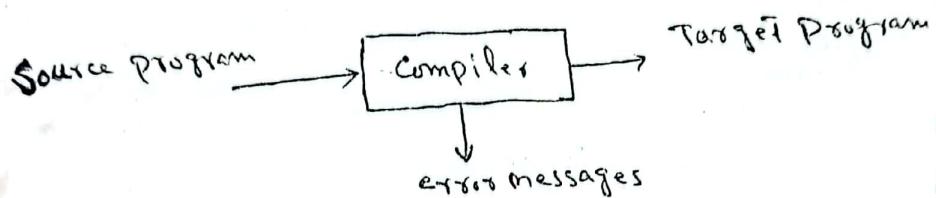


COMPILER THEORY

(1)

Compiler:- A compiler is a program that reads a program written in one language (i.e source program) and translates it into an equivalent program in another language (i.e the target language). Besides this, the compiler reports to its user the presence of errors in the source program.



There are thousands of source program languages, ranging from traditional programming languages such as Fortran and Pascal to specialized languages that can be used in every area of computer application.

The target languages are equally varied; a target language may be another programming language or the machine language of any computer.

Compilers are sometimes classified as single-pass, multipass, debugging or optimizing, depending on how they have been constructed or on what function they are supposed to perform. Despite this apparent complexity, the basic functions that any compiler must perform are essentially the same.

Mr.

The Analysis - Synthesis Model of Compilation

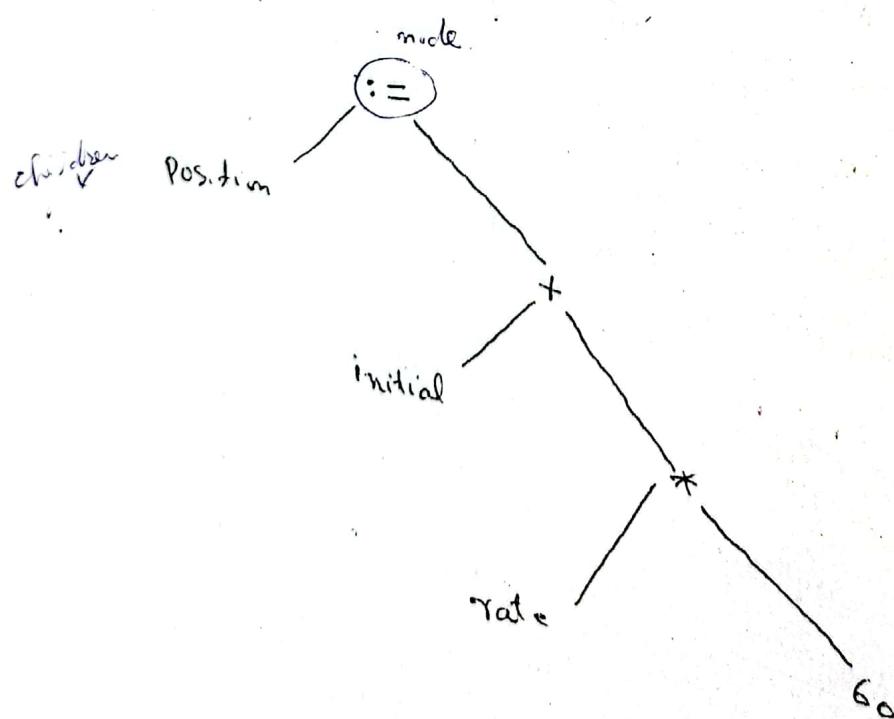
There are two parts of compilation: analysis and synthesis.

The analysis part breaks up the source program into pieces and creates an intermediate representation of the source program.

The synthesis part constructs the desired target program from the intermediate representation.

During analysis, the operations implied by the source program are determined and recorded in a hierarchical form called a tree. Often a special kind of tree called a syntax tree is used, in which each node represents an operation and the children of a node represent the arguments of the operation.

For example, a syntax tree for an assignment statement is as follows:



many software tools that manipulate source programs, first perform some kind of analysis.
Some examples of such tools include:

1. Structure Editors : A structure editor takes as input a sequence of commands to build a source program. The structure editor not only performs the text-creation and modification functions of an ordinary text editor, but it also analyzes the program text, putting an appropriate hierarchical structure on the source program. Thus, the structure editor can perform additional tasks that are useful in the preparation of programs.

for example, it can check that the input is correctly formed, can supply keywords automatically (e.g. when the user types while, the editor supplies the matching do and reminds the user that a conditional must come between them), and can jump from a begin or left parenthesis to its matching end or right parenthesis.

further, the output of such an editor is often similar to the output of the analysis phase of a Compiler.

2. Pretty Printers : A pretty printer analyzes a program and prints it in such a way that the structure of the program becomes clearly visible. for example, comments may appear in a special font, and statements may appear with an amount of indentation proportional to the depth of their nesting in the hierarchical organization of the statements.

3- static checkers :- A static checker reads a program, analyzes it, and attempts to discover bugs without running the program. For example, a static checker may detect the parts of the source program can never be executed, or that a certain variable might be used before being defined.

4. interpreters :- Instead of producing a large program as a translation, an interpreter performs the operations implied by the source program. For an assignment statement, for example, an interpreter might build a tree and then carry out the operations at the nodes as it walks the tree.

Consider the tree on page 2.
At the root it would discover it (interpreter) had an assignment to perform, so it would call a routine to evaluate the expression on the right, and then store the resulting value in the location associated with the identifier position. At the right side of the root the routine would discover it had to compute the sum of two expressions. It would call itself recursively to compute the value of the expression $rate * 60$. It would then add that value to the value of the variable $initial$.

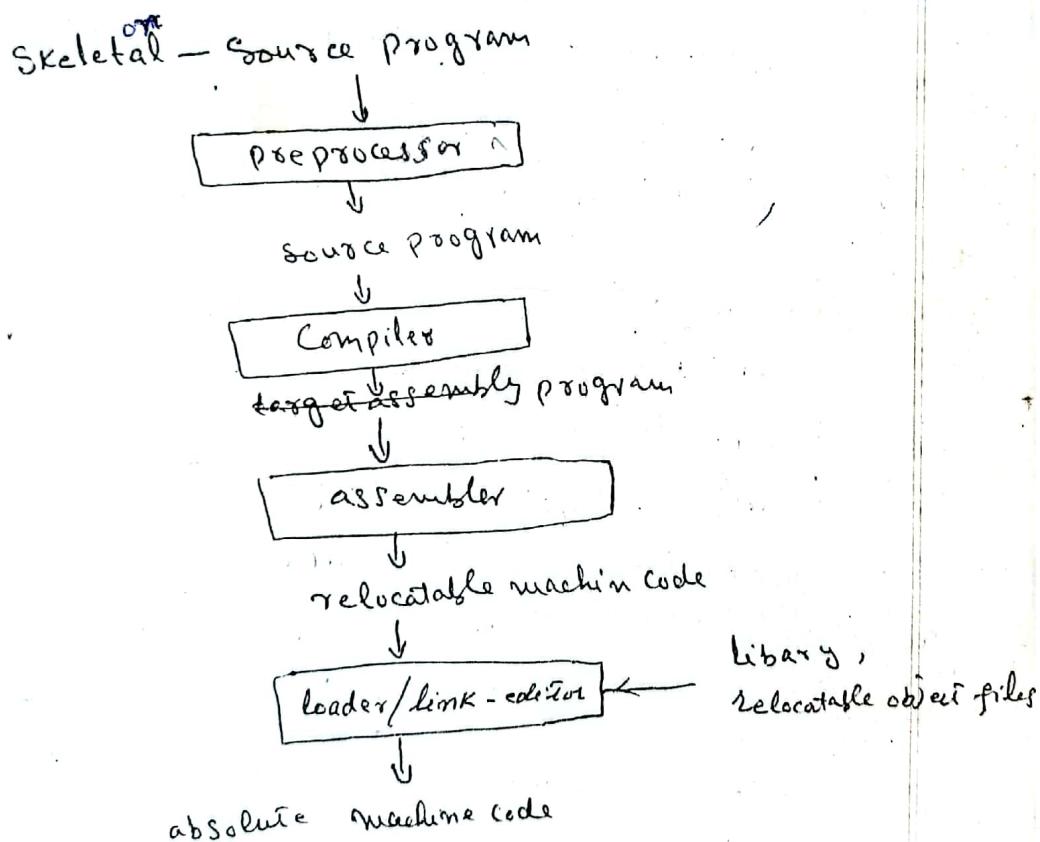
Interpreters are frequently used to execute command languages. Since each operator executed in a command language is usually an invocation of a complex routine such as a compiler.

The Content of a Compiler:

(3)

In addition to a compiler, several other programs may be required to create an executable target program.

A source program may be divided in modules stored in separate files. The task of collecting the source program is sometimes given to a distinct program, called a preprocessor. The preprocessor may also expand shorthands called macros, into source language statements. The following Fig shows a typical "Compilation".



The target program created by the compiler may require further processing before it can be run. The compiler creates assembly code that is translated by an assembler into machine code and then linked together with some library routines into the code that actually runs on the machine.

Analysis of the Source Program:

In Compiling, analysis consists of three phases

- 1- Linear analysis, in which the stream of characters make the source program is read from left to right and given into tokens that are sequences of characters have a collective meaning.
2. Hierarchical analysis, in which characters or tokens are grouped hierarchically into nested collections with collective meaning.
- 3- Semantic analysis; in which certain checks are performed to ensure that the components of a program fit together meaningfully.

Lexical Analysis:

In a compiler, linear analysis is called lexical analysis or scanning. For example, in lexical analysis the characters in the assignment statement

position := initial + rate * 60

would be grouped into the following tokens:

- The identifier position
- The assignment symbol :=
- The identifier initial
- The plus sign
- The identifier rate
- The multiplication sign
- The number 60

