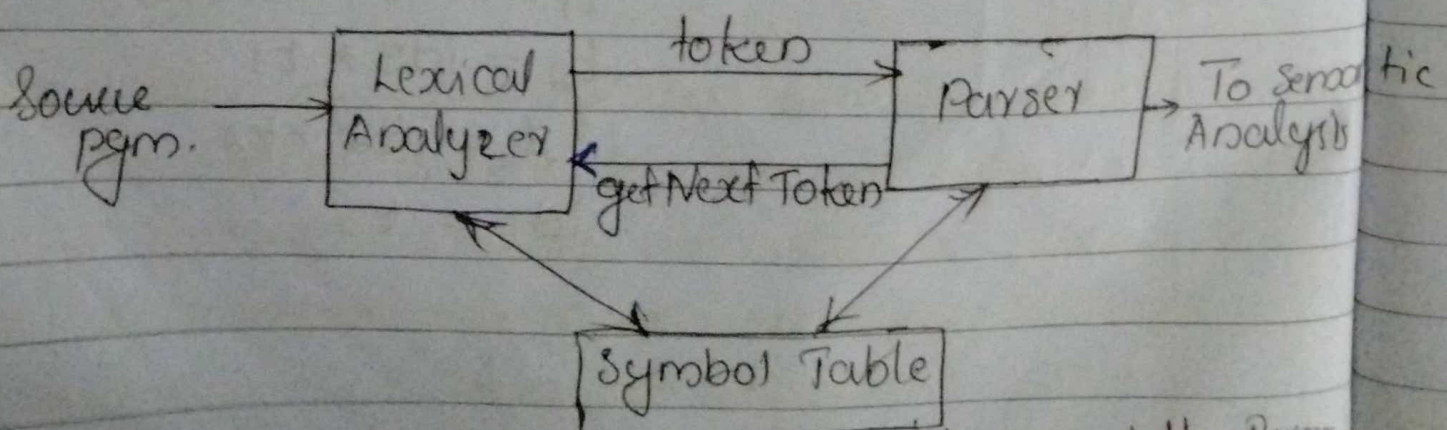
Figure 1.7: Translation of an assignment statement

1.3 LEXICAL ANALYSIS & ITS ROLE

- First phase of a compiler.
- It reads the i/p characters of the source program, group them into lexemes, and produce as o/p a sequence of tokens for each lexeme in the source program.
- The stream of tokens is sent to the parser for syntax analysis.



- It also interact with the symbol table.
- When the lexical analyzer discovers a lexeme constituting an identifier, it needs to enter that lexeme into the symbol table.
- In some cases, the information regarding the kind of identifier may be read from the symbol table by the lexical analyzer to assist it in determining the proper tokens it must pass.



Interactions b/w the lexical analyzer and the Parser

Other tasks of Lexical Analyzer:

1. Stripping out comments and whitespaces (blank, new line, tab and perhaps other characters that are used to separate tokens in the i/p).
2. Correlating error messages generated by the compiler with the source prog. For instance, the lexical analyzer may keep track of the no. of newline characters seen, so it can associate a line no. with each error message.
3. If the source prog uses a macro preprocessor, the expansion of macros may also be performed by the lexical analyzer.

→ Lexical analyzers are divided into 2 processes

1. Scanning - that donot require tokenization of i/p such as deletion of comments and compaction of consecutive whitespaces characters into one.
2. Lexical Analysis - It produce tokens from the o/p of the scanner.

1.3.1 Tokens, Patterns and Lexemes

Token: - is a pair consisting of a token name and an optional attribute value.

- The token name is an abstract symbol representing a kind of lexical unit. e.g: a particular keyword, or a sequence of i/p

characters denoting an identifier.

- The token names are the I/P symbols that the parser processes.

Patterns: — It is a description of the forms that the lexemes of a token may take. In the case of a keyword as a token, the pattern is just the sequence of characters that forms the keyword.

Lexeme: — Is a sequence of characters in the source program that matches the pattern for a token and is identified by the lexical analyzer as an instance of that token.

- The word generated by the lexical analysis may be of different kinds:

- * identifier

- * keyword (if, while, ...)

- * punctuation

- * numbers

- * literals

- Such a kind is called a **TOKEN** and an element of this kind is called **LEXEME**.

