

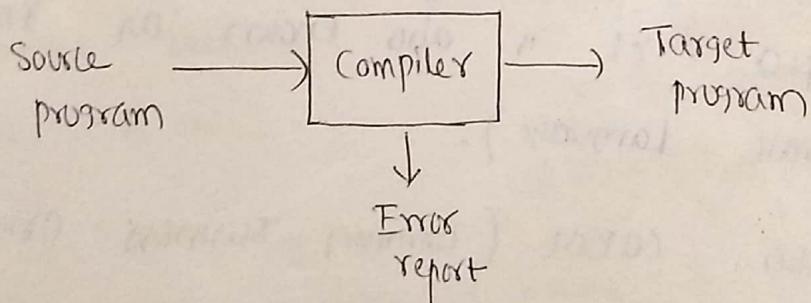
Introduction to Compiling

Compiler :

- * A Compiler is a program that reads a source program written in one language (i.e) Source language (high level language) and translates it into an equivalent program in another language (i.e) Target language (low level language).

(or)

A compiler is a piece of code that translates the high level language into machine language (or) low level language.



- * After compilation, compiler generates an error report whenever an error occurs in the source program.

History of the Compiler :-

- * The term "Compiler" was coined by "Grace Murray Hopper".

- * In earlier days, computers have very limited memory

capacity and also created many technical problems when implementing a compiler.

- * The first compiler was written by "Grace Murray Hopper" in 1952 for A-O programming language (Arithmetic Language version -0).
- * In 1957, the FORTRAN (Formula Translating System) Compiler was developed by IBM, lead by John Backus Team .
- * The first compiler which was written by Grace Murray Hopper was "Incomplete Compiler".
- * FORTRAN was the first complete compiler.
- * In 1958, ALGOL 58 (Algorithmic Language) was developed . It is also known as IAL (International Algebraic Language).
- * In 1960, COBOL (common Business Oriented language) compiler was developed .

Differences

between Compiler and Interpreter :-

Compiler

Interpreter

- (i) It takes the entire program as input.
- (ii) It takes a single instruction as input.

- (iii) It generates an error report after compilation, whenever an error occurs.
 - (iv) By using Compiler, control statements are executed faster.
 - (v) It takes a large amount of time in analysing & processing the source program.
- (ii) It stops the translation after it gets the first error.
 - (vi) By using interpreter, control statements are executed slower.
 - (vii) It takes lesser time in the same process.

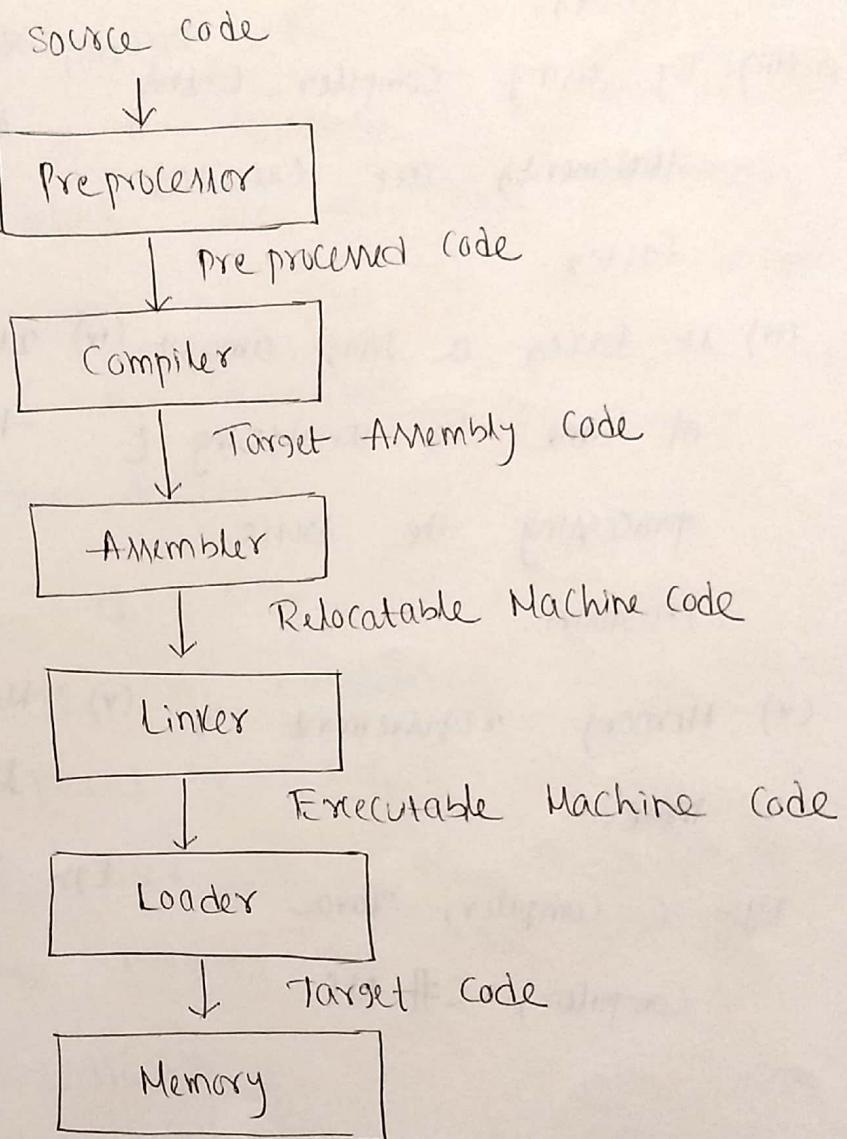
Eg:- C Compiler, Java Compiler, C# etc.

(v) Memory requirement is less.

Eg:- BASIC (Beginners All purpose Symbolic Instruction Code), LISP (List Processing), etc.

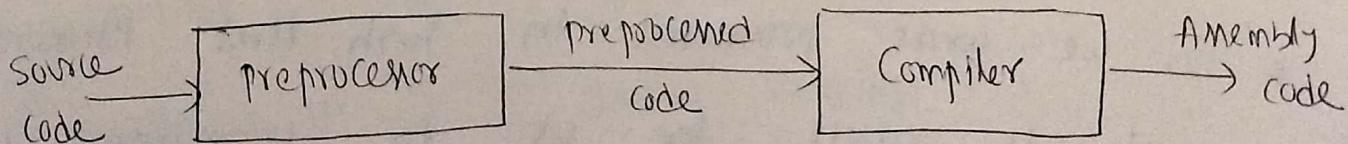
- Language Processing System :- (columns of compiler)
- * Computer system is made of hardware and software.
- * The hardware understands a language, which humans cannot understand.
- * So we write programs in high-level language, which is easier for us to understand and remember.

* These programs are then fed into a series of tools and OS components to get the desired code that can be used by the machine.



(i) Preprocessor :-

* A preprocessor is considered as a part of compiler, that produces input for the compiler.



* The output of preprocessor given as input to the compiler. (3)

* The tasks performed by the preprocessor are :-

- (a) It allows the user to use "Macros" in the program. A Macro is a set of instructions which can be used repeatedly in the program. Whenever in the program, "Macro name" is identified then it is replaced by its "Macro definition".

Eg:- `#define PI 3.14`

whenever 'PI' is encountered in a program, it is replaced by '3.14'.

- (b) It allows the user to use the "Header files", which may be required by the program.

Eg:- `#include <stdio.h>`

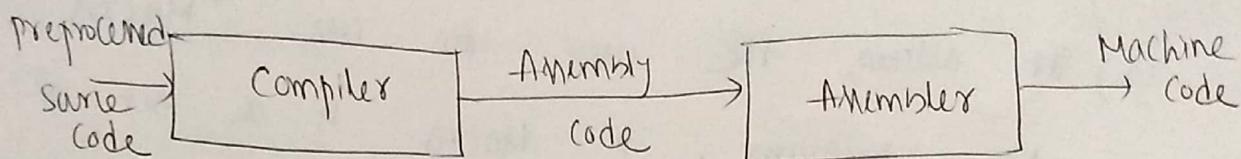
The header file contains different predefined functions. Those can be included by using preprocessor.

(ii) Compiler :-

- A Compiler is a program that converts high-level language in to assembly language.

(iii) Assembler :-

- * An assembler translates the assembly language into machine code.



- * Assembly code is a mnemonic version of machine code.
- * The output of an assembler is an "Object file", which contains machine instructions as well as the data required to place the instructions in memory.

(iv) Linker :-

- * Linker is a computer program that links and merges various object files together in order to make an executable file.
- * The major task of a linker is to search and locate referenced routines in a program and to determine the memory location where these codes will be loaded.

(v) Loader :-

- * Loader is a part of operating system and it is responsible for loading executable files into memory and execute them.

* It calculates the size of program and creates memory space for it. (h)

Phases of Compiler :-

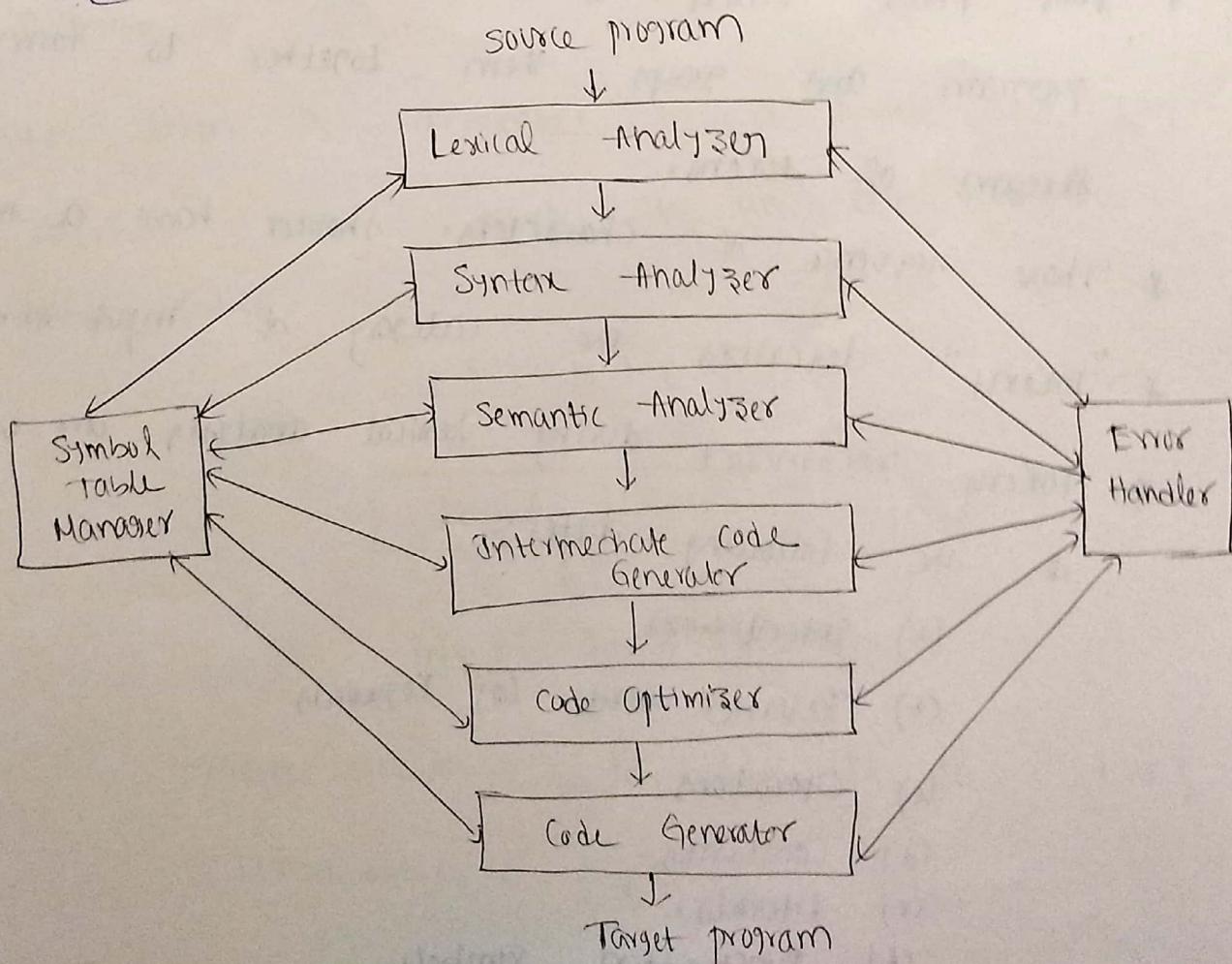
* The compilation can be done in to 2 parts :-

(i) Analysis.