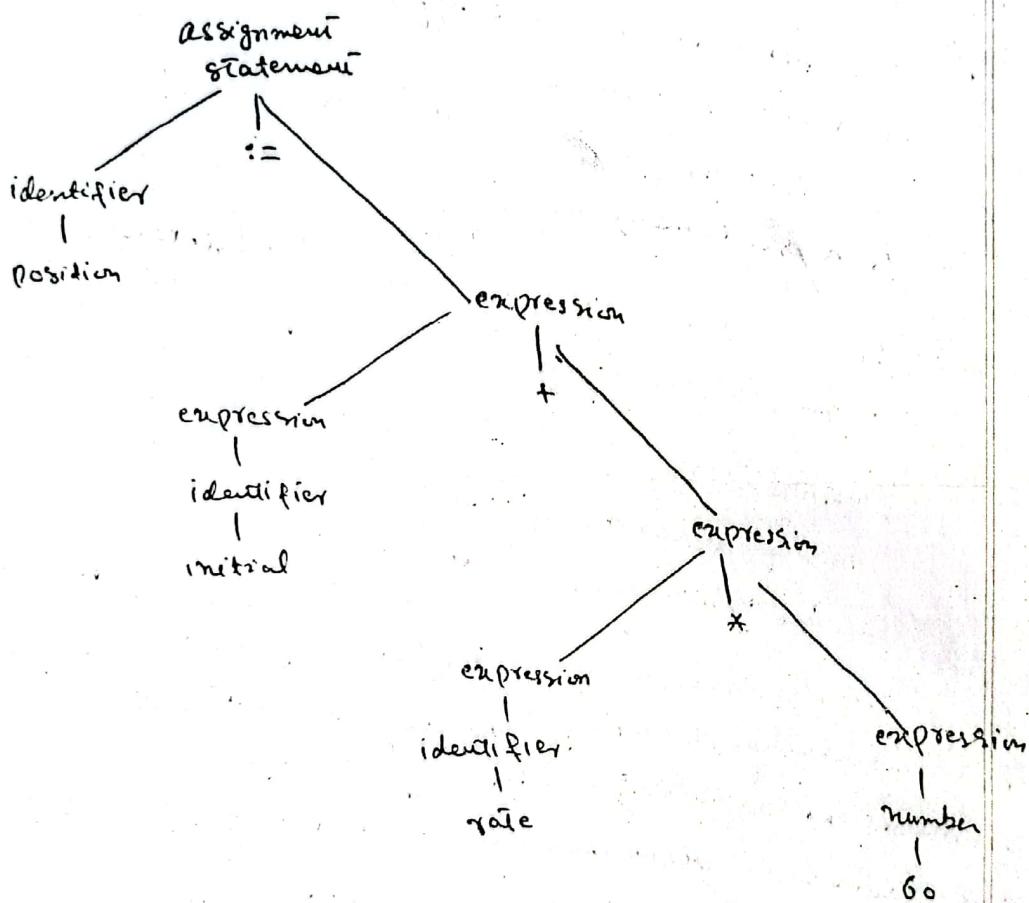
The blanks separating the characters of these tokens would normally be eliminated during lexical analysis

Syntax Analysis

(4)

A hierarchical analysis is called parsing or syntax analysis. It involves grouping the tokens of the source program into grammatical phrases that are used by the compiler to synthesize the output. Usually, the grammatical phrase of the source program are represented by a parse tree such as in following fig, the parse tree of

position := initial + rate * 60 is,



The hierarchical structure of a program is usually expressed by recursive rules. For example, we might have the following rules as part of the definition of expressions:

1 - Any identifier is an expression.

2 - Any number is an expression.

3 - If expression₁ and expression₂ are expressions,

Then so are

expression₁, + expression₂

expression₁, * expression₂

(expression₁)

Rules (1) and (2) are (non-recursive) basis rules, while (3) defines expressions in terms of operators applied to other expressions. Thus by rule (1), initial and rate are EXPRESSIONS. By rule (2), 60 is an expression, while by rule (3), we can first infer that $rate * 60$ is an expression and finally that $initial + rate * 60$ is an expression.

Similarly, many languages define statements

recursively by rules such as:

1 - If identifier₁ is an identifier, and expression₁ is an expression, then

identifier₁ := expression₂
is a statement

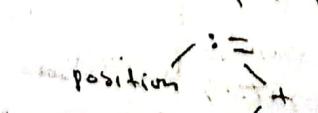
2 - If expression₁ is an expression and statement₁ is a statement, then:

while (expression₁) do statement₂
if (expression₁) Then statement₂
are statements.

The division between lexical and syntactic analysis is somewhat arbitrary - we usually choose a division that simplifies the overall task of analysis. One factor in determining the division is whether source language construct is inherently recursive or not. Lexical constructs do not require recursion while syntactic constructs often do.

A Syntax tree is a compressed representation of the parse tree in which the operators appear as interior nodes, and the operands of an operator are children of the node for that operator.

e.g.



Semantic Analysis: (gt find logical Errors)

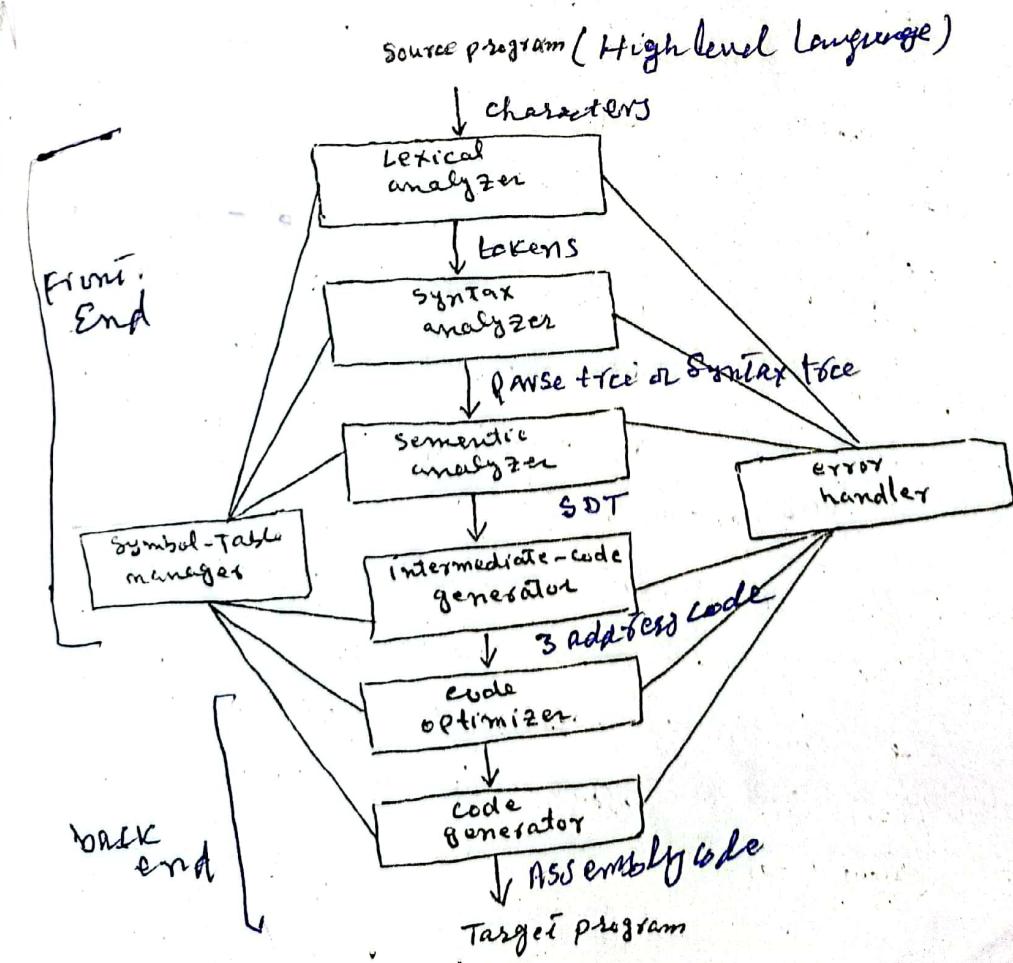
(5)

The semantic analysis phase checks the source program for semantic errors and gathers type information for the subsequent code-generation phase. It uses the hierarchical structure determined by the syntax analysis phase to identify the operators and operands of expressions and statements.

An important component of semantic analysis is type checking. Hence the compiler checks that each operator has ^{variables identified, number} operands that are permitted by the source language specification. For example, many programming language definitions require a compiler to report an error every time a real number is used to index an array.

93/ The Phases of a Compiler :

A compiler operates in phases, each of which transforms the source program from one representation to another. A typical decomposition of a compiler is:



The first three phases, forming the bulk of the analysis part, two other activities, symbol-table management and error handling, are shown interacting with the six phases of lexical analysis, syntax analysis, semantic analysis, intermediate-code generation, code optimization, and code generation.

Symbol-Table Management:-

An essential function of a compiler is to record the identifiers used in the source program and all information about the various attributes of each identifier. These attributes may provide information about the storage allocated for an identifier, its type, its scope (where in the program it is valid), and in the case of procedure names, such things as the number and types of its arguments.

A symbol table is a data structure containing a record for each identifier, with fields for the attributes. The identifier. The data structure allows us to find the record for each identifier quickly and to store or retrieve data from that record quickly.

When an identifier in the source program is detected by the lexical analyzer, the identifier is entered onto the symbol table. However, the attributes of an identifier cannot normally be determined during lexical analysis. For example, in a Pascal declaration like

Var position, initial, rate : real;

The type real is not known when position, initial, "rate" are seen by the lexical analyzer.

The remaining phases enter information about identifiers into the symbol table and then use this information in various ways.

Error Detection and Reporting:

(6) ⑥

Each phase can encounter errors. However, after detecting an error, a phase must somehow deal with that error, so that compilation can proceed, allowing further errors in the source program to be detected.

The syntax and semantic analysis phases usually handle

a large fraction of the errors detectable by the compiler. The lexical phase can detect errors where the characters remaining in the input do not form any tokens of the language. Errors where the token stream violates the structure rule (syntax) of the language are determined by the syntax analysis phase. During semantic analysis the compiler tries to detect constructs that have the right syntactic structure but no meaning to the operations involved, e.g., if we try to add two ~~structure~~ identifiers, one of which is the name of an array, and the other the name of a procedure.

The Analysis Phases:

As translation progresses, the compiler's internal representation of the source program changes. We show these representations by considering the translation of the statement.

$$\text{position} := \text{initial} + \text{rate} * 60 \quad \text{--- (1)}$$

The following fig shows the representation of this statement after each phase.

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



Lexical analyzer



$$id_1 := id_2 + id_3 * 60$$

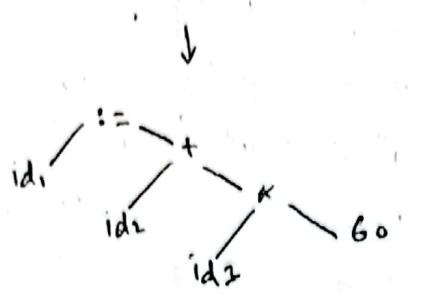
~~id~~



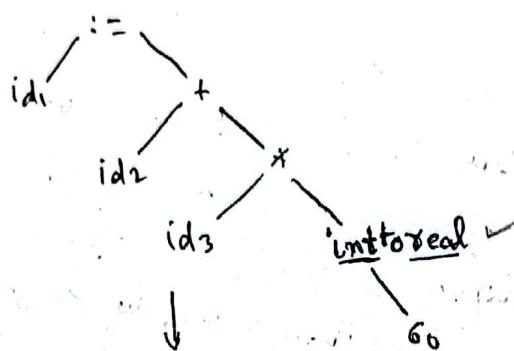
Syntax analyzer

symbol table

| | | |
|---|----------|-----|
| 1 | Position | --- |
| 2 | initial | --- |
| 3 | rate | --- |
| 4 | | |



semantic analyzer



intermediate code generator

temp1 := inttoreal (60)

temp2 := id3 + temp1

temp3 := id2 + temp2

id1 := temp3

code optimizer

temp1 := id3 * 60.0

id1 := id2 + temp1



code generator

MOVF R2, id3
~~MULF #60.0~~
MULF R2, #60.0
MOVF R1, id2
ADDF R1, R2
MOVF id1, R1
MOVF id1, R1

MOVF id3, R2
MULF #60.0, R2
MOVF id2, R1
ADDF R2, R1
MOVF R1, id1

m.s.
 lariy
 m.g.
 lariy

The lexical analysis phase reads the characters in the source program and groups them into a stream of tokens in which each token represents a logically sequence of characters, such as an identifier, a keyword (if, while etc), a punctuation character, a multi-character operator like `==`. The character sequence forming a token is called the lexeme for the token.

