



**Note Junction**  
Best Note Provider

Note By: Roshan BiSt



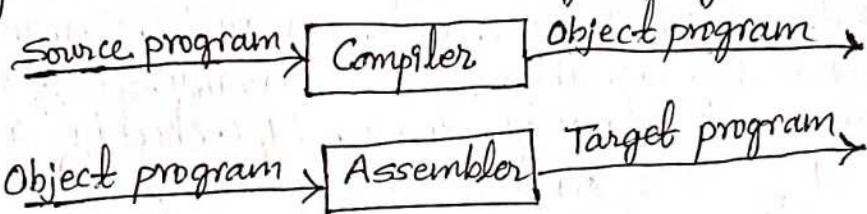
## UNIT-1

### Introduction to Compiler

A compiler is a translator software program that takes its input in the form of program written in one particular programming language and produce the output in the form of program in another language.

#### Features of Compiler:

- A compiler is a translator that converts the high-level language into the machine language.
- High-level language is written by a developer and machine language can be understood by the processor.
- Compiler is used to show errors to the programmer.
- Compiler executes a program into two parts: In first, source program is compiled into object program, In second, object program is translated into target program.



#### Compiler Structure:

A compiler operates in phases. A phase is a logically interrelated operation that takes source program in one representation and produces output in another representation. There are two phases of compilation:

1) Analysis phase: In analysis part, an intermediate representation is created from the given source program. This part is called front end of the compiler. This part consists of mainly four phases:

↳ Lexical Analysis: Lexical analysis or scanning is the process where the source program is read from left-to-right and grouped into tokens. Tokens are sequences of characters with a collective meaning. In any programming language tokens may be constants, operators, reserved words etc. The Lexical Analyzer takes a source program as input, and produces a stream of tokens as output.

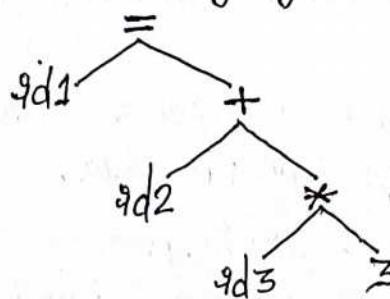
### Example:

Input string:  $c = a + b * 3$

Tokens:  $id1 = id2 + id3 * 3$

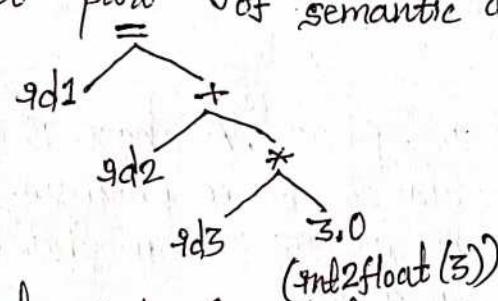
**Syntax Analysis:** In this phase, the syntax analyzer takes the token produced by lexical analyzer as input and generates a parse tree as output. In syntax analysis phase, the parser checks that the expression made by the token is syntactically correct or not, according to the rules that define the syntax of the source language.

### Example:



**Semantic Analysis:** In this phase, the semantic analyzer checks the source program for semantic errors and collects the type information for the code generation. Semantic analyzer checks whether they form a sensible set of instructions in the programming language or not. Type-checking is an important part of semantic analyzer.

### Example:



**Intermediate Code Generation:** If the program syntactically and semantically correct then intermediate code generator generates a simple machine independent intermediate language. The intermediate code should be generated in such a way that it can easily be translated into the target machine code.

### Example:

$t1 = 3.0;$

$t2 = id3 * t1;$

$t3 = id2 + t2;$

$id1 = t3;$

2) Synthesis phase: In synthesis part, the equivalent target program is created from intermediate representation of the program created by analysis part. This part is also called back end of the compiler. This part consists of mainly two phases:

i) Code Optimization: It is used to improve the intermediate code so that the output of the program could run faster and takes less space. It removes the unnecessary lines of the code and arranges the sequence of statements in order to speed up the program execution without wasting resources.

Example:  $t2 = 9d3 * 3.0;$   
 $9d1 = 9d2 + t2;$

ii) Code Generation: Code generation is the final stage of the compilation process. It takes the optimized intermediate code as input and maps it to the target machine language.

Example: MOV R1, 9d3  
MUL R1, #3.0  
MOV R2, 9d2  
ADD R1, R2  
MOV 9d1, R1.

④ Symbol Table: Symbol tables are data structures that are used by compilers to hold information about source-program constructs. The information is collected incrementally by the analysis phase of compiler and used by the synthesis phases to generate the target code. Entries in the symbol table contain information about an identifier such as its type, its position in storage, and any other relevant information.

⑤ Error Handling: Whenever an error is encountered during the compilation of the source program, an error handler is invoked. Error handler generates a suitable error reporting message regarding the error encountered. Errors can be encountered at any phase as:

Lexical analysis phase: due to misspelled tokens, unrecognized characters etc.

Syntax analysis phase: due to syntactic violation of language.

Intermediate code generation: due to incompatibility of operands type for operator.

Code optimization phase: due to some unreachable statements.

