

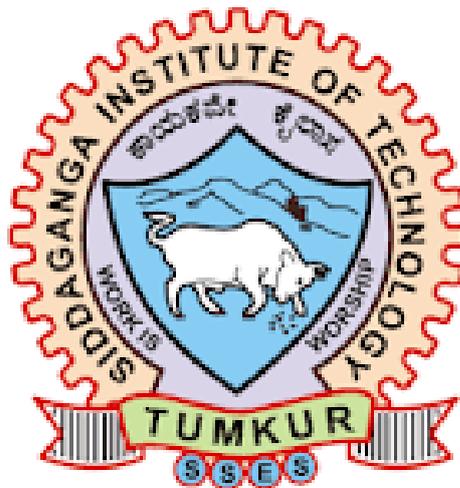
# **DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

## **Lecture Notes**

**Course: System Software & Compiler Design**

**Course Code:6CS01**

**Faculty: Dr. Shreenath KN**



## **SIDDAGANGA INSTITUTE OF TECHNOLOGY TUMKUR-3**

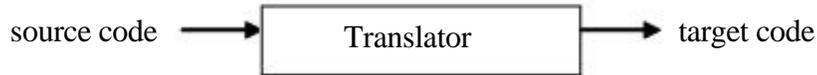
An Autonomous Institution, Affiliated to VTU, Belagavi & Recognised by AICTE  
and Accredited by NBA, New Delhi

## UNIT I- LEXICAL ANALYSIS

### INTRODUCTION TO COMPILING

Translator:

It is a program that translates one language to another.



Types of Translator:

- 1.Interpreter
- 2.Compiler
- 3.Assembler

1.Interpreter:

It is one of the translators that translate high level language to low level language.



During execution, it checks line by line for errors.

Example: Basic, Lower version of Pascal.

2.Assembler:

It translates assembly level language to machine code.



Example: Microprocessor 8085, 8086.

3.Compiler:

It is a program that translates one language(source code) to another language (target code).



It executes the whole program and then displays the errors.

Example: C, C++, COBOL, higher version of Pascal.

Difference between compiler and interpreter:

<b>Compiler</b>	<b>Interpreter</b>
It is a translator that translates high level to low level language	It is a translator that translates high level to low level language
It displays the errors after the whole program is executed.	It checks line by line for errors.
Examples: Basic, lower version of Pascal.	Examples: C, C++, Cobol, higher version of Pascal.

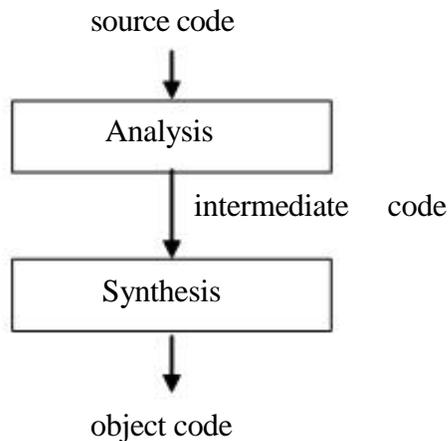
## PARTS OF COMPILATION

There are 2 parts to compilation:

1. Analysis
2. Synthesis

Analysis part breaks down the source program into constituent pieces and creates an intermediate representation of the source program.

Synthesis part constructs the desired target program from the intermediate representation.



Software tools used in Analysis part:

1) 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.
- For example , it can supply key words automatically - while .... do and begin..... end.

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.

3) Static checkers :

- A static checker reads a program, analyzes it, and attempts to discover potential bugs without running the program.
- For example, a static checker may detect that parts of the source program can never be executed.

4) Interpreters :

- Translates from high level language ( BASIC, FORTRAN, etc..) into machine language.  An interpreter might build a syntax tree and then carry out the operations at the nodes as it walks the tree.
- Interpreters are frequently used to execute command language since each operator executed in a command language is usually an invocation of a complex routine such as an editor or compiler.

## ANALYSIS OF THE SOURCE PROGRAM

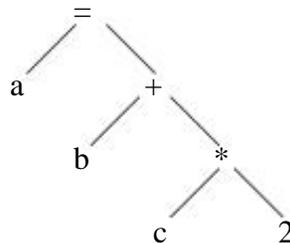
Analysis consists of 3 phases:

Linear/Lexical Analysis :

- It is also called scanning. It is the process of reading the characters from left to right and grouping into tokens having a collective meaning.
- For example, in the assignment statement  $a=b+c*2$ , the characters would be grouped into the following tokens:
  - i) The identifier1 'a'
  - ii) The assignment symbol (=)
  - iii) The identifier2 'b'
  - iv) The plus sign (+)
  - v) The identifier3 'c'
  - vi) The multiplication sign (\*)
  - vii) The constant '2'

Syntax Analysis :

- It is called parsing or hierarchical analysis. It involves grouping the tokens of the source program into grammatical phrases that are used by the compiler to synthesize output.
- They are represented using a syntax tree as shown below:



- A syntax tree is the tree generated as a result of syntax analysis in which the interior nodes are the operators and the exterior nodes are the operands.
- This analysis shows an error when the syntax is incorrect.

Semantic Analysis :

- It checks the source programs for semantic errors and gathers type information for the subsequent code generation phase. It uses the syntax tree to identify the operators and operands of statements.
- An important component of semantic analysis is type checking. Here the compiler checks that each operator has operands that are permitted by the source language specification.

## PHASES OF COMPILER

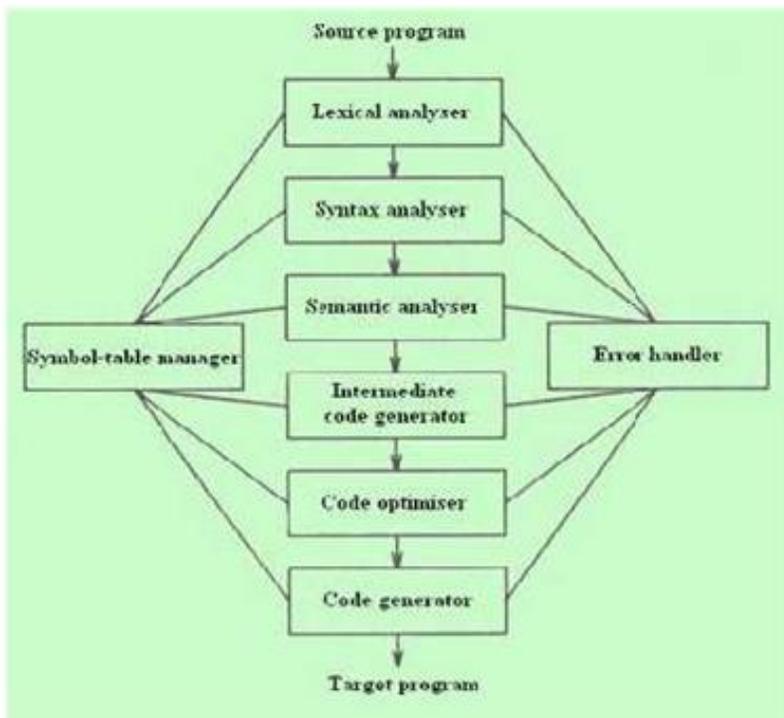
A Compiler operates in phases, each of which transforms the source program from one representation into another. The following are the phases of the compiler:

Main phases:

- 1) Lexical analysis
- 2) Syntax analysis
- 3) Semantic analysis
- 4) Intermediate code generation
- 5) Code optimization
- 6) Code generation

Sub-Phases:

- 1) Symbol table management
- 2) Error handling



### LEXICAL ANALYSIS:

- It is the first phase of the compiler. It gets input from the source program and produces tokens as output.
- It reads the characters one by one, starting from left to right and forms the tokens.

Token : It represents a logically cohesive sequence of characters such as keywords, operators, identifiers, special symbols etc.

Example:  $a + b = 20$

Here,  $a, b, +, =, 20$  are all separate tokens.

Group of characters forming a token is called the Lexeme.

- The lexical analyser not only generates a token but also enters the lexeme into the symbol table if it is not already there.

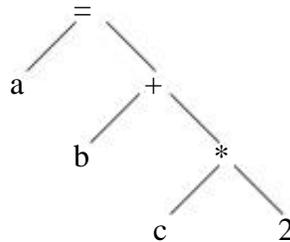
## SYNTAX ANALYSIS:

- It is the second phase of the compiler. It is also known as parser.
- It gets the token stream as input from the lexical analyser of the compiler and generates syntax tree as the output.

### Syntax tree:

It is a tree in which interior nodes are operators and exterior nodes are operands.

Example: For  $a=b+c*2$ , syntax tree is



## SEMANTIC ANALYSIS:

- It is the third phase of the compiler.
- It gets input from the syntax analysis as parse tree and checks whether the given syntax is correct or not.
- It performs type conversion of all the data types into real data types.

## INTERMEDIATE CODE GENERATION:

- It is the fourth phase of the compiler.
- It gets input from the semantic analysis and converts the input into output as intermediate code such as three-address code.
- The three-address code consists of a sequence of instructions, each of which has at most three operands.

Example:  $t1=t2+t3$

## CODE OPTIMIZATION:

- It is the fifth phase of the compiler.
- It gets the intermediate code as input and produces optimized intermediate code as output.
- This phase reduces the redundant code and attempts to improve the intermediate code so that faster-running machine code will result.
- During the code optimization, the result of the program is not affected.

To improve the code generation, the optimization involves

- deduction and removal of dead code (unreachable code).
- calculation of constants in expressions and terms.
- collapsing of repeated expression into temporary string.
- loop unrolling.
- moving code outside the loop.
- removal of unwanted temporary variables.

## CODE GENERATION:

- It is the final phase of the compiler.
- It gets input from code optimization phase and produces the target code or object code as result.
- Intermediate instructions are translated into a sequence of machine instructions that perform the same task.
- The code generation involves
  - allocation of register and memory
  - generation of correct references
  - generation of correct data types
  - generation of missing code

## SYMBOL TABLE MANAGEMENT:

- Symbol table is used to store all the information about identifiers used in the program.
- It is a data structure containing a record for each identifier, with fields for the attributes of the identifier.
- It allows to find the record for each identifier quickly and to store or retrieve data from that record.
- Whenever an identifier is detected in any of the phases, it is stored in the symbol table.

## ERROR HANDLING:

- Each phase can encounter errors. After detecting an error, a phase must handle the error so that compilation can proceed.
- In lexical analysis, errors occur in separation of tokens.
- In syntax analysis, errors occur during construction of syntax tree.
- In semantic analysis, errors occur when the compiler detects constructs with right syntactic structure but no meaning and during type conversion.
- In code optimization, errors occur when the result is affected by the optimization.  In code generation, it shows error when code is missing etc.