- A word is recognized to be a lexeme for a certain token by **PATTERN MATCHING**.

eg: letter followed by letters and digits is a pattern that matches a word like x or y with the token id (= identifier)

Token	Lexemes	Pattern
if	if	if
else	else	else
ID	x, y, no	letter followed by letters & digits
number	1, 3.14, 2.02e23	any numeric constant
comparison	<=, !=	comparison operators
literal	"Hello", "Hi"	any string of characters between "and"

- During the analysis, the compiler manages a symbol table by:

- * recognizing the identifier at the source program
- * collecting information (called ATTRIBUTES) about them: storage allocation, type, scope (for bns) signature.

- when the identifier x is bound by lexical analyzer

- * It generates token id
- * enters lexeme x into symbol table (if not present)

* associates to the generated tokens a pointer to the symbol-table entry x. This pointer is called

The lexical value of the tokens.

1.3.2 Attributes for Tokens.

When more than one lexeme can match a pattern, the lexical analyzer must provide the subsequent compiler phases additional information about the particular lexeme that matched.

eg: Pattern for tokens number matches both 0 & 1 but it is important for code generator to know which lexeme was bound in the source prog.

So the lexical analyzer returns to the parser not only a token-name, but also an attribute value that describes the lexeme.

eg: For the token id, we need to associate with the token a great deal of info such as its lexeme, type, location (where it is bound). This is stored in the symbol table.

Thus, the appropriate attribute value for an identifier is a pointer to the symbol-table entry for that identifier.

Example

$$E = M * C ** 2$$

Tokens generated for the above line of code
<id, pointer to the symbol-table entry for E>
<assign-op>

< id, pointer to symbol-table entry for M >

< mult-op >

< id, pointer to symbol-table entry for c >

< exp-op >

< number, integer value 2 >

or
pointer to symbol table entry for 2.

Note:

- In certain pairs, especially operators, punctuation^{ns}, and keywords, there is no need for attribute value.

- Token number has been given an integer-valued attribute.

1.4 INPUT BUFFERING.

It is the technique to improve the execution speed of the lexical analysis phase by speeding up the reading of source program.

- Here we use two buffers that are alternately reloaded.

- Each buffer is of the same size N , and N is usually the size of a disk block, e.g. 4096 bytes.

- Using one system read command we can read N characters into a buffer, rather than using one system call per character.

- If fewer than N characters remain in the I/O file, then a special character, represented by eof, marks

The end of the source file.

- Two pointers to the i/p are maintained.
- 1) Pointer `lexemeBegin`, marks the beginning of current lexeme, whose extent we are attempting to determine.
- 2) Pointer `forward` scans ahead until a pattern match is found. If the forward has passed to the end of the lexeme, it must be retracted one pos to the

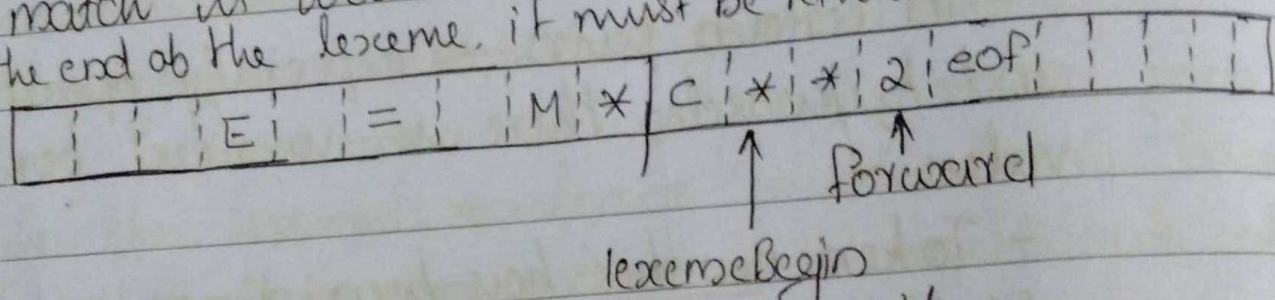


Fig: using a pair of i/p buffers.

Working

- Initially both pointers point to the first character of the next lexeme.
 - Forward pointer scans. If a lexeme is found.
 - It is set to the last character of the lexeme found.
 - Thus for each character read, two tests are made one for the end of buffer and the other to determine the character that is read.
 - Both tests can be combined by extending the buffer to hold a sentinel character at the end.
- ⇒ Sentinel is a special character that cannot be part of source prog, such as eof.