(ii) Synthesis.

* In "Analysis" part, the source program is read and broken down into pieces and the syntax & meaning of the source program is determined and then intermediate code is created.

* In "synthesis" part, the intermediate code is taken and converted into target code.



* The analysis part consists of 3 phases:-

- (i) Lexical Analysis
- (ii) Syntax Analysis
- (iii) Semantic Analysis

* The synthesis part consists of 3 phases:-

- (i) Intermediate code generator
- (ii) Code optimizer.
- (iii) Code generator.

(i) Lexical Analysis :-

* This phase is also called as "Scanner" (or)
"Tokeniser".

* This phase reads the characters in the source program and groups them together to form a stream of tokens.

* Those sequence of characters should have a meaning.

* "Tokens" describes the category of input string.

* Tokens recognised during lexical analysis are usually

* Tokens of the following types:-

- (a) Identifiers.
- (b) Reserved words (a) Keywords
- (c) Operators
- (d) Constants.
- (e) Literals.
- (f) Punctuation symbols.

Eg: int a = 100; (5)

There are '5' tokens in the above statements.

- * A "lexeme" is a sequence of characters that is matched by the pattern to represent a token.
- * A "pattern" is a rule which describes a set of lexemes that can represent a token.

<u>Eg:</u>	<u>TOKEN</u>	<u>Lexeme</u>	<u>Pattern</u>
	=	=	letter followed by letters
	id	area,	
		pi, b2	and digits

* The lexical analyzer removes white spaces (blank space, new line, tab) and comments from the source program.

* Each token is augmented by a "lexical value" or "attribute value", which is an entry into the symbol table.

Eg:

a := b + c * 60



Lexical Analyzer



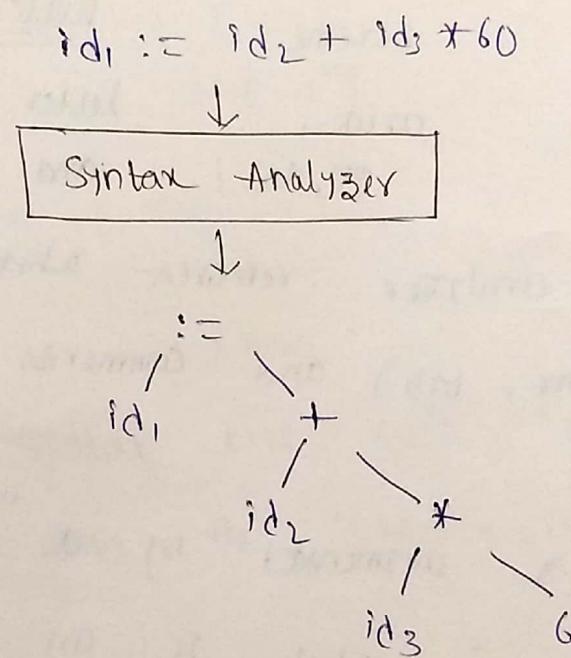
id₁ := id₂ + id₃ * 60

Token Stream = { id₁, id₂, id₃, GO, =, +, * }

Lexical values = { 1, 2, 3 }

(ii) Syntax Analysis :-

- * This phase is also called as "parser".
- * This phase takes the token stream as input and shows the structural representation of a program as output called "parse tree" (or) "syntax tree".
- * This phase is responsible for identifying syntax errors in the source program.



- * In the expression, " $a = b + c * 60$ ", multiplication is to be performed before addition because of "operator precedence".

- * The hierarchical structure is formed by following recursive rules (or) production rules.
- * Those rules are expressed by using CFG (Context free grammar)

* Content Free Grammar is a combination of (6)

$G = (V, P, T, S)$ where,

V = Variables (or) Non-Terminals

P = Productions

T = Terminals

S = Start symbol

For the expression "a = b + c * 60", the CFG will be

$S \rightarrow E = E$

| $E + E$

| $E * E$

$E \rightarrow id$

| num

* In the syntax tree, each "interior node" represents "Operators" and the "children nodes" represents "Operands".

(iii) Semantic Analysis :-

* This phase uses the syntax tree produced by the syntax analysis and checks the program consistency.

(ie) → correct no. of arguments to a function call.

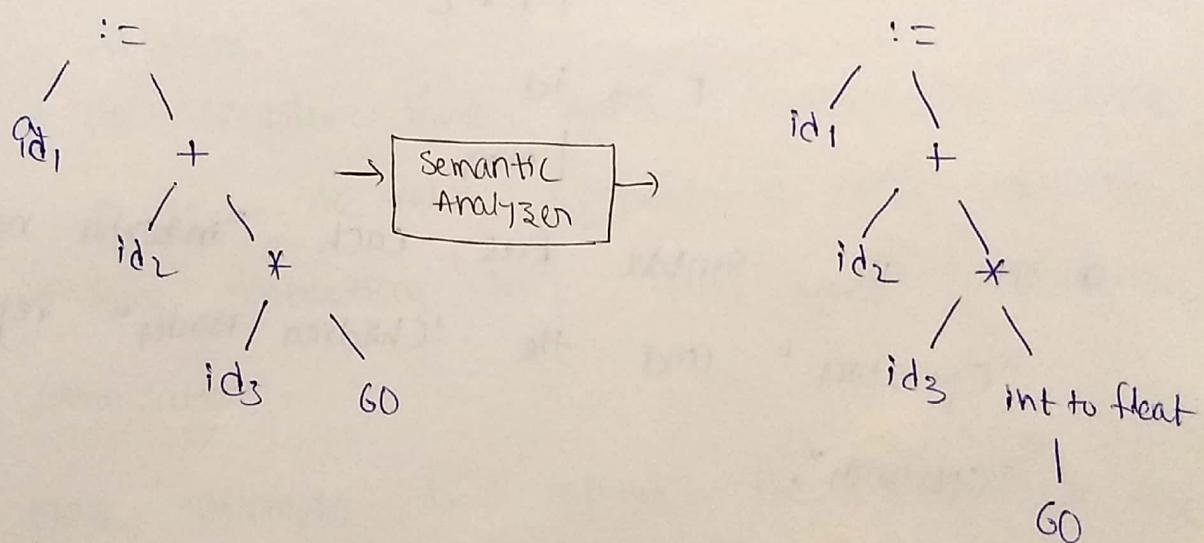
→ variables should be declared before used.

* This phase checks the source program for semantic errors and gathers information for next phase.

* "Type checking" is done by this phase.

Eg:- $a := b + c * 60 \Rightarrow id_1 := id_2 + id_3 * 60$

We assume that all identifiers are declared as "float" but '60' is an integer value. So, we need to convert the integer into float. This can be done by creating an extra node, converts integer to float value.



(iv) Intermediate Code Generator :-

* The intermediate code have 2 important properties:-

→ It is easy to produce.

→ Easy to translate into the target code.

* The intermediate code can be in different forms.

such as Three address code, quadruple, triple, RPN etc.

- * We consider an intermediate code or "Three address code", which is like assembly language. (7)
- * Every memory location act as a register.
- * The "three address code" consists of instructions, each of which has atmost 3 operands.
- * Properties of Three address code:-
- Each instruction has atmost one operator other than assignment operator.
 - The compiler has to decide the order of operations.
 - The compiler must generate a temporary variable to hold the value computed by each instruction.
 - Some instructions have fewer than 3 operands.

(v) Code Optimization :-

- * This phase attempts to improve the intermediate code so that "Faster running machine code" will be produced.
- * By optimizing the code, overall running time of the target program can be improved and less consumption of memory.

$$id_1 = id_2 + id_3 * 60$$

The intermediate code for the above statements is

$\text{temp1} = \text{int to float (60)}$

$\text{temp2} = id_3 * \text{temp1}$

$\text{temp3} = id_2 + \text{temp2}$

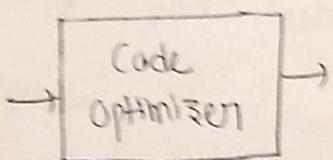
$id_1 = \text{temp3}$

$\text{temp1} = \text{int to float (60)}$

$\text{temp2} = id_3 * \text{temp1}$

$\text{temp3} = id_2 + \text{temp2}$

$id_1 = \text{temp3}$



$\text{temp1} = id_3 * 60.0$

$id_1 = id_2 + \text{temp1}$

(vi) Code Generator:

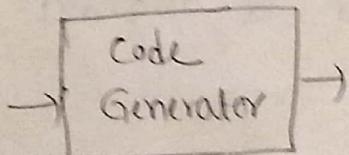
* This Phase contains machine code (or) Assembly code.

* Memory location are selected for each variable used in the program.

* The intermediate code is translated into a sequence of machine instructions which perform the same task.

$\text{temp1} = id_3 * 60.0$

$id_1 = id_2 + \text{temp1}$



MOVF id3, R1

MULF #600, R1

MOVF id2, R2

ADDF R2, R1

MOVF R1, id1

(8)

- * The first and second operands in each instruction represent source & destination respectively.
- * "F" in each instruction tells us those instructions deal with floating point numbers.
- * The "#" indicates that "60.0" is treated as a constant value.

Symbol Table Management :-

- * A symbol table is a data structure containing a record for each identifier and constant.
- * This data structure allows us to find the record for each identifier quickly and to store (or) retrieve data from the record quickly.
- * The lexical analyser identifies the identifiers & constants in the source program and that information is entered into the symbol table.

Compiler Construction Tools :-

- * Writing a Compiler is tedious and time consuming task.
- * There are some specialized tools which helps in implementation of various phases of a compiler.
- * They are :-
 - (i) Scanner generator.
 - (ii) Parser generator.
 - (iii) Syntax directed translation engine.
 - (iv) Automatic code generator.
 - (v) Data flow engine.

(i) Scanner Generator :-

- * This tool automatically generates lexical analyser phase based on "Regular Expression".
- * From a specification which describes a set of strings for a specific language over the input set.

Eg: Input set (Σ) = {0, 1}

A set of strings which ends with '10'

$$L = \{10, 010, 110, 0010, 0110, \dots\}$$

Then, the regular expression for the above language

$$\text{is } (0|1)^* 10$$

* UNIX operating system has the utility for generating a lexical analyzer called "LEX". (9)

(iii) Parser Generator :-

* This tool generates "syntax analyzer" phase.

* The specification given to this tool is in the form of CFG (Context Free Grammar)

* Context Free Grammar is :-

- A set of terminal symbols.
- A set of non-terminal symbols.
- A set of productions. Non-terminal symbols are replaced with other non-terminal (or) terminal symbols.
- A start symbol, which is a non-terminal symbol.