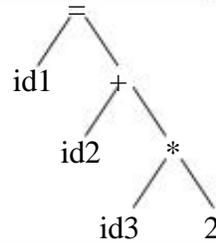
To illustrate the translation of source code through each phase, consider the statement  $a=b+c*2$ . The figure shows the representation of this statement after each phase:

a=b+c\*2

Lexical analyser

id1=id2+id3\*2

Syntax analyser

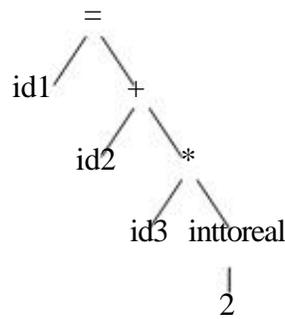


Symbol Table

a	id1
b	id2
c	id3

Semantic analyser

=



Intermediate code generator

temp1=inttoreal(2)  
temp2=id3\*temp1  
temp3=id2+temp2  
id1=temp3

Code optimizer

temp1=id3\*2.0  
id1=id2+temp1

Code generator

MOVF id3,R2  
MULF #2.0,R2  
MOVF id2,R1  
ADDF R2,R1  
MOVF R1,id1

## COUSINS OF COMPILER

1. Preprocessor
2. Assembler
3. Loader and Link-editor

### PREPROCESSOR

A preprocessor is a program that processes its input data to produce output that is used as input to another program. The output is said to be a preprocessed form of the input data, which is often used by some subsequent programs like compilers.

They may perform the following functions :

1. Macro processing
2. File Inclusion
3. Rational Preprocessors
4. Language extension

#### 1. Macro processing:

A macro is a rule or pattern that specifies how a certain input sequence should be mapped to an output sequence according to a defined procedure. The mapping process that instantiates a macro into a specific output sequence is known as macro expansion.

#### 2. File Inclusion:

Preprocessor includes header files into the program text. When the preprocessor finds an #include directive it replaces it by the entire content of the specified file.

#### 3. Rational Preprocessors:

These processors change older languages with more modern flow-of-control and datastructuring facilities.

#### 4. Language extension :

These processors attempt to add capabilities to the language by what amounts to built-in macros. For example, the language Equel is a database query language embedded in C.

### ASSEMBLER

Assembler creates object code by translating assembly instruction mnemonics into machine code. There are two types of assemblers:

- One-pass assemblers go through the source code once and assume that all symbols will be defined before any instruction that references them.
- Two-pass assemblers create a table with all symbols and their values in the first pass, and then use the table in a second pass to generate code.

### LINKER AND LOADER

A linker or link editor is a program that takes one or more objects generated by a compiler and combines them into a single executable program.

Three tasks of the linker are :

1. Searches the program to find library routines used by program, e.g. printf(), math routines.
2. Determines the memory locations that code from each module will occupy and relocates its instructions by adjusting absolute references
3. Resolves references among files.

A loader is the part of an operating system that is responsible for loading programs in memory, one of the essential stages in the process of starting a program.

## GROUPING OF THE PHASES

Compiler can be grouped into front and back ends:

- Front end: analysis (machine independent)

These normally include lexical and syntactic analysis, the creation of the symbol table, semantic analysis and the generation of intermediate code. It also includes error handling that goes along with each of these phases.

- Back end: synthesis (machine dependent)

It includes code optimization phase and code generation along with the necessary error handling and symbol table operations.

Compiler passes

A collection of phases is done only once (single pass) or multiple times (multi pass)

- Single pass: usually requires everything to be defined before being used in source program.
- Multi pass: compiler may have to keep entire program representation in memory.

Several phases can be grouped into one single pass and the activities of these phases are interleaved during the pass. For example, lexical analysis, syntax analysis, semantic analysis and intermediate code generation might be grouped into one pass.

COMPILER CONSTRUCTION TOOLS These are specialized tools that have been developed for helping implement various phases of a compiler. The following are the compiler construction tools:

1) Parser Generators:

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

-It consumes a large fraction of the running time of a compiler. -

Example-YACC (Yet Another Compiler-Compiler).

2) Scanner Generator:

-These generate lexical analyzers, normally from a specification based on regular expressions. -The basic organization of lexical analyzers is based on finite automation.

3) Syntax-Directed Translation:

-These produce routines that walk the parse tree and as a result generate intermediate code. - Each translation is defined in terms of translations at its neighbor nodes in the tree.

4) Automatic Code Generators:

-It takes a collection of rules to translate intermediate language into machine language. The rules must include sufficient details to handle different possible access methods for data.

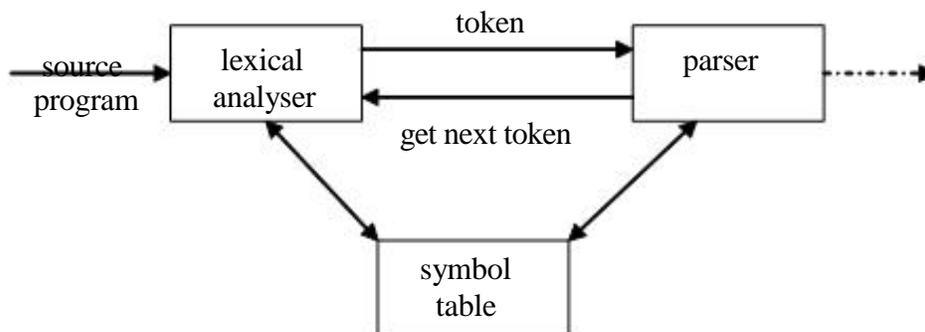
5) Data-Flow Engines:-It does code optimization using data-flow analysis, that is, the gathering of information about how values are transmitted from one part of a program to each other part.

## LEXICAL ANALYSIS

Lexical analysis is the process of converting a sequence of characters into a sequence of tokens. A program or function which performs lexical analysis is called a lexical analyzer or scanner. A lexer often exists as a single function which is called by a parser or another function.

### 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 as output a sequence of tokens that the parser uses for syntax analysis.



- Upon receiving a “get next token” command from the parser, the lexical analyzer reads input characters until it can identify the next token.

ISSUES OF LEXICAL ANALYZER There are three issues in lexical analysis:

- To make the design simpler.
- To improve the efficiency of the compiler.
- To enhance the computer portability.

### TOKENS

A token is a string of characters, categorized according to the rules as a symbol (e.g., IDENTIFIER, NUMBER, COMMA). The process of forming tokens from an input stream of characters is called tokenization.

A token can look like anything that is useful for processing an input text stream or text file. Consider this expression in the C programming language: `sum=3+2;`

Lexeme	Token type
sum	Identifier
=	Assignment operator
3	Number
+	Addition operator
2	Number

## LEXEME:

Collection or group of characters forming tokens is called Lexeme.

## PATTERN:

A pattern is a description of the form that the lexemes of a token may take.

In the case of a keyword as a token, the pattern is just the sequence of characters that form the keyword. For identifiers and some other tokens, the pattern is a more complex structure that is matched by many strings.

## Attributes for Tokens

Some tokens have attributes that can be passed back to the parser. The lexical analyzer collects information about tokens into their associated attributes. The attributes influence the translation of tokens.

- i) Constant : value of the constant
- ii) Identifiers: pointer to the corresponding symbol table entry.

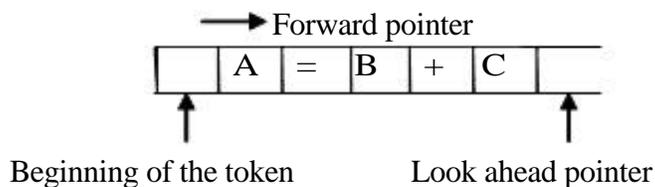
## ERROR RECOVERY STRATEGIES IN LEXICAL ANALYSIS:

The following are the error-recovery actions in lexical analysis:

- 1) Deleting an extraneous character.
- 2) Inserting a missing character.
- 3) Replacing an incorrect character by a correct character.
- 4) Transforming two adjacent characters.
- 5) Panic mode recovery: Deletion of successive characters from the token until error is resolved.

## INPUT BUFFERING

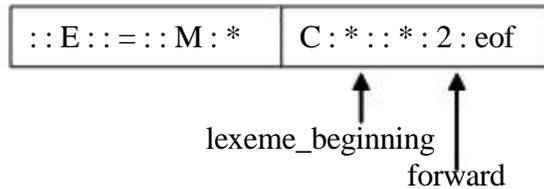
We often have to look one or more characters beyond the next lexeme before we can be sure we have the right lexeme. As characters are read from left to right, each character is stored in the buffer to form a meaningful token as shown below:



We introduce a two-buffer scheme that handles large look aheads safely. We then consider an improvement involving "sentinels" that saves time checking for the ends of buffers.

## BUFFER PAIRS

- A buffer is divided into two N-character halves, as shown below



- Each buffer is of the same size N, and N is usually the number of characters on one disk block. E.g., 1024 or 4096 bytes.
- Using one system read command we can read N characters into a buffer.
- If fewer than N characters remain in the input file, then a special character, represented by eof, marks the end of the source file.
- Two pointers to the input are maintained:
  1. Pointer lexeme\_beginning, marks the beginning of the current lexeme, whose extent we are attempting to determine.
  2. Pointer forward scans ahead until a pattern match is found.
    - Once the next lexeme is determined, forward is set to the character at its right end.
- The string of characters between the two pointers is the current lexeme. After the lexeme is recorded as an attribute value of a token returned to the parser, lexeme\_beginning is set to the character immediately after the lexeme just found.

Advancing forward pointer: Advancing forward pointer requires that we first test whether we have reached the end of one of the buffers, and if so, we must reload the other buffer from the input, and move forward to the beginning of the newly loaded buffer. If the end of second buffer is reached, we must again reload the first buffer with input and the pointer wraps to the beginning of the buffer.

Code to advance forward pointer:

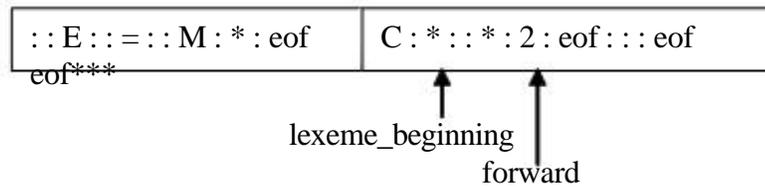
```

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 second half;
    move forward to beginning of first half
end
else forward := forward + 1;
```

## SENTINELS

- For each character read, we make two tests: one for the end of the buffer, and one to determine what character is read. We can combine the buffer-end test with the test for the current character if we extend each buffer to hold a sentinel character at the end.
- The sentinel is a special character that cannot be part of the source program, and a natural choice is the character eof.

- The sentinel arrangement is as shown below:



Note that eof retains its use as a marker for the end of the entire input. Any eof that appears other than at the end of a buffer means that the input is at an end.

Code to advance forward pointer:

```

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

```

## SPECIFICATION OF TOKENS

There are 3 specifications of tokens:

- 1) Strings
- 2) Language
- 3) Regular expression

### Strings and Languages

An alphabet or character class is a finite set of symbols.

A string over an alphabet is a finite sequence of symbols drawn from that alphabet. A language is any countable set of strings over some fixed alphabet.

In language theory, the terms "sentence" and "word" are often used as synonyms for

"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  $\epsilon$ , is the string of length zero.