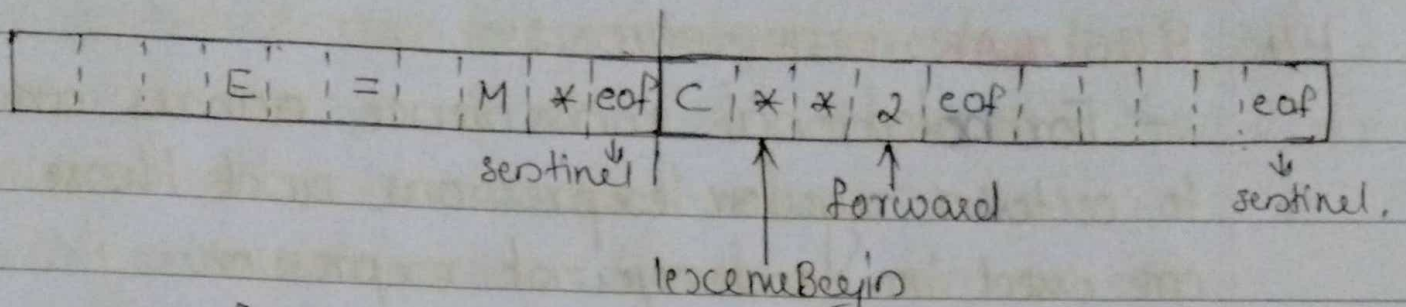


Fig: sentinels at the end of each bubble.

Fig. shows Algorithm for advancing forward.

```
Switch (* forward++)
{
```

```
  case eof:
```

```
    if (forward is at end of first bubble)
    {
```

```
      reload second bubble;
```

```
      forward = beginning of second bubble;
```

```
    }
```

```
    else if (forward is at end of second bubble)
    {
```

```
      reload first bubble;
```

```
      forward = beginning of first bubble;
```

```
    }
```

```
    else /* eof within bubble marks the end of input */
      terminate lexical analysis;
```

```
    break;
```

```
  case for other characters;
```

```
  }
```


1.4 REGULAR DEFINITIONS:

- For notational convenience, names are given to certain regular expressions, and those names are used in subsequent expressions. (These names are also symbols).

- If Σ is an alphabet of basic symbols, then a regular definition is a sequence of definitions of the form:

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

...

$$d_n \rightarrow r_n.$$

where:

- 1) Each d_i is a new symbol, not in Σ and not the same as any other of the d 's, and
- 2) Each r_i is a regular expression over the alphabet $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$

Eg: C identifiers are strings of letters, digits and underscores. The regular definition for the language of C identifiers are;

$$\text{letter}_- \rightarrow A/B/\dots/2/a/b/\dots/z/$$

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

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

Unsigned numbers (integer or floating point) are strings such as 5280, 0.01234, 6.36E+07 or 1.89E-4.

Regular definition for unsigned numbers

digit \rightarrow 0|1|...|9

digits \rightarrow digit digit*

Optional fraction \rightarrow .digits/e

Optional Exponent \rightarrow (E(+|-|e) digits)/e

number \rightarrow digits optional fraction optional Exponent

- i.e. an Optional Fraction is either a decimal point (dot) followed by one or more digits, or it is missing (the empty string)

- An Optional Exponent, if not missing, is the letter E followed by an optional + or - sign, followed by one or more digits

- Note

At least one digit must follow the dot.

so the number does not match 1., but match 1.0

1.4.1 Extensions as Regular Expressions.

1. One or More Instances

The unary, postfix operator + represents the positive closure of a regular expression and its language. That is, if r is a R.E, then $(r)^+$ denotes the language $(L(r))^+$. The operator + has

The same precedence and associativity as the operator $*$

The following algebraic laws relates the Kleene closure and positive closure

$$\gamma^* = \gamma^+ / \epsilon$$

$$\gamma^+ = \gamma\gamma^* = \gamma^*\gamma$$

2. Zero or One Instance

The unary postfix operator '?' means zero or one occurrence

ie,

$$\gamma? = \gamma / \epsilon$$

$$L(\gamma?) = L(\gamma) \cup \{\epsilon\}$$

3. Character classes

- A r.e $a_1 a_2 \dots a_n$, where the a_i 's are each symbols of the alphabet, can be replaced by the shorthand $[a_1 a_2 \dots a_n]$.