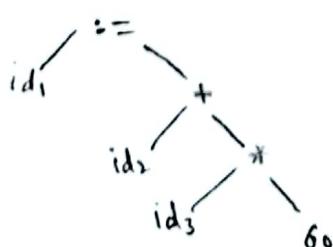
Certain tokens will be augmented by a "lexical value". For example, when an identifier like rate is found, the lexical analyzer not only generates a token, say id, but also enters the lexeme rate into the symbol table. The lexical value associated with occurrence of id points to the symbol-table entry for rate.

We shall use id_1 , id_2 and id_3 for position, initial, and rate respectively, to emphasize that the internal representation of an identifier is different from the character sequence forming the identifier. The representation of (1) after lexical analysis is therefore:

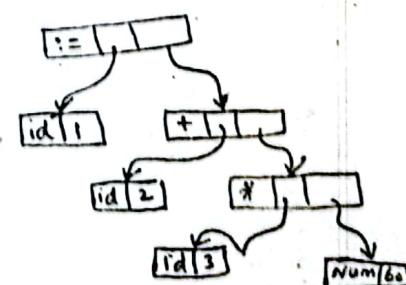
$$id_1 == id_2 + id_3 * 60 \quad - \quad (1)$$

Syntax analysis imposes a hierarchical structure on the token stream, which we shall show by syntax tree in fig (a). A typical data structure for the tree is shown in fig (b).

In which an interior node is a record with a field for the operator and two fields containing pointers to the records for the left and right children. A leaf is a record with two or more fields, one to identify the token at the leaf, and the other to record information about the token.



(a)



(b)

The data structure in (b) for the syntax tree in (a)

Intermediate Code Generation:

After syntax and semantic analysis, some compilers generate an explicit intermediate representation of the source program. This intermediate representation should have two important properties; it should be easy to produce, and easy to translate into the target program.

We consider an intermediate form called "three-address code", which is like the assembly language for a machine in which every memory location can act like a register. Three-address code consists of a sequence of instructions, each of which has at most three operands. The source program in (1) in three-address code is:

$$\begin{aligned} \text{temp1} &= \text{inttoreal(6)} \\ \text{temp2} &= \text{id3} * \text{temp1} \\ \text{temp3} &= \text{id2} + \text{temp2} \\ \text{id1} &= \text{temp3} \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \quad \text{--- (3)} \quad \text{my copy}$$

This intermediate form has several properties:-

- i) Each three-address instruction has at most one operator in addition to the assignment. Thus when generating these instructions, the compiler has to decide the order in which operations are to be done; the multiplication precedes the addition in source program (1).
- ii) The compiler must generate a temporary name to hold the value computed by each instruction.

- iii) Some "three address" instructions have fewer than three operands, e.g., the first and last instructions in (3)
- id2, temp2 etc.

Code Optimization:

The code optimization phase attempts to improve the intermediate code, so that faster-running machine code will result. Some optimizations are trivial. For example, a natural algorithm generates the intermediate code (3), using an instruction for each operator in the tree representation after semantic analysis - even though there is a better way to perform the same calculation, using the two instructions ^{floating}

$$\left. \begin{array}{l} \text{temp1 := id3 * 60.0} \\ \text{id1 := id2 + temp1} \end{array} \right\} \quad (4)$$

Code Generation:

The final phase of the compiler is the generation of target code, consisting normally of relocatable m/c code or assembly code. Memory locations are selected for each of the variables used by the program. A crucial aspect is the assignment of variables to registers. For example, using registers 1 and 2, the translation of the code in (4) will become

^{relocatable code} ^{MovF R2,id3} ^{MovF R2,R2} ^{805D6}
^{Code} MOVF id3 , R2 } (5)
^{MULF #60.0, R2}
^{MovF id2 , R1}
^{ADD F R2 , R1}
^{MovF R1 , id1}

The first and the 2nd operands of each instruction specify a source and destination, respectively. The F in each instruction tells us that instructions deal with floating-point numbers.

Handwritten

COUSINS OF THE COMPILER :-

Preprocessors :- preprocessors produce input to compiler
They may perform the following functions :-

- 1 - macro processing: A preprocessor may allow a user to define macros that are shorthands for longer constructs.
- 2 - File inclusion: A preprocessor may include header files into the program text.

3 - Rational preprocessors :- These processors attempt to bring older languages with more modern flow-of-control and data-structuring facilities. For example, such a preprocessor might provide the user with built-in macros for constructs like while-loops or if-statements, where none exist in the programming language itself.

4 - language extensions: These processors attempt to add capabilities to the language by what amounts to built-in macros. For example, the language SQL is a database query language embedded in C. Statements beginning with # are taken by the preprocessor to be database statements, unrelated to C, and are translated into procedure calls on routines that perform the database access.

Assemblers :-

Some compilers produce assembly code as intermediate output that is passed to an assembler for further processing. Other compilers perform the job of the assembler, producing relocatable m/c code that can be passed directly to the loader/link-editor.

Assembly code is a mnemonic version of m/c code ,
 in which names are used instead of binary codes for operations
 and names are also given to memory addresses.
 A typical sequence of assembly instructions might be

| | | | |
|-----------------------|------------------------|---|-----|
| R ₁ , a | MOV a, R ₁ | } | (6) |
| D R ₁ , #2 | ADD #2, R ₁ | | |
| b, R ₁ | MOV R ₁ , b | | |

This code moves the contents of the address a into register R₁. Then adds the constant 2 to it and finally stores the result in the location named b . Thus, it computes
 $b := a + 2$ ✓

TWO-PASS ASSEMBLY

The simplest form of assembler makes two passes over the input, where a pass consists of reading an input file once. In the 1st phase pass, all the identifiers that denote storage locations are found and stored in a symbol table . Identifiers are assigned storage locations as they are encountered for the first time , so after reading (6) , for example, the symbol table might contain the entries as follows

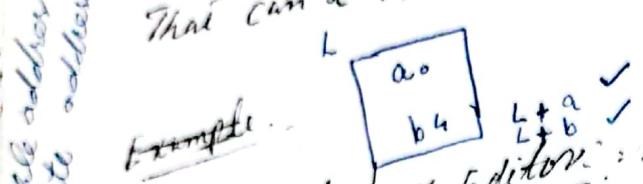
| IDENTIFIER ✓ | ADDRESS ✓ |
|--------------|-----------|
| a | 0 ✓ |
| b | 4 ✓ 0.23 |

Here , we have assumed that a word , consisting of four bytes , is set for each identifier , and their addresses are assigned starting from byte 0 .

In the second pass , the assembler scans the input again . This time , it translates each operation code into the sequence of bits representing that operation in m/c language , and it translates each identifier representing a location into the address given for that identifier , in the symbol table .

The output of the 2nd pass is usually relocatable m/c code , meaning that it can be loaded starting at

any location L in memory i.e., if L is moved to all addresses in the code, then all references will be correct. Thus, the output of the assembler must define those positions of instructions that refer to addresses that can be relocated.



Loaders and Link-Editor :-

usually a program called a loader performs a two functions of loading and link editing. The process of link editing consists of taking relocatable M/C code, altering the relocatable addresses and placing the altered instructions and data in memory at the proper locations.

The Link-editor allows us to make a single program from several files of relocatable m/c code. These files may have been the result of several different compilations and one or more may be library files of routines provided by the system and available to any program that needs them.

THE GROUPING OF PHASES:-

Implementation activities from more than one phase are often grouped together.

Front and Back Ends:-

Often, the phases are collected into a front end and a back end. The front end consists of those phases dependent primarily on the source program (or source language) and largely independent of the target m/c. These normally include lexical and syntactic analysis, the creation of symbolic table, semantic analysis, and the generation of intermediate code.

(10)

A certain amount of code optimization can be done by front end as well. The front end also includes the error handling that goes along with each of these phases.

The back end includes those portions of compiler that depend on the target rule, and generally these portions do not depend on the source language. In the back end, we find aspects of the code optimization phase, and we find code generation, along with the necessary error handling and symbol-table operations.

PASSES:- Several phases of compilation are usually implemented in a single pass consisting of reading an input file and writing an output file.
It is common for several phases to be grouped into one pass, and for the activity of these phases to be interleaved during the pass. For example, lexical analysis, syntax analysis, Semantic analysis and intermediate code generation might be grouped into one pass. If so, the token stream after lexical analysis may be translated directly into intermediate code.

In an implementation of a compiler, portions of one or more phases are combined into a module called a PASS. A pass reads the source program or output of the previous pass, makes the transformations specified by its phases, and writes output into an intermediate file, which may then be read by a subsequent pass. If several phases are grouped into one pass, then the operation of the phases may be interleaved, with control alternating among several phases.

A multi-pass compiler can be made to use less space than a single-pass compiler, since the space occupied by the compiler program for one pass can be used by the following pass. A multi-pass compiler is slower than a single-pass compiler, because each pass reads and writes an intermediate file. Thus, running on computers with small memory would normally use several passes while, on a computer with a large random memory, a compiler with fewer passes would be possible.

Compiler - Construction Tools:

A number of tools have been developed specifically to help construct compilers. These range from scanners and parser generators to complete systems, variously called compile-compilers, compile-generators or translator-writing systems, which produce a compiler from some form of specification of a source language and target language. The input specification for these systems may contain:

- a description of the lexical and syntactic structures of the source language.
- a description of what output is to be generated for each source language construct.
- a description of the target machine language.

The following is the list of some useful compiler construction tools or aids provided by existing compilers.

Compilers are:

1- parser generator: These produce syntax analyzers, normally from input that is based on a context-free grammar.

2- Scanner generator: These automatically generate lexical analyzers, normally from a specification based on regular expressions.

3- Syntax-directed translation engines:

These produce collections of routines that walk the parse tree, generating intermediate code.

4- Automatic code generators:

Such a tool takes a collection of rules that define the translation of each operation of the intermediate language into the mic language for the target machine. The rules must include sufficient detail that we can handle the different possible access methods for data: e.g., variables may be in registers or a fixed (scatter) location in memory.

5- Data-flow engines:

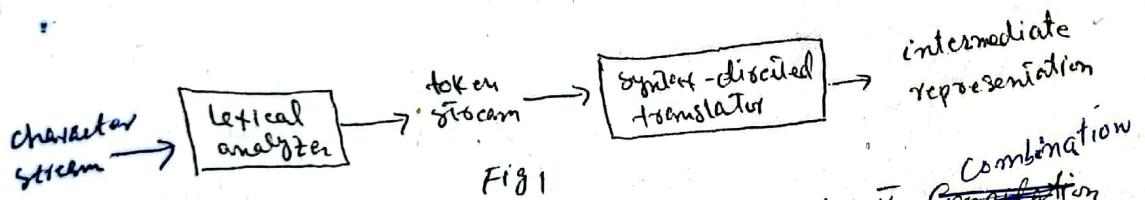
Much of the information needed to perform good code optimization involves "data-flow analysis". The gathering of information about how values are transmitted from one part of the program to each other part.



Overview :- A programming language can be defined by describing what its programs looks like (the syntax of the language) and what its programs mean (the semantics of the language). For specifying the syntax of a language, we use a notation called Context-free grammars or BNF (Backus-Naur Form).

A grammar-oriented compiling technique, known as syntax-directed translation, is very helpful for organizing compiler front end (before machine code).