## Block Diagram:

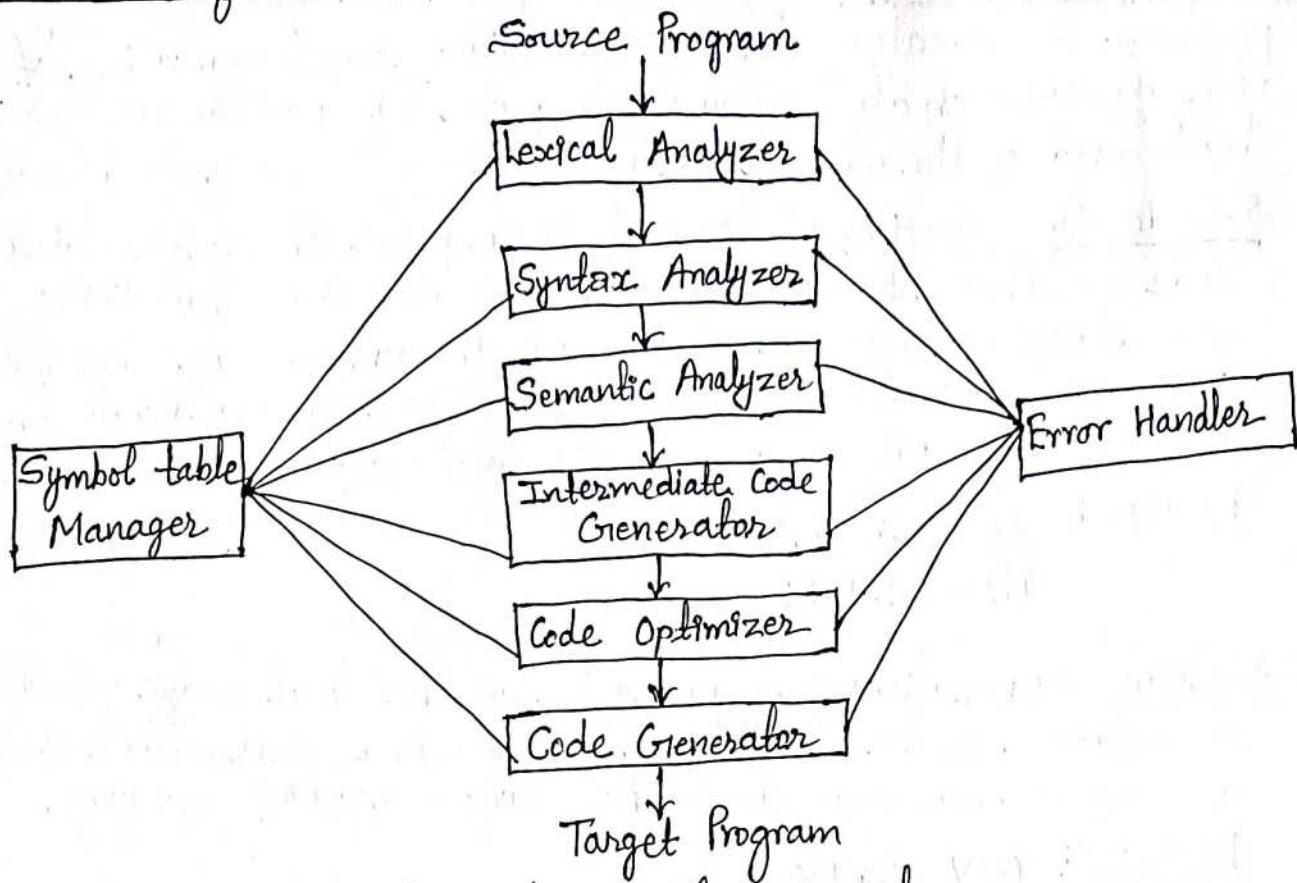


Fig: Phases of a Compiler

## ④ Compiler vs Interpreter:

Compiler	Interpreter
<ul style="list-style-type: none"> <li>⇒ Compiler is a translator which takes input i.e., high-level language, and produces an output of low-level language.</li> <li>⇒ Compiler scans the whole program in one go.</li> <li>⇒ Errors are shown at the end together as it scans full code in one go.</li> <li>⇒ It does not require source code for later execution.</li> <li>⇒ It converts source code into object code.</li> </ul>	<ul style="list-style-type: none"> <li>⇒ An Interpreter is a program that translates a programming language into a comprehensible language.</li> <li>⇒ Interpreter translates one statement at a time.</li> <li>⇒ Errors are shown line by line as it scans code one line at a time.</li> <li>⇒ It requires source code for later execution.</li> <li>⇒ It does not convert source code into object code instead it scans it line by line.</li> </ul>

## ④ Simple or One pass Compiler vs. Multi-pass Compiler:

(Simple) One pass compiler	Multi-pass compiler
i) In a one pass compiler all the phases are combined into one pass.	i) In multi-pass compiler different phases of compiler are grouped into multiple phases.
ii) Here intermediate representation of source program is not created.	ii) Here intermediate representation of source program is created.
iii) It is faster than multi-pass compiler.	iii) It is slightly slower than one pass compiler.
iv) It is also called narrow compiler.	iv) It is also called wide compiler.
v) Pascal's compiler is an example of one pass compiler.	v) C++ compiler is an example of multi-pass compiler.

⑤ Preprocessor: A preprocessor, generally considered as a part of compiler, is a tool that produces input for compilers. It deals with macro-processing, argumentation, file inclusion etc.

### Functions:

- Macro processing: A preprocessor may allow a user to define macros that are short hands for longer constructs.
- File inclusion: A preprocessor may include header files onto the program text.
- Rational preprocessor: These preprocessors augment older languages with more modern flow-of-control and data structuring facilities.
- Language Extensions: These preprocessor attempts to add capabilities to the language by certain amounts to build-in macro.

④ Macros: A macro stands for macroinstruction is a programmable pattern which translates a certain sequence of input into a present sequence of output. Macros can make tasks less repetitive by representing a complicated sequence of keystrokes, mouse movements, commands, or other types of input.

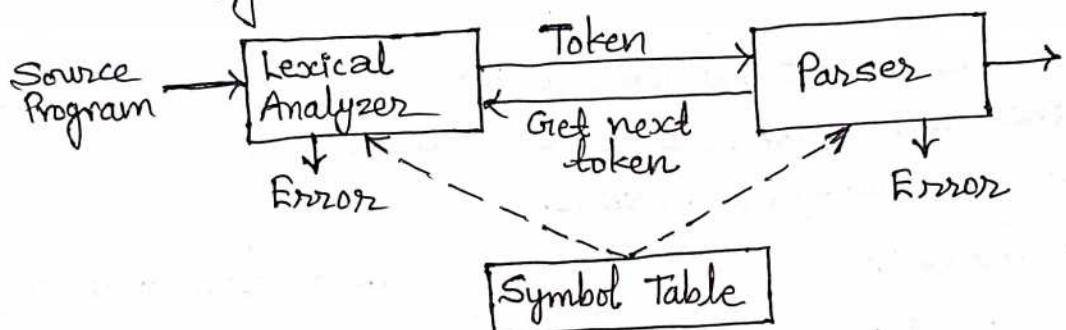
#### Features of macro processor:

- It represents a group of commonly used statements in the source programming language.
- Using macro instructions programmer can leave the mechanical details to be handled by macro processor.
- Macro processor designs is not directly related to the computer architecture on which it runs.
- Macro processor involves definition, invocation and expansion.

UNIT-2Lexical Analyzer

lexical analysis is  
done for validation  
of tokens

④ Lexical Analysis: The lexical analysis is the first phase of a compiler where a lexical analyzer acts as an interface between the source program and the rest of the phases of compiler. It reads input characters of the source program, groups them into lexemes, and produces a sequence of tokens for each lexeme. The tokens are then sent to the parser for syntax analysis.



Example: Example of Lexical Analysis, Tokens, Non-Tokens  
Consider the following code that is fed to Lexical Analyzer.

```

#include <stdio.h>
int largest (int x, int y)
{
    if (x>y) //This will compare two numbers
        return x;
    else
        return y;
}
    
```

Examples of Tokens created:

Lexeme	Token
int	Keyword
largest	Identifier
(	Operator
int	Keyword
x	Identifier
,	Operator
int	Keyword
y	Identifier
)	Operator
{	Operator
if	Keyword

## Examples of Non-Tokens:

Type	Examples
Comment	//this will compare 2 numbers.
Pre-processor directive	#include <stdio.h>

## ②. Role of Lexical Analyzer:

- Lexical analyzer helps to identify token into the symbol table.
- It can either work as a separate module or as a sub-module.
- It is responsible for eliminating comments and white spaces from the source program.
- It reads the input character and produces output sequence of tokens that the parser uses for syntax analysis.
- It also generates lexical errors.
- Lexical analyzer is used by web browsers to format and display a web page with the help of parsed data from Javascript, HTML, CSS.