- $[abc]$ is shorthand for $a|b|c$. and

$[a-z]$ is shorthand for $a|b|\dots|z$

- Using shorthand, the regular defn. for c- identifiers are as follows:

letter_ $\rightarrow [A-Za-z_]$

digit $\rightarrow [0-9]$

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

Regular expr. for unsigned nos:

digit $\rightarrow [0-9]$
 digits $\rightarrow \text{digit}^+$
 number $\rightarrow \text{digits}(\cdot \text{digits})?(E[+-]? \text{digits})?$

1.5 Specification of Tokens using Regular Expressions.

A grammar for branching stmts,

stmt \rightarrow if expr then stmt
 | if expr then stmt else stmt
 | ϵ

expr \rightarrow term relop term
 | term

term \rightarrow id
 | number

The terminals of the grammar, which are if, then, else, relop, id and number, are the names of tokens as far as lexical analyzer is concerned.

- The patterns for these tokens are described using regular exprs as follows:

(patterns for id & number are same as before)

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

number \rightarrow digits (\cdot digits)? ($[E [+ -]]?$ digits)?
 letter \rightarrow [A-Za-z]
 id \rightarrow letter (letter | digit)*
 if \rightarrow if
 then \rightarrow then
 else \rightarrow else
 relop \rightarrow <|> |<=|>=| |= |<>

- For the language, lexical analyzer will recognize the keywords 'if', 'then', and 'else' as well as lexemes that match the patterns for relop, id and number.

- In addition, lexical analyzer is given the job of stripping out white space, by recognizing the tokens, ws obtained by

ws \rightarrow (blank | tab | newline)+

\rightarrow Here blank, tab and newline are abstract symbols that is used to express the ASCII characters of the same name.

\rightarrow Token ws is different from other tokens. As when ws is recognized, it is not returned to parser, but instead lexical analyzer restarts from the character that follows the ws (white space)

- After performing lexical analysis, a set of tokens and attribute value (if present) are passed to the parser.

The following table shows, for each lexeme or family of lexemes, which token name is returned to parser and what attribute value.

Lexemes	Token names	Attribute value
Any ws	-	-
if	if	-
then	then	-
else	else	-
Any id	id	Pointer to table entry
Any number	number	Pointer to table entry
<	relop	LT
<=	relop	LE
=	relop	EQ
<>	relop	NE
>	relop	GT
>=	relop	GE

For the 6 relational operators, symbolic constants LT, LE etc are used as attribute value, in order to indicate the instance of the tokens relop bound.

1.6 Implementing Lexical Analyser using FA.

Transition diagrams:

- Transition diagrams have a collection of nodes or circles, called states. Each state represents a condition that could occur during the process of scanning the I/P looking for a lexeme that matches one of several patterns.

- Edges are directed from one state of the transition diagram to another. Each edge is labeled by a symbol or set of symbols.

- If we are in some state S , and the next I/P symbol is a , we look for an edge out of state S labeled by a . If we find such an edge, advances the forward pointer and enter the state of the transition diagram to which that edge leads.