Eg: Consider an expression

$$a+b*c-d$$

The CFG for the above expression is

$$S \rightarrow a \mid b \mid c \mid S+S \mid S-S \mid S*S$$

check the grammar, whether the given grammar generates the given expression (or) not.

$$S \rightarrow S + \underline{S}$$

$$\rightarrow S + S * \underline{S}$$

$$\rightarrow S + S * S - \underline{S}$$

$$\rightarrow S + S * S - d$$

$$\rightarrow S + \underline{S} * \overline{c} - d$$

$$\rightarrow \underline{s} + b * c - d$$

$$\rightarrow a + b * c - d$$

* UNIX operating system has the utility for generating the syntax analyser called "YACC" (Yet Another Compiler Compiler)

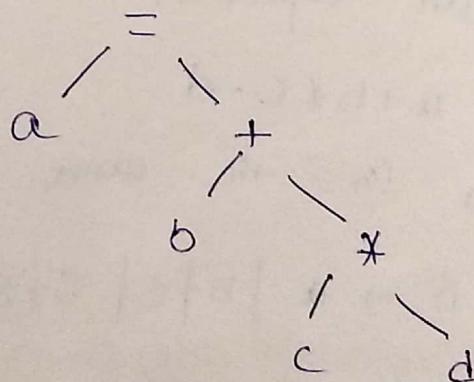
(iii) Syntax Directed Translation engine :-

- * By using this tool, the parse tree is scanned completely to generate "Intermediate Code".
- * The translation is done for each node of the tree.

Eg:- Consider a statement

$$a = b + c * d$$

The syntax tree for the above statement is



The intermediate code for the above syntax tree ,

$$\text{temp1} = c * d$$

$$\text{temp2} = b + \text{temp1}$$

$$a = \text{temp2}$$

(iv) Automatic Code Generator:- (10)

- * This tool takes an intermediate code as input and converts it into machine code.
- * In this tool, "Template Matching" technique is used.
- * The intermediate code statements are replaced by the templates, which represents machine instructions.

(v) Data Flow Engine:-

- * This tool is used for "Code optimization".
- * To perform good code optimization, we need data flow analysis.
- * This tool gathers information about how values are transmitted from one part of the program to the other.

Types of Compiler:

(i) Incremental Compiler:

This type of compiler performs recompilation of only modified source code rather than compiling the whole source program.

(ii) Cross Compiler:

A compiler which runs on one machine and produces the target code for another machine.

A Simple One-pass Compiler

- * In general, compiler is a program that reads a program written in one language, which is called the "Source language" and translates it into another language, which is called the "Target language".
- * Traditionally, source language was a "high-level language" and target language was a "low level language".
- * "Pass" and "phase" are 2 terms often used with compilers.
- * A "pass" refers to the traversal of a compiler through the entire program.
- * No. of "passes" of a compiler is the no. of times it goes over the source program.
- * A "phase" is a single-independent part of a compiler.
- * A standard way to classify compilers is by the no. of passes.
- * Initially, computers did not have enough memory to hold such a program that did the complete job.
- * Due to this limitation, compilers were broken down into smaller sub programs that did its partial job over the source code.

- * "One-pass Compiler" compiles the entire source program in a single pass. (11)
- * It is easier to write a one-pass compiler and they perform faster than "Multi-pass compilers".
- * A multi-pass compiler is made up of several stages. Each stage is considered as a "phase".
- * The advantage of having different phases is that the development of compiler can be distributed.
- * It improves the modularity and reuse by allowing phases to be replaced by improved ones or additional phases to the compiler.

Context Free Grammar (CFG):-

- * A grammar is used to specify the syntax of a language.
- * A grammar describes the hierarchical structure of any programming language.
- * The Context Free Grammar has 4 components:-
 - A set of terminal symbols.
 - A set of non-terminal symbols.
 - A set of productions, where each production consists of :-

- a non-terminal on left-hand side of a production.
- an arrow and
- a sequence of terminals (or) non-terminals on right-hand side of a production.

Eg: Consider a statement

$$a + b - c$$

The CFG for the above statement

$$E \rightarrow E+E$$

$$| E-E$$

$$E \rightarrow a | b | c$$

check whether the above statement is generated from the grammar or not.

$$E \rightarrow E+E$$

$$\rightarrow E+E-E$$

$$\rightarrow E+E-\underline{E}$$

$$\rightarrow \underline{E}+b-c$$

$$\rightarrow a+b-c$$

Derivation :-

- * It is the process of deriving an input string from the given grammar.
- * A grammar derives strings by beginning with the start symbol and repeatedly replacing a non-terminal by the body of a production.
- * The geometrical representation of derivation is called "Derivation Tree".
- * Derivations are of 2 types:-
 - (i) Left-most derivation.
 - (ii) Right-most derivation.

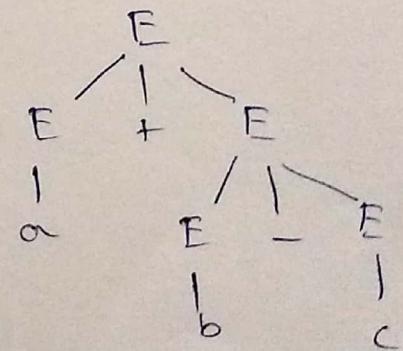
(i) Left-most derivation:-

The process of constructing a parse tree by expanding the left most non-terminal first in the derivation process is called "Left most derivation".

Eg: $\text{CFG} : E \rightarrow E + E$ input: a+b-c
 $|\underline{E-E}$
 $|a|b|c$

Derivation: $F \rightarrow \underline{E+E}$
 $\rightarrow a+\underline{E}$
 $\rightarrow a+\underline{E-E}$
 $\rightarrow a+b-\underline{E}$
 $\rightarrow a+b-c$

Parse tree:-



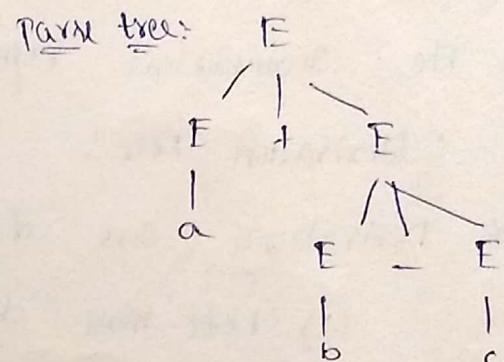
(ii) Right-most Derivation:

The process of constructing a parse tree by expanding the right-most non-terminal first in the derivation process is called "Right-most Derivation".

$$\text{Ex: } \text{CFG: } E \rightarrow E+E \quad | \quad E-E \quad | \quad a \quad | \quad b \quad | \quad c$$

input: a+b-c

$$\begin{aligned} \text{Derivation: } & E \rightarrow E+E \\ & \rightarrow E+E-E \\ & \rightarrow E+E-c \\ & \rightarrow E+b-c \\ & \rightarrow a+b-c \end{aligned}$$



$$\text{Ex: } E \rightarrow E+E \quad | \quad E+E \quad | \quad E-E \quad | \quad a \quad | \quad b \quad | \quad c \quad | \quad d$$

input: a+b*c-d

$$\begin{aligned} \text{RMD: } & E \rightarrow E+E \\ & \rightarrow E+E*E \\ & \rightarrow E+E*E-E \\ & \rightarrow E+E*E-d \\ & \rightarrow E+E*c-d \\ & \rightarrow E+b*c-d \\ & \rightarrow a+b*c-d \end{aligned}$$

$$\begin{aligned} \text{LMD: } & E \rightarrow E+E \\ & \rightarrow a+E \\ & \rightarrow a+E+E \\ & \rightarrow a+b+E \\ & \rightarrow a+b*c-E \\ & \rightarrow a+b*c-d \end{aligned}$$

Ambiguous Grammar :-

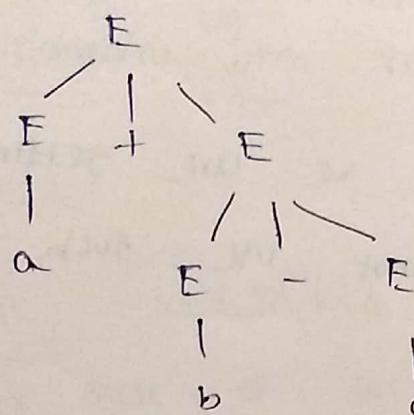
A grammar is said to be ambiguous, if it generates more than one parse tree.

Eg:- $E \rightarrow E + E$

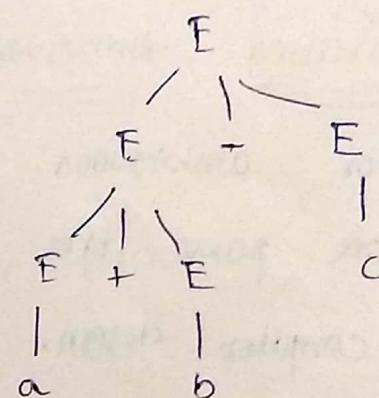
$$| \\ E - E$$

| a | b | c

input:- a + b - c



(i)



(ii)

The above grammar generates 2 parse trees for the given input string. Hence, the given grammar is Ambiguous.

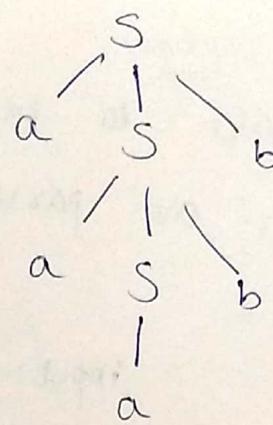
Unambiguous grammar :-

A grammar is said to be unambiguous, if it generates only one parse tree.

Eg:- $S \rightarrow aSb$

| a

input:- aaabb



The above grammar generates only one parse tree. Hence,
the given grammar is unambiguous.

Converting Ambiguous grammar into unambiguous grammar:-

- * For ambiguous grammar, we are getting more than one parse tree. We cannot use such grammar in compiler design.
- * Ambiguous grammar can be converted into unambiguous grammar by following "Disambiguity Rule".
- * If we see what are the reasons for ambiguity and if we work on that root cause we can have unambiguous grammar.

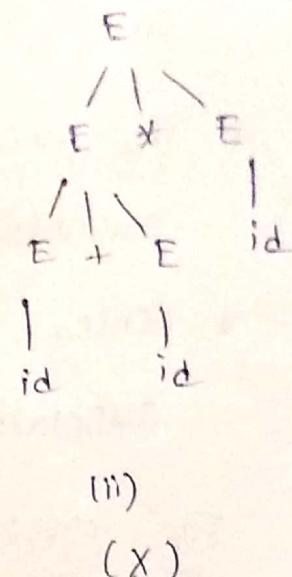
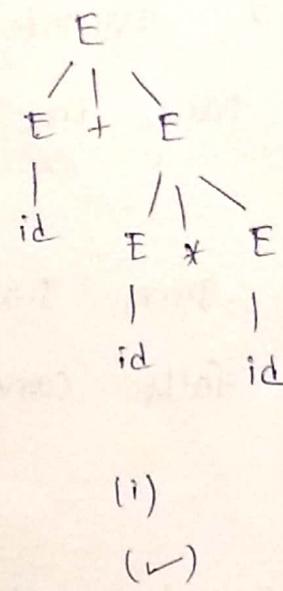
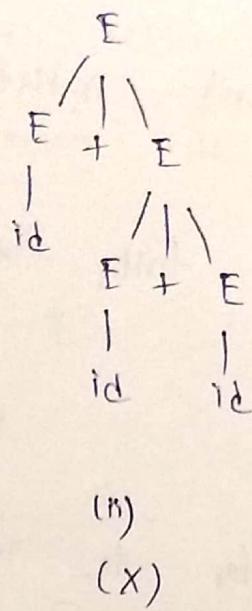
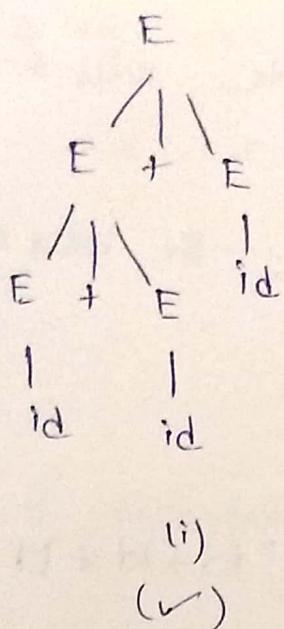
$$\text{Eg: } E \rightarrow E + E \\ | E * E \\ | id$$

input 1 :- $id + id + id$

input 2 :- $id + id * id$

Input = id + id + id

Input = id + id * id



For input:

= "id + id + id", whenever there are 2 operators
* we have "id + id + id", whenever there are 2 operators
on either side of an operand, which operator should
associate with the operand.