## ③. Lexemes, Patterns, Tokens:

Lexemes: A lexeme is a sequence of alphanumeric characters that is matched against the pattern for token. A sequence of input characters that make up a single token is called a lexeme. A token can represent more than one lexeme.

Example: The token "string constant" may have a number of lexemes such as "bh", "sum", "area", "name" etc.

Patterns: Patterns are the rules for describing whether a given lexeme belongs to a token or not. The rule associated with each set of string is called pattern. Lexeme is matched against pattern to generate token. Regular expressions are widely used to specify patterns.

Token: Token is word, which describes the lexeme in source program. It is generated when lexeme is matched against pattern. A token is a logical building block of language. They are the sequence of characters having a collective meaning.

Example 1: Example showing lexeme, token and pattern for variables.

- Lexeme: A1, Sum, Total
- Pattern: Starting with a letter and followed by letter or digit but not a keyword.
- Token: ID

Example 2: Example showing lexeme, token and pattern for floating number.

- Lexeme: 123.45
- Pattern: Starting with digit followed by a digit or optional fraction and or optional exponent.
- Token: NUM

### Specifications of Tokens:

Regular Expression is a way to specify tokens. The regular expression represents the regular languages. The language is the set of strings and string is the set of alphabets. Thus, following terminologies are used to specify tokens:

1) Alphabets: The set of symbols is called alphabets. Example any finite set of symbols  $\Sigma = \{0, 1\}$  is a set of binary alphabets.

2) Strings: Any finite sequence of alphabets is called a string. Length of string is the total number of occurrence of alphabets. e.g., the length of string 'Kanchanpur' is 10. A string of zero length is known as empty string and is denoted by  $\epsilon$  (epsilon).

3) Language: A language is considered as a finite set of strings over some finite set of alphabets. Computer languages are considered as finite sets, and mathematically set operations can be performed on them.

4) Union: It is defined as the set that consists of all elements belonging to either set A or set B or both. In regular expression + symbol is used to represent union operation.

$$A = \{\text{dog, ba, na}\} \text{ and } B = \{\text{house, ba}\}$$

$$\text{then, } A \cup B = A + B = \{\text{dog, ba, na, house}\}$$

v) Concatenation: It is the operation of joining character strings end-to-end. In regular expression dot operator (.) is used to represent concatenation operation.

$$A \cdot B = \{\text{doghouse, dogba, bahouse, baba, nahouse, naba}\}$$

w) Kleene Closure: The Kleene closure  $\Sigma^*$  gives the set of all possible strings of all possible lengths over  $\Sigma$  including  $\{\emptyset\}$ .

x) Positive Closure: The set  $\Sigma^+$  is the infinite set of all possible strings of all possible lengths over  $\Sigma$  excluding  $\{\emptyset\}$ .

# Regular Expressions: Regular expressions are the algebraic expressions that are used to describe tokens of a programming language. It uses three regular operations called union, concatenation and star.

Examples: Given the alphabet  $A = \{0, 1\}$ .

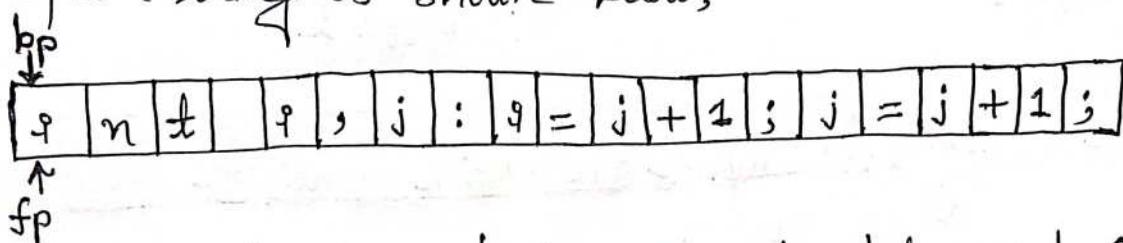
1)  $1(1+0)^*0$  denotes the language of all strings that begins with a '1' and ends with a '0'.

2)  $(01)^* + (10)^*$  denotes the set of all strings that describe alternating 1s and 0s. [Note: Practice R.E for identifiers in C, one dimensional array, two dimensional array & floating numbers. (6, 10, 11, 12 no. in rec book under RE topic.)]

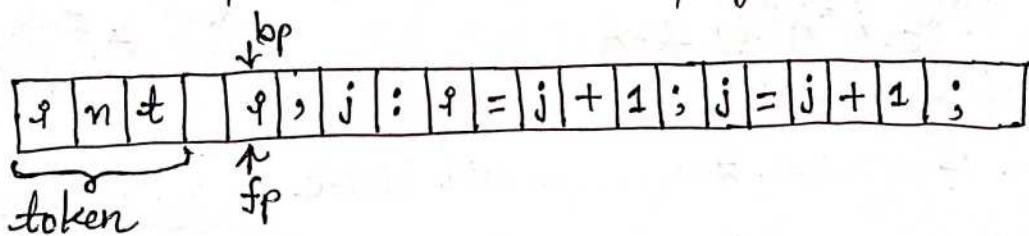
④ Recognition of Tokens: A recognizer for a language is a program that takes a string  $w$ , and answers "Yes" if  $w$  is a sentence of that language, otherwise "No." The tokens that are specified using RE are recognized by using transition diagram or finite automata (FA). Starting from the start state we follow the transition defined. If the transition leads to the accepting state, then the token is matched and hence the lexeme is returned, otherwise other transition diagrams are tried out until we process all the transition diagram or the failure is detected. Recognizer of tokens takes the language  $L$  and the string  $s$  as input and try to verify whether  $s \in L$  or not. There are two types of Finite Automata.

- 1). Deterministic Finite Automata (DFA).
- 2). Non Deterministic Finite Automata (NFA).

④ Input Buffering: Reading character by character from secondary storage is slow process and time consuming as well so, we use buffer technique to eliminate this problem and increase efficiency. The lexical analyzer scans the input from left to right one character at a time. It uses two pointers, begin ptr (bp) and forward ptr (fp) to keep track of the pointer of input scanned. Initially both the pointers point to the first character of the input string as shown below;



The forward ptr moves ahead to search for end of lexeme. As soon as the blank space is encountered, it indicates end of lexeme. In above example as soon as fp encounters a blank space the lexeme 'int' is identified. Then both bp and fp are set at next token as shown below and this process will be repeated for the whole program.



One Buffer Scheme: In this scheme, only one buffer is used to store the input string but the problem with this scheme is that if lexeme is very long then it crosses the buffer boundary, to scan rest of the lexeme the buffer has to be refilled.

Two Buffer Scheme: In this scheme, two buffers are used to store input string and they are scanned alternately. When end of current buffer is reached the other buffer is filled. The only problem with this method is that if length of the lexeme is longer than length of the buffer then scanning input cannot be scanned completely.