### Operations on strings

The following string-related terms are commonly used:

1. A prefix of string  $s$  is any string obtained by removing zero or more symbols from the end of string  $s$ .

For example, ban is a prefix of banana.

2. A suffix of string  $s$  is any string obtained by removing zero or more symbols from the beginning of  $s$ .

For example, nana is a suffix of banana.

3. A substring of  $s$  is obtained by deleting any prefix and any suffix from  $s$ .

For example, nan is a substring of banana.

4. The proper prefixes, suffixes, and substrings of a string  $s$  are those prefixes, suffixes, and substrings, respectively of  $s$  that are not  $\epsilon$  or not equal to  $s$  itself.

5. A subsequence of  $s$  is any string formed by deleting zero or more not necessarily consecutive positions of  $s$ .

For example, baan is a subsequence of banana.

### Operations on languages:

The following are the operations that can be applied to languages:

1. Union
2. Concatenation
3. Kleene closure
4. Positive closure

The following example shows the operations on strings:

Let  $L=\{0,1\}$  and  $S=\{a,b,c\}$

1. Union :  $L \cup S=\{0,1,a,b,c\}$
2. Concatenation :  $L.S=\{0a,1a,0b,1b,0c,1c\}$
3. Kleene closure :  $L^*=\{ \epsilon,0,1,00,\dots\}$
4. Positive closure :  $L^+=\{0,1,00,\dots\}$

Regular Expressions Each regular expression  $r$  denotes a language  $L(r)$ .

Here are the rules that define the regular expressions over some alphabet  $\Sigma$  and the languages that those expressions denote:

1.  $\epsilon$  is a regular expression, and  $L(\epsilon)$  is  $\{ \epsilon \}$ , that is, the language whose sole member is the empty string.
2. If 'a' is a symbol in  $\Sigma$ , then 'a' is a regular expression, and  $L(a) = \{a\}$ , that is, the language with one string, of length one, with 'a' in its one position.
3. Suppose  $r$  and  $s$  are regular expressions denoting the languages  $L(r)$  and  $L(s)$ . Then,
  - a)  $(r)|(s)$  is a regular expression denoting the language  $L(r) \cup L(s)$ .
  - b)  $(r)(s)$  is a regular expression denoting the language  $L(r)L(s)$ .
  - c)  $(r)^*$  is a regular expression denoting  $(L(r))^*$ .
  - d)  $(r)$  is a regular expression denoting  $L(r)$ .
4. The unary operator  $*$  has highest precedence and is left associative.
5. Concatenation has second highest precedence and is left associative.
6.  $|$  has lowest precedence and is left associative.

Regular set

A language that can be defined by a regular expression is called a regular set. If two regular expressions  $r$  and  $s$  denote the same regular set, we say they are equivalent and write  $r = s$ .

There are a number of algebraic laws for regular expressions that can be used to manipulate into equivalent forms. For instance,  $r|s = s|r$  is commutative;  $r|(s|t)=(r|s)|t$  is associative.

Regular Definitions

Giving names to regular expressions is referred to as a Regular definition. If  $\Sigma$  is an alphabet of basic symbols, then a regular definition is a sequence of definitions of the form

$d_1 \rightarrow r_1$   
 $d_2 \rightarrow r_2$   
 .....  
 $d_n \rightarrow r_n$

1. Each  $d_i$  is a distinct name.
2. Each  $r_i$  is a regular expression over the alphabet  $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$ .

Example: Identifiers is the set of strings of letters and digits beginning with a letter. Regular definition for this set:

letter  $\rightarrow A | B | \dots | Z | a | b | \dots | z |$

digit  $\rightarrow 0 | 1 | \dots | 9$

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

### 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  $^*$ .

2. Zero or one instance (?):

- The unary postfix operator ? means

“zero or one instance of”. - The notation

$r?$  is a shorthand for  $r | \epsilon$ .

- If ‘ $r$ ’ is a regular expression, then  $(r)?$  is a regular expression that denotes the language  $L(r) \cup \{ \epsilon \}$ .

3. Character Classes:

- The notation  $[abc]$  where  $a$ ,  $b$  and  $c$  are alphabet symbols denotes the regular expression  $a | b | c$ .

- Character class such as  $[a - z]$  denotes the regular expression  $a | b | c | d | \dots | z$ .
- We can describe identifiers as being strings generated by the regular expression,  $[A-Za-z][A-Za-z0-9]^*$

Non-regular Set A language which cannot be described by any regular expression is a non-regular set.

Example: The set of all strings of balanced parentheses and repeating strings cannot be described by a regular expression. This set can be specified by a context-free grammar.

## RECOGNITION OF TOKENS

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+)?
```

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.

## Transition diagram

A transition diagram is similar to a flowchart for (a part of) the lexer. We draw one for each possible token. It shows the decisions that must be made based on the input seen. The two main components are circles representing states (think of them as decision

points of the lexer) and arrows representing edges (think of them as the decisions made).

The transition diagram for relop is shown below.

1. The double circles represent accepting or final states at which point a lexeme has been found. There is often an action to be done (e.g., returning the token), which is written to the right of the double circle.
2. If we have moved one (or more) characters too far in finding the token, one (or more) stars are drawn.
3. An imaginary start state exists and has an arrow coming from it to indicate where to begin the process.

It is fairly clear how to write code corresponding to this diagram. You look at the first character, if it is <, you look at the next character. If that character is =, you return (relop,LE) to the parser. If instead that character is >, you return (relop,NE). If it is another character, return (relop,LT) and adjust the input buffer so that you will read this character again since you have not used it for the current lexeme. If the first character was =, you return (relop,EQ).

### Transition diagram for relational operators

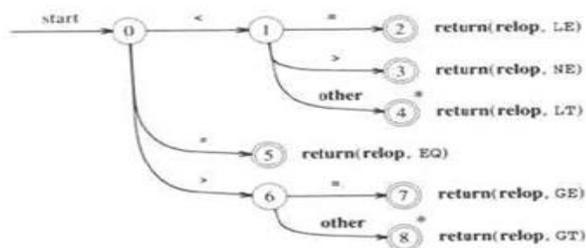
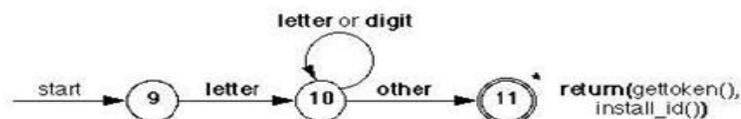


Fig. 3.12. Transition diagram for relational operators.

### Transition diagram for identifiers and keywords



## Recognition of Reserved Words and Identifiers

Recognizing keywords and identifiers presents a problem. Usually, keywords like `if` or `then` are reserved, so they are not identifiers even though they look like identifiers.

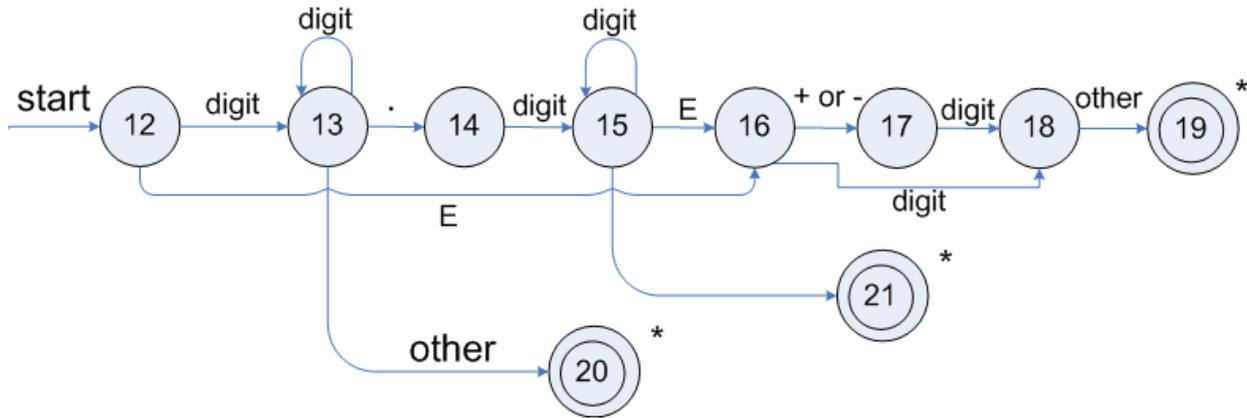
There are two ways that we can handle reserved words that look like identifiers:

1. Install the reserved words in the symbol table initially. A field of the symbol-table entry indicates that these strings are never ordinary identifiers, and tells which token they represent. When we find an identifier, a call to `installID` places it in the symbol table if it is not already there and returns a pointer to the symbol-table entry for the lexeme found. Of course, any identifier not in the symbol table during lexical analysis cannot be a reserved word, so its token is `id`. The function `getToken` examines the symbol table entry for the lexeme found, and returns whatever token name the symbol table says.
2. Create separate transition diagrams for each keyword. Note that such a transition diagram consists of states representing the situation after each successive letter of the keyword is seen.

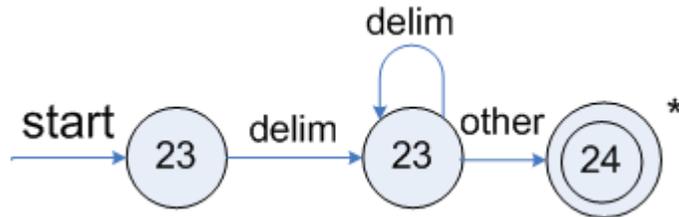
## Recognition of Numbers

The transition diagram for token `number` is shown in Fig. and is so far the most complex diagram we have seen. Beginning in state 12, if we see a digit, we go to state 13. In that state, we can read any number of additional digits. However, if we see anything but a digit or a dot, we have seen a number in the form of an integer; 123 is an example. That case is handled by entering state 20, where we return token `number` and a pointer to a table of constants where the found lexeme is entered. These mechanics are not shown on the diagram but are analogous to the way we handled identifiers. If we instead see a dot in state 13, then we have an "optional fraction." State 14 is entered, and we look for one or more additional digits; state 15 is used for that purpose. If we see an E, then we have an "optional exponent," whose recognition is the job of states 16 through 19. Should we, in state 15, instead see anything but E or a digit, then we have come to the end of the fraction, there is no exponent,

and we return the lexeme found , via state 21.



The transition diagram , shown in Fig , is for whitespace . In that diagram , we look for one or more "whitespace" characters , represented by delim in that diagram — typically these characters would be blank , tab , newline.



### Architecture of a Transition -Diagram -Based Lexical Analyzer

There are several ways that a collection of transition diagrams can be used to build a lexical analyzer . Regardless of the overall strategy , each state is represented by a piece of code . We may imagine a variable `state` holding the number of the current state for a transition diagram . A switch based on the value of `state` takes us to code for each of the possible states , where we find the action of that state . Often , the code for a state is itself a switch statement or multiway branch that determines the next state by reading and examining the next input character .

Example 3.10: In Fig . 3.18 we see a sketch of `getRe1opO` , a `C++` function whose job is to simulate the transition

diagram of Fig. 3.13 and return an object of type TOKEN, that is, a pair consisting of the token name (which must be `relop` in this case) and an attribute value (the code for one of the six comparison operators in this case). `getRelop` first creates a new object `retToken` and initializes its first component to `RELOP`, the symbolic code for token `relop`.

We see the typical behavior of a state in case 0, the case where the current state is 0. A function `nextChar()` obtains the next character from the input and assigns it to local variable `c`. We then check `c` for the three characters we expect to find, making the state transition dictated by the transition diagram of Fig. 3.13 in each case. For example, if the next input character is `=`, we go to state 5.

If the next input character is not one that can begin a comparison operator, then a function `fail()` is called. What `fail()` does depends on the global error-recovery strategy of the lexical analyzer. It should reset the forward pointer to `lexemeBegin`, in order to allow another transition diagram to be applied to the true beginning of the unprocessed input. It might then change the value of state to be the start state for another transition diagram, which will search for another token. Alternatively, if there is no other transition diagram that remains unused, `fail()` could initiate an error-correction phase that will try to repair the input and find a lexeme, as discussed in Section 3.1.4.

We also show the action for state 8 in Fig. 3.18. Because state 8 bears a \*, we must retract the input pointer one position (i.e., put `c` back on the input stream). That task is accomplished by the function `retract()`. Since state 8 represents the recognition of lexeme `>=`, we set the second component of the returned object, which we suppose is named `attribute`, to `GT`, the code for this operator.

To place the simulation of one transition diagram in perspective, let us consider the way code like Fig. 3.18 could fit into the entire lexical analyzer.

1. We could arrange for the transition diagrams for each token to be tried sequentially. Then, the function `fail()` of Example 3.10 resets the pointer forward and starts the next transition diagram, each time it is called.

This method allows us to use transition diagrams for the individual

key

words, like the one suggested in Fig. 3.15. We have only to use these before we use the diagram for `id`, in order for the keywords to be reserved words.

2. We could run the various transition diagrams "in parallel," feeding the next input character to all of them and allowing each one to make what

ever transitions it required. If we use this strategy, we must be careful

to resolve the case where one diagram finds a lexeme that matches its pattern, while one or more other diagrams are still able to process input.

The normal strategy is to take the longest prefix of the input that matches any pattern. That rule allows us to prefer identifier then next to keyword

then, or the operator `->` to `-`, for example.

3. The preferred approach, and the one we shall take up in the following

sections, is to combine all the transition diagrams into one. We allow the

transition diagram to read input until there is no possible next state, and

then take the longest lexeme that matched any pattern, as we discussed

in item (2) above. In our running example, this combination is easy,

because no two tokens can start with the same character; i.e., the first

character immediately tells us which token we are looking for. Thus, we

could simply combine states 0, 9, 12, and 22 into one start state, leaving other transitions intact. However, in general, the problem

of combining transition diagrams for several tokens is more complex, as we shall

see

shortly.