$(id + (id) + id)$
 ↙ ↘
 left Right
 associative associative

* In case of "+", the final output doesn't depend
on what the actual associativity is.
* But, if we consider some other operators other than '+',
which have equal precedence then it should be "left
associative".

* In the first parse tree, the left most '+' is
evaluated first. In the second parse tree, the right
most '+' is evaluated first.

- * The first parse tree is the correct form tree based on rules of associativity.
- * The second parse tree doesn't satisfy the rules of associativity.
- * Hence, the given grammar fails because the rules of associativity fail completely.

For input:-

- * we got 2 parse trees for the input "id + id * id".
- * When we have different operators of different precedence present in an expression, the highest precedence operator should be evaluated first.
- * In the first parse tree, the '*' is evaluated first and then '+' is evaluated.
- * In the second parse tree, the '+' is evaluated first and then '*' is evaluated.
- * Among those 2 parse trees, the first parse tree is the correct parse tree which satisfies the "operator precedence rules".
- * Hence, the given grammar fails because the operator precedence is not taken care.
- * If we take care of both operator precedence & associativity then the grammar is unambiguous.

(15)

Solution for Associativity problem :-

* Why we are getting both associativity ?

We are defining the grammar without any order.

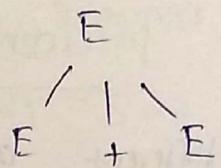
$$E \rightarrow E+E$$

$$\quad | \quad E+E$$

1 id

Input:- id + id + id

* According to the grammar,



If we want one more '+', we could go either left side (or) right side.

* If we grow the parse tree in left direction,
then it is left associativity.

* If we grow the parse tree in right direction, then
it is right associativity.

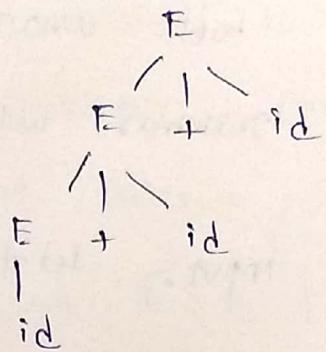
* So, we need to restrict the grammar to grow the
parse tree in only one direction (i.e.) left

* Rewrite the grammar to have left associativity. Then,

$$E \rightarrow E+id$$

1 id

input : id + id * id



- * The given grammar is defined as left recursive, which grows the parse tree in left direction only.

Solution for operator precedence problem :-

= = = = = different

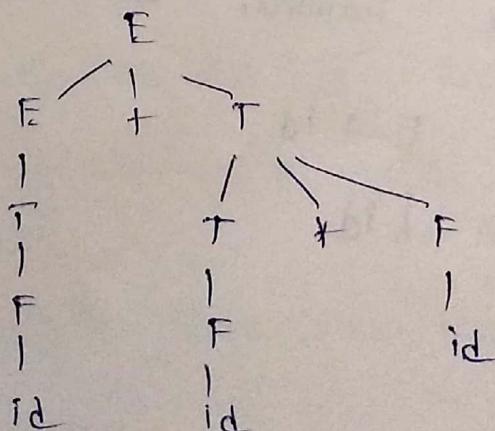
- * When there are different operators of different operator precedences, the highest precedence operator should be at bottom (or) least level. Then only it is evaluated first.

- * To fix this precedence problem, we need to define the grammar to be of different levels.

$$E \rightarrow E+E \\ | E+E \\ | id$$

$$E \rightarrow E+E|T \\ T \rightarrow T*T|F \\ F \rightarrow id$$

input : id + id * id



Recursive Grammar :-

A grammar is said to be recursive, if atleast one of the production contains same non-terminal symbols on both sides.

Eg:-

$$(i) \quad S \rightarrow ASB \mid \epsilon$$

$$A \rightarrow a$$

$$B \rightarrow b$$

The given grammar is recursive because the production $S \rightarrow ASB$ contains same non-terminal symbol 'S' on both sides.

$$(ii) \quad S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

The given grammar is not recursive. Because, no production contains the same non-terminal symbols on both sides.

$$(iii) \quad S \rightarrow AB$$

$$A \rightarrow aB \mid a$$

$$B \rightarrow aA \mid b$$

The given grammar is "Indirect recursive" because the production $A \rightarrow aB$
 $\rightarrow aaA$

derives the same non-terminal 'A' indirectly.

Non - recursive grammar :-

A grammar is said to be non - recursive , if no production contains same non-terminal symbols on both sides.

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

The given grammar is non - recursive because no production contains same non-terminal symbols on both sides.

Types of Recursive Grammar:-

Recursive grammar are of 2 types:-

(i) Left recursive

(ii) Right recursive

(i) Left Recursive:-

A grammar is said to be left recursive , if the left most non-terminal symbol on right hand side of a production is equal to the non-terminal symbol on LHS of production .

$$\begin{array}{l} S \rightarrow Sa \\ | b \end{array}$$

The given grammar is left recursive

(iii) Right Recursion:-

* grammar is said to be right recursive, if the right most non-terminal symbol on RHS of a production is equal to the non-terminal symbol on LHS of a production.

$$\text{Ex: } S \rightarrow aS \\ | a$$

Left Recursion:-

* Because of left recursive grammar, the top down parser can enter in to "Infinite loop".

$$\text{Ex: } E \rightarrow E + T$$

Every non-terminal have its own procedure which describes its production. In the above production, the procedure of ' E ' calls its own procedure recursively. Then, infinite loop will be formed.

* A left recursive production can be eliminated by rewriting the production.

If the production is in the form

$$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m | \beta_1 | \beta_2 | \dots | \beta_n$$

where α_i & β_j are the sequence of terminals and non-terminals that doesn't start with ' A '.

We can eliminate the left recursion by rewriting the production rule as:-

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

Ex: $E \rightarrow E + T$
 | T

The production $E \rightarrow E + T$ is in the form of $A \rightarrow A\alpha \mid \beta$. where, $A = E$, $\alpha = +T$, $\beta = T$.

Then, the left recursion can be eliminated

as follows:-

$$\begin{array}{ccc} E \rightarrow E + T & & E \rightarrow TE' \\ | T & \Rightarrow & E' \rightarrow +TE' \mid \epsilon \end{array}$$

Ex: $T \rightarrow T * F$
 | F

$$A = T, \alpha = *F, \beta = F$$

$$\begin{array}{ccc} T \rightarrow T * F & & T \rightarrow FT' \\ | F & \Rightarrow & T' \rightarrow *FT' \mid \epsilon \end{array}$$

Ex: $S \rightarrow S_0 S_1 S$
 | 01

$$A = S, \alpha = 0S_1S, \beta = 01$$

$$\begin{array}{ccc} S \rightarrow S_0 S_1 S & & S \rightarrow 0S_1 \\ | 01 & \Rightarrow & S' \rightarrow 0S_1SS_1 \\ & & | \epsilon \end{array}$$

Eg:-

$$S \rightarrow (L)$$

| x

$$L \rightarrow L, S$$

| S

$$A = L, \alpha = , S, \beta = S$$

$$S \rightarrow (L)$$

| x

$$L \rightarrow L, S$$

| S

$$S \rightarrow (L)$$

| x

$$L \rightarrow SL'$$

$$L' \rightarrow , SL'$$

| E

NON-Deterministic Grammar:-

A grammar is said to be non-deterministic, if a production having multiple forms which contains a common prefix.

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_n$$

where α' is the common prefix in multiple forms of a production

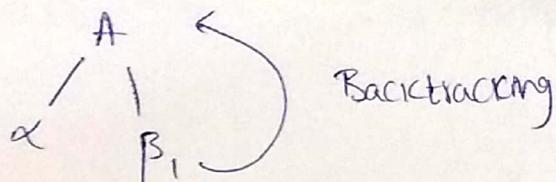
Problem of non-deterministic grammar:-

Consider an example

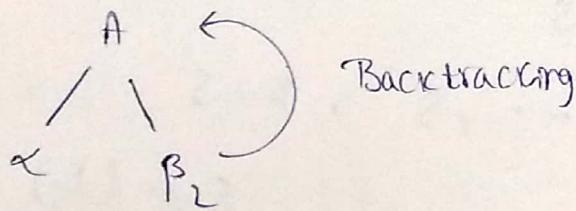
$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \alpha\beta_3$$

Input :- $\alpha\beta_3$

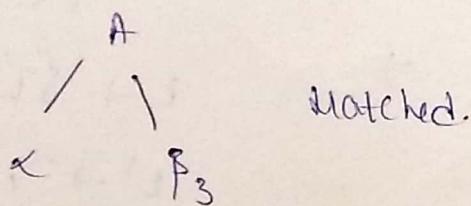
(i)



(ii)



(iii)



* Backtracking is happening because of common prefix.

* Backtracking is happening because we are actually making a decision of which production to choose without seeing enough of the input

* Given input ' $\alpha\beta_3$ ', we are making decision by seeing only ' α '.

* We have to postpone the decision making, this can be done by using "Left Factoring".

Left Factoring :-

* This is the process of Converting non-deterministic grammar to deterministic grammar.

* Common prefixes are to be written separately.

* If the production is in the form :-

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n$$

↓,

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

Ex:- $S \rightarrow iE + S$

$$\mid iE + S$$

where $\alpha = iE + S$

↓ a

$$E \rightarrow b$$

By applying left factoring the grammar is converted

Q1 :-

$$S \rightarrow iE + S S'$$

↓ a

$$S' \rightarrow eS$$

↓ e

$$E \rightarrow b$$

Ex:- ~~S → iE + S~~ $S \rightarrow aSSbS$

$$\mid iE + S$$

↓ a

~~E → S~~

$$\mid aSaSb$$

↓ abb

↓ b

Here $\alpha = a$

After applying left factoring then

$$S \rightarrow a S'$$

| b

$$S' \rightarrow S S b S$$

| S a S b

| b b

Apply left Factoring again. Then,

$$S \rightarrow a S'$$

| b

$$S' \rightarrow S S''$$

| b b

$$S'' \rightarrow S b S$$

| a S b

Eg: $S \rightarrow b S S a a S$

| b S S a a S b Here $\alpha = b S$, then

| b S b

| a

$$S \rightarrow b S S' | a$$

$$S' \rightarrow S a S$$

Here $\alpha = S a$, then

| S a S b

| b

$$S \rightarrow b S S' | a$$

$$S' \rightarrow S a S'' | b$$

$$S'' \rightarrow a S | S b$$

Parsing:

- * It is the process of determining if a string token can be generated by a grammar.
- * The parser scans the input string from left to right and identifies that derivation is leftmost or rightmost.
- * The parser use production rule for choosing the appropriate derivation. Finally a parse tree is constructed.
- * When the parse tree can be constructed from "null" and expanded to leaves. such type of parsing is called "Top down parsing".
- * When the parse tree can be constructed from leaves to root, such type of parsing is called "Bottom up parsing".

2.5 A TRANSLATOR FOR SIMPLE EXPRESSIONS

Using the techniques of the last three sections, we now construct a syntax-directed translator, in the form of a working C program, that translates arithmetic expressions into postfix form. To keep the initial program manageable small, we start off with expressions consisting of digits separated by plus and minus signs. The language is extended in the next two sections to include numbers, identifiers, and other operators. Since expressions appear as a construct in so many languages, it is worth studying their translation in detail.