In Compiler, the lexical analyzer converts the stream of input characters into a stream of tokens that becomes the input to the following phase as in fig.



The "syntax-directed translator" on the fig is the combination of a syntax analyzer and an intermediate code generator.

Syntax Definition: we introduce a notation, called a context-free grammar (or grammar), for specifying the syntax of a language.

A grammar naturally describes the hierarchical structure of many programming language constructs. For example, an if-else statement in C has the form

if (expression) statement else statement

That is, the statement is the sequence of the keyword if, an opening parenthesis, an expression, a closing parenthesis, a statement, the keyword else, and another statement.

Using the variable expr to denote an expression and the variable stmt to denote a statement, this sequencing rule can be expressed as

$\text{if } \text{expr} \rightarrow \text{stmt}$

in which the arrow may be read as "can have the form". Such a rule is called a production. In a production lexical elements like the keyword if and the parenthesis are called tokens. Variables like expr and stmt represent sequences of tokens and called nonterminals.

A context-free grammar has four components:

- 1- A set of tokens, known as terminal symbols.
- 2- A set of nonterminals
- 3- A set of productions where each production consists of a nonterminal, called the left side, of the production, an arrow, and a sequence of tokens and/or nonterminals, called the right side of the production.
- 4- A designation of one of the nonterminals as the start symbol

We follow the convention of specifying grammars by their production, with the productions for the start symbol listed first. We assume that digits, signs such as <, and underline strings such as while are called terminals.

Example 1 :- We use expressions consisting of digits and plus and minus signs, e.g., $9 - 5 + 2 \cdot 3 - 1$ and γ . Since a plus or minus sign must appear between two digits, we refer to such expressions as "lists of digits separated by plus or minus signs". The following grammar describes the syntax of these expressions. The productions are:

$$\begin{array}{l} \text{terminal} \\ \rightarrow \text{list} \rightarrow \text{list} + \text{digit} \quad \text{terminal} \quad (1) \\ \rightarrow \text{list} \rightarrow \text{list} - \text{digit} \quad \text{terminal} \quad (2) \\ \rightarrow \text{list} \rightarrow \text{digit} \quad \text{terminal} \quad (3) \\ \text{digit} \rightarrow , 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 \quad (4) \end{array}$$

The right sides of the productions with nonterminal list on the left side can equivalently be grouped:

$$\text{list} \rightarrow \text{list} + \text{digit} \mid \text{list} - \text{digit} \mid \text{digit}$$

* The tokens of the grammar are the symbols

$$+ - 0 1 2 3 4 5 6 7 8 9 \quad \text{special character}$$

The nonterminals are the list and digit, with list being the starting nonterminal because its productions are given first.

We say a production is for a nonterminal if the nonterminal appears on the left side of the production. A string of tokens is

a sequence of zero or more tokens. The string containing a zero tokens, written as ϵ , is called the empty string.

A grammar derives strings by beginning with the

start symbol and repeatedly replacing a nonterminal

by the right side of a production for that nonterminal.

The token strings that can be derived from the start

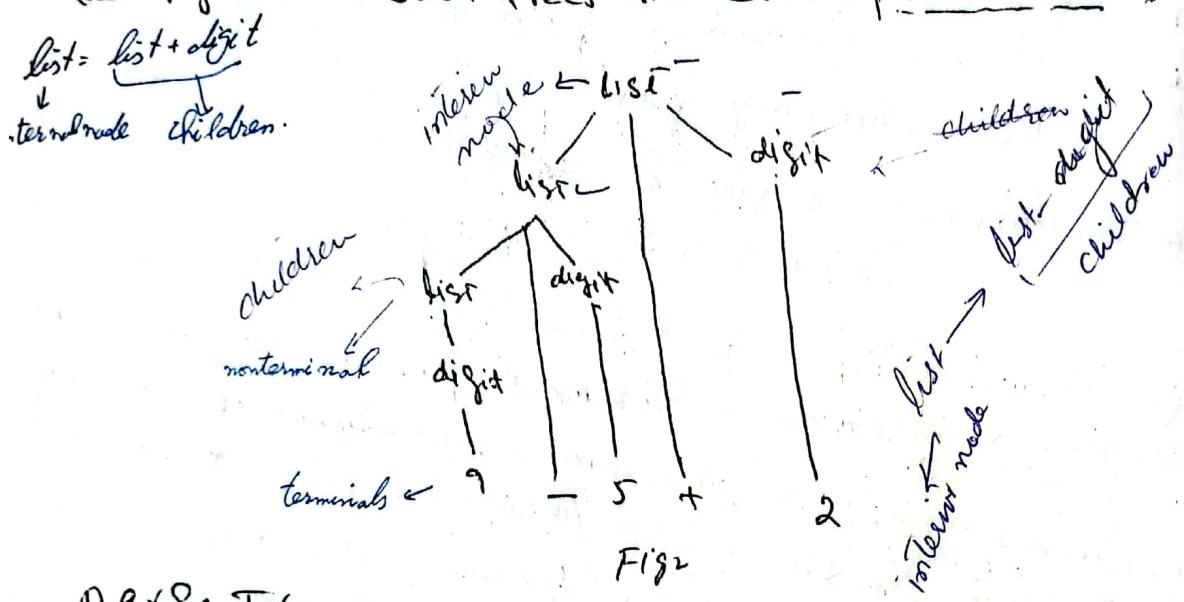
symbol form the language defined by the grammar.

Example 2 :- The language defined by the grammar of example 1, consists of lists of digits separated by plus and minus signs.

We can deduce that $9 - 5 + 2$ is a list as follows

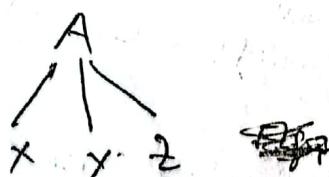
- a) q is a list by production(3), since q is a digit.
- b) $q-5$ is a list by production(2), since q is a list and 5 is a digit.
- c) $q-5+2$ is a list by production(1), since $q-5$ is a list and 2 is a digit.

This reasoning is shown by the tree. Each node in the tree is labeled by a grammar symbol. An interior node and its children correspond to a production; The interior node corresponds to the left side of the production, the children of the right side. Such trees are called parse trees.



Parse Trees

A parse tree often shows how the symbol of a grammar derives a string in the language. If nonterminal A has a production $A \rightarrow xyz$, then a parse tree may have an interior node labeled A with three children labeled x , y and z , from left to right.



Formally, given a context free grammar, a parse tree is a tree with the following properties:

- The root is labeled by the start symbol
 - Each leaf is labeled by a token or by ϵ
 - Each interior node is labeled by a nonterminal.
 - If A is the nonterminal labeling some interior node and x_1, x_2, \dots, x_n are the labels of the children of that node from left to right, then $A \rightarrow x_1 x_2 \dots x_n$ is a production. Here x_1, x_2, \dots, x_n stand for a symbol that is either a terminal or a nonterminal. As a special case, if $A \rightarrow \epsilon$ then a node labeled A may have a single child labeled ϵ .
- list of
nonterminals
are
digits, +, and digit

Example 3 :- On fig 2, the root is labeled $list$, the start symbol on example(1). The children of the root are labeled from left to right, $list$, $,$, and $digit$.

The leaves of a parse tree read from left to right form the yield of the tree, which is the string generated from the nonterminal at the root of the parse tree. On fig 2, the generated string is $q-5+2$. On their fig, all the leaves are shown at the bottom level.

#. (1) Another definition of the language generated by a grammar is as "the set of strings that can be generated by some parse tree". The process of finding a parse tree for a given string of tokens is called parsing that string. //

Ambiguity :- If a grammar can have more than one parse tree generating a given string of tokens, then such a grammar is said to be ambiguous.

To show that a grammar is ambiguous, all we need

To do is find a token string that has more than one parse tree. Since a string with more than one meaning usually has more than one parse tree, for compiling applications we need to design unambiguous grammar.

EXAMPLE 4: suppose we did not distinguish between digits and lists as in example(1). We could have written the grammar

$$\text{string} \rightarrow \text{string} + \text{string} \mid \text{string} - \text{string} \mid 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

However, fig (3) shows that an expression like $9-5+2$ now has more than one parse tree. The two trees for $9-5+2$ correspond to the two ways of parenthesizing the expression: $(9-5)+2$ and ~~$9-(5+2)$~~

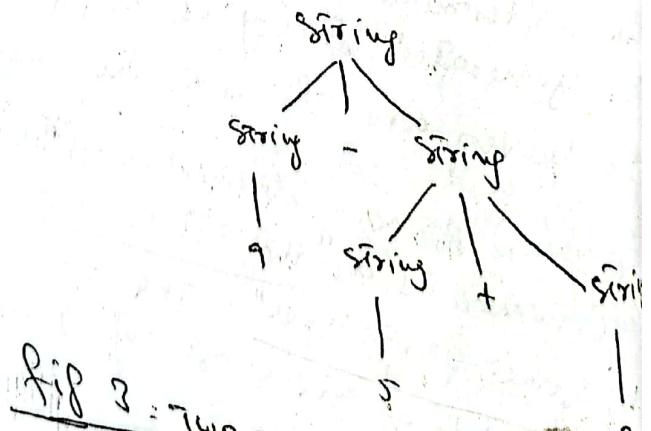
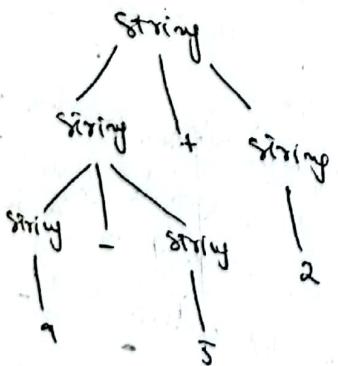


fig 3: Two parse trees for $9-5+2$.

Associativity of operators.

By Rule, $9+5+2$ is equivalent to $(9+5)+2$ and $9-5-2$ is equivalent to $(9-5)-2$. When an operand like 5 has operators to its left and right, conventions are needed for deciding which operator takes that operand. We say the operator + associates to the left because an operand with + signs on both sides of it is taken by the operator to its left. In most programming languages, the four arithmetic operators addition, subtraction, multiplication and division are left associative.

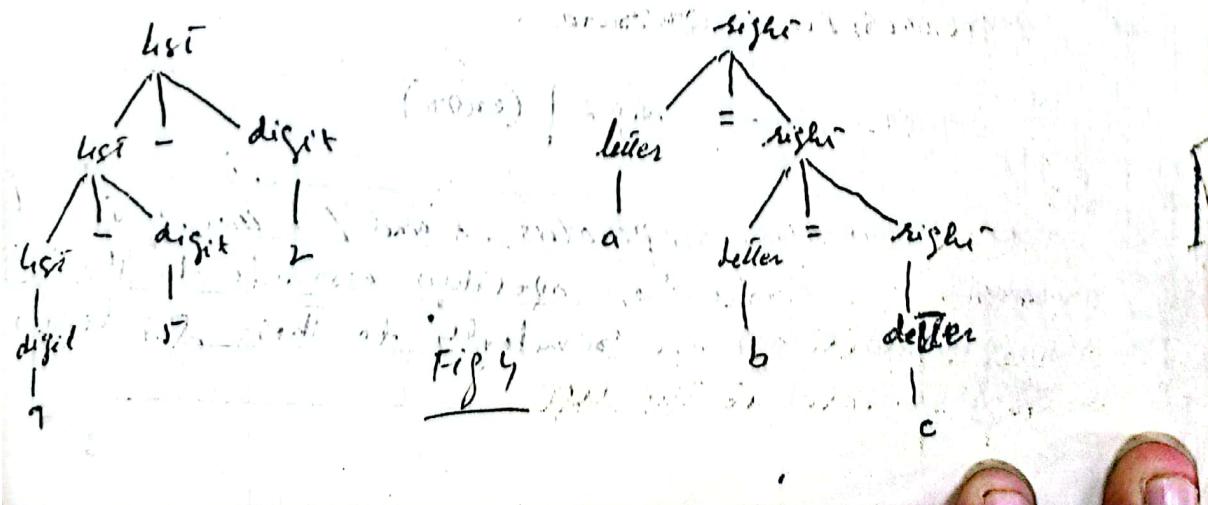
Some common operators such as exponentiation are right associative. As another example, the assignment operator = in C is right associative in C, the expression $a=b=c$ is treated in the same way as the expression $a=(b=c)$.