Initially both the bp and fp are pointing to the first character of first buffer. Then the fp moves forward in search of

end of lexeme. As soon as blank character is recognized, the string between bp and fp is identified as corresponding token. To identify, the boundary of first buffer end of buffer character (eof) should be placed at the end of first buffer, similarly for second buffer also. Alternatively both the buffers can be filled up until end of the input program and stream of tokens is identified.

	a	n	t		=		i	+ 1		eof	
--	---	---	---	--	---	--	---	-----	--	-----	--

Buffer 1

	;	j	=		j		+		1		;		eof	
--	---	---	---	--	---	--	---	--	---	--	---	--	-----	--

Buffer 2.

⊗ Read Following basic topics from 4<sup>th</sup> sem note of TOC (Unit-2 and Unit-3):

- 1) Finite Automata.
  - 2) Deterministic Finite Automata (DFA).
  - 3) Non Deterministic Finite Automata (NFA).
  - 4) Epsilon NFA ( $\epsilon$ -NFA).
  - 5) Minimization of DFA. (state partition method)
  - 6) Equivalence of Regular Expression and Finite Automata.
  - 7) Reduction of Regular Expression to  $\epsilon$ -NFA.
  - 8) Converting  $\epsilon$ -NFA to its equivalent DFA.
- ⊗ Design of a lexical Analyzer: [Imp]
- Lexical Analyzer can be designed using following two algorithms:

Algorithm 1: Regular Expression  $\rightarrow$  NFA  $\rightarrow$  DFA (two steps: first RE to NFA, then NFA to DFA).

Algorithm 2: Regular Expression  $\rightarrow$  DFA (directly convert a regular expression into a DFA).

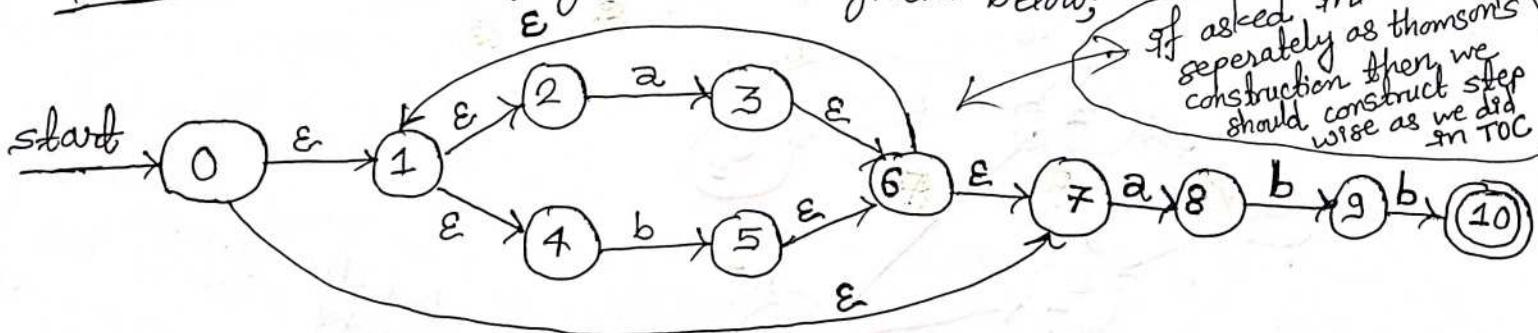
### #Algorithm 1:

This consists of following two steps:

- 1) Regular Expression to NFA (Thomson's Construction):
- 2) Conversion from NFA to DFA (Subset. Construction Algorithm):

Example: For Regular Expression  $(a+b)^*abb$ , first convert this RE to NFA, then convert resulting NFA to DFA.

Solution: The NFA of given RE is given below;



Now we convert above NFA to DFA as,

The starting state of DFA =  $S_0 = \epsilon\text{-closure}(S_0)$

Mark  $S_0$ ,

Since we have two input symbols a and b so we check for each input

$$S_0 = \{0, 1, 2, 4, 7\}$$

starting state  $S_0$

start state वाले ए

लिए कुन कुन state

सम पूर्णका है यहांका

list

जो नंबर

सेट अप्पा

अप्पा State

जो denote

JIT

like here

$S_1$

For a:  $\epsilon\text{-closure}(S(S_0, a)) = \epsilon\text{-closure}(3, 8) = \{1, 2, 3, 4, 6, 7, 8\} \rightarrow S_1$

For b:  $\epsilon\text{-closure}(S(S_0, b)) = \epsilon\text{-closure}(5) = \{1, 2, 4, 5, 6, 7\} \rightarrow S_2$

For a:  $\epsilon\text{-closure}(S(S_1, a)) = \epsilon\text{-closure}(3, 8) = \{1, 2, 3, 4, 6, 7, 8\} \rightarrow S_1$

For b:  $\epsilon\text{-closure}(S(S_1, b)) = \epsilon\text{-closure}(5) = \{1, 2, 4, 5, 6, 7, 9\} \rightarrow S_3$

Marks  $S_1$ ,

For a:  $\epsilon\text{-closure}(S(S_1, a)) = \epsilon\text{-closure}(3, 8) = \{1, 2, 3, 4, 6, 7, 8\} \rightarrow S_1$

For b:  $\epsilon\text{-closure}(S(S_1, b)) = \epsilon\text{-closure}(5) = \{1, 2, 4, 5, 6, 7\} \rightarrow S_2$

Marks  $S_2$ ,

For a:  $\epsilon\text{-closure}(S(S_2, a)) = \epsilon\text{-closure}(3, 8) = \{1, 2, 3, 4, 6, 7, 8\} \rightarrow S_1$

For b:  $\epsilon\text{-closure}(S(S_2, b)) = \epsilon\text{-closure}(5, 10) = \{1, 2, 4, 5, 6, 7, 10\} \rightarrow S_4$

Marks  $S_3$ ,

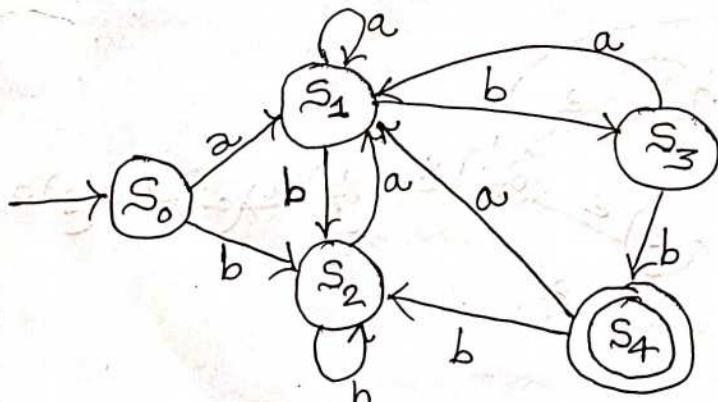
For a:  $\epsilon\text{-closure}(S(S_3, a)) = \epsilon\text{-closure}(3, 8) = \{1, 2, 3, 4, 6, 7, 8\} \rightarrow S_1$

For b:  $\epsilon\text{-closure}(S(S_3, b)) = \epsilon\text{-closure}(5) = \{1, 2, 4, 5, 6, 7\} \rightarrow S_2$

$S_0$  is the starting state.  $S_4$  is an accepting state of DFA since final state of NFA is 10 i.e., states containing 10 are member of final state of DFA. (Since  $S_4 = \{1, 2, 4, 5, 6, 7, 10\}$ ).