```

TOKEN  getRelopO
{
    TOKEN  retToken = new (RELOP ) ;
    while ( 1 ) { /* repeat character processing until
a return
                or failure occurs*/
        switch (state ) {
            case 0: c = nextCharQ ;
                    if ( c == '<>' ) state = 1;
                    else if      ( c == '=' )
state = 5;
                    else if      ( c == '>' )
state = 6;
                    else fail ( ) ; /* lexeme is not a
relop */
                break ;
            case 1:
                case 8 : retract ( ) ;
                        retToken .attribute = GT ;
                        return (retToken ) ;

```

**Figure 3.18: Sketch of implementation of relop transition diagram**

## FINITE AUTOMATA

Finite Automata is one of the mathematical models that consist of a number of states and edges. It is a transition diagram that recognizes a regular expression or grammar.

### Types of Finite Automata

There are two types of Finite Automata :

- Non-deterministic Finite Automata (NFA)
- Deterministic Finite Automata (DFA)

### Non-deterministic Finite Automata

NFA is a mathematical model that consists of five tuples denoted by  $M = \{Q_n, \Sigma, \delta, q_0, f_n\}$

$Q_n$  - finite set of states

$\Sigma$  - finite set of input symbols

$\delta$  - transition function that maps state-symbol pairs to set of states

$q_0$  - starting state

$f_n$  - final state

### Deterministic Finite Automata

DFA is a special case of a NFA in which

- i) no state has an  $\epsilon$ -transition.
- ii) there is at most one transition from each state on any input.

DFA has five tuples denoted by

$$M = \{Q_d, \Sigma, \delta, q_0, f_d\}$$

$Q_d$  - finite set of states

$\Sigma$  - finite set of input symbols

$\delta$  - transition function that maps state-symbol pairs to set of states

$q_0$  - starting state

$f_d$  - final state

Construction of DFA from regular expression

The following steps are involved in the construction of DFA from regular expression:

- i) Convert RE to NFA using Thomson's rules
- ii) Convert NFA to DFA
- iii) Construct minimized DFA

## SYNTAX ANALYSIS

### SYNTAX ANALYSIS

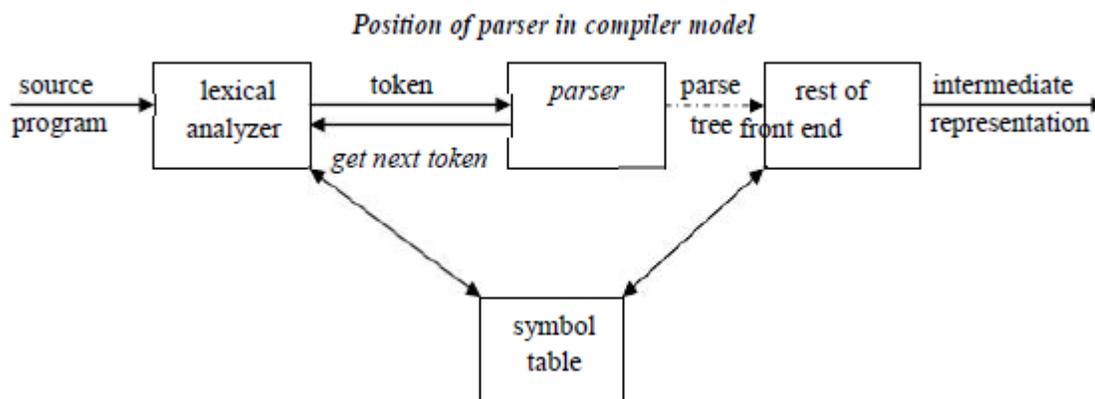
Syntax analysis is the second phase of the compiler. It gets the input from the tokens and generates a syntax tree or parse tree.

#### Advantages of grammar for syntactic specification :

1. A grammar gives a precise and easy-to-understand syntactic specification of a programming language.
2. An efficient parser can be constructed automatically from a properly designed grammar.
3. A grammar imparts a structure to a source program that is useful for its translation into object code and for the detection of errors.
4. New constructs can be added to a language more easily when there is a grammatical description of the language.

#### THE ROLE OF PARSER

The parser or syntactic analyzer obtains a string of tokens from the lexical analyzer and verifies that the string can be generated by the grammar for the source language. It reports any syntax errors in the program. It also recovers from commonly occurring errors so that it can continue processing its input.



### **Functions of the parser :**

1. It verifies the structure generated by the tokens based on the grammar.
2. It constructs the parse tree.
3. It reports the errors.
4. It performs error recovery.

### **Issues :**

Parser cannot detect errors such as:

1. Variable re-declaration
2. Variable initialization before use.
3. Data type mismatch for an operation.

The above issues are handled by Semantic Analysis phase.

### **Syntax error handling :**

Programs can contain errors at many different levels. For example :

1. Lexical, such as misspelling a keyword.
2. Syntactic, such as an arithmetic expression with unbalanced parentheses.
3. Semantic, such as an operator applied to an incompatible operand.
4. Logical, such as an infinitely recursive call.

Functions of error handler :

1. It should report the presence of errors clearly and accurately.
2. It should recover from each error quickly enough to be able to detect subsequent errors.
3. It should not significantly slow down the processing of correct programs.

### **Error recovery strategies :**

The different strategies that a parser uses to recover from a syntactic error are:

1. Panic mode
2. Phrase level
3. Error productions
4. Global correction

### **Panic mode recovery:**

On discovering an error, the parser discards input symbols one at a time until a synchronizing token is found. The synchronizing tokens are usually delimiters, such as semicolon or **end**. It has the advantage of simplicity and does not go into an infinite loop. When multiple errors in the same statement are rare, this method is quite useful.

### **Phrase level recovery:**

On discovering an error, the parser performs local correction on the remaining input that allows it to continue. Example: Insert a missing semicolon or delete an extraneous semicolon etc.

**Error productions:**

The parser is constructed using augmented grammar with error productions. If an error production is used by the parser, appropriate error diagnostics can be generated to indicate the erroneous constructs recognized by the input.

**Global correction:**

Given an incorrect input string  $x$  and grammar  $G$ , certain algorithms can be used to find a parse tree for a string  $y$ , such that the number of insertions, deletions and changes of tokens is as small as possible. However, these methods are in general too costly in terms of time and space.

## CONTEXT-FREE GRAMMARS

A Context-Free Grammar is a quadruple that consists of **terminals**, **non-terminals**, **start symbol** and **productions**.

**Terminals :** These are the basic symbols from which strings are formed.

**Non-Terminals :** These are the syntactic variables that denote a set of strings. These help to define the language generated by the grammar.

**Start Symbol :** One non-terminal in the grammar is denoted as the “Start-symbol” and the set of strings it denotes is the language defined by the grammar.

**Productions :** It specifies the manner in which terminals and non-terminals can be combined to form strings. Each production consists of a non-terminal, followed by an arrow, followed by a string of non-terminals and terminals.

**Example of context-free grammar:** The following grammar defines simple arithmetic expressions:

$$expr \rightarrow expr \ op \ expr$$
$$expr \rightarrow (expr)$$
$$expr \rightarrow - \ expr$$
$$expr \rightarrow \mathbf{id}$$
$$op \rightarrow +$$
$$op \rightarrow -$$
$$op \rightarrow *$$
$$op \rightarrow /$$
$$op \rightarrow \uparrow$$

In this grammar,

- **id** + - \* /  $\uparrow$  ( ) are terminals.
- *expr* , *op* are non-terminals.
- *expr* is the start symbol.
- Each line is a production.

### Derivations:

Two basic requirements for a grammar are :

1. To generate a valid string.
2. To recognize a valid string.

**Derivation** is a process that generates a valid string with the help of grammar by replacing the non-terminals on the left with the string on the right side of the production.

**Example :** Consider the following grammar for arithmetic expressions :

$$E \rightarrow E+E \mid E * E \mid ( E ) \mid - E \mid \mathbf{id}$$

To generate a valid string - ( id+id ) from the grammar the steps are

1.  $E \rightarrow - E$
2.  $E \rightarrow - ( E )$
3.  $E \rightarrow - ( E+E )$
4.  $E \rightarrow - ( id+E )$
5.  $E \rightarrow - ( id+id )$

In the above derivation,

- E is the start symbol.
- - (id+id) is the required sentence (only terminals).
- Strings such as E, -E, -(E), . . . are called sentinel forms.

**Types of derivations:**

The two types of derivation are:

1. Left most derivation
2. Right most derivation.

- In leftmost derivations, the leftmost non-terminal in each sentinel is always chosen first for replacement.
- In rightmost derivations, the rightmost non-terminal in each sentinel is always chosen first for replacement.

**Example:**

Given grammar  $G : E \rightarrow E+E \mid E * E \mid ( E ) \mid - E \mid id$

Sentence to be derived : - (id+id)

LEFTMOST DERIVATION

- $E \rightarrow - E$
- $E \rightarrow - ( E )$
- $E \rightarrow - ( E+E )$
- $E \rightarrow - ( id+E )$
- $E \rightarrow - ( id+id )$

RIGHTMOST DERIVATION

- $E \rightarrow - E$
- $E \rightarrow - ( E )$
- $E \rightarrow - ( E+E )$
- $E \rightarrow - ( E+id )$
- $E \rightarrow - ( id+id )$

- String that appear in leftmost derivation are called **left sentinel forms**.
- String that appear in rightmost derivation are called **right sentinel forms**.

**Sentinels:**

Given a grammar G with start symbol S, if  $S \rightarrow \alpha$ , where  $\alpha$  may contain non-terminals or terminals, then  $\alpha$  is called the sentinel form of G.

### Yield or frontier of tree:

Each interior node of a parse tree is a non-terminal. The children of node can be a terminal or non-terminal of the sentinel forms that are read from left to right. The sentinel form in the parse tree is called **yield** or **frontier** of the tree.

### Ambiguity:

A grammar that produces more than one parse for some sentence is said to be **ambiguous grammar**.

Example : Given grammar  $G : E \rightarrow E+E \mid E * E \mid ( E ) \mid - E \mid id$

The sentence  $id+id*id$  has the following two distinct leftmost derivations:

$$E \rightarrow E + E$$

$$E \rightarrow id + E$$

$$E \rightarrow id + E * \underline{E}$$

$$E \rightarrow id + id * \underline{E}$$

$$E \rightarrow id + id * \underline{id}$$

$$E \rightarrow E * E$$

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

$$\underline{E} \rightarrow id + E * \underline{E}$$