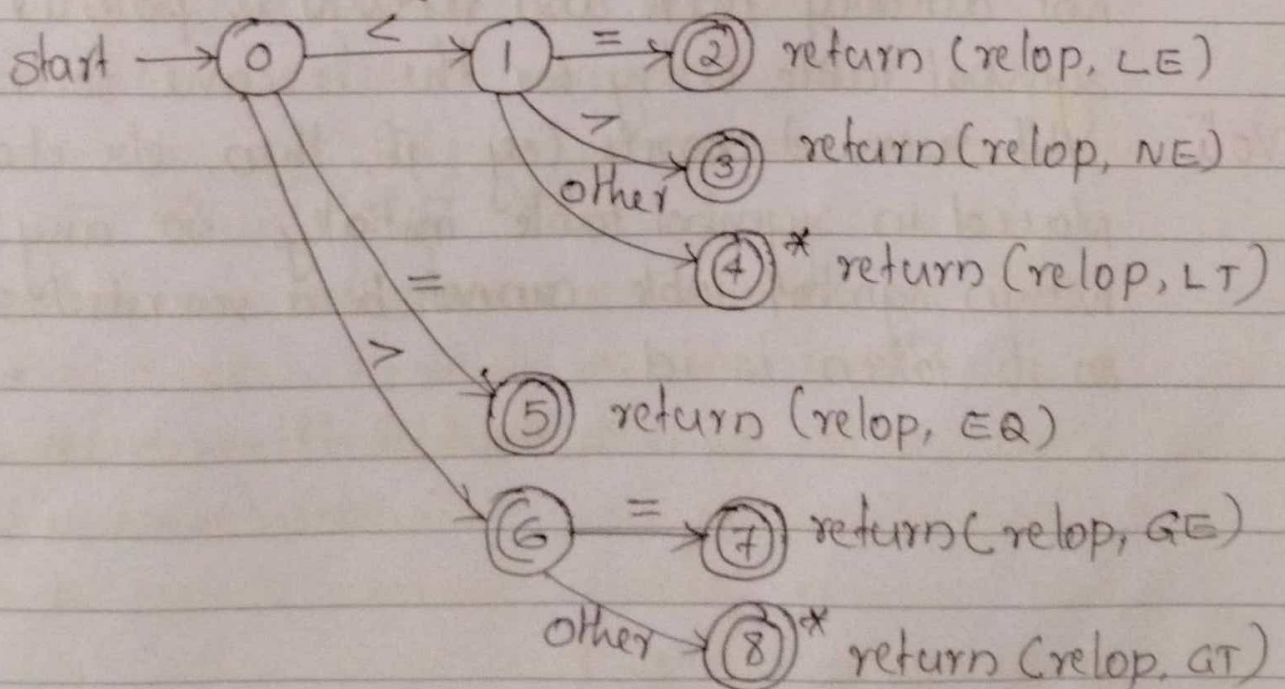
Some important conventions about transition diagrams are:

1. Certain states are said to be accepting, or final. These states indicate that a lexeme has been found, although the actual lexeme may not consist of all positions between the lexeme's begin and forward pointers. Indicate an accepting state by a double circle and if there is an action to be taken-

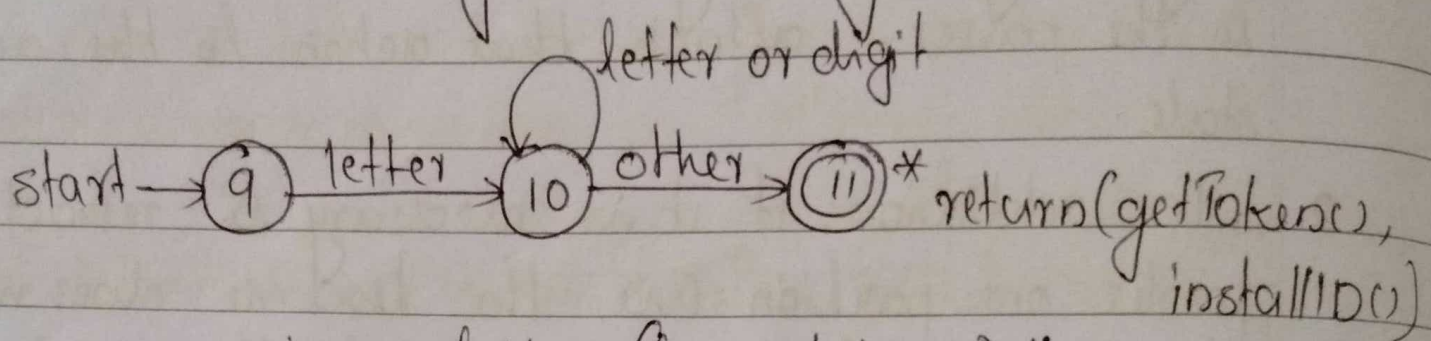
typically returning a token and an attribute value to the parser - attach that action to the accepting state.