A syntax-directed translation scheme can often serve as the specification for a translator. We use the scheme in Fig. 2.19 (repeated from Fig. 2.13) as the definition of the translation to be performed. As is often the case, the underlying grammar of a given scheme has to be modified before it can be parsed with a predictive parser. In particular, the grammar underlying the scheme in Fig. 2.19 is left-recursive, and as we saw in the last section, a predictive parser cannot handle a left-recursive grammar. By eliminating the left-recursion, we can obtain a grammar suitable for use in a predictive recursive-descent translator.

$expr \rightarrow expr + term$	{ print('+') }
$expr \rightarrow expr - term$	{ print('−') }
$expr \rightarrow term$	
$term \rightarrow 0$	{ print('0') }
$term \rightarrow 1$	{ print('1') }
...	
$term \rightarrow 9$	{ print('9') }

Fig. 2.19. Initial specification of infix-to-postfix translator.

Abstract and Concrete Syntax

A useful starting point for thinking about the translation of an input string is an *abstract syntax tree* in which each node represents an operator and the children of the node represent the operands. By contrast, a parse tree is called a *concrete syntax tree*, and the underlying grammar is called a *concrete syntax* for the language. Abstract syntax trees, or simply *syntax trees*, differ from parse trees because superficial distinctions of form, unimportant for translation, do not appear in syntax trees.

For example, the syntax tree for $9 - 5 + 2$ is shown in Fig. 2.20. Since + and - have the same precedence level, and operators at the same precedence level are evaluated left to right, the tree shows $9 - 5$ grouped as a subexpression. Comparing Fig. 2.20 with the corresponding parse tree of Fig. 2.2, we note that the syntax tree associates an operator with an interior node, rather than making the operator be one of the children.

It is desirable for a translation scheme to be based on a grammar whose parse trees are as close to syntax trees as possible. The grouping of subexpressions by the grammar in Fig. 2.19 is similar to their grouping in syntax trees. Unfortunately, the grammar of Fig. 2.19 is left-recursive, and hence not suitable for predictive

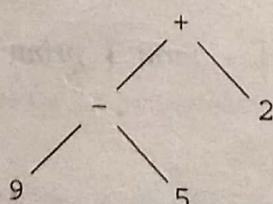


Fig. 2.20. Syntax tree for $9 - 5 + 2$.

Adapting the Translation Scheme

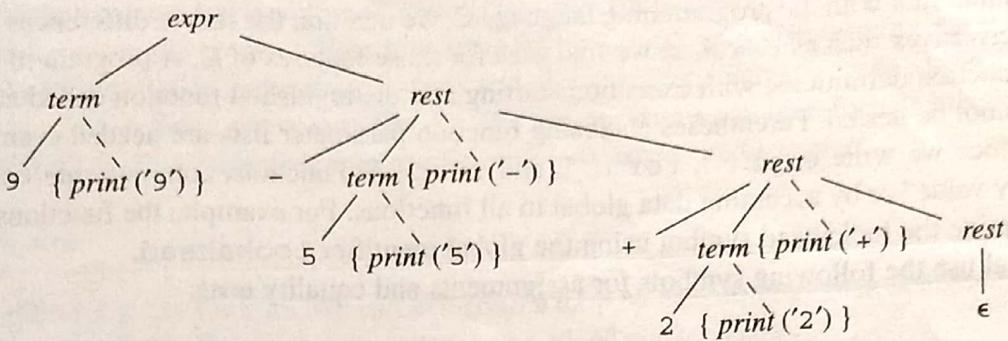
The left-recursion elimination technique sketched in Fig. 2.18 can also be applied to productions containing semantic actions. We extend the transformation in Section 5.5 to take synthesized attributes into account. The technique transforms the productions $A \rightarrow A\alpha | A\beta | \gamma$ into

$$\begin{aligned} A &\rightarrow \gamma R \\ R &\rightarrow \alpha R | \beta R | \epsilon \end{aligned}$$

When semantic actions are embedded in the productions, we carry them along in the transformation. Here, if we let $A = \text{expr}$, $\alpha = + \text{term} \{ \text{print}('+) \}$, $\beta = - \text{term} \{ \text{print}('-) \}$, and $\gamma = \text{term}$, the transformation above produces the translation scheme (2.14). The expr productions in Fig. 2.19 have been transformed into the productions for expr and the new nonterminal rest in (2.14). The productions for term are repeated from Fig. 2.19. Notice that the underlying grammar is different from the one in Example 2.9 and the difference makes the desired translation possible.

$$\begin{aligned} \text{expr} &\rightarrow \text{term rest} \\ \text{rest} &\rightarrow + \text{term} \{ \text{print}('+) \} \text{ rest} | - \text{term} \{ \text{print}('-) \} \text{ rest} | \epsilon \\ \text{term} &\rightarrow 0 \{ \text{print}('0') \} \\ \text{term} &\rightarrow 1 \{ \text{print}('1') \} \\ \text{term} &\rightarrow 9 \{ \text{print}('9') \} \end{aligned} \tag{2.14}$$

Figure 2.21 shows how $9-5+2$ is translated using the above grammar.

Fig. 2.21. Translation of $9 - 5 + 2$ into $95 - 2 + \epsilon$.

Procedures for the Nonterminals *expr*, *term*, and *rest*

We now implement a translator in C using the syntax-directed translation scheme (2.14). The essence of the translator is the C code in Fig. 2.22 for the functions *expr*, *term*, and *rest*. These functions implement the corresponding nonterminals in (2.14).

```

expr()
{
    term(); rest();
}

rest()
{
    if (lookahead == '+') {
        match('+'); term(); putchar('+'); rest();
    }
    else if (lookahead == '-') {
        match('-'); term(); putchar('-'); rest();
    }
    else ;
}

term()
{
    if (isdigit(lookahead)) {
        putchar(lookahead); match(lookahead);
    }
    else error();
}

```

Fig. 2.22. Functions for the nonterminals *expr*, *rest*, and *term*.

The function *match*, presented later, is the C counterpart of the code in Fig. 2.17 to match a token with the lookahead symbol and advance through the input. Since each token is a single character in our language, *match* can be implemented by comparing and reading characters.

3. Lexical Analysis

The role of lexical analyzer:-

- * Lexical analyzer is the first phase of a compiler.
- * It reads the input source program from left to right one character at a time and generates the sequence of tokens.
- * Tokens recognized during lexical analysis are usually of the following types:-
 - (i) Identifiers.
 - (ii) keywords.
 - (iii) constants.
 - (iv) operators.
 - (v) symbols.

- * The lexical analyzer interacts with the syntal analyzer and symbol table & Error report handler.

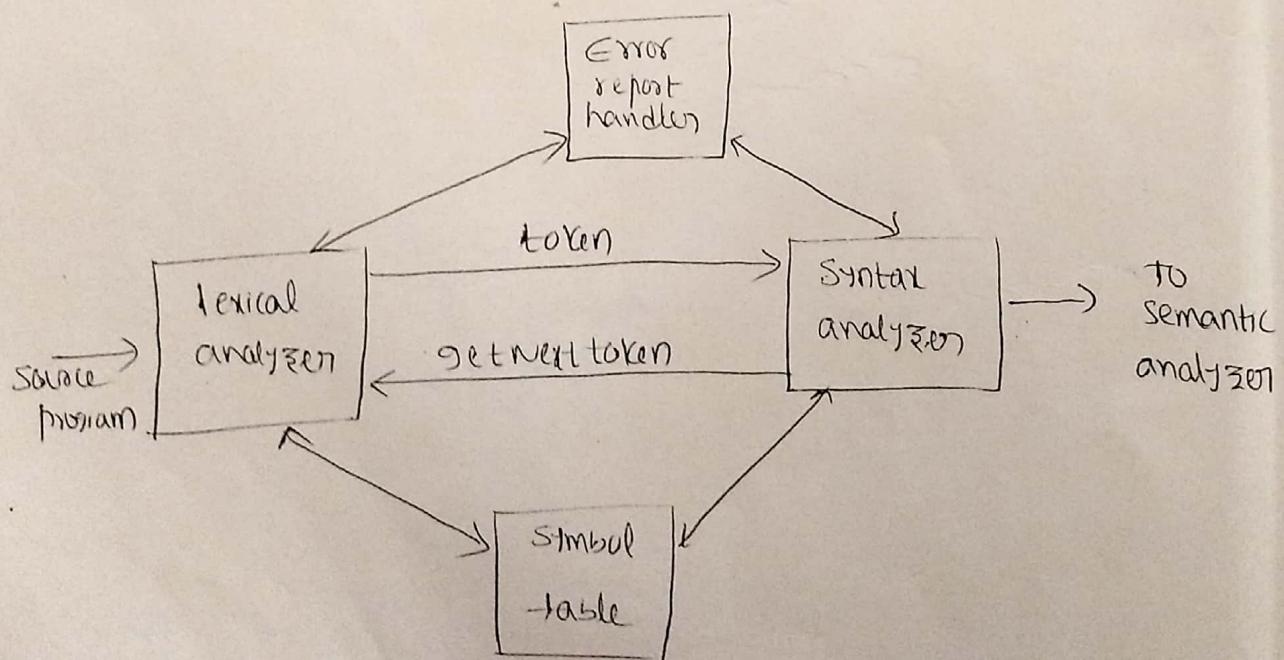


fig: Lexical analyzer interaction with other phases.

The stream of tokens is sent to the syntax analyzer. 27

40

A token is a pair consisting of token name & attribute value (or) lexical value, which describes the category of input string.

Lexical value (or) attribute value is an entry into the symbol table.

The character sequence for forming a token is called the "lexeme" for a token.

The lexical analyzer generates one token at a time & whenever the remaining tokens get generated to the lexical syntax analyzer makes a request.

"getNextToken" command causes the lexical analyzer to generate next token which returns to the syntax analyzer.

The lexical analyzer is responsible for comments and whitespaces.

→ removing stream of tokens.

→ produce

Lexical analyzer keeps track of line numbers wherever a new line character (\n) is seen. This could be helpful in identifying the error in the program.

Ex:
int max(^③ int ^④ x, ^⑤ int ^⑥ y)

⑩ 1

^⑩ return (^⑪ ^⑫ ^⑬ ^⑭ x > ^⑮ y ? ^⑯ ^⑰ ^⑱ ^⑲ x : ^⑳ y);

⑲ }
⑳ }

Mark tolerance = 22

lexeme

int

max

(

int

x

,

int

y

)

{

return

(

x

>

y

?

x

:

TOKEN

Keyword

identifier

symbol

Keyword

identifier

operator

Keyword

identifier

symbol

symbol

Keyword

symbol

identifier

operator

identifier

operator

identifier

operator

lexeme

TOKEN

identifier

Symbol

Symbol

Symbol

Symbol

3.2 INPUT BUFFERING

This section covers some efficiency issues concerned with the buffering of input. We first mention a two-buffer input scheme that is useful when look-ahead on the input is necessary to identify tokens. Then we introduce some useful techniques for speeding up the lexical analyzer, such as the use of "sentinels" to mark the buffer end.

There are three general approaches to the implementation of a lexical analyzer.

1. Use a lexical-analyzer generator, such as the Lex compiler discussed in Section 3.5, to produce the lexical analyzer from a regular-expression-based specification. In this case, the generator provides routines for reading and buffering the input.
2. Write the lexical analyzer in a conventional systems-programming language, using the I/O facilities of that language to read the input.
3. Write the lexical analyzer in assembly language and explicitly manage the reading of input.

The three choices are listed in order of increasing difficulty for the implementor. Unfortunately, the harder-to-implement approaches often yield faster lexical analyzers. Since the lexical analyzer is the only phase of the compiler that reads the source program 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 lexical analysis is a concern in compiler design. While the bulk of the chapter is devoted to the first approach, the design and use of an automatic generator, we also consider techniques that are helpful in manual design. Section 3.4 discusses transition diagrams, which are a useful concept for the organization of a hand-designed lexical analyzer.