Strings like $a=b=c$ with a right associative operator is generated by the following grammar:

$$\text{right} \rightarrow \text{letter} = \text{right} \mid \text{letter}$$

$$\text{letter} \rightarrow a \mid b \mid c \mid \dots \mid z$$

The contrast between a parse tree for a left-associative operator like - and a parse tree for a right-associative operator like = is shown in fig 4. Note that the parse tree for $9-5-2$ grows down towards the left, whereas the parse tree for $a=b=c$ grows down towards the right.



precedence of operators:

Consider the expression $9+5*2$. There are two possible interpretations of this expression: $(9+5)*2$ or $9+(5*2)$. The associativity of + and * does not resolve this ambiguity. For this reason, we need to know the precedence of operators.

We say that * has higher precedence than + as * takes its operands before + does. In ordinary arithmetic multiplication and division have higher precedence than addition and subtraction. Therefore, 5 is taken by * in both $9+5*2$ and $9*(5+2)$; i.e., the expressions are equivalent to $9+(5*2)$ and $(9*5)+2$, respectively.

8

Syntax of expressions:

A grammar for arithmetic expressions can be constructed from a table showing the associativity and precedence of operators. We start with the four common arithmetic operators and a precedence table, showing the operators in order of increasing precedence.

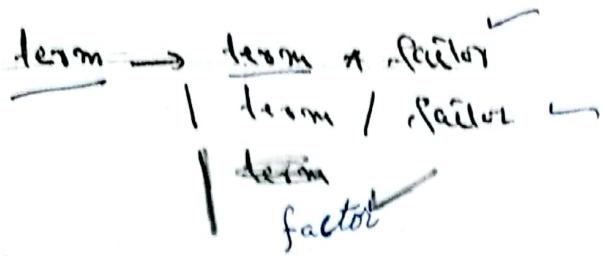
left associative: + -

left associative: * /

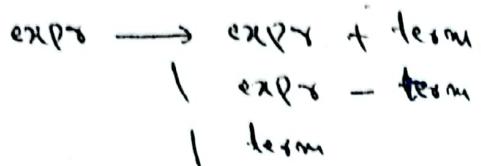
We create two nonterminals `expr` and `term` for the two levels of precedence, and an extra nonterminal factor for generating basic units in expressions. The basic units in expressions are presently digits and parenthesized expressions.

factor \rightarrow digit | (expr)

Now consider the operators, * and /, that have the high precedence. Since these operators associate to the left, the production are similarly to those for `term`.



Similarly, expr generates lists of terms separated by the additive operators.



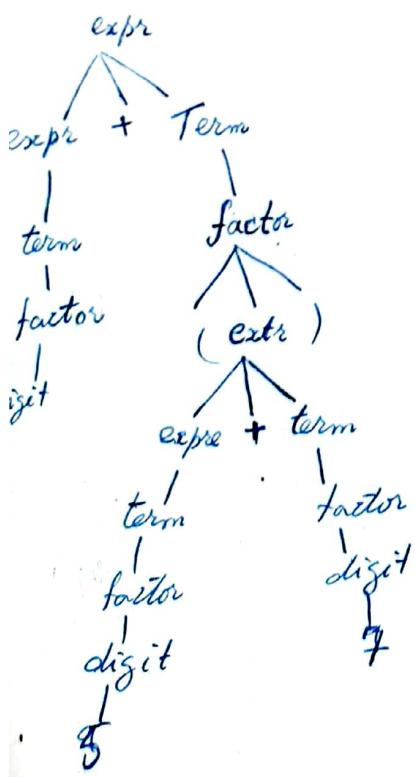
~~so far~~

The reducing grammar is therefore,

$$\begin{array}{l} \text{expr} \rightarrow \text{expr} + \text{term} \mid \text{expr} - \text{term} \mid \text{term} \\ \text{term} \rightarrow \text{term} * \text{factor} \mid \text{term} / \text{factor} \mid \text{factor} \\ \text{factor} \rightarrow \text{digit} \mid (\text{expr}) \end{array}$$

This grammar treats an expression as a list of terms separated by either + or - signs, and a term as a list of factors separated by * or / signs.

$$9 + (5 + 7)$$



SYNTAX-DIRECTED TRANSLATION

To translate a programming language, a compiler may need to track of many quantities besides the code generated for the construct. For example, the compiler may need to know the type of construct or location of the JST instruction in its target code or the number of instructions generated. We therefore talk about attributes associated with constructs. An attribute may represent any quantity, e.g., a string, a type, a memory location etc.

A syntax-directed definition specifies the formulation of a construct in terms of attributes associated with its syntactic components.

Postfix Notation:

The Postfix notation of an expression E can be defined as follows:

- 1 - If E is a variable or constant, then the postfix notation for E is E itself.
 - 2 - If E is an expression of the form $E_1 \text{ op } E_2$, where op is any binary operation, then the postfix notation for E is $E_1 E_2 \text{ op}$, where E_1 and E_2 are the postfix notations of E_1 and E_2 , respectively.
 - 3 - If E is an expression of the form (E_1) , then the postfix notation for E_1 is also the postfix notation for E .
- No parentheses are needed in postfix notation because the position and arity (number of arguments) of the operators permits only one decoding of a postfix notation. For example, the postfix notation for $(9-5)+2$ is $95-2+$ and the postfix notation for $9-(5+2)$ is $952+-$.

Syntax-Directed Definitions:

A syntax-directed definition uses a context-free grammar to specify the syntactic structure of

The input . With each grammar symbol , it associates the input . With each production , a set of semantic rules , and with each production , a set of attributes , and with each production , a set of attributes associated with rules for computing values of the attributes appearing in that production . The grammar and the set of semantic rules constitute the syntax-directed definition .

A translation is an input-output mapping .

The output of each input x is specified in the following manner . First , construct a parse tree for x . Suppose a node n in the parse tree is labeled by the grammar symbol X . We write $X.a$ to denote the value of attribute a of X at that node . The value of $X.a$ at n is computed using the semantic rule for attribute a associated with the X -production used at node n .

A parse tree showing the attribute values at each node called a annotated parse tree .

Synthesized Attributes

An attribute is said to be synthesized if its value at a parse-tree node is determined from attribute values at the children of the node . Synthesized attributes have the desirable property that they can be evaluated during a single bottom-up traversal of the parse tree .

Example 5:

A syntax-directed definition for translating expressions consisting of digits separated by plus or minus signs into postfix notation is shown in fig 5. Associated with each nonterminal is a string-valued attribute that represents the postfix notation for the expression generated by that nonterminal in a parse tree .

The value of attribute is a string .

production

$\text{expr} \rightarrow \text{expr}_1 + \text{term}$
 $\text{expr} \rightarrow \text{expr}_1 - \text{term}$
 $\text{expr} \rightarrow \text{term}$
 $\text{term} \rightarrow 0$
 $\text{term} \rightarrow 1$
 \vdots
 $\text{term} \rightarrow 9$

Semantic Rule

$\text{expr.t} := \text{expr}_1.t \parallel \text{term.t} \parallel '+'$
 $\text{expr.t} := \text{expr}_1.t \parallel \text{term.t} \parallel '-'$
 $\text{expr.t} := \text{term.t}$
 $\text{term.t} := '0'$
 $\text{term.t} := '1'$
 \vdots
 $\text{term.t} := '9'$

Fig 5: - Syntax-directed definition for infix to postfix translation

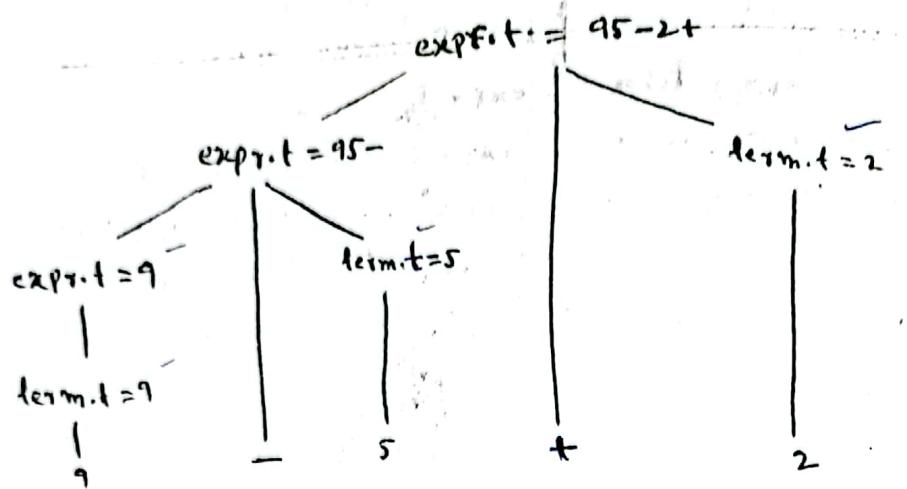
The Postfix form of a digit is the digit itself; e.g.; the semantic rule associated with the production $\text{term} \rightarrow 9$ defines term.t to be 9 whenever this production is used at a node in a parse tree. When the production $\text{expr} \rightarrow \text{term}$ is applied, the value of term.t becomes the value of expr.t.

The production $\text{expr} \rightarrow \text{expr}_1 + \text{term}$ derives an expression containing a plus operator. The left operand of the plus operator is given by expr_1 and the right operand by term . The semantic rule:

$\text{expr.t} := \text{expr}_1.t \parallel \text{term.t} \parallel '+'$

associated with this production defines the value of attribute expr.t by concatenating the postfix forms expr.t and term.t of the left and right operands, respectively, and then appending the plus sign. The operator \parallel in semantic rules represents string concatenation.

Fig 6 contains the annotated parse tree corresponding to the tree of Fig 2. The value of the t-attribute at each node has been computed using the semantic rule associated with the production used at that node. The value of the attribute at the root is the Postfix notation for the string generated by the parse tree.