2. In addition, if it is necessary to retract the forward pointer one position (i.e. the lexeme does not include the symbol that got us to the accepting state) then we shall additionally place a * near that accepting state.
3. One state is designated as the start state, or initial state. It is indicated by an edge, labeled "start" entering from nowhere. The transition diagram always begins in the start state before any i/p symbol have been read.

eg. Transition diagram for relop :



Transition diagram for keywords & identifiers



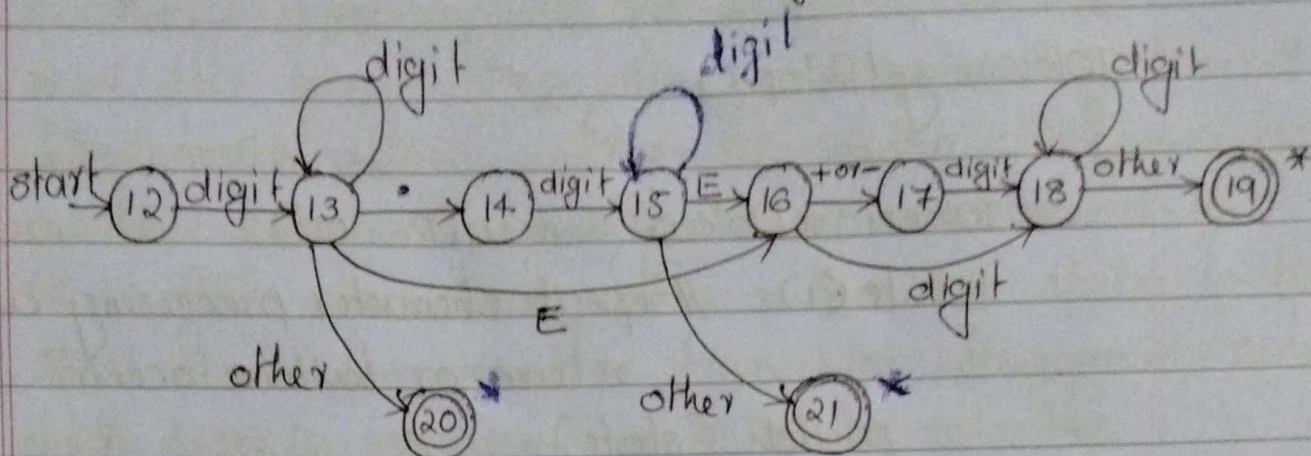
$id = \text{letter} (\text{letter} | \text{digit})^*$

getTokens() → examines the symbol table entry for the lexeme found, and returns the token name that the symbol table entry says this lexeme represents. i.e., either id or one of the keyword tokens that was initially installed in the table.

installID() → when we bind an identifier, a call to installID places it in the symbol table if it is not already there and returns a pointer to the symbol table entry for the lexeme found.

Note: All reserved words (e.g., if, then, else etc) are placed in symbol table initially. So any identifier not in symbol table cannot be a reserved word, so its token is id.

Transition diagram for unsigned numbers



digit \rightarrow 0 | 1 | ... | 9

digits \rightarrow digit digit*

optional fraction \rightarrow . digits / e

optional exponent \rightarrow (E (+|-) digits) / e

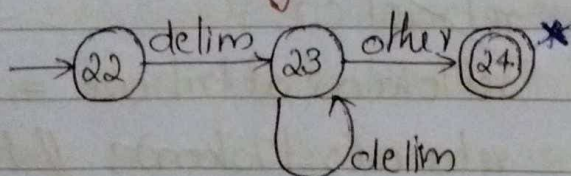
number \rightarrow digit . optional fraction . optional exponent

Eg: 12.3

123

123.45E-6

Transition diagram for whitespace



delims \rightarrow blank / tab / newline

ws \rightarrow delims*

Compiler - Construction Tools

There are specialized tools that help to implement various phases of a compiler. Some commonly used compiler - construction tools include:

1. Parser generators \rightarrow automatically produce syntax analyzers from a grammatical description of a programming language.
2. Scanner generators \rightarrow produce lexical analyzers from regular expression description of the tokens of a language.
3. Syntax - directed translation engine \rightarrow produce collections of routines for walking a parse tree and generating intermediate code.
4. Code generator generators \rightarrow produce a code generator from a collection of rules for translating each operation of the intermediate language into the machine language for a target machine.
5. Data-flow analysis engine \rightarrow gathering of information about how values are transmitted from one part of a program to each other part. Data-flow analysis 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.