Now, constructing DFA using above information;

1. No new states  
here  $S_1, S_2$  are already checked  
so we stop here. We stop when there are no any new states to check



Note: Practice more examples from KEC book

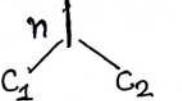
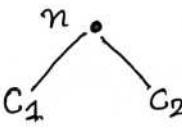
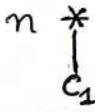
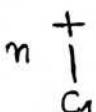
## # Algorithm 2 (Conversion from RE to DFA Directly):

Conversion Steps:

for exam mostly numericals are asked theory for understanding only

1. Augment the given regular expression by concatenating it with special symbol #.
2. Create the syntax tree for this augmented regular expression.
3. Then each alphabet symbol (including #) will be numbered as position numbers.
4. Compute functions nullable, firstpos, lastpos, and followpos.
5. Finally construct DFA directly from a regular expression by computing the functions nullable( $n$ ), firstpos( $n$ ), lastpos( $n$ ), and followpos( $i$ ) from the syntax tree.

Rules for calculating nullable, firstpos and lastpos:

Node $n$	nullable ( $n$ )	firstpos ( $n$ )	lastpos ( $n$ )
A leaf labelled $\epsilon$	True	$\emptyset$	$\emptyset$
A leaf with position $i$	False	$\{i\}$ (position of leaf node)	$\{i\}$
An or node 	Nullable ( $c_1$ ) or Nullable ( $c_2$ ).	firstpos ( $c_1$ ) $\cup$ firstpos ( $c_2$ )	lastpos ( $c_1$ ) $\cup$ lastpos ( $c_2$ )
A concatenation node 	Nullable ( $c_1$ ) and Nullable ( $c_2$ ).	If Nullable ( $c_1$ ) firstpos ( $c_1$ ) $\cup$ firstpos ( $c_2$ ). else firstpos ( $c_1$ )	If Nullable ( $c_2$ ) lastpos ( $c_1$ ) $\cup$ lastpos ( $c_2$ ). else lastpos ( $c_2$ ).
A star node 	True	firstpos ( $c_1$ )	lastpos ( $c_1$ )
A +ve closure node 	False	firstpos ( $c_1$ )	lastpos ( $c_1$ )

Now we calculate followpos for each  $i$  we numbered. Now we can construct DFA with starting state  $S_1 = \text{firstpos}(\text{root})$  from syntax tree and marking each state for each input until no new unmarked states occur.

Example: Convert the regular expression  $(a|b)^* \cdot a \#$  into equivalent DFA by direct method.

Solution:

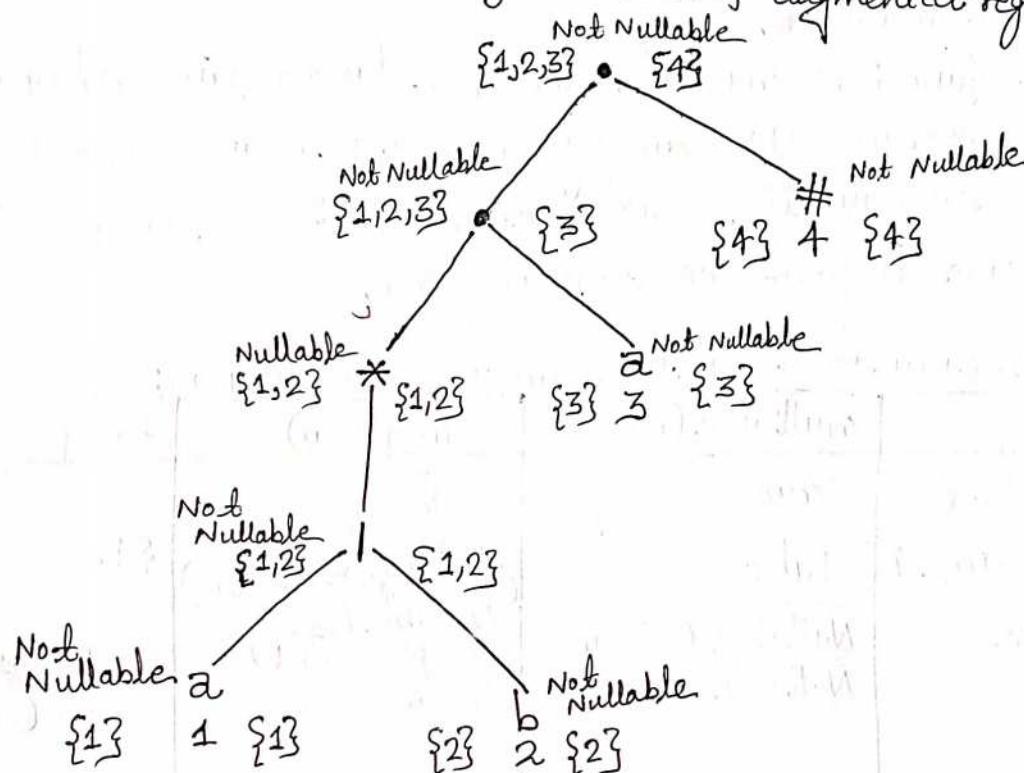
Step 1: At first we augment the given regular expression as,

alphabets including # numbered  
 $(a|b)^* \cdot a \#$

1 2 3 4

last  $\#$  या augment के लिए  
 कीन operator द्वारा product form में है  
 नीचे असली dot operator जैसा हो

Step 2: Now we construct syntax tree of augmented regular expression as,



Step 3: Now we compute followpos as,

$$\text{followpos}(1) = \{1, 2, 3\}$$

$$\text{followpos}(2) = \{1, 2, 3\}$$

$$\text{followpos}(3) = \{4\}$$

$$\text{followpos}(4) = \{\emptyset\}$$

Step 4: Now we start with starting state as,

$$S_1 = \text{firstpos (root node of syntax tree)} = \{1, 2, 3\}$$

Mark  $S_1$ ,

$$\text{For } a: \text{followpos}(1) \cup \text{followpos}(3) = \{1, 2, 3, 4\} \rightarrow S_2$$

$$\text{for } b: \text{followpos}(2) = \{1, 2, 3\} \rightarrow S_1$$

Mark  $S_2$ ,

$$\text{For } a: \text{followpos}(1) \cup \text{followpos}(3) = \{1, 2, 3, 4\} \rightarrow S_2$$

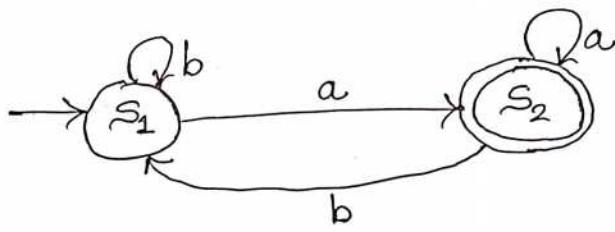
$$\text{for } b: \text{followpos}(2) = \{1, 2, 3\} \rightarrow S_1$$

Since we are marking for a  
 for  $S_1 = \{1, 2, 3\}$ . Now look  
 at step 1 where where  
 a occurs in 1, 2, 3. Here  
 a occurs in 1 & 3.  
 So Union of 1 & 3

marked as  
 new state  
 or state name  $S_2$

No new states occur so we stop marking here. Accepting state means state containing position of  $\#$  (i.e., 4). So, here  $S_2 = \{1, 2, 3, 4\}$  is accepting state.