~~Fig 6~~ : Attribute values at nodes in a parse tree

Depth-First Traversal

A traversal of a tree starts at the root in some order. Here, Semantic rules will be evaluated using the depth-first traversal defined in Fig 7. It starts at the root and recursively visits the children of each node in left-to-right order. The semantic rules at a given node are evaluated once all descendents of that node have been visited. It is called "depth-first" because it visits an unvisited child of a node whenever it can, so it tries to visit nodes as far away from the root as quickly as it can.

procedure visit (n : node);
 begin
 for each child m of n , from left to right do
 visit (m);

evaluate semantic rules at node n

end

~~Fig 8~~ :- A depth-first-traversal of a tree

Translation Schemes

A translation scheme is a context free grammar in which program fragments called semantic actions are added within the right sides of productions. A translation scheme is like a syntax-directed definition, except that the order of evaluation of the semantic rules is explicitly shown. The position at which an action is to be executed is shown by enclosing it between braces and writing it within the right side of a production, as in the following example:

$\text{nest} \rightarrow + \text{term} \{ \text{print}('+' \}) \text{rest}$

A translation scheme generates an output for each sentence it generates by the underlying grammar by executing the actions in the order they appear during the depth-first traversal of a parse tree. For example, consider a parse tree with a node labeled rest, representing a production. The action $\{ \text{print}('+' \})$ will be performed after the subtree for term is traversed but before the child for rest is visited.

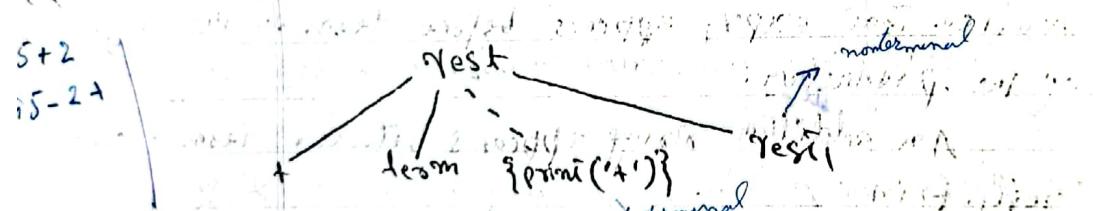


Fig 8: An extra leaf is constructed for a semantic action when drawing a parse tree for a translation scheme.

When drawing a parse tree for a translation scheme we indicate an action by constructing for it an extra child, connected by a dashed line to the node for its production. For example, the portion of the parse tree for the above production and action is drawn in the above fig 8. The node for a semantic action has no children, so the action is performed when that node is first seen.

Emitting a Translation:

The Semantic action in translation schemes write the output of a translation into a string. When the output is emitted, in this fashion, the order in which the characters are printed is important.

The string representing the translation of the nonterminal on the left side of each production is the concatenation of the translations of the nonterminals on the right, in the same order as in the production, with some additional strings (perhaps none) interleaved. A syntax-directed definition with this property is termed Simple. For example, consider the first production and semantic rule from the syntax-directed definition of expr given earlier.

Production Rule & SEMANTIC RULE

$$\text{expr} \rightarrow \text{expr}_1 + \text{term} \quad \text{expr}_1 := \text{expr}_1, t \parallel \text{term} \parallel '+'$$

Here the translation expr_1, t is the concatenation of the translations of expr_1 and term , followed by the symbol $+$. Notice that expr_1 appears before term on the right side of the production.

An additional string appears between term and $+$ in expr_1, t .

Produced Rule & SEMANTIC RULE

$$\text{rest} \rightarrow + \text{term} \text{ rest} \quad \text{rest} := \text{term} \parallel '+' \parallel \text{rest}, t$$

but again, the nonterminal term appears before rest , on the right side of both rules.

Simple syntax-directed definition can be implemented with translation schemes in which actions print the additional strings in the order they appear in the definition. The actions in the following productions print the additional strings in (5) and (6), respectively,



$\text{expr} \rightarrow \text{expr} + \text{term} \{ \text{print}('+' \})$

$\text{rest} \rightarrow + \text{term} \{ \text{print}('+' \}) \text{ rest}$

Example 6 :- Fig 5 contained a simple definition of translating expressions into postfix form. A translation scheme derived from this definition is given in fig 9 and a parse tree with actions for $9-5+2$ is shown in fig 10. Note that although figures 6 and 10 represent the same input-output mapping, the translation in the two cases is conducted differently; Fig 6 attaches the output to the root of the parse tree, while fig 10 prints the output incrementally.

$\text{expr} \rightarrow \text{expr} + \text{term} \{ \text{print}('+' \})$

$\text{expr} \rightarrow \text{expr} - \text{term} \{ \text{print}('-' \})$

$\text{expr} \rightarrow \text{term}$

$\text{term} \rightarrow 0 \{ \text{print}('0' \})$

$\text{term} \rightarrow 1 \{ \text{print}('1' \})$

:

$\text{term} \rightarrow 9 \{ \text{print}('9' \})$

Fig 9:- Actions translating expressions into postfix notation.

The root of fig 10 represents the ~~product~~ of first production in fig 9. In a depth-first traversal, we first perform all the actions in the subtree for the left operand expr we use when we traverse the leftmost subtree of the root. We then visit the leaf + at which there is no action. We next perform the actions in the subtree for the right operand term and finally, the semantic action $\{ \text{print}('+' \})$ at the ultra node.

Since the productions for term have only a digit on the right side, that digit is printed by the

actions for the productions. No output is necessary for the production $\text{expr} \rightarrow \text{term}$, and only the operator for needs to be printed in the action for the first three productions. When executed during a depth-first traversal of the parse tree, the actions in fig. 10 print 95-2+.

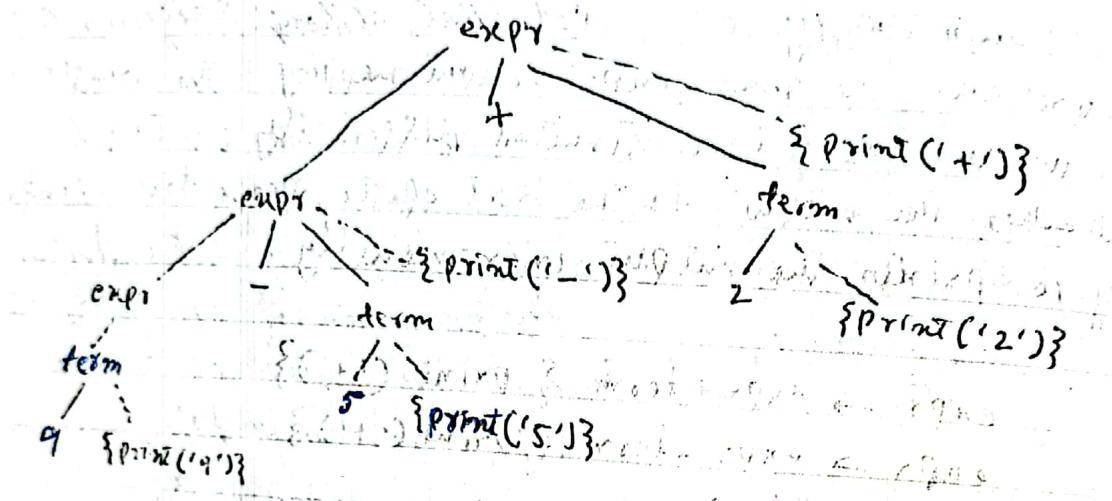


Fig 10:- Actions translating 9-5+2 into 95-2+.

As a general rule, most parsing methods processes their input from left to right in a "greedy" fashion; that is they consume as much as of a parse tree as possible before reading the next input token. In a simple translation scheme, actions are also done in a left-to-right order.



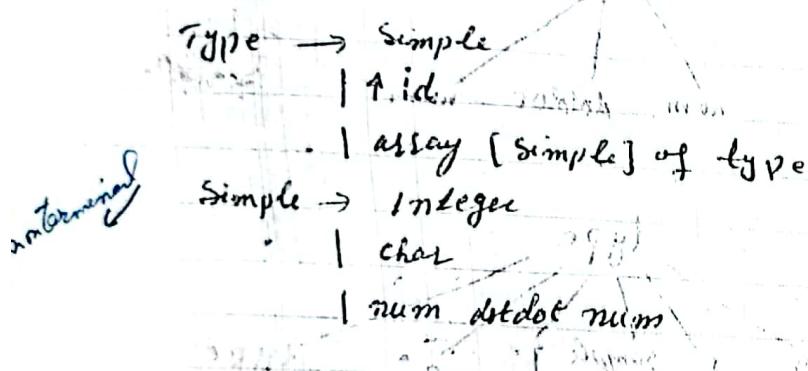
PARSING

Parsing is the process of determining if a string of tokens can be generated by a grammar.

Most parsing methods have two classes, called the top-down & bottom-up methods. These terms refer to the order which nodes in the parse tree are constructed. In the former, construction starts at the root and proceeds towards the leaves, while in the latter, construction starts at the leaves and proceeds towards the root.

Top-down Parsing :-

The following grammar generates a subset of the types of Pascal. We use the token dotdot for "..."



The top-down construction of a parse tree is done by starting with the root, labeled with the starting non-terminal, and repeatedly performing the following two steps:

- 1 - At node n , labeled with non-terminal A , select one of the productions for A , and construct children of n for the symbols on the right side of the production.
- 2 - find the next node at which a subtree is to be constructed.

The top-down construction of the parse tree
for the above grammar is :

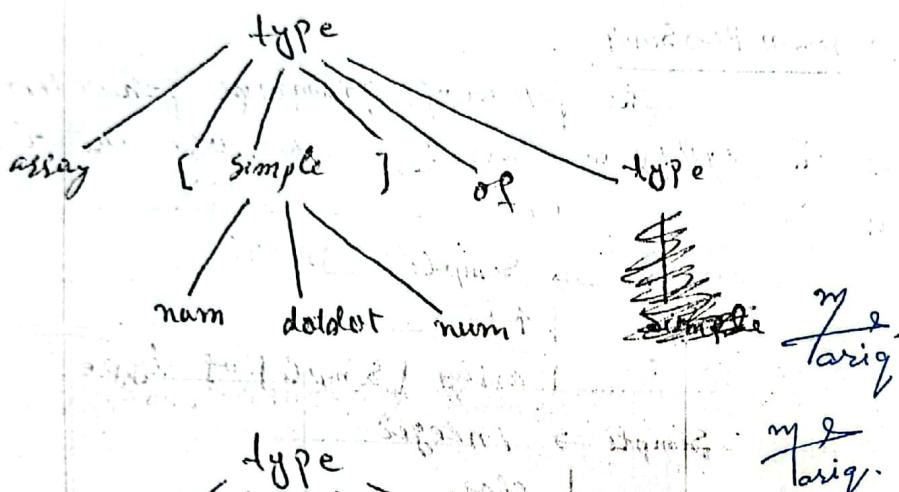
(a).

type

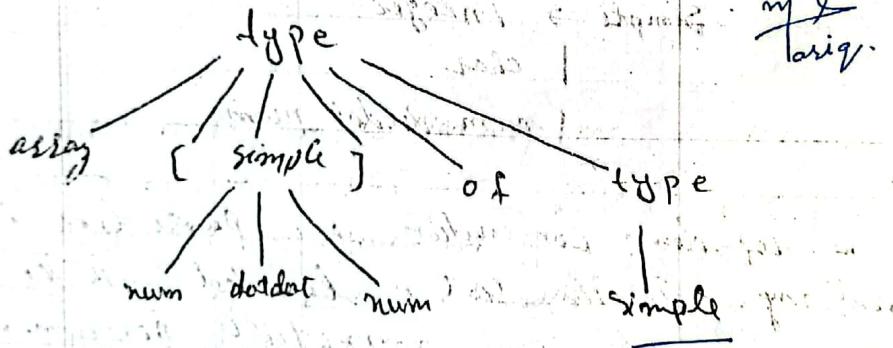
(b)