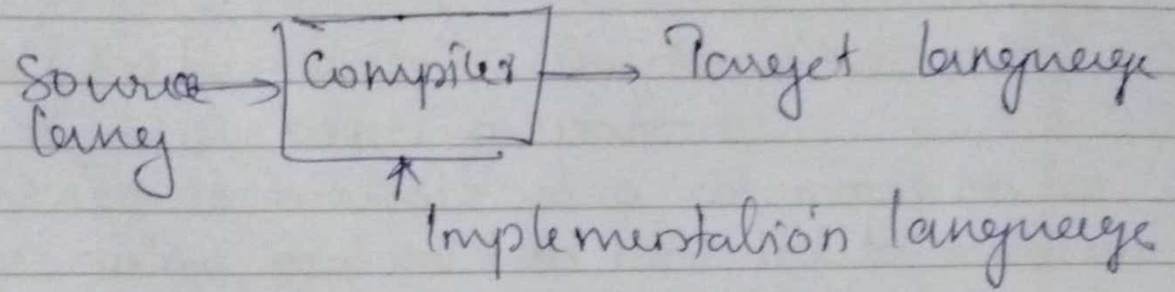
Bootstrapping

Process
- Is used to design a compiler

Generally a compiler can be represented using three languages

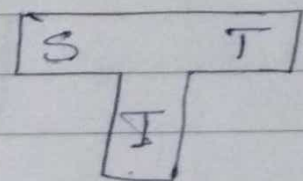
- 1) Source lang
- 2) Target lang
- 3) Implementation lang

- The language in which the compiler is written.



The compilation process can be explained with the help of T-diagram

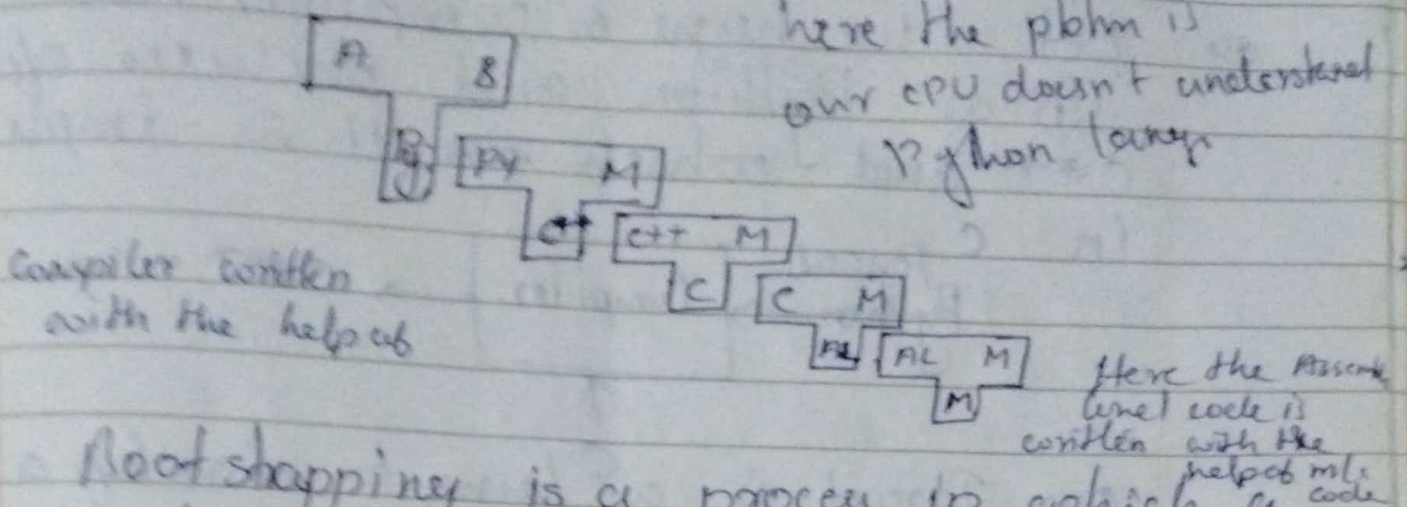
T-diagram



- S - Source lang
- T - Target
- I - Implementation lang

Q4. Let, we have a compiler with A as input B as o/p and the comp^r is written in Python lang.

here the pbm is our CPU doesn't understand Python lang.



Bootstrapping is a process in which a simple language is used in order to translate a complicated program. But this complicated program produces even more complicated program and likewise the process will continue and this process is known as bootstrapping.