Buffer Pairs

For many source languages, there are times when the lexical analyzer needs to look ahead several characters beyond the lexeme for a pattern before a match can be announced. The lexical analyzers in Chapter 2 use a function `ungetc` to push lookahead characters back into the input stream. Because a large amount of time can be consumed moving characters, specialized buffering techniques have been developed to reduce the amount of overhead required to process an input character. Many buffering schemes can be used, but since the techniques are somewhat dependent on system parameters, we shall only outline the principle behind one class of schemes here.

We use a buffer divided into two N -character halves, as shown in Fig. 3.3. Typically, N is the number of characters on one disk block, e.g., 1024 or 4096.

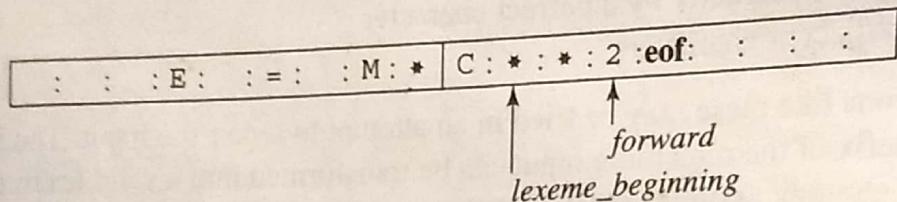


Fig. 3.3. An input buffer in two halves.

We read N input characters into each half of the buffer with one system read command, rather than invoking a read command for each input character. If fewer than N characters remain in the input, then a special character `eof` is read into the buffer after the input characters, as in Fig. 3.3. That is, `eof` marks the end of the source file and is different from any input character.)

Two pointers to the input buffer are maintained. The string of characters between the two pointers is the current lexeme. Initially, both pointers point to the first character of the next lexeme to be found. One, called the forward pointer, scans ahead until a match for a pattern is found. Once the next lexeme is determined, the forward pointer is set to the character at its right end. After the lexeme is processed, both pointers are set to the character immediately past the lexeme. With this scheme, comments and white space can be treated as patterns that yield no token.

If the forward pointer is about to move past the halfway mark, the right half is filled with N new input characters. If the forward pointer is about to move past the right end of the buffer, the left half is filled with N new characters and the forward pointer 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, and this limited lookahead may make it impossible to recognize tokens in situations where the distance that the forward pointer must travel is more than the length of the buffer. For example, if we see

```
DECLARE ( ARG1, ARG2, ..., ARGn )
```

in a PL/I program, we cannot determine whether `DECLARE` is a keyword or an array name until we see the character that follows the right parenthesis. In either case, the lexeme ends at the second `E`, but the amount of lookahead needed is proportional to the number of arguments, which in principle is unbounded.

Sentinels

If we use the scheme of Fig. 3.3 exactly as shown, we must check each time we move the forward pointer that we have not moved off one half of the buffer; if we do, then we must reload the other half. That is, our code for advancing the forward pointer performs tests like those shown in Fig. 3.4.

```

if forward at end of first half then begin
    reload second half;
    forward := forward + 1
end
else if forward at end of second half then begin
    reload first half;
    move forward to beginning of first half
end
else forward := forward + 1;

```

Fig. 3.4. Code to advance forward pointer.

Except at the ends of the buffer halves, the code in Fig. 3.4 requires two tests for each advance of the forward pointer. We can reduce the two tests to one if we extend each buffer half to hold a *sentinel* character at the end. The *sentinel* is a special character that cannot be part of the source program. A natural choice is *eof*; Fig. 3.5 shows the same buffer arrangement as Fig. 3.3, with the sentinels added.

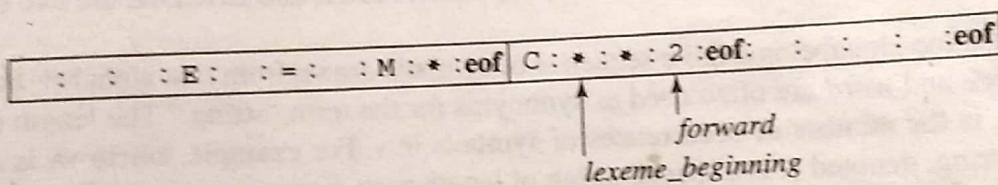


Fig. 3.5. Sentinels at end of each buffer half.

With the arrangement of Fig. 3.5, we can use the code shown in Fig. 3.6 to advance the forward pointer (and test for the end of the source file). Most of the time the code performs only one test to see whether *forward* points to an *eof*. Only when we reach the end of a buffer half or the end of the file do we perform more tests. Since *N* input characters are encountered between *eof*'s, the average number of tests per input character is very close to 1.

```

forward := forward + 1;
if forward = eof then begin
    if forward at end of first half then begin
        reload second half;
        forward := forward + 1
    end
    else if forward at end of second half then begin
        reload first half;
        move forward to beginning of first half
    end
    else /* eof within a buffer signifying end of input */
        terminate lexical analysis
end

```

Fig. 3.6. Lookahead code with sentinels.

3.3 SPECIFICATION OF TOKENS

Regular expressions are an important notation for specifying patterns. Each pattern matches a set of strings, so regular expressions will serve as names for sets of strings. Section 3.5 extends this notation into a pattern-directed language for lexical analysis.

Strings and Languages

The term *alphabet* or *character class* denotes any finite set of symbols. Typical examples of symbols are letters and characters. The set $\{0, 1\}$ is the *binary alphabet*. ASCII and EBCDIC are two examples of computer alphabets.

A *string* over some alphabet is a finite sequence of symbols drawn from that alphabet. In language theory, the terms *sentence* and *word* are often used as synonyms for the term “string.” The length of a string s , usually written $|s|$, is the number of occurrences of symbols in s . For example, banana is a string of length six. The *empty string*, denoted ϵ , is a special string of length zero. Some common terms associated with parts of a string are summarized in Fig. 3.7.

The term *language* denotes any set of strings over some fixed alphabet. This definition is very broad. Abstract languages like \emptyset , the *empty set*, or $\{\epsilon\}$, the set containing only the empty string, are languages under this definition. So too are the set of all syntactically well-formed Pascal programs and the set of all grammatically correct English sentences, although the latter two sets are much more difficult to specify. Also note that this definition does not ascribe any meaning to the strings in a language. Methods for ascribing meanings to strings are discussed in Chapter 5.

TERM	DEFINITION
<i>prefix of s</i>	A string obtained by removing zero or more trailing symbols of string s ; e.g., ban is a prefix of banana.
<i>suffix of s</i>	A string formed by deleting zero or more of the leading symbols of s ; e.g., nana is a suffix of banana.
<i>substring of s</i>	A string obtained by deleting a prefix and a suffix from s ; e.g., nan is a substring of banana. Every prefix and every suffix of s is a substring of s , but not every substring of s is a prefix or a suffix of s . For every string s , both s and ϵ are prefixes, suffixes, and substrings of s .
<i>proper prefix, suffix, or substring of s</i>	Any nonempty string x that is, respectively, a prefix, suffix, or substring of s such that $s \neq x$.
<i>subsequence of s</i>	Any string formed by deleting zero or more not necessarily contiguous symbols from s ; e.g., baaa is a subsequence of banana.

Fig. 3.7. Terms for parts of a string.

If x and y are strings, then the *concatenation* of x and y , written xy , is the string formed by appending y to x . For example, if $x = \text{dog}$ and $y = \text{house}$, then $xy = \text{doghouse}$. The empty string is the identity element under concatenation. That is, $s\epsilon = \epsilon s = s$.

If we think of concatenation as a "product", we can define string "exponentiation" as follows. Define s^0 to be ϵ , and for $i > 0$ define s^i to be $s^{i-1}s$. Since ϵs is s itself, $s^1 = s$. Then, $s^2 = ss$, $s^3 = sss$, and so on.

Operations on Languages

There are several important operations that can be applied to languages. For lexical analysis, we are interested primarily in union, concatenation, and closure, which are defined in Fig. 3.8. We can also generalize the "exponentiation" operator to languages by defining L^0 to be $\{\epsilon\}$, and L^i to be $L^{i-1}L$. Thus, L^i is L concatenated with itself $i - 1$ times.

Example 3.2. Let L be the set $\{A, B, \dots, Z, a, b, \dots, z\}$ and D the set $\{0, 1, \dots, 9\}$. We can think of L and D in two ways. We can think of L as the alphabet consisting of the set of upper and lower case letters, and D as the alphabet consisting of the set of the ten decimal digits. Alternatively, since a symbol can be regarded as a string of length one, the sets L and D are each finite languages. Here are some examples of new languages created from L and D by applying the operators defined in Fig. 3.8.

OPERATION	DEFINITION
<i>union</i> of L and M written $L \cup M$	$L \cup M = \{ s \mid s \text{ is in } L \text{ or } s \text{ is in } M \}$
<i>concatenation</i> of L and M written LM	$LM = \{ st \mid s \text{ is in } L \text{ and } t \text{ is in } M \}$
<i>Kleene closure</i> of L written L^*	$L^* = \bigcup_{i=0}^{\infty} L^i$ L^* denotes "zero or more concatenations of" L .
<i>positive closure</i> of L written L^+	$L^+ = \bigcup_{i=1}^{\infty} L^i$ L^+ denotes "one or more concatenations of" L .

Fig. 3.8. Definitions of operations on languages.

1. $L \cup D$ is the set of letters and digits.
2. LD is the set of strings consisting of a letter followed by a digit.
3. L^4 is the set of all four-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 and digits beginning with a letter.
6. D^+ is the set of all strings of one or more digits.

□

Regular Expressions

In Pascal, an identifier is a letter followed by zero or more letters or digits; that is, an identifier is a member of the set defined in part (5) of Example 3.2. In this section, we present a notation, called regular expressions,

Regular Definitions

For notational convenience, we may wish to give names to regular expressions and to define regular expressions using these names as if they were symbols. If Σ is an alphabet of basic symbols, then a *regular definition* is a sequence of definitions of the form

$$\begin{aligned}d_1 &\rightarrow r_1 \\d_2 &\rightarrow r_2 \\\dots \\d_n &\rightarrow r_n\end{aligned}$$

where each d_i is a distinct name, and each r_i is a regular expression over the symbols in $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$, i.e., the basic symbols and the previously defined names. By restricting each r_i to symbols of Σ and the previously defined names, we can construct a regular expression over Σ for any r_i by repeatedly replacing regular-expression names by the expressions they denote. If r_i used d_j for some $j \geq i$, then r_i might be recursively defined, and this substitution process would not terminate.

To distinguish names from symbols, we print the names in regular definitions in boldface.

Example 3.4. As we have stated, the set of Pascal identifiers is the set of strings of letters and digits beginning with a letter. Here is a regular definition for this set.

$$\begin{aligned}\text{letter} &\rightarrow \text{A} \mid \text{B} \mid \dots \mid \text{z} \mid \text{a} \mid \text{b} \mid \dots \mid \text{z} \\\text{digit} &\rightarrow 0 \mid 1 \mid \dots \mid 9 \\\text{id} &\rightarrow \text{letter} (\text{letter} \mid \text{digit})^*\end{aligned}$$

Example 3.5. Unsigned numbers in Pascal are strings such as 5280, 39.37, 6.336E4, or 1.894E-4. The following regular definition provides a precise specification for this class of strings:

$$\begin{aligned}
 \text{digit} &\rightarrow 0 \mid 1 \mid \dots \mid 9 \\
 \text{digits} &\rightarrow \text{digit} \text{ digit}^* \\
 \text{optional_fraction} &\rightarrow . \text{ digits} \mid \epsilon \\
 \text{optional_exponent} &\rightarrow (E (+ \mid - \mid \epsilon) \text{ digits}) \mid \epsilon \\
 \text{num} &\rightarrow \text{digits optional_fraction optional_exponent}
 \end{aligned}$$

This definition says that an **optional_fraction** is either a decimal point followed by one or more digits, or it is missing (the empty string). An **optional_exponent**, if it is not missing, is an E followed by an optional + or - sign, followed by one or more digits. Note that at least one digit must follow the period, so **num** does not match 1, but it does match 1.0. \square

Notational Shorthands

Certain constructs occur so frequently in regular expressions that it is convenient to introduce notational shorthands for them.