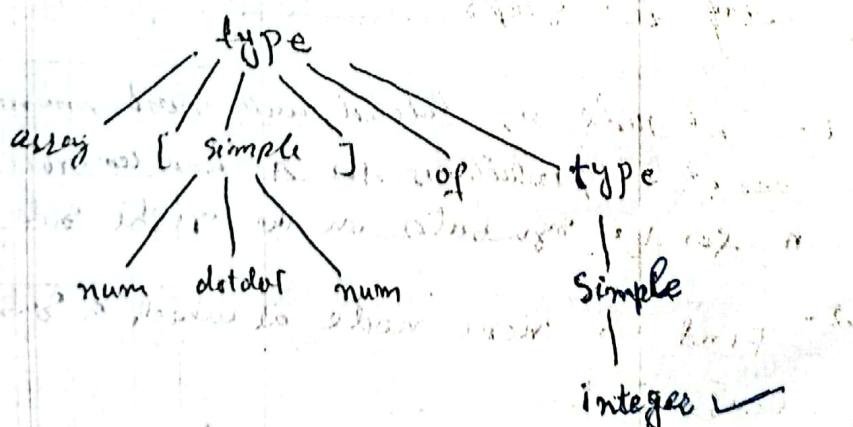
(c)



(d)



(e)



LEXICAL ANALYSIS

23

22

The Role of THE LEXICAL ANALYZER:

The lexical analyzer is the first phase of a compiler. Its main task is to read the input characters, and produce an output in sequence of tokens that the parser the parser uses for syntax analysis. The interaction of the lexical analyzer with parser is shown in fig 1.

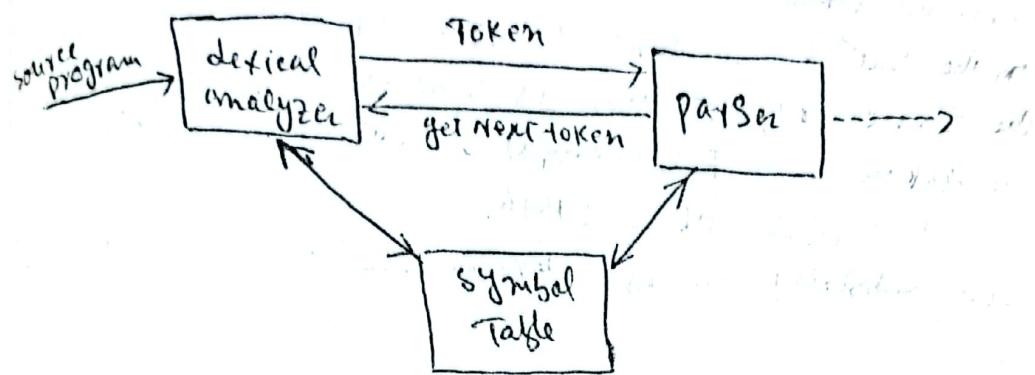


FIG 1:- Interaction of Lexical Analyzer with Parser

Since the lexical analyzer is the part of the compiler that reads the source text, it may also perform certain secondary tasks. One such task is stripping out from the source program comments and white space in the form of blank, tab, and newline characters. Another is correlating error messages from the compiler with the source program. For example, the lexical analyzer may keep track of the number of newline characters seen, so that a line number can be associated with an error message.

Sometimes, lexical analyzers are divided into two phases, the first called "scanning" and the second "lexical analysis". The scanning is responsible for doing simple tasks, while the lexical analyzer does the more complex operations. For example, a C program might use a scanner to eliminate blanks from the input.

Tokens, patterns, lexemes

When talking about lexical analysis, we use the terms token, pattern and lexeme with specific meaning. Examples of their use are shown in fig. 2.13.

In general, there is a set of strings in the input which the same token is produced as output. This set of strings is described by a rule called a Pattern, associated with the token. The pattern is said to match each string in the set. A lexeme is a sequence of characters in the source program that is matched by the pattern of a token. For example, in the Pascal statement

const PI = 3.1416;

The substring PI is a lexeme for the token "Identifier".

| TOKEN | SAMPLE LEXIMES | INFORMAL DESCRIPTION OF PATTERN |
|----------|----------------------|---------------------------------------|
| const | const | const |
| if | if | if |
| relation | <, <=, =, >, >= | < or <= or = or > or >= |
| id | PI, count, D2 | letter followed by letters and digits |
| num | 3.1416, 0.6, 0.02E23 | any numeric constant |
| literal | "A\$alarm" | Any character between " and " |

Fig 2:- Examples of tokens!

We treat tokens as terminal symbols in the grammar of the source language, using underline names to represent tokens. (The lexemes matched by the pattern for the token represent strings of characters in the source program that can be treated together as a lexical unit.)

In most programming languages, the following constructs, treated as tokens: keywords, operators, identifiers, units, literal strings, and punctuation symbols, in parentheses, commas, and semicolons. In the example above, when the character sequence pi appears in source program, a token representing an identifier is passed to the parser.

A pattern is a rule describing the set of elements that can represent a particular token in source form. It consists of individual tokens, their attributes, and attributes for tokens.

The lexical analyzer collects information about tokens and their associated attributes. As a practical matter, token has usually only a single attribute — a pointer to symbolic-table entry in which the information about token is kept; the pointer becomes the attribute for token.

Example 1 :- The tokens and associated attribute-values in the FORTRAN statement

$E = M * C * T / 2$ are written below as a sequence of pairs:

< id, > Pointer to 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
< div-op, >
< num, integer value 2 >

Lexical Errors - If the string "fi" is encountered in a C program for the first time in the context

fi (a = f(n);)

a lexical analyzer cannot tell whether fi is a misspelling of the keyword if, or an undeclared function identifier. Since fi is valid identifier, the lexical analyzer may return the token for an identifier and let some other phase of the compiler handle any error.

But suppose a situation does arise in which the lexical is unable to proceed because none of the patterns for tokens matches a prefix of the remaining input. The simplest recovery strategy is "panic mode" recovery. We delete the successive character from the remaining input until the lexical analyzer can find a well-formed token. Other possible error-recovery actions are:

- 1 - deleting an extraneous character
- 2 - inserting a missing character
- 3 - replacing an incorrect character by a correct character
- 4 - transposing two adjacent characters

INPUT BUFFERING :

Buffer Pairs:-

For many source languages, the lexical analyzer needs to look ahead several characters beyond the lexeme for a pattern before a match can be announced.

We use pair buffers divided into two N-character halves as shown in Fig 3.

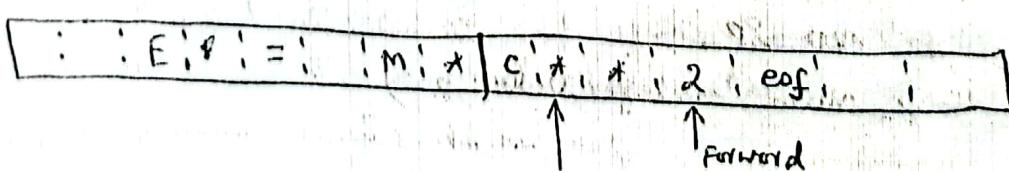


Fig 3:- An input buffer in two halves

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; 25 fewer than N characters remain in the input. Then a special character `eof` is read into the buffer after the input characters. That is, `eof` marks the end of the source file.

Two pointers to the input buffer are maintained. The string of characters between the two pointers is the current item. Initially, both pointers point to the first character of the next item to be found. One, called the forward pointer, scans ahead until a match for the pattern is found. Once the next item is determined, a forward pointer is set to the character at its right end. After the item is processed, both pointers are set to the character immediately past the item.

If the forward pointer is about to move past the next item 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.

Sentinel?

If we use the scheme of fig 3, 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 a test shown in fig 4. (Note the double `if` test, which handles both ends of the buffer.)

```

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 the beginning of first half
end
else forward := forward + 1

```

FIG 4:- Code to advance forward pointer.

Except at the ends of the buffer halves, the code in fig 4 requires two tests for each advance of a 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 5 shows the same buffer arrangement as fig 3 with the sentinels added.

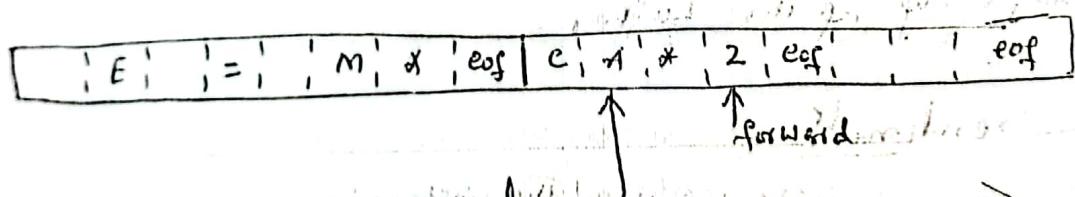


FIG 5:- Sentinels at end of each buffer half.

With the arrangement of fig 5, we can use the code shown in fig 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.

```

forward := forward + 1;
if forward  $\uparrow$  = 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
        terminate lexical analysis
end

```

fig6:- lookahead code with Sentinels

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.

String and Languages:

The term alphabet or character class denotes a finite set of symbols. Typical examples of symbols, letters and characters. The set $\{0, 1\}$ is the binary alphabet.

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 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.

5' 2
2 2 vs
2 2

The term language denotes any set of strings over some fixed alphabet. Languages like \emptyset , the empty set or $\{\epsilon\}$, the set containing only the empty string, are languages under this definition.

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, $\underline{S}\epsilon = \epsilon S = S$.

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 S^0 is S itself, $S^1 = S$, then $S^2 = \underline{SS}$, $S^3 = \underline{SSS}$ and so on.

Some common terms associated with parts of string are shown in fig 7.

| <u>TERM</u> | <u>DEFINITION</u> |
|------------------|---|
| prefix of s | A string obtained by removing zero or more trailing symbols of string s . e.g., <u>ba</u> is a prefix of <u>banana</u> . |
| suffix of s | A string formed by deleting zero or more of the leading symbols of s . e.g., <u>na</u> is a suffix of <u>banana</u> . |
| substring of s | A string obtained by deleting a prefix and a suffix from s e.g., <u>na</u> is a substring of <u>banana</u> . Every prefix and every suffix of s is a substring of s , but not every string 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 . |

6.5

S.t

"prefix, suffix, a substring of S"

Any nonempty string x , that is, respectively, a prefix, suffix, or substring of S such that $S \neq x$.

Sequence of S

Any string formed by deleting zero or more not necessarily contiguous symbols from S , e.g; banana is a subsequence of banana

FIG 8:- Terms for Parts of a String

Operations and Languages:-

There are several important operations that can be applied to languages, which are union, concatenation, and closure as shown in fig 8. We can also parallelize the exponentiation operator to languages by defining L^0 to be $\{\epsilon\}$ and L^i to be $L^{i-1}L$.