$$\underline{E} \rightarrow id + id * \underline{E}$$

$$\underline{E} \rightarrow id + id * \underline{id}$$

The two corresponding parse trees are :



### WRITING A GRAMMAR

There are four categories in writing a grammar :

1. Regular Expression Vs Context Free Grammar
2. Eliminating ambiguous grammar.
3. Eliminating left-recursion
4. Left-factoring.

Each parsing method can handle grammars only of a certain form hence, the initial grammar may have to be rewritten to make it parsable.

## Regular Expressions vs. Context-Free Grammars:

REGULAR EXPRESSION	CONTEXT-FREE GRAMMAR
It is used to describe the tokens of programming languages.	It consists of a quadruple where $S \rightarrow$ start symbol, $P \rightarrow$ production, $T \rightarrow$ terminal, $V \rightarrow$ variable or non-terminal.
It is used to check whether the given input is valid or not using <b>transition diagram</b> .	It is used to check whether the given input is valid or not using <b>derivation</b> .
The transition diagram has set of states and edges.	The context-free grammar has set of productions.
It has no start symbol.	It has start symbol.
It is useful for describing the structure of lexical constructs such as identifiers, constants, keywords, and so <u>forth</u> .	It is useful in describing nested structures such as balanced parentheses, matching <u>begin-end's</u> and so <u>on</u> .

- The lexical rules of a language are simple and RE is used to describe them.
- Regular expressions provide a more concise and easier to understand notation for tokens than grammars.
- Efficient lexical analyzers can be constructed automatically from RE than from grammars.
- Separating the syntactic structure of a language into lexical and nonlexical parts provides a convenient way of modularizing the front end into two manageable-sized components.

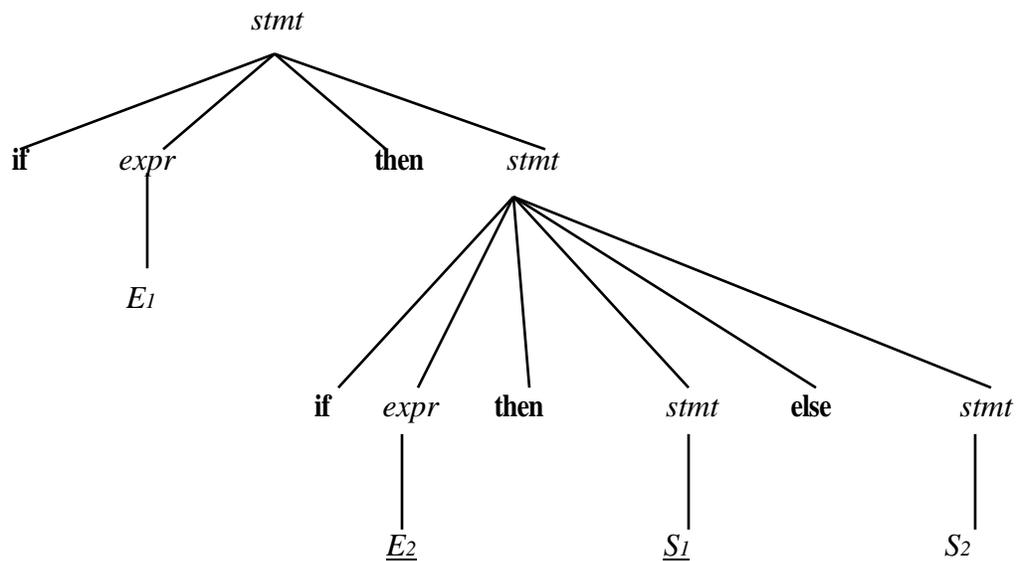
### Eliminating ambiguity:

Ambiguity of the grammar that produces more than one parse tree for leftmost or rightmost derivation can be eliminated by re-writing the grammar.

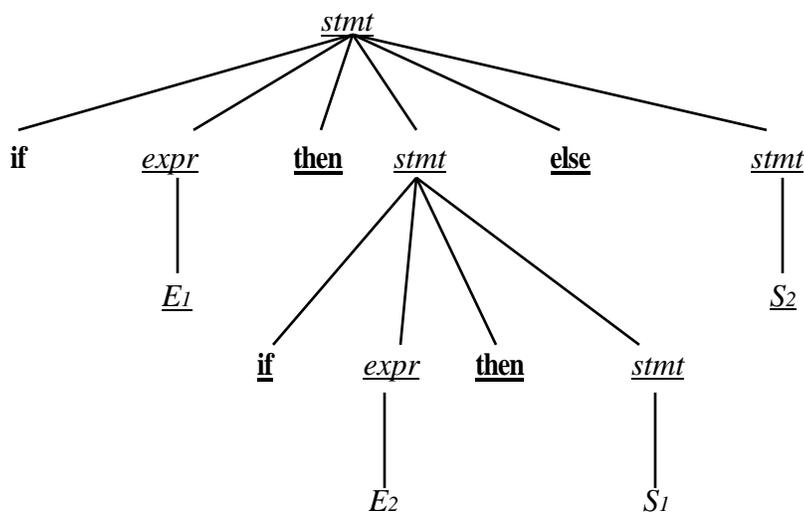
Consider this example,  $G: stmt \rightarrow \text{if } expr \text{ then } stmt \mid \text{if } expr \text{ then } stmt \text{ else } stmt \mid \text{other}$

This grammar is ambiguous since the string **if E<sub>1</sub> then if E<sub>2</sub> then S<sub>1</sub> else S<sub>2</sub>** has the following two parse trees for leftmost derivation :

1.



2.



To eliminate ambiguity, the following grammar may be used:

$$stmt \rightarrow matched\_stmt \mid unmatched\_stmt$$

$$matched\_stmt \rightarrow \text{if } expr \text{ then } matched\_stmt \text{ else } matched\_stmt \mid \text{other}$$

$$unmatched\_stmt \rightarrow \text{if } expr \text{ then } stmt \mid \text{if } expr \text{ then } matched\_stmt \text{ else } unmatched\_stmt$$

### Eliminating Left Recursion:

A grammar is said to be *left recursive* if it has a non-terminal  $A$  such that there is a derivation  $A \Rightarrow A\alpha$  for some string  $\alpha$ . Top-down parsing methods cannot handle left-recursive grammars. Hence, left recursion can be eliminated as follows:

If there is a production  $A \rightarrow A\alpha \mid \beta$  it can be replaced with a sequence of two productions

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

without changing the set of strings derivable from  $A$ .

**Example :** Consider the following grammar for arithmetic expressions:

$$E \rightarrow E+T \mid T$$

$$T \rightarrow T*F \mid F$$

$$F \rightarrow (E) \mid \text{id}$$

First eliminate the left recursion for  $E$  as

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

Then eliminate for  $T$  as

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

Thus the obtained grammar after eliminating left recursion is

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \epsilon$$

$$F \rightarrow (E) \mid \text{id}$$

**Algorithm to eliminate left recursion:**

1. Arrange the non-terminals in some order  $A_1, A_2 \dots A_n$ .
2. **for**  $i := 1$  **to**  $n$  **do begin**
  - for**  $j := 1$  **to**  $i-1$  **do begin**
    - replace each production of the form  $A_i \rightarrow A_j \gamma$   
by the productions  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$   
where  $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all the current  $A_j$ -productions;
  - end**
  - eliminate the immediate left recursion among the  $A_i$ -productions
- end**

### Left factoring:

Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive parsing. When it is not clear which of two alternative productions to use to expand a non-terminal A, we can rewrite the A-productions to defer the decision until we have seen enough of the input to make the right choice.

If there is any production  $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ , it can be rewritten as

$$A \rightarrow \alpha A'$$

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

Consider the grammar,  $G : S \rightarrow iEtS \mid iEtSeS \mid a$   
 $E \rightarrow b$

Left factored, this grammar becomes

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \varepsilon$$

$$E \rightarrow b$$

### PARSING

It is the process of analyzing a continuous stream of input in order to determine its grammatical structure with respect to a given formal grammar.

#### Parse tree:

Graphical representation of a derivation or deduction is called a parse tree. Each interior node of the parse tree is a non-terminal; the children of the node can be terminals or non-terminals.

#### Types of parsing:

1. Top down parsing
  2. Bottom up parsing
- Top-down parsing : A parser can start with the start symbol and try to transform it to the input string.  
Example : LL Parsers.
  - Bottom-up parsing : A parser can start with input and attempt to rewrite it into the start symbol.  
Example : LR Parsers.

### TOP-DOWN PARSING

It can be viewed as an attempt to find a left-most derivation for an input string or an attempt to construct a parse tree for the input starting from the root to the leaves.

## Types of top-down parsing :

1. Recursive descent parsing
2. Predictive parsing

### 1. RECURSIVE DESCENT PARSING

- Recursive descent parsing is one of the top-down parsing techniques that uses a set of recursive procedures to scan its input.
- This parsing method may involve **backtracking**, that is, making repeated scans of the input.

#### Example for backtracking :

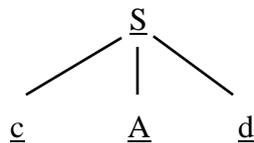
Consider the grammar  $G : S \rightarrow cAd$   
 $A \rightarrow ab \mid a$

and the input string  $w=cad$ .

The parse tree can be constructed using the following top-down approach :

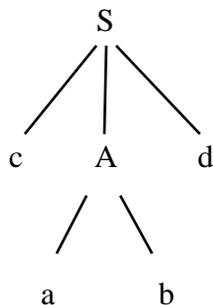
#### Step1:

Initially create a tree with single node labeled S. An input pointer points to 'c', the first symbol of w. Expand the tree with the production of S.



#### Step2:

The leftmost leaf 'c' matches the first symbol of w, so advance the input pointer to the second symbol of w 'a' and consider the next leaf 'A'. Expand A using the first alternative.



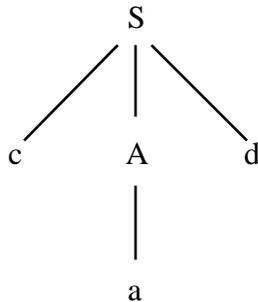
#### Step3:

The second symbol 'a' of w also matches with second leaf of tree. So advance the input pointer to third symbol of w 'd'. But the third leaf of tree is b which does not match with the input symbol **d**.

Hence discard the chosen production and reset the pointer to second position. This is called backtracking.

**Step4:**

Now try the second alternative for A.

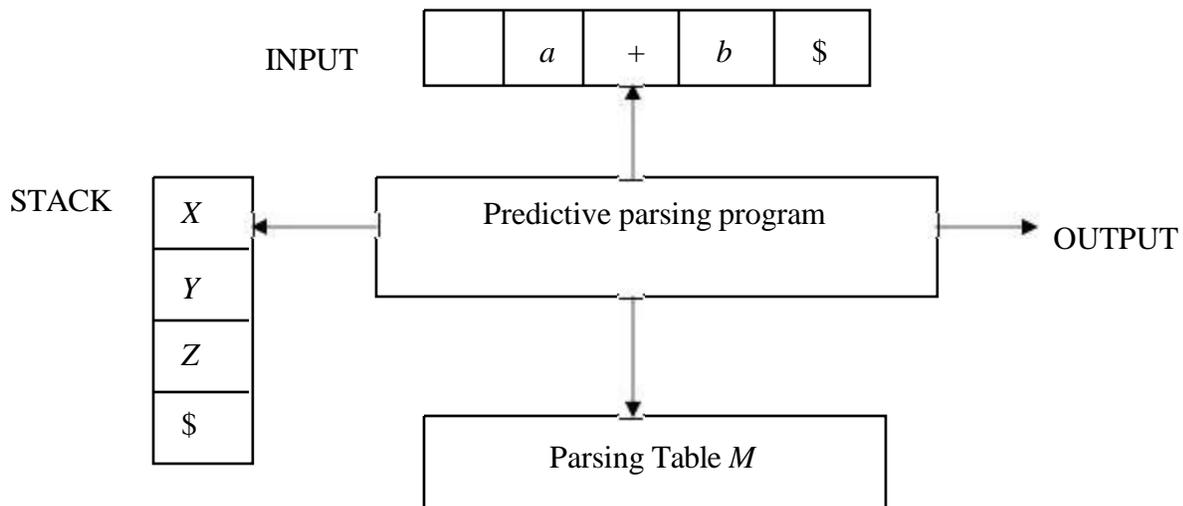


Now we can halt and announce the successful completion of parsing.

## 2. PREDICTIVE PARSING

- Predictive parsing is a special case of recursive descent parsing where no backtracking is required.
- The key problem of predictive parsing is to determine the production to be applied for a non-terminal in case of alternatives.

### Non-recursive predictive parser



The table-driven predictive parser has an input buffer, stack, a parsing table and an output stream.

**Input buffer:**

It consists of strings to be parsed, followed by \$ to indicate the end of the input string.

**Stack:**

It contains a sequence of grammar symbols preceded by \$ to indicate the bottom of the stack. Initially, the stack contains the start symbol on top of \$.

**Parsing table:**

It is a two-dimensional array  $M[A, a]$ , where 'A' is a non-terminal and 'a' is a terminal.

**Predictive parsing program:**

The parser is controlled by a program that considers  $X$ , the symbol on top of stack, and  $a$ , the current input symbol. These two symbols determine the parser action. There are three possibilities:

1. If  $X = a = \$$ , the parser halts and announces successful completion of parsing.
2. If  $X = a \neq \$$ , the parser pops  $X$  off the stack and advances the input pointer to the next input symbol.
3. If  $X$  is a non-terminal, the program consults entry  $M[X, a]$  of the parsing table  $M$ . This entry will either be an  $X$ -production of the grammar or an error entry.  
If  $M[X, a] = \{X \rightarrow UVW\}$ , the parser replaces  $X$  on top of the stack by WVU.  
If  $M[X, a] = \text{error}$ , the parser calls an error recovery routine.

**Algorithm for nonrecursive predictive parsing:**

**Input :** A string  $w$  and a parsing table  $M$  for grammar  $G$ .

**Output :** If  $w$  is in  $L(G)$ , a leftmost derivation of  $w$ ; otherwise, an error indication.

**Method :** Initially, the parser has  $\$S$  on the stack with  $S$ , the start symbol of  $G$  on top, and  $w\$$  in the input buffer. The program that utilizes the predictive parsing table  $M$  to produce a parse for the input is as follows:

```
set  $ip$  to point to the first symbol of  $w\$$ ;  
repeat  
    let  $X$  be the top stack symbol and  $a$  the symbol pointed to by  $ip$ ;  
    if  $X$  is a terminal or $ then  
        if  $X = a$  then  
            pop  $X$  from the stack and advance  $ip$   
        else  $error()$   
    else /*  $X$  is a non-terminal */  
        if  $M[X, a] = X \rightarrow Y_1 Y_2 \dots Y_k$  then begin
```

```

        pop X from the stack;
        push  $Y_k, Y_{k-1}, \dots, Y_1$  onto the stack, with  $Y_1$  on top;
        output the production  $X \rightarrow Y_1 Y_2 \dots Y_k$ 
    end
    else error()
until X = $          /* stack is empty */

```

### **Predictive parsing table construction:**

The construction of a predictive parser is aided by two functions associated with a grammar  $G$  :

1. FIRST
2. FOLLOW

#### **Rules for first( ):**

1. If  $X$  is terminal, then  $\text{FIRST}(X)$  is  $\{\underline{X}\}$ .
2. If  $X \rightarrow \varepsilon$  is a production, then add  $\varepsilon$  to  $\text{FIRST}(X)$ .
3. If  $X$  is non-terminal and  $X \rightarrow a\alpha$  is a production then add  $a$  to  $\underline{\text{FIRST}(X)}$ .
4. If  $X$  is non-terminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then place  $a$  in  $\text{FIRST}(X)$  if for some  $i$ ,  $a$  is in  $\text{FIRST}(Y_i)$ , and  $\varepsilon$  is in all of  $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$ ; that is,  $Y_1, \dots, Y_{i-1} \Rightarrow \varepsilon$ . If  $\varepsilon$  is in  $\text{FIRST}(Y_j)$  for all  $j=1, 2, \dots, k$ , then add  $\varepsilon$  to  $\text{FIRST}(X)$ .

#### **Rules for follow( ):**

1. If  $S$  is a start symbol, then  $\text{FOLLOW}(S)$  contains  $\underline{\$}$ .
2. If there is a production  $A \rightarrow \alpha B \beta$ , then everything in  $\text{FIRST}(\beta)$  except  $\varepsilon$  is placed in  $\text{follow}(B)$ .
3. If there is a production  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha B \beta$  where  $\text{FIRST}(\beta)$  contains  $\varepsilon$ , then everything in  $\text{FOLLOW}(A)$  is in  $\text{FOLLOW}(B)$ .

### **Algorithm for construction of predictive parsing table:**

**Input :** Grammar  $G$

**Output :** Parsing table  $M$

**Method :**

1. For each production  $A \rightarrow \alpha$  of the grammar, do steps 2 and 3.
2. For each terminal  $a$  in  $\text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$ .
3. If  $\varepsilon$  is in  $\text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$  for each terminal  $b$  in  $\text{FOLLOW}(A)$ . If  $\varepsilon$  is in  $\text{FIRST}(\alpha)$  and  $\$$  is in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \$]$ .
4. Make each undefined entry of  $M$  be **error**.

**Example:**

Consider the following grammar :

$$E \rightarrow E+T \mid T$$

$$T \rightarrow T*F \mid F$$

$$F \rightarrow (E) \mid id$$

After eliminating left-recursion the grammar is

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \varepsilon$$

$$F \rightarrow (E) \mid id$$

**First() :**

$$FIRST(E) = \{ (, id \}$$

$$FIRST(E') = \{ +, \varepsilon \}$$

$$FIRST(T) = \{ (, id \}$$

$$FIRST(T') = \{ *, \varepsilon \}$$

$$FIRST(F) = \{ (, id \}$$

**Follow() :**

$$FOLLOW(E) = \{ \$, ) \}$$

$$FOLLOW(E') = \{ \$, ) \}$$

$$FOLLOW(T) = \{ +, \$, ) \}$$

$$FOLLOW(T') = \{ +, \$, ) \}$$

$$FOLLOW(F) = \{ +, *, \$, ) \}$$

**Predictive parsing table :**

NON-TERMINAL	id	+	*	(	)	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \varepsilon$	$E' \rightarrow \varepsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \varepsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \varepsilon$	$T' \rightarrow \varepsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$		

**Stack implementation:**

stack	Input	Output
\$E	id+id*id \$	
\$E'T	id+id*id \$	E → TE'
\$E'T'F	id+id*id \$	T → FT'
\$E'T'id	id+id*id \$	F → id
\$E'T'	+id*id \$	
\$E'	+id*id \$	T' → ε
\$E'T+	+id*id \$	E' → +TE'
\$E'T	id*id \$	
\$E'T'F	id*id \$	T → FT'
\$E'T'id	id*id \$	F → id
\$E'T'	*id \$	
\$E'T'F*	<u>*id</u> \$	<u>T'</u> → <u>*FT'</u>
\$E'T'F	<u>id</u> \$	
\$E'T'id	<u>id</u> \$	<u>F</u> → <u>id</u>
\$E'T'	\$	
\$E'	\$	<u>T'</u> → ε
\$	\$	<u>E'</u> → ε

**LL(1) grammar:**

The parsing table entries are single entries. So each location has not more than one entry. This type of grammar is called LL(1) grammar.

Consider this following grammar:

$$S \rightarrow iEtS \mid iEtSeS \mid a$$

$$E \rightarrow b$$

After eliminating left factoring, we have

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \varepsilon$$

$$E \rightarrow b$$

To construct a parsing table, we need FIRST() and FOLLOW() for all the non-terminals.

$$\text{FIRST}(S) = \{ i, a \}$$

$$\text{FIRST}(S') = \{ e, \varepsilon \}$$

$$\text{FIRST}(E) = \{ b \}$$

$$\text{FOLLOW}(S) = \{ \$, e \}$$

$FOLLOW(S') = \{ \$, e \}$

$FOLLOW(E) = \{ t \}$

**Parsing table:**

NON-TERMINAL	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow eS$ $S' \rightarrow \epsilon$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

Since there are more than one production, the grammar is not LL(1) grammar.

**Actions performed in predictive parsing:**

1. Shift
2. Reduce
3. Accept
4. Error

**Implementation of predictive parser:**

1. Elimination of left recursion, left factoring and ambiguous grammar.
2. Construct  $FIRST()$  and  $FOLLOW()$  for all non-terminals.
3. Construct predictive parsing table.
4. Parse the given input string using stack and parsing table.

## **BOTTOM-UP PARSING**

Constructing a parse tree for an input string beginning at the leaves and going towards the root is called bottom-up parsing.

A general type of bottom-up parser is a **shift-reduce parser**.

### **SHIFT-REDUCE PARSING**

Shift-reduce parsing is a type of bottom-up parsing that attempts to construct a parse tree for an input string beginning at the leaves (the bottom) and working up towards the root (the top).

**Example:**

Consider the grammar:

$S \rightarrow aABe$

$A \rightarrow Abc \mid b$

$B \rightarrow d$

The sentence to be recognized is **abbcd**.

## REDUCTION (LEFTMOST)

abcde (A → b)  
aAbcde (A → Abc)  
aAde (B → d)  
aABe (S → aABe)

S

The reductions trace out the right-most derivation in reverse.