1. *One or more instances.* The unary postfix operator $^+$ means “one or more instances of.” If r is a regular expression that denotes the language $L(r)$, then $(r)^+$ is a regular expression that denotes the language $(L(r))^+$. Thus, the regular expression a^+ denotes the set of all strings of one or more a ’s. The operator $^+$ has the same precedence and associativity as the operator $*$. The two algebraic identities $r^* = r^+ \mid \epsilon$ and $r^+ = rr^*$ relate the Kleene and positive closure operators.
2. *Zero or one instance.* The unary postfix operator $?$ means “zero or one instance of.” The notation $r?$ is a shorthand for $r \mid \epsilon$. If r is a regular expression, then $(r)?$ is a regular expression that denotes the language $L(r) \cup \{\epsilon\}$. For example, using the $^+$ and $?$ operators, we can rewrite the regular definition for **num** in Example 3.5 as

$$\begin{aligned}
 \text{digit} &\rightarrow 0 \mid 1 \mid \dots \mid 9 \\
 \text{digits} &\rightarrow \text{digit}^+ \\
 \text{optional_fraction} &\rightarrow (. \text{ digits})? \\
 \text{optional_exponent} &\rightarrow (E (+ \mid -)?) \text{ digits } ? \\
 \text{num} &\rightarrow \text{digits optional_fraction optional_exponent}
 \end{aligned}$$

3. *Character classes.* The notation $[abc]$ where a , b , and c are alphabet symbols denotes the regular expression $a \mid b \mid c$. An abbreviated character class such as $[a-z]$ denotes the regular expression $a \mid b \mid \dots \mid z$. Using character classes, we can describe identifiers as being strings generated by the regular expression

$[A-Za-z][A-Za-z0-9]^*$

3.4 RECOGNITION OF TOKENS

In the previous section, we considered the problem of how to specify tokens. In this section, we address the question of how to recognize them. Throughout this section, we use the language generated by the following grammar as a running example.

Example 3.6. Consider the following grammar fragment:

```
stmt → if expr then stmt  
      | if expr then stmt else stmt  
      | ε  
expr → term relop term  
      | term  
term → id  
      | num
```

where the terminals **if**, **then**, **else**, **relop**, **id**, and **num** generate sets of strings given by the following regular definitions:

```
if → if  
then → then  
else → else  
relop → < | <= | = | < > | > | >=  
id → letter ( letter | digit )*  
num → digit+ ( . digit+ )? ( E( + | - )? digit+ )?
```

where **letter** and **digit** are as defined previously.

For this language fragment the lexical analyzer will recognize the keywords if, then, else, as well as the lexemes denoted by **relop**, **id**, and **num**. To simplify matters, we assume keywords are reserved; that is, they cannot be used as identifiers. As in Example 3.5, **num** represents the unsigned integer and real numbers of Pascal.

In addition, we assume lexemes are separated by white space, consisting of nonnull sequences of blanks, tabs, and newlines. Our lexical analyzer will strip out white space. It will do so by comparing a string against the regular definition **ws**, below.

```
delim → blank | tab | newline  
ws → delim+
```

If a match for **ws** is found, the lexical analyzer does not return a token to the parser. Rather, it proceeds to find a token following the white space and returns that to the parser.

REGULAR EXPRESSION	TOKEN	ATTRIBUTE-VALUE
ws	-	-
if	if	-
then	then	-
else	else	-
id	id	pointer to table entry
num	num	pointer to table entry
<	relop	LT
<=	relop	LE
=	relop	EQ
< >	relop	NE
>	relop	GT
>=	relop	GE

Fig. 3.10. Regular-expression patterns for tokens.

Our goal is to construct a lexical analyzer that will isolate the lexeme for the next token in the input buffer and produce as output a pair consisting of the appropriate token and attribute-value, using the translation table given in Fig. 3.10. The attribute-values for the relational operators are given by the symbolic constants LT, LE, EQ, NE, GT, GE. □

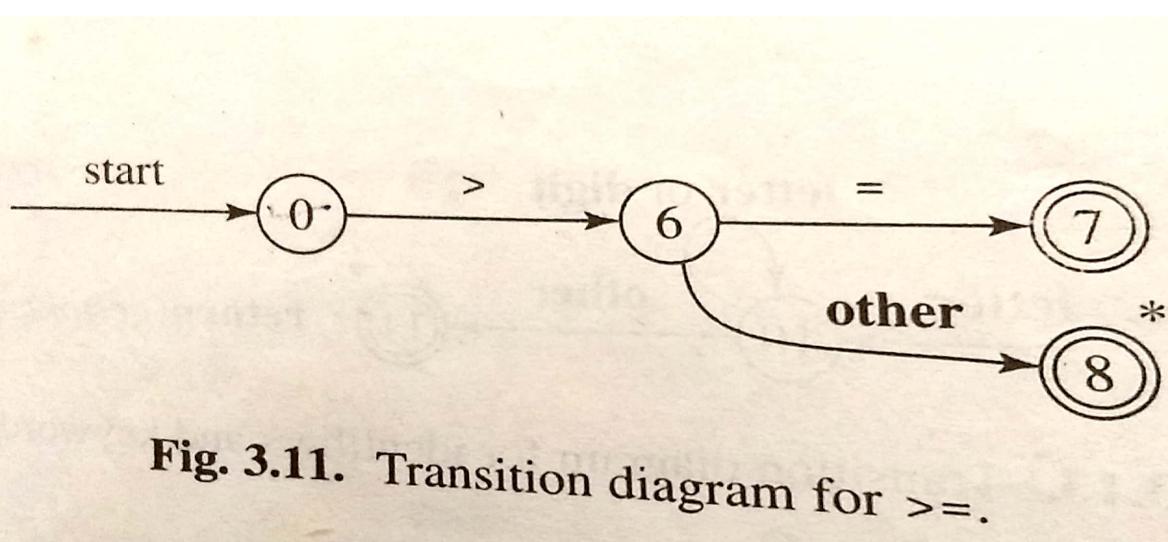


Fig. 3.11. Transition diagram for \geq .

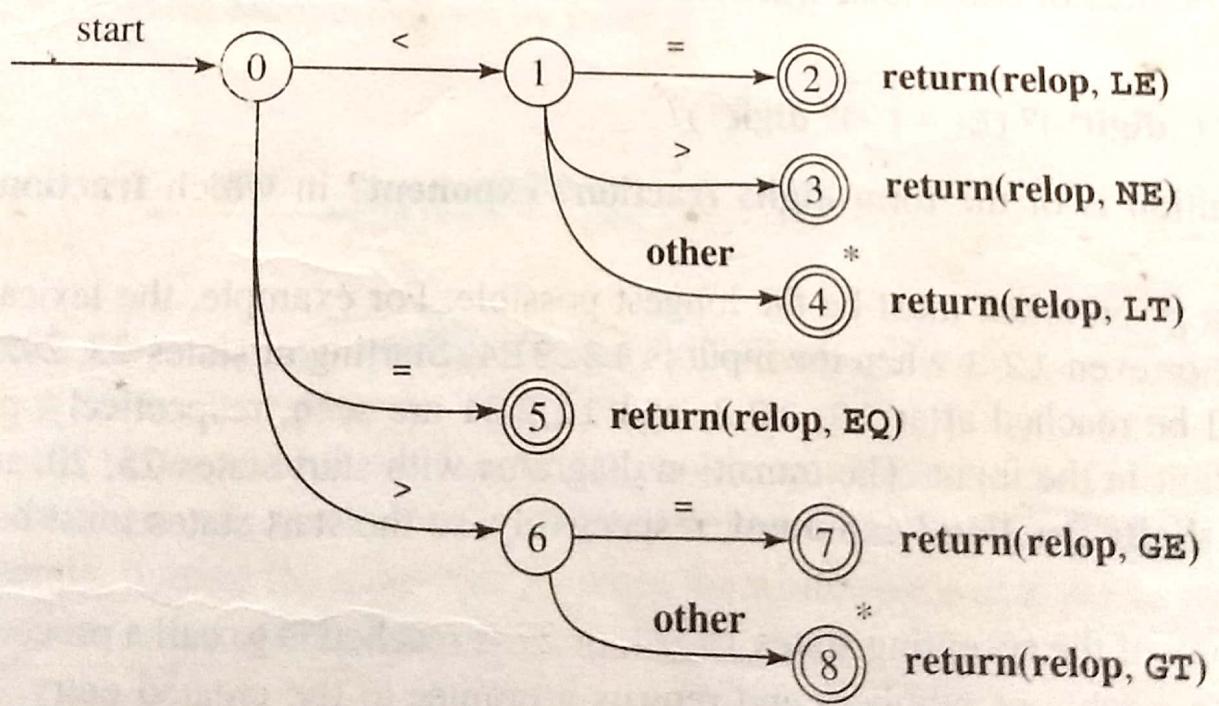


Fig. 3.12. Transition diagram for relational operators.

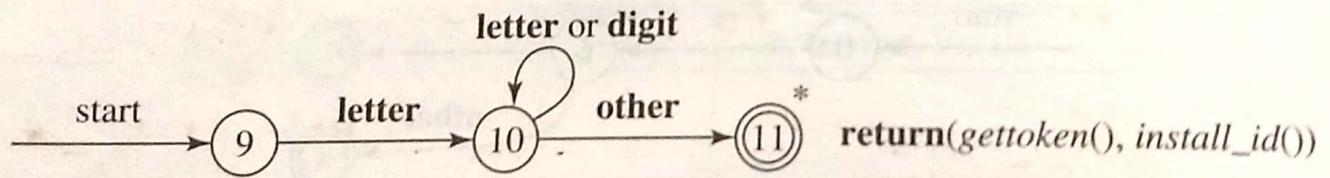


Fig. 3.13. Transition diagram for identifiers and keywords.