Now, based on the above information, resulting DFA of given regular expression is as follows:



Practice these questions also if any confusion or want to try to match answer then refer tec book for solution page no 37.

Q. Convert the regular expression  $(a|\varepsilon)bc^*$  into equivalent DFA by direct method.

Q. Convert the regular expression  $ba(a+b)^*ab$  into equivalent DFA by direct method.

### Flex: An Introduction

possibility  
of asking in short note  
for 2.5 marks

Flex is a tool for generating scanners. A scanner is a program which recognizes lexical patterns in text. The flex program reads the given input files, or its standard input if no file names are given, for a description of a scanner to generate. The description is in the form of pairs of regular expressions and C code, called rules. Flex generates as output a C source file, 'lex.yy.c' by default, which defines a routine `yylex()`. This file can be compiled and linked with the flex runtime library to produce an executable. When the executable is run, it analyzes its input for occurrences of the regular expressions. Whenever it finds one, it executes the corresponding C code.

Flex Specification: A flex specification consists of three parts:

Regular definitions, C declarations in % { % }

Translation rules

This is syntax. % % are opening and closing tags.

% %

User defined auxiliary procedures

The translation rules are of the form:

P<sub>1</sub>      {action 1}

P<sub>2</sub>      {action 2}

P<sub>n</sub>      {action n}

In all parts of the specification comments of the form /\*comment text\*/ are permitted.

## Syntax Analyzer

about 20 marks की  
स्टोरेज ये बात

④ Syntax Analysis: All programming languages have certain syntactic structures. We need to verify the source code written for a language is syntactically valid. The validity of the syntax is checked by the syntax analysis. Syntaxes are represented using context free grammar (CFG), or Backus Naur Form (BNF). Parsing is the act of performing syntax analysis to verify an input program's compliance with the source language. The purpose of syntax analysis or parsing is to check that we have a valid sequence of tokens.

### ④ Role of Syntax Analyzer/Parser:

syntax analyzer  
and Parser are same  
thing

The parser 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 has to report any syntax errors if occurs.

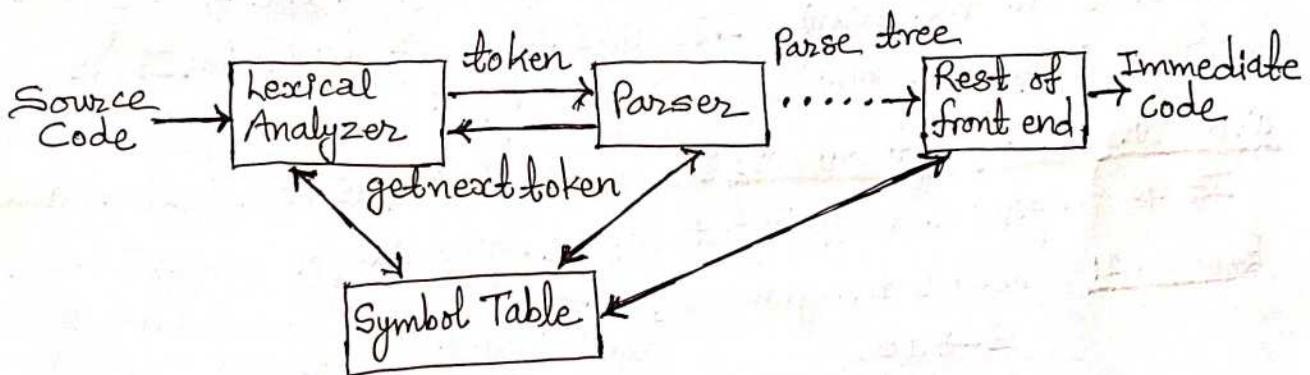


Fig: Role of Parser in a compiler model.

The tasks of parser can be expressed as;

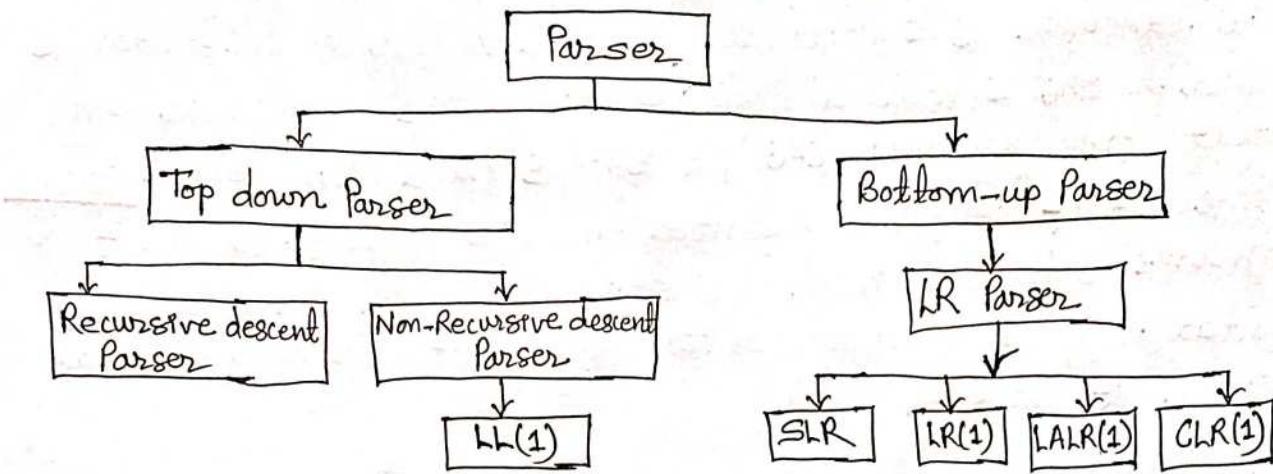
- Analyzes the context free syntax.
- Generates the parse tree.
- Provides the mechanism for context sensitive analysis.
- Determine the errors and tries to handle them.

### ④ Basic Topics to understand from TOC Unit-4 or Kec Book:

1. Concept of context free grammar including production rules of language.
2. Derivations (leftmost and rightmost)
3. Parse Trees
4. Ambiguity of a Grammar.

These topics are not asked in exam  
maybe rarely parse tree asked but  
concept is needed for upcoming  
important topic Parsing.

\* Parsing: Parser is a compiler that is used to break the data into smaller elements coming from lexical analysis phase. A parser takes input in the form of sequence of tokens and produces output in the form of parse tree.



### 1) Top down Parser:

Top-down parser is a parser which generates parse for the given input string with the help of grammar productions by expanding the non-terminals i.e., it starts from the start symbol and ends on the terminals. It uses left most derivation. It is further classified into following two types:

a) Recursive descent parser: It is also known as backtracking parser. It is a general inefficient technique, only used for small production rules.