## RIGHTMOST DERIVATION

S → aABe  
→ aAde  
→ aAbcde  
→ abcde

### Handles:

A handle of a string is a substring that matches the right side of a production, and whose reduction to the non-terminal on the left side of the production represents one step along the reverse of a rightmost derivation.

### **Example:**

Consider the grammar:

$E \rightarrow E+E$   
 $E \rightarrow E^*E$   
 $E \rightarrow (E)$   
 $E \rightarrow id$

And the input string  $id_1+id_2*id_3$

The rightmost derivation is :

$E \rightarrow \underline{E+E}$   
→  $E+\underline{E^*E}$   
→  $E+E^*\underline{id_3}$   
→  $E+\underline{id_2}*id_3$   
→  $\underline{id_1}+id_2*id_3$

In the above derivation the underlined substrings are called **handles**.

### Handle pruning:

A rightmost derivation in reverse can be obtained by “**handle pruning**”.

(i.e.) if  $w$  is a sentence or string of the grammar at hand, then  $w = \gamma_n$ , where  $\gamma_n$  is the  $n^{\text{th}}$  right-sentinel form of some rightmost derivation.

## Stack implementation of shift-reduce parsing :

	Stack	Input	Action
1.	\$	id <sub>1</sub> +id <sub>2</sub> *id <sub>3</sub> \$	shift
2.	\$ id <sub>1</sub>	+id <sub>2</sub> *id <sub>3</sub> \$	reduce by E→id
	\$E	+id <sub>2</sub> *id <sub>3</sub> \$	shift
	\$ E+	id <sub>2</sub> *id <sub>3</sub> \$	shift
	\$ E+id <sub>2</sub>	*id <sub>3</sub> \$	reduce by E→id
	\$ E+E	*id <sub>3</sub> \$	shift
	\$ E+E*	id <sub>3</sub> \$	shift
	\$ E+E*id <sub>3</sub>	\$	<u>reduce</u> by E→id
	\$ E+E*E	\$	<u>reduce</u> by E→ E *E
	\$ E+E	\$	<u>reduce</u> by E→ E+E
	\$E	\$	<u>accept</u>

### Actions in shift-reduce parser:

- shift – The next input symbol is shifted onto the top of the stack.
- reduce – The parser replaces the handle within a stack with a non-terminal.
- accept – The parser announces successful completion of parsing.
- error – The parser discovers that a syntax error has occurred and calls an error recovery routine.

### Conflicts in shift-reduce parsing:

There are two conflicts that occur in shift shift-reduce parsing:

1. **Shift-reduce conflict:** The parser cannot decide whether to shift or to reduce.
2. **Reduce-reduce conflict:** The parser cannot decide which of several reductions to make.

#### 1. Shift-reduce conflict:

**Example:**

Consider the grammar:

$E \rightarrow E+E \mid E * E \mid id$  and input  $id+id*id$

Stack	Input	Action	Stack	Input	Action
\$ E+E	*id \$	Reduce by $E \rightarrow E+E$	\$E+E	*id \$	Shift
\$E	*id \$	Shift	\$E+E*	id \$	Shift
\$ E*	id \$	Shift	\$E+E*id	\$	Reduce by $E \rightarrow id$
\$ E*id	\$	Reduce by $E \rightarrow id$	\$E+E*E	\$	Reduce by $E \rightarrow E*E$
\$ E*E	\$	Reduce by $E \rightarrow E*E$	\$E+E	\$	Reduce by $E \rightarrow E*E$
\$E			\$E		

## 2. Reduce-reduce conflict:

Consider the grammar:

$M \rightarrow R+R \mid R+c \mid \underline{R}$

$R \rightarrow c$

and input  $c+c$

Stack	Input	Action	Stack	Input	Action
\$	<u>c+c</u> \$	<u>Shift</u>	\$	<u>c+c</u> \$	Shift
\$c	<u>+c</u> \$	Reduce by $R \rightarrow c$	\$c	<u>+c</u> \$	Reduce by $R \rightarrow c$
\$R	+c \$	Shift	\$R	+c \$	Shift
\$ R+	c \$	Shift	\$ R+	c \$	Shift
\$ R+c	\$	Reduce by $R \rightarrow c$	\$ R+c	\$	Reduce by $M \rightarrow R+c$
\$ R+R	\$	Reduce by $M \rightarrow R+R$	\$M	\$	
\$M	\$				

### Viable prefixes:

- $\alpha$  is a viable prefix of the grammar if there is  $w$  such that  $\alpha w$  is a right sentinel form.
- The set of prefixes of right sentinel forms that can appear on the stack of a shift-reduce parser are called viable prefixes.
- The set of viable prefixes is a regular language.

### **LR PARSERS**

An efficient bottom-up syntax analysis technique that can be used to parse a large class of CFG is called LR( $k$ ) parsing. The 'L' is for left-to-right scanning of the input, the 'R' for constructing a rightmost derivation in reverse, and the ' $k$ ' for the number of input symbols. When ' $k$ ' is omitted, it is assumed to be 1.

### Advantages of LR parsing:

- ✓ It recognizes virtually all programming language constructs for which CFG can be written.
- ✓ It is an efficient non-backtracking shift-reduce parsing method.
- ✓ A grammar that can be parsed using LR method is a proper superset of a grammar that can be parsed with predictive parser.
- ✓ It detects a syntactic error as soon as possible.

### Drawbacks of LR method:

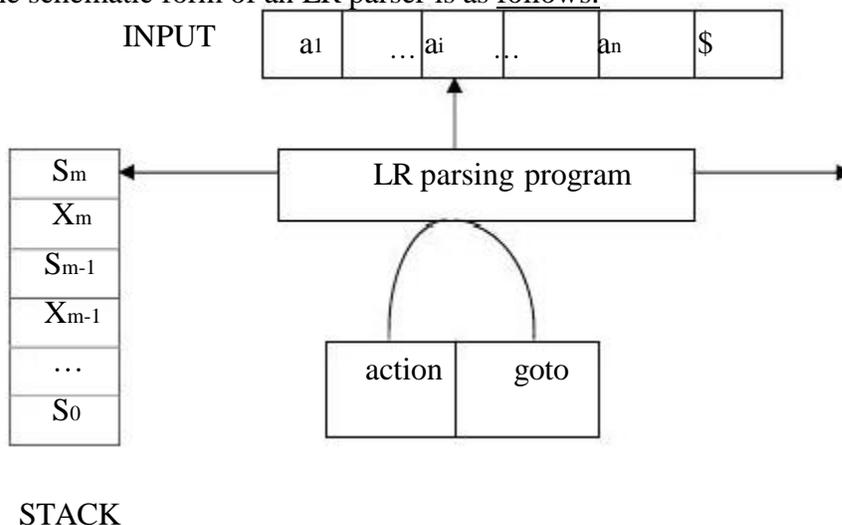
It is too much of work to construct a LR parser by hand for a programming language grammar. A specialized tool, called a LR parser generator, is needed. Example: YACC.

### Types of LR parsing method:

1. SLR- Simple LR
  - Easiest to implement, least powerful.
2. CLR- Canonical LR
  - Most powerful, most expensive.
3. LALR- Look-Ahead LR
  - Intermediate in size and cost between the other two methods.

### The LR parsing algorithm:

The schematic form of an LR parser is as follows:



It consists of : an input, an output, a stack, a driver program, and a parsing table that has two parts (*action* and *goto*).

- The driver program is the same for all LR parser.
- The parsing program reads characters from an input buffer one at a time.
- The program uses a stack to store a string of the form  $s_0X_1s_1X_2s_2\dots X_ms_m$ , where  $s_m$  is on top. Each  $X_i$  is a grammar symbol and each  $s_i$  is a state.
- The parsing table consists of two parts : *action* and *goto* functions.

**Action** : The parsing program determines  $s_m$ , the state currently on top of stack, and  $a_i$ , the current input symbol. It then consults  $action[s_m, a_i]$  in the action table which can have one of four values :

1. shift  $s$ , where  $s$  is a state,
2. reduce by a grammar production  $A \rightarrow \beta$ ,
3. accept, and
4. error.

**Goto** : The function *goto* takes a state and grammar symbol as arguments and produces a state.

### **LR Parsing algorithm:**

**Input:** An input string  $w$  and an LR parsing table with functions *action* and *goto* for grammar  $G$ .

**Output:** If  $w$  is in  $L(G)$ , a bottom-up-parse for  $w$ ; otherwise, an error indication.

**Method:** Initially, the parser has  $s_0$  on its stack, where  $s_0$  is the initial state, and  $w\$$  in the input buffer. The parser then executes the following program :