OPERATIONS

union of L and M
written LUM

DEFINITION

$$LUM = \{s \mid s \text{ is in } L \text{ or } s \text{ is in } M\}$$

concatenation of L and M
written LM

$$LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$$

Kleene closure of L
written L^*

$$L^* = \bigcup_{i=0}^{\infty} L^i \quad \begin{matrix} \text{Concatenation} \\ \text{not zero} \end{matrix}$$

$\{ L^* \text{ denotes "zero or more concatenations of } L \text{"}$

positive closure of L
written L^+

~~$L^+ = \bigcup_{i=1}^{\infty} L^i$~~

$\{ L^+ \text{ denotes "one or more concatenations of } L \text{"}$

FIG 8:- Definitions of Operations on Languages

Example 2: Let L be the set $\{A, B, \dots, Z, a, b, c, \dots, z\}$ and D the set $\{0, 1, \dots, 9\}$.
The new languages, created from L and D by applying the above operations are as follows:

- 1- LUD 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 .
- 5- $L(LUD)^*$ 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

On Pascal, an identifier is a letter followed by zero or more letters or digits. But then, here, with a notation, called regular expression, we might define Pascal identifiers as
 $\text{letter} \underset{\text{or}}{(}\text{letter} \mid \text{digit})^*$ containing at least one

The vertical bar here means "or". The parentheses are used to group subexpressions, the star means zero or more instances of the parenthesized expressions.

Each regular expression γ denotes a language $L(\gamma)$. The defining rules specify how $L(\gamma)$ is formed by combining in various ways the languages denoted by the subexpressions of γ .

These are the rules that define the regular expression over alphabet Σ . Associated with each rule is a specification of the language denoted by the regular expression being defined.

\emptyset is a regular expression that denotes $\{\emptyset\}$, i.e., the set containing the empty string.

If a is a symbol in Σ , then a is a regular expression that denotes $\{a\}$ i.e., the set containing the string a .

Suppose r and s are regular expressions denoting the languages $L(r)$ and $L(s)$. Then,

(r) $|$ (s) is a regular expression denoting $\underline{L(r) \cup L(s)}$

N.B:

- b) $\underline{r} \underline{s}$ is a regular expression denoting $\underline{L(r)} \underline{L(s)}$
- c) \underline{r}^* is a regular expression denoting $(L(r))^*$
- d) \underline{r}^2 is a regular expression denoting $L(r)^2$

Language denoted by a regular expression is said to be regular set.

unnecessary parentheses can be avoided in regular expressions if we adopt the conventions that:

1. The unary operator $*$ has the highest precedence.
2. Concatenation has the second highest precedence.
3. $|$ has the lowest precedence.

Under these conventions, $(a) | ((b)^* c)$ is equivalent to $\underline{ab^*c}$. Both expressions denote the set of strings that are in a single a or zero or more b 's followed by one c .

Example:- Let $\Sigma = \{a, b\}$

1. The regular expression $a|b$ denotes the set $\{a, b\}$

2. The regular expression $(a|b)(a|b)$ denotes $\{a^2, ab, ba, bb\}$,
The set of all strings of a^2 and b^2 of length two.

3. The regular expression a^* denotes the set of all strings of zero or more a 's i.e., $\{ \epsilon, aa, aaa, \dots \}$

4 - The regular expression $(a|b)^*$ denotes the set of all strings containing zero or more instances of an a or b, i.e., the set of all strings of a's and b's. Another regular expression for this set is $(a^*b^*)^*$.

5 - The regular expression $a|ab$ denotes the set containing the string a and all strings consisting of zero or more a's followed by a b.

If two regular expressions r and s denote the same language, we say r and s are equivalent and write r = s. For example, $(a|b) = (b|a)$

Fig 9 shows the some algebraic laws that hold for regular expressions. r, s, and t

| <u>Axiom</u> | <u>Description</u> |
|-------------------------------------|--|
| $r s = s r$ | $ $ is commutative |
| $r(s t) = (r s) t$ | $ $ is associative |
| $(rs)t = r(st)$ | Concatenation is associative |
| $r(s(t)) = rs(t)$ | Concatenation distributes over $ $ |
| $(s t)r = sr tr$ | |
| $\epsilon r = r$ $r\epsilon = r$ | ϵ is the identity element for concatenation |
| $r^* = (r \epsilon)^*$ | relation between $*$ and ϵ |
| $r^* = r^{**}$ | |

H → for a note take

Fig 9: Algebraic Properties of Regular Expressions.

Regular Definitions:

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

$$d_1 \rightarrow t_1$$

$$d_2 \rightarrow t_2$$

$$\vdots$$

$$d_n \rightarrow t_n$$

in which each d_i is a distinct name over the symbols in $\{d_1, d_2, \dots, d_{i-1}\}$, i.e., the basic symbols and the previously defined names.

EXAMPLE 4: - The set of Pascal identifiers is the set of strings of letters and digits beginning with a letter.

The regular definition for this set is

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

$$\underline{\text{digit}} \rightarrow 0 | 1 | 2 | \dots | 9$$

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

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

EXAMPLE 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:

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

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

$$\underline{\text{optional fraction}} \rightarrow \cdot \underline{\text{digits}} | \epsilon$$

$$\underline{\text{optional exponent}} \rightarrow (E (+|-|) \underline{\text{digits}}) | \epsilon$$

$$\underline{\text{num}} \rightarrow \underline{\text{digits}} \underline{\text{optional fraction}} \underline{\text{optional exponent}}$$

Notational Shorthands:-

Notational shorthands for components in regular expressions are shown in instances

- ~~$\star = \star \circ \star$~~
- one or more instances: The unary operator \star means "one or more instances of". If γ is a regular expression that denotes the language $L(\gamma)$, then $(\gamma)^\star$ is a regular expression that denotes the language $(L(\gamma))^\star$. Thus, the regular expression a^\star denotes the set of all strings of one or more a 's. The two algebraic identities $\gamma^\star = \gamma^\star \mid \epsilon$ and $\gamma^\star = \gamma \gamma^\star$ relate the Kleene and positive closure operators.

- zero or one instance: The unary operator $?$ means "zero or one instance of". The notation $\gamma ?$ is a shorthand for $\gamma \mid \epsilon$. If γ is a regular expression, then $(\gamma) ?$ is a regular expression that denotes the language $L(\gamma) \cup \{\epsilon\}$. For example, using the \star and $?$ operators, we can rewrite the regular definition of num in Example 5 as

Digit Shorthand

digit $\rightarrow b / 1 / - \rightarrow 19$
digit⁺ \rightarrow digit⁺ \rightarrow 1 or more instances not zero

optional fraction $\rightarrow (\cdot \text{ digits}) ?$

optional exponent $\rightarrow (E (+|-) ? \text{ digits}) ?$

num \rightarrow digits optional fraction optional exponent

- 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 $a \mid b \mid c \mid \dots \mid z$.

RECOGNITION OF TOKENS :-

EXAMPLE 6: Consider the following grammar fragment:

$$\text{stmt} \rightarrow \text{if } \text{expr} \cdot \underline{\text{then}} \text{ stmt}$$

$$\quad | \quad \underline{\text{if }} \text{expr} \cdot \underline{\text{then}} \text{ stmt } \underline{\text{else}} \text{ stmt}$$

$$\quad | \quad \epsilon$$

$$\text{expr} \rightarrow \text{term} \cdot \underline{\text{relop}} \text{ term}$$

$$\quad | \quad \text{term}$$

$$\text{term} \rightarrow \underline{\text{id}} \cdot \underline{\text{num}}$$

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

$$\underline{\text{if}} \rightarrow \text{if}^*$$

$$\underline{\text{Then}} \rightarrow \text{Then}^*$$

$$\underline{\text{else}} \rightarrow \text{else}^*$$

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

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

$$\underline{\text{num}} \rightarrow \underline{\text{digit}}^+ (\cdot \underline{\text{digit}}^+)? ((\underline{E} (+|-))? \underline{\text{digit}}^+)?$$

where letter and digit are as defined in example 4.

For this language, the lexical analyzer will recognize the keywords if, Then, else as well as the lexemes denoted by relop, id and num. For simplicity we assume keywords are reserved; that is, they cannot be used as identifiers.

In addition, we assume lexemes are separated by white space, consisting of sequences of blanks, tabs, and New Lines. Our lexical analyzer will skip over white space. It will do so by comparing a string against the regular expression ws - below

delim → balance | 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 whitespace and returns that to the parser.

The lexical analyzer 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 in Fig 10. The attribute-values for the relational operators are given by the symbolic constants LT, LE, EQ, NE, GT, GE.

| REGULAR EXPRESSION | TOKEN | ATTRIBUTE-VALUE |
|--------------------|-------|--------------------------|
| WS | - | - |
| if | if | - |
| Then | Then | - |
| else | else | - |
| token (| - | - |
| id | id | - |
| num | num | Pointer to Table entry ✓ |
| < | relOp | Pointer to Table entry ✓ |
| <= | relOp | LT |
| = | relOp | LE |
| > | relOp | EQ |
| >= | relOp | NE |
| | | GT |
| | | GE |

Fig 10:- Regular-expression Patterns for TOKENS

Transition Diagrams: As an intermediate step in the construction of a lexical analyzer, we first produce a stylized flowchart, called a transition diagram. Transition diagrams illustrate the actions that take place when a lexical analyzer is fed by the parser to get the next token. Suppose that the input buffer is as in Fig 3 and the lexeme beginning pointer points to the character following the last lexeme found. We use a transition diagram to keep track of information about lexemes that are seen, as forward pointers scan the input. We do so by moving from position to position in a diagram as characters are read.

Positions in a transition diagram are drawn as circles and are called states. The states are connected by arrows, called edges. Edges leaving state s have labels indicating the input characters that can next appear after the transition diagram has reached state s . The label other refers to any character that is not indicated by any of the other edges leaving s .

One state is labeled the start state; it is the initial state of the transition diagram where control resides when we begin to try to recognize a token. If one entering a state, we read the next input character. If there is an edge from the current state whose label matches this input character, we then go to the state pointed to by the edge. Otherwise, we indicate failure.

Fig 11 shows a transition diagram for the tokens \geq and $>$. The transition diagram works as follows:-

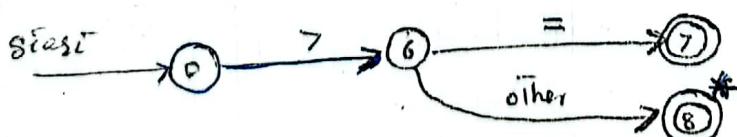


Fig:- 11. Transition diagram for $>=$

as start state is state 0. In state 0, we read the next input character. The edge labeled $>$ from state 0 is to be followed to state 6 if this input character is $>$. Otherwise we have failed to recognize either $> \circ \cdot 2 \geq$.

On reaching state 6 we read the next input character. The edge labeled $=$ from state 6 is to be followed to state 7 if this input character is $=$. Otherwise the edge labeled other indicates that we are to go to state 8. The double circle on state 7 indicates that it is an accepting state, a state on which the token \geq has been found.

Notice that the character $>$ and another extra character are read as we follow the sequence of edges from the start state to the accepting state 8. Since the extra character is not a part of the relational operator \geq , we must reset the forward pointer one character. We use \ast to indicate states on which this input retraction must take place.

In general, there may be several transition diagrams, each specifying a group of tokens. If failure occurs while we are following one transition diagram, then we must reset the forward pointer to where it was in the next state of this diagram, and reactivate the next transition diagram. If failure occurs in all transition diagrams, then a lexical error has been detected.

The End. of the News.