Example 1: Consider the grammar,

$$S \rightarrow aBc$$

$$B \rightarrow bc/b$$

Input string: abc

Solution:

For solving this every time we use three columns input, output & rule.  
Input मा दिएको input string राख्न  
output मा production rule मा starting symbol  
जुन हो राख्ने / र rule मा starting symbol  
मा दिएको production rule राख्ने

Input	Output	Rule used
abc	S	use $S \rightarrow aBc$
abc	aBc	Match symbol a
bc	Bc	use $B \rightarrow bc$
bc	bcc	Match symbol b
c	cc	Match symbol c
∅	c	Dead end, back-track.
bc	Bc	use $B \rightarrow b$
bc	bc	Match symbol b
c	c	Match symbol c
∅	∅	Accepted

above step मा  $S \rightarrow aBc$   
rule use गरेकोले S  
लाई abc को replace  
गर्ने

left to right corresponding  
symbol match जस्तै symbol  
match जस्तै c जस्तै symbol  
delete जस्तै here a प्राप्त  
deleted

now we have capital  
B which is non terminal  
so it should be reduced  
to terminal using production  
rule accept  
backtrack जस्तै

backtrack जस्तै

Example 2: Consider the grammar,

$$S \rightarrow abc \mid aab$$

$$B \rightarrow bc \mid a$$

Input string: aaa

Solution:

Input	Output	Rule used
aaa	S	use $S \rightarrow abc$
aaa	aBc	Match symbol a
aa	Bc	use $B \rightarrow bc$
aa	bcc	Dead end, backtrack
aa	Bc	use $B \rightarrow a$
aa	ac	Match symbol a
a	c	Dead end, backtrack
aa	Bc	Dead end, backtrack
aa	Bc	Dead end, backtrack
aaa	S	use $S \rightarrow aab$
aaa	aAb	Match symbol a
aa	aB	Match symbol a
a	B	use $B \rightarrow bc$
a	bc	Dead end, backtrack
a	B	use $B \rightarrow a$
a	a	Match symbol a
∅	∅	Accepted

(1) मुख्य को दूर कर  
symbol terminal  
(i.e., small) मात्र  
जो match होता है  
जो dead end होता है  
recently used rule  
मात्र backtrack  
होता है

(2) back track गेरे input  
in output copy करें  
recently used production  
rule की

(3) B replaced by a by  
rule and symbol a matched  
so will get deleted in next step

(4) dead end आएर recent  
production rule B → a मात्र  
जारी जसमा input व output  
aa व Bc दूर होने हामिले B  
लाई reduce जाने पर्ने terminal  
symbol (i.e., small) मात्र but  
B → bc व B → a already  
check जारीकरें एवं आएर  
अब rule में B gives so  
फिर backtrack नहीं मानिए  
production S → abc मात्र

(\*) Left Recursion: A grammar becomes left-recursive if it has any non-terminal 'A' whose derivation contains 'A' itself as the left-most symbol. Left recursion becomes problem for top-down parsers because left recursion leads to infinite loop so we reduce or eliminate it before solving top-down parsers. A grammar is left recursive if it has a non-terminal A such that there is a derivation:

$$A \rightarrow A\alpha \text{ for some string } \alpha,$$

Let we have production as follows:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid A\alpha_3 \mid \dots \mid A\alpha_n \mid B_1 \mid B_2 \mid B_3 \mid \dots \mid B_n$$

where,  $B_1, B_2, \dots, B_n$  do not start with A

Now we eliminate immediate left recursion as;

$$A \rightarrow B_1 A' \mid B_2 A' \mid \dots \mid B_n A'$$

→ general form  
of left recursion  
eliminate it

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

A<sub>1</sub> to A<sub>n</sub> left recursion नहीं होती  
दूरी पर A' से होती है

B means it can  
be terminal or  
non-terminal symbol

→ माध्यिक left recursion जारी होती है  
पर A' से होती है

Example 1: Eliminate left recursion from following grammar,

$$S \rightarrow SBC | Sab | ab | ba | a | b$$

$$B \rightarrow bc/a$$

Solution:

$$S \rightarrow abs' | bas' | as' | bs'$$

$$S' \rightarrow BCS' | ABS' | \epsilon$$

$$B \rightarrow bc/a$$

### Non-Immediate left recursion:

Example: Let's take a grammar with non-immediate left recursion.

$$S \rightarrow Aa | b$$

$$A \rightarrow Sc | d$$

This grammar is not immediately left-recursive but is still left recursive, so, at first we make immediate recursive as,

$$S \rightarrow Sca | da | b$$

$$A \rightarrow Sc | d$$

Now, eliminate left recursion as,

$$S \rightarrow das' | bs'$$

$$S' \rightarrow cas' | \epsilon$$

$$A \rightarrow Sc | d$$

⊗. Left-factoring: If more than one production rules of a grammar have a common prefix string, then the top-down parser cannot make a choice as to which of the production it should take to parse the string. This grammar is called left factoring grammar.

Let's take an example;

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \alpha\beta_3 | \dots | \alpha\beta_n | \gamma$$

Now eliminate left factoring as;

$$A \rightarrow \alpha(\beta_1 | \beta_2 | \beta_3 | \dots | \beta_n) | \gamma$$

$$A' \rightarrow \beta_1 | \beta_2 | \beta_3 | \dots | \beta_n$$

General form

repeat इन start symbol  
common factor left-  
factoring करें

Example 2: Eliminate left recursion from following grammar,

$$E \rightarrow E + T | T$$

$$T \rightarrow T * F | F$$

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

Solution:

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

Recursive E+T फॉर्म  
E नहीं है तो, E skip +T  
taken and E' added with E,  
as in general form

non-recursice part with T'

directly left recursion नहीं है लेकिन अब production रखेर check जाएँ left recursion आउएँ, यस्ते लाइ non-immediate left recursion नहीं है, असली solve जाएँ अब production rule रखेर पहिला left recursion करनावें अवधि को यस्ते solve जाएँ।

production की right side rule E<sub>1</sub>  
में same symbol और start symbol  
हो जाएँ (more than 1 rule) तो, it  
is called left-factorial grammar

Example 1: Eliminate left factorial from following grammar:

$$S \rightarrow \#E\$S \mid \#E\$S\$ \mid a$$

$$B \rightarrow b$$

Solution:

$$S \rightarrow \#E\$S(E \mid \$) \mid a$$

$$B \rightarrow b$$

The resulting grammar with left factorial free is,

$$S \rightarrow \#E\$S' \mid a$$

$$S' \rightarrow E \mid \$$$

$$B \rightarrow b$$

Example 2: Eliminate left factorial from following grammar:

$$S \rightarrow bSSaas \mid bSSaSb \mid bSb \mid a$$

Solution:

$$S \rightarrow bSS' \mid a$$

$$S' \rightarrow Saas \mid Sasb \mid b$$

The resulting grammar with left factorial free is,

$$S \rightarrow bSS' \mid a$$

$$S' \rightarrow SaS''$$

$$S'' \rightarrow as \mid sb \mid b$$

### b) Non-Recursive Descent Parser:

A form of recursive-descent parsing that does not require any back-tracking is known as predictive parsing. It is also called LL(1) parsing table technique since we would be building a table for string to be parsed. It has capability to predict which production is to be used to replace input string.