Cross Compiler

- A compiler which runs on one m/c produces an object code for another m/c.

For a source language statement $a = b * c - 2$ where a, b and c are floating point variables $*$ and $-$ represent multiplication and subtraction on the same data type. Show the i/p and o/p at each of the complete phases.

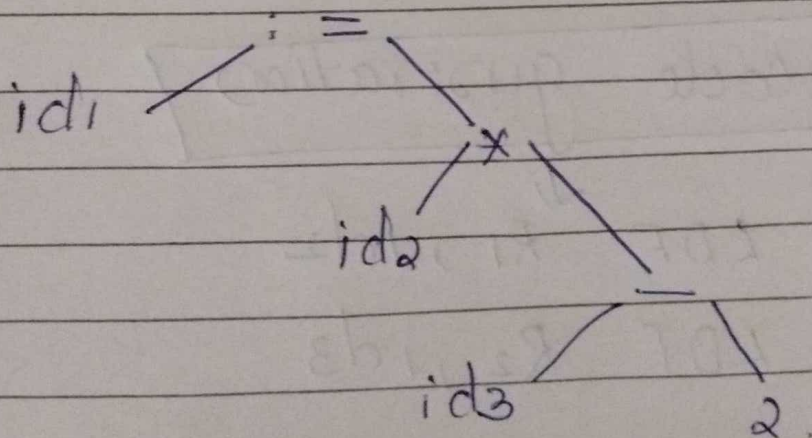
$$a = b * c - 2$$

Lexical Analyzer

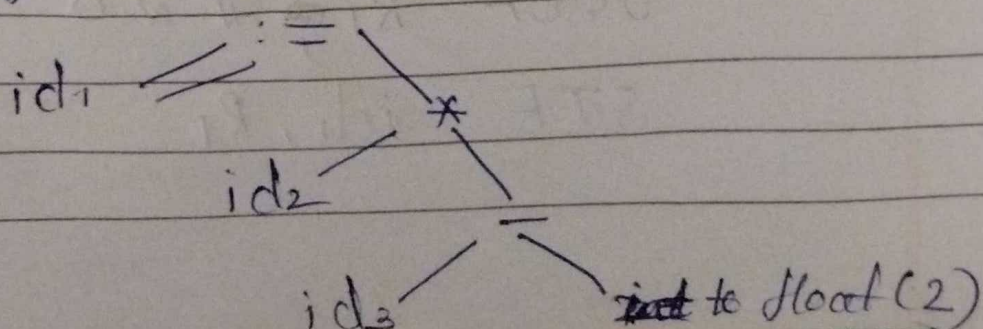
$$id_1 = id_2 * id_3 - 2$$

$\langle id_1 \rangle \langle = \rangle \langle id_2 \rangle \langle * \rangle \langle id_3 \rangle \langle - \rangle \langle 2 \rangle$

Syntax Analyzer



Semantic Analyzer



↓
[intermediate code representation]

↓
 $t_1 = id_2 \times id_3$
 $t_2 = \text{int-to-float}(2)$
 $t_3 = t_1 - t_2$
 $id_1 = t_3$

↓
[code optimizer]

↓
 $t_1 = id_2 \times id_3$
 $id_1 = t_1 - 2.0$

↓
[code generation]

↓
LDF R_1, id_2
LDF R_2, id_3
MUL R_1, R_2
SUBF $R_1, \#2.0$
STF id_1, R_1

**Explain how the regular expressions and finite state automata can be used
for the specification and recognition of tokens**

The lexical analyzer needs to scan and identify only a finite set of valid string/token/lexeme that belong to the language in hand. It searches for the pattern defined by the language rules.

Regular expressions have the capability to express finite languages by defining a pattern for finite strings of symbols. Regular expression is an important notation for specifying patterns. Each pattern matches a set of strings, so regular expressions serve as names for a set of strings.

The only problem left with the lexical analyzer is how to verify the validity of a regular expression used in specifying the patterns of keywords of a language. A well-accepted solution is to use finite automata for verification

Finite automata is a state machine that takes a string of symbols as input and changes its state accordingly. Finite automata is a recognizer for regular expressions. When a regular expression string is fed into finite automata, it changes its state for each literal. If the input string is successfully processed and the automata reaches its final state, it is accepted