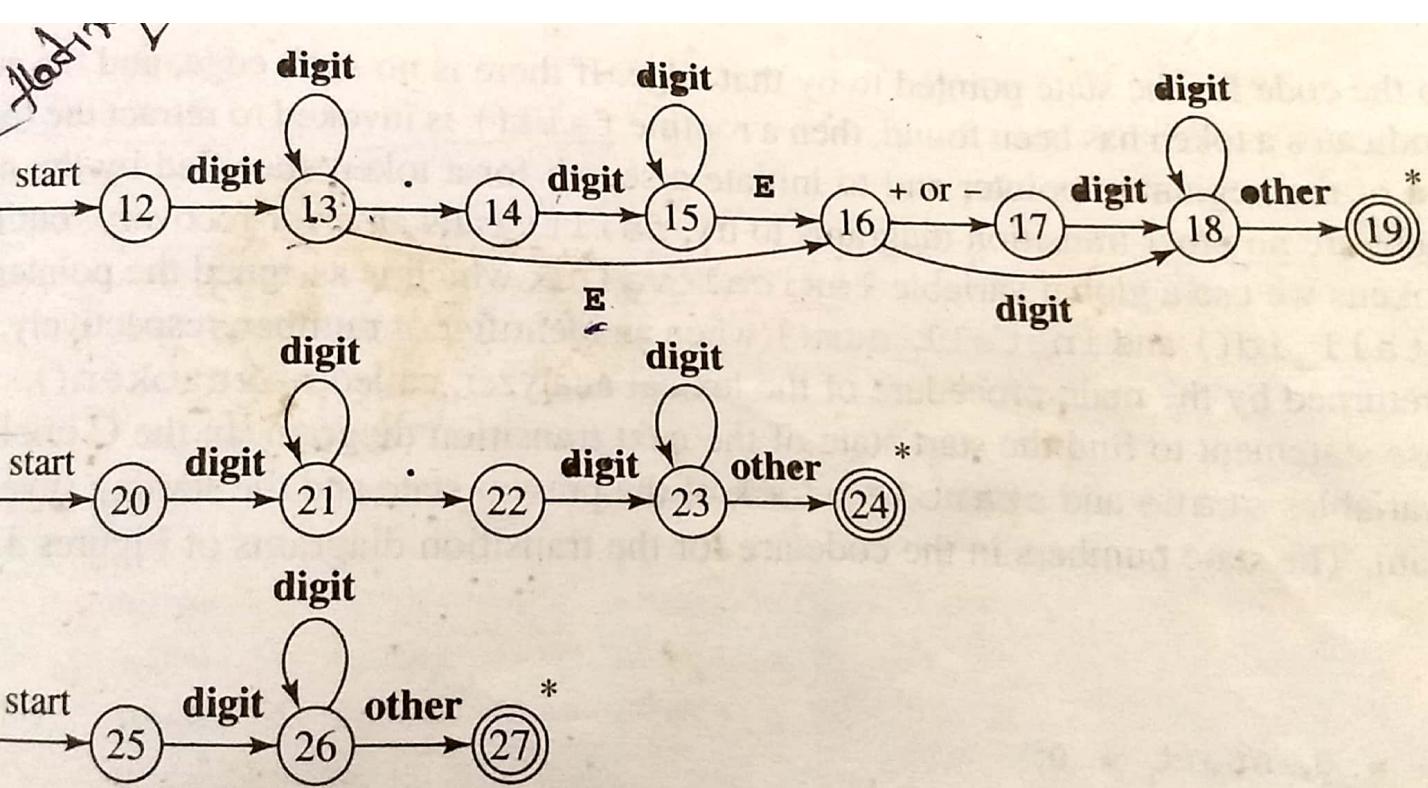
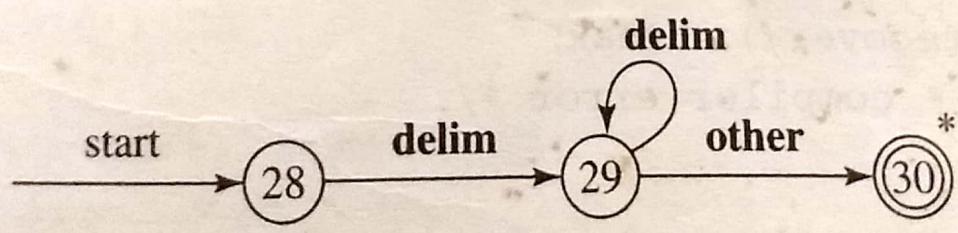


Fig. 3.14. Transition diagrams for unsigned numbers in Pascal.



```
int state = 0, start = 0;
int lexical_value;
    /* to "return" second component of token */
int fail( )
{
    forward = token_beginning;
    switch (start) {
        case 0:    start = 9; break;
        case 9:    start = 12; break;
        case 12:   start = 20; break;
        case 20:   start = 25; break;
        case 25:   recover(); break;
        default:   /* compiler error */
    }
    return start;
}
```

Fig. 3.15. C code to find next start state.

```

token nexttoken()
{
    while(1) {
        switch (state) {
            case 0: c = nextchar();
                /* c is lookahead character */
                if (c==blank || c==tab || c==newline) {
                    state = 0;
                    lexeme_beginning++;
                    /* advance beginning of lexeme */
                }
                else if (c == '<') state = 1;
                else if (c == '=') state = 5;
                else if (c == '>') state = 6;
                else state = fail();
                break;
                ... /* cases 1-8 here */

            case 9: c = nextchar();
                if (isletter(c)) state = 10;
                else state = fail();
                break;
            case 10: c = nextchar();
                if (isletter(c)) state = 10;
                else if (isdigit(c)) state = 10;
                else state = 11;
                break;
            case 11: retract(1); install_id();
                return (gettoken());
                ... /* cases 12-24 here */

            case 25: c = nextchar();
                if (isdigit(c)) state = 26;
                else state = fail();
                break;
            case 26: c = nextchar();
                if (isdigit(c)) state = 26;
                else state = 27;
                break;
            case 27: retract(1); install_num();
                return (NUM);
        }
    }
}

```

Fig. 3.16. C code for lexical analyzer.

3.5 A LANGUAGE FOR SPECIFYING LEXICAL ANALYZERS

Several tools have been built for constructing lexical analyzers from special-purpose notations based on regular expressions. We have already seen the use of regular expressions for specifying token patterns. Before we consider algorithms for compiling regular expressions into pattern-matching programs, we will see an example of a tool that might use such an algorithm.

In this section, we describe a particular tool, called Lex, that has been widely used to specify lexical analyzers for a variety of languages. We refer to the tool as the *Lex compiler*, and to its input specification as the *Lex language*. Discussion of an existing tool will allow us to show how the specification of patterns using regular expressions can be combined with actions, e.g., making entries into a symbol table, that a lexical analyzer may be required to perform. Lex-like specifications can be used even if a Lex compiler is not available; the specifications can be manually transcribed into a working program using the transition diagram techniques of the previous section.

Lex is generally used in the manner depicted in Fig. 3.17. First, a specification of a lexical analyzer is prepared by creating a program `lex.1` in the Lex language. Then, `lex.1` is run through the Lex compiler to produce a C program `lex.yy.c`. The program `lex.yy.c` consists of a tabular representation of a transition diagram.

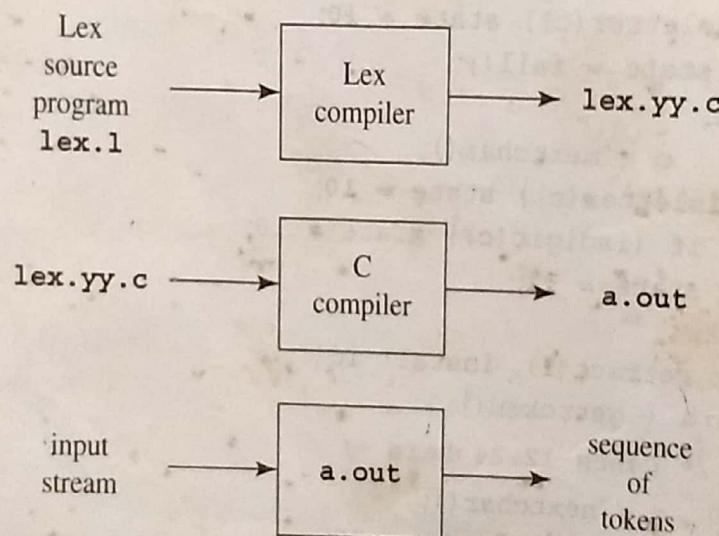


Fig. 3.17. Creating a lexical analyzer with Lex.

constructed from the regular expressions of `lex.1`, together with a standard routine that uses the table to recognize lexemes. The actions associated with regular expressions in `lex.1` are pieces of C code and are carried over directly to `lex.yy.c`. Finally, `lex.yy.c` is run through the C compiler to produce an object program `a.out`, which is the lexical analyzer that transforms an input stream into a sequence of tokens.

Lex Specifications

A Lex program consists of three parts:

declarations
%%
translation rules
%%
auxiliary procedures

The declarations section includes declarations of variables, manifest constants, and regular definitions. (A manifest constant is an identifier that is declared to represent a constant.) The regular definitions are statements similar to those given in Section 3.3 and are used as components of the regular expressions appearing in the translation rules.

The translation rules of a Lex program are statements of the form

$$\begin{array}{ll} p_1 & \{ \text{action}_1 \} \\ p_2 & \{ \text{action}_2 \} \\ \dots & \dots \\ p_n & \{ \text{action}_n \} \end{array}$$

where each p_i is a regular expression and each action_i is a program fragment describing what action the lexical analyzer should take when pattern p_i matches a lexeme. In Lex, the actions are written in C; in general, however, they can be in any implementation language.

The third section holds whatever auxiliary procedures are needed by the actions. Alternatively, these procedures can be compiled separately and loaded with the lexical analyzer.

A lexical analyzer created by Lex behaves in concert with a parser in the following manner. When activated by the parser, the lexical analyzer begins reading its remaining input, one character at a time, until it has found the longest prefix of the input that is matched by one of the regular expressions p_i . Then, it executes action_i . Typically, action_i will return control to the parser. However, if it does not, then the lexical analyzer proceeds to find more lexemes, until an action causes control to return to the parser. The repeated search for lexemes until an explicit return allows the lexical analyzer to process white space and comments conveniently.

The lexical analyzer returns a single quantity, the token, to the parser. To pass an attribute value with information about the lexeme, we can set a global variable called `yyval`.

Example 3.11. Figure 3.18 is a Lex program that recognizes the tokens of Fig. 3.10 and returns the token found. A few observations about the code will introduce us to many of the important features of Lex.

In the declarations section, we see (a place for) the declaration of certain manifest constants used by the translation rules.⁴ These declarations are surrounded by the special brackets `%{` and `%}`. Anything appearing between these brackets is copied directly into the lexical analyzer `lex.yy.c`, and is not treated as part of the regular definitions or the translation rules. Exactly the same treatment is accorded the auxiliary procedures in the third section. In Fig. 3.18, there are two procedures, `install_id` and `install_num`, that are used by the translation rules; these procedures will be copied into `lex.yy.c` verbatim.

Also included in the definitions section are some regular definitions. Each such definition consists of a name and a regular expression denoted by that name. For example, the first name defined is `delim`; it stands for the character class `[\t\n]`, that is, any of the three symbols blank, tab (represented by `\t`), or newline (represented by `\n`). The second definition is of white space, denoted by the name `ws`. White space is any sequence of one or more delimiter characters. Notice that the word `delim` must be surrounded by braces in Lex, to distinguish it from the pattern consisting of the five letters `delim`.

In the definition of `letter`, we see the use of a character class. The shorthand `[A-Za-z]` means all of the capital letters A through Z or the lowercase letters a through z. The fifth definition, of `id`, uses parentheses, which are metasymbols in Lex, with their natural meaning as groupers. Similarly, the vertical bar is a Lex metasymbol representing union.

```

%{

    /* definitions of manifest constants
    LT, LE, EQ, NE, GT, GE,
    IF, THEN, ELSE, ID, NUMBER, RELOP */
}

/* regular definitions */
delim      [ \t\n]
ws         {delim}+
letter     [A-Za-z]
digit      [0-9]
id         {letter}({letter} | {digit})*
number    {digit}+(\.{digit}+)?(E[+\-]?{digit}+)?

%%

{ws}        {/* no action and no return */}
if          {return(IF);}
then        {return(THEN);}
else        {return(ELSE);}
{id}        {yyval = install_id(); return(ID);}
{number}    {yyval = install_num(); return(NUMBER);}
"<"        {yyval = LT; return(RELOP);}
"<="       {yyval = LE; return(RELOP);}
"="         {yyval = EQ; return(RELOP);}
"<>"      {yyval = NE; return(RELOP);}
">"        {yyval = GT; return(RELOP);}
">="       {yyval = GE; return(RELOP);}

%%

install_id() {
    /* procedure to install the lexeme, whose
    first character is pointed to by yytext and
    whose length is yylen, into the symbol table
    and return a pointer thereto */
}

install_num() {
    /* similar procedure to install a lexeme that
    is a number */
}

```

Fig. 3.18. Lex program for the tokens of Fig. 3.10.