```
set ip to point to the first input symbol of  $w\$$ ;  
repeat forever begin  
    let  $s$  be the state on top of the stack and  
     $a$  the symbol pointed to by ip;  
    if  $action[s, a] = \text{shift } s'$  then begin  
        push  $a$  then  $s'$  on top of the stack;  
        advance ip to the next input symbol  
    end  
    else if  $action[s, a] = \text{reduce } A \rightarrow \beta$  then begin  
        pop  $2 * |\beta|$  symbols off the stack;  
        let  $s'$  be the state now on top of the stack;  
        push  $A$  then  $goto[s', A]$  on top of the stack;  
        output the production  $A \rightarrow \beta$   
    end  
    else if  $action[s, a] = \text{accept}$  then  
        return  
    else  $error()$   
end
```

## CONSTRUCTING SLR(1) PARSING TABLE:

To perform SLR parsing, take grammar as input and do the following:

1. Find LR(0) items.
2. Completing the closure.
3. Compute  $goto(I,X)$ , where, I is set of items and X is grammar symbol.

### LR(0) items:

An  $LR(0)$  item of a grammar G is a production of G with a dot at some position of the right side. For example, production  $A \rightarrow XYZ$  yields the four items :

$A \rightarrow \cdot XYZ$   
 $A \rightarrow X \cdot YZ$   
 $A \rightarrow XY \cdot Z$   
 $A \rightarrow XYZ \cdot$

### Closure operation:

If I is a set of items for a grammar G, then  $closure(I)$  is the set of items constructed from I by the two rules:

1. Initially, every item in I is added to  $closure(I)$ .
2. If  $A \rightarrow \alpha \cdot B\beta$  is in  $closure(I)$  and  $B \rightarrow \gamma$  is a production, then add the item  $B \rightarrow \cdot \gamma$  to I, if it is not already there. We apply this rule until no more new items can be added to  $closure(I)$ .

### Goto operation:

$Goto(I, X)$  is defined to be the closure of the set of all items  $[A \rightarrow \alpha X \cdot \beta]$  such that  $[A \rightarrow \alpha \cdot X\beta]$  is in  $I$ .

Steps to construct SLR parsing table for grammar G are:

1. Augment G and produce  $G'$
2. Construct the canonical collection of set of items C for  $G'$
3. Construct the parsing action function *action* and *goto* using the following algorithm that requires FOLLOW(A) for each non-terminal of grammar.

### Algorithm for construction of SLR parsing table:

**Input** : An augmented grammar  $G'$

**Output** : The SLR parsing table functions *action* and *goto* for  $G'$

**Method** :

1. Construct  $C = \{I_0, I_1, \dots, I_n\}$ , the collection of sets of LR(0) items for  $G'$ .
2. State  $i$  is constructed from  $I_i$ . The parsing functions for state  $i$  are determined as follows:
  - (a) If  $[A \rightarrow \alpha \cdot a\beta]$  is in  $I_i$  and  $goto(I_i, a) = I_j$ , then set  $action[i, a]$  to "shift j". Here  $a$  must be terminal.
  - (b) If  $[A \rightarrow \alpha \cdot]$  is in  $I_i$ , then set  $action[i, a]$  to "reduce  $A \rightarrow \alpha$ " for all  $a$  in FOLLOW(A).
  - (c) If  $[S' \rightarrow S \cdot]$  is in  $I_i$ , then set  $action[i, \$]$  to "accept".

If any conflicting actions are generated by the above rules, we say grammar is not SLR(1).

3. The *goto* transitions for state  $i$  are constructed for all non-terminals  $A$  using the rule:  
If  $goto(I_i, A) = I_j$ , then  $goto[i, A] = j$ .
4. All entries not defined by rules (2) and (3) are made “error”
5. The initial state of the parser is the one constructed from the set of items containing  $[S' \rightarrow .S]$ .

**Example for SLR parsing:**

Construct SLR parsing for the following grammar :

G :  $E \rightarrow E + T \mid T$   
 $T \rightarrow T * F \mid F$   
 $F \rightarrow (E) \mid id$

The given grammar is :

G :  $E \rightarrow E + T$  ----- (1)  
 $E \rightarrow T$  ----- (2)  
 $T \rightarrow T * F$  ----- (3)  
 $T \rightarrow F$  ----- (4)  
 $F \rightarrow (E)$  ----- (5)  
 $F \rightarrow id$  ----- (6)

**Step 1 :** Convert given grammar into augmented grammar.

**Augmented grammar :**

$E' \rightarrow E$   
 $E \rightarrow E + T$   
 $E \rightarrow T$   
 $T \rightarrow T * F$   
 $T \rightarrow F$   
 $F \rightarrow (E)$   
 $F \rightarrow id$

**Step 2 :** Find LR (0) items.

$I_0 : E' \rightarrow .E$   
 $E \rightarrow .E + T$   
 $E \rightarrow .T$   
 $T \rightarrow .T * F$   
 $T \rightarrow .F$   
 $F \rightarrow .(E)$   
 $F \rightarrow .id$

GOTO (  $I_0$  , E)

$I_1 : E' \rightarrow E .$   
 ~~$E \rightarrow E . + T$~~

GOTO (  $I_4$  , id )

$I_5 : F \rightarrow id .$   


---

GOTO (I<sub>0</sub>, T)  
I<sub>2</sub> : E → T.  
T → T.\*F

GOTO (I<sub>0</sub>, F)  
I<sub>3</sub> : T → F.

GOTO (I<sub>0</sub>, (  
I<sub>4</sub> : F → (.E)  
E → .E+T  
E → .T  
T → .T\*F  
T → .F  
F → .(E)  
F → .id

GOTO (I<sub>0</sub>, id)  
I<sub>5</sub> : F → id.

GOTO (I<sub>1</sub>, +)  
I<sub>6</sub> : E → E + .T  
T → .T\*F  
T → .F  
F → .(E)  
F → .id

GOTO (I<sub>2</sub>, \*)  
I<sub>7</sub> : T → T\* .F  
F → .(E)  
F → .id

GOTO (I<sub>4</sub>, E)  
I<sub>8</sub> : F → (E.)  
E → E.+T

GOTO (I<sub>4</sub>, T)  
I<sub>2</sub> : E → T.  
T → T.\*F

GOTO (I<sub>4</sub>, F)  
I<sub>3</sub> : T → F.

GOTO (I<sub>6</sub>, T)  
I<sub>9</sub> : E → E + T.  
T → T.\*F

GOTO (I<sub>6</sub>, F)  
I<sub>3</sub> : T → F.

GOTO (I<sub>6</sub>, (  
I<sub>4</sub> : F → (.E)

GOTO (I<sub>6</sub>, id)  
I<sub>5</sub> : F → id.

GOTO (I<sub>7</sub>, F)  
I<sub>10</sub> : T → T\*F.

GOTO (I<sub>7</sub>, (  
I<sub>4</sub> : F → (.E)  
E → .E+T  
E → .T  
T → .T\*F  
T → .F  
F → .(E)  
F → .id

GOTO (I<sub>7</sub>, id)  
I<sub>5</sub> : F → id.

GOTO (I<sub>8</sub>, )  
I<sub>11</sub> : F → (E).

GOTO (I<sub>8</sub>, +)  
I<sub>6</sub> : E → E + .T  
T → .T\*F  
T → .F  
F → .(E)  
F → .id

GOTO (I<sub>9</sub>, \*)  
I<sub>7</sub> : T → T\* .F  
F → .(E)  
F → .id

GOTO (I<sub>4</sub>, (

I<sub>4</sub> : F → ( . E)

E → . E + T

E → . T

T → . T \* F

T → . F

F → . ( E )

F → id

FOLLOW (E) = { \$ , ) , + }

FOLLOW (T) = { \$ , + , ) , \* }

FOLLOW (F) = { \* , + , ) , \$ }

**SLR parsing table:**

	<u>ACTION</u>						<u>GOTO</u>		
	<u>id</u>	<u>±</u>	<u>*</u>	<u>(</u>	<u>)</u>	<u>\$</u>	<u>E</u>	<u>T</u>	<u>F</u>
I <sub>0</sub>	<u>s5</u>			<u>s4</u>			<u>1</u>	2	3
I <sub>1</sub>		<u>s6</u>				<u>ACC</u>			
I <sub>2</sub>		<u>r2</u>	<u>s7</u>		<u>r2</u>	<u>r2</u>			
I <sub>3</sub>		<u>r4</u>	<u>r4</u>		<u>r4</u>	<u>r4</u>			
I <sub>4</sub>	<u>s5</u>			<u>s4</u>			<u>8</u>	2	3
I <sub>5</sub>		<u>r6</u>	<u>r6</u>		<u>r6</u>	<u>r6</u>			
I <sub>6</sub>	<u>s5</u>			<u>s4</u>				9	3
I <sub>7</sub>	s5			s4					10
I <sub>8</sub>		s6			s11				
I <sub>9</sub>		r1	s7		r1	r1			
I <sub>10</sub>		r3	r3		r3	r3			
I <sub>11</sub>		r5	r5		r5	r5			

Blank entries are error entries.

**Stack implementation:**

Check whether the input **id + id \* id** is valid or not.

STACK	INPUT	ACTION
0	id + id * id \$	GOTO ( I <sub>0</sub> , id ) = s5 ; <b>shift</b>
0 id 5	+ id * id \$	GOTO ( I <sub>5</sub> , + ) = r6 ; <b>reduce</b> by F → id
OF3	+ id * id \$	GOTO ( I <sub>0</sub> , F ) = 3 GOTO ( I <sub>3</sub> , + ) = r4 ; <b>reduce</b> by T → F
OT2	+ id * id \$	GOTO ( I <sub>0</sub> , T ) = 2 GOTO ( I <sub>2</sub> , + ) = r2 ; <b>reduce</b> by E → T
OE1	+ id * id \$	GOTO ( I <sub>0</sub> , E ) = 1 GOTO ( I <sub>1</sub> , + ) = s6 ; <b>shift</b>
OE1+6	id * id \$	GOTO ( I <sub>6</sub> , id ) = s5 ; <b>shift</b>
0 E 1 + 6 id 5	* id \$	GOTO ( I <sub>5</sub> , * ) = r6 ; <b>reduce</b> by F → id
OE1+6F3	<u>*</u> id \$	<u>GOTO</u> ( I <sub>6</sub> , F ) = <u>3</u> <u>GOTO</u> ( I <sub>3</sub> , * ) = r4 ; <b>reduce</b> by T → F
OE1+6T9	<u>*</u> id \$	<u>GOTO</u> ( I <sub>6</sub> , T ) = <u>9</u> <u>GOTO</u> ( I <sub>9</sub> , * ) = s7 ; <b>shift</b>
OE1+6T9*7	<u>id</u> \$	<u>GOTO</u> ( I <sub>7</sub> , id ) = s5 ; <b>shift</b>
0 E 1 + 6 T 9 * 7 id <u>5</u>	<u>\$</u>	<u>GOTO</u> ( I <sub>5</sub> , \$ ) = r6 ; <b>reduce</b> by F → id
0 E 1 + 6 T 9 * 7 F <u>10</u>	<u>\$</u>	<u>GOTO</u> ( I <sub>7</sub> , F ) = <u>10</u> <u>GOTO</u> ( I <sub>10</sub> , \$ ) = r3 ; <b>reduce</b> by T → T * F
OE1+6T9	<u>\$</u>	<u>GOTO</u> ( I <sub>6</sub> , T ) = <u>9</u> <u>GOTO</u> ( I <sub>9</sub> , \$ ) = r1 ; <b>reduce</b> by E → E + T
OE1	<u>\$</u>	GOTO ( I <sub>0</sub> , E ) = 1 GOTO ( I <sub>1</sub> , \$ ) = <b>accept</b>