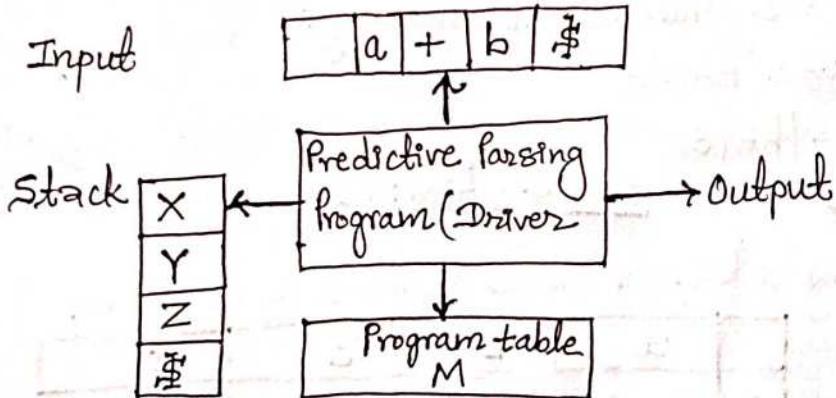


Fig: Model of a non-recursive predictive parser.

- Input Buffer: It contains the string to be parsed followed by a symbol \$.
- Stack: It contains sequence of grammar symbols with \$ at bottom.

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

Output stream: A production rule representing a step of the derivation sequence of the string in the input buffer.

### Recursive predictive descent parser vs. Non-recursive predictive descent parser:

Recursive predictive descent parser	Non-recursive predictive descent parser
It is a technique which may require backtracking process.	It is a technique that does not require any kind of backtracking.
It uses procedures for every non-terminal entity to parse strings.	It finds out productions to use by replacing input string.
It is a type of top-down parsing built from a set of mutually recursive procedures where each procedure implements one of non-terminals of grammar.	It is a type of top-down approach, which is also a type of recursive parsing that does not use technique of backtracking.
It contains several small functions one for each non-terminals in grammar.	The predictive parser uses a look ahead pointer which points to next input symbols to make it parser backtracking free, predictive parser puts some constraints on grammar.
It accepts all kinds of grammars.	It accepts only a class of grammars known as LL(k) grammar.

Example: (To demonstrate Non-recursive descent parser):

Consider the grammar  $G_1$  given by:

$$S \rightarrow aAa \mid BAA \mid \epsilon$$

$$A \rightarrow cA \mid bA \mid \epsilon$$

$$B \rightarrow b$$

Input string: bcba

Solution:

Parsing table for above grammar is as follows:

	a	b	c	\$
S	$S \rightarrow aAa$	$S \rightarrow BAA$		$S \rightarrow \epsilon$
A	$A \rightarrow \epsilon$	$A \rightarrow bA$	$A \rightarrow cA$	
B		$B \rightarrow b$		

DT Table 3(B) Topic  
LL(1) DT construct STT  
DT DT | For now  
consider as we  
know to draw this  
table or this table  
is given in question

Now do not  
focus on how  
this table  
is constructed  
just use this  
table

Now parsing of input string  $w = bcba$  using non-recursive descent parser is as follows:

Stack की सुरक्षा  $\rightarrow$   
अंतिम सимвол तो Start symbol of grammar लेंदूँ

②  
Input को सुरक्षा  
Input String र last मा  
\\$ symbol लेराहूँ।

Stack में last character  
आया तो input में first  
character का जावा करें।

Stack	Remaining Input	Action
$\$S$	$bcba\$$	choose $S \rightarrow BAa$
$\$aAB$	$bcba\$$	choose $B \rightarrow b$
$\$aAb$	$bcba\$$	match $b$
$\$aA$	$cba\$$	choose $A \rightarrow cA$
$\$aAc$	$cba\$$	match $c$
$\$aA$	$ba\$$	choose $A \rightarrow bA$
$\$aAb$	$ba\$$	match $b$
$\$aA$	$a\$$	choose $A \rightarrow \epsilon$
$\$a$	$a\$$	match $a$
$\$$	$\$$	Accept

③  
प्राप्ति table  
मात्र stack की  
अंतिम symbol  
अंतिम input  
symbol की first  
rule देता है।  
Here we have  
stack top = \$  
1st input = b  
So, S & b in  
table gives  
BAa.

④  
We have  $S$  in stack top  
but we choose  $S \rightarrow BAa$   
So,  $S$  is replaced by  
 $BAa$  but as we know stack is  
last in first out (LIFO) so we have to  
store it in reverse order as  $aAB$ .

Now top of stack contains  $B$   
and start of input contains  $b$   
so according to table rule is  $B \rightarrow b$   
So we choose  $B \rightarrow b$

⑤  
Now top of stack contains  $b$   
and first input also  $b$   
so both matched and  $b$   
is deleted

⑥  
Similarly we  
proceed and finally  
if we get both  
stack and input as  
\\$ symbol only then  
we accept string  
otherwise reject

### Constructing LL(1) Parsing Table: [V.Imp]

The parse table construction requires two functions: FIRST and FOLLOW.

A grammar  $G_1$  is suitable for LL(1) parsing table if the grammar  
is free from left recursion and left factoring.

A grammar  $\xrightarrow{\text{eliminate left recursion}} \xrightarrow{\text{eliminate left factor}}$

Grammar is  
suitable now for  
predictive parsing  
(LL(1) grammar)

To compute LL(1) parsing table, at first we need to compute  
FIRST and FOLLOW functions.

Compute FIRST:  $\text{FIRST}(\alpha)$  is a set of terminal symbols which occur as first symbols in strings derived from  $\alpha$  where  $\alpha$  is any string of grammar symbols.

Rules:

1. If 'a' is terminal, then  $\text{FIRST}(a) = \{a\}$ .

2. If  $A \rightarrow E$  is a production, then  $\text{FIRST}(A) = E$ .

3. For any non-terminal A with production rules  $A \rightarrow \alpha_1 | \alpha_2 | \alpha_3 | \dots | \alpha_n$  then  $\text{FIRST}(A) = \text{FIRST}(\alpha_1) \cup \text{FIRST}(\alpha_2) \cup \dots \cup \text{FIRST}(\alpha_n)$ .

4. If the production rule of the form,  $A \rightarrow \beta_1 \beta_2 \beta_3 \dots \beta_n$  then,  
 $\text{FIRST}(A) = \text{FIRST}(\beta_1 \beta_2 \beta_3 \dots \beta_n)$ .

Example 1: Find FIRST of following grammar symbols,

$$R \rightarrow aS | (R)S$$

$$S \rightarrow +RS | aRS | *S | \epsilon$$

Solution:

$$\text{FIRST}(aS) = \text{FIRST}(a) = \{a\}$$

$$\text{FIRST}(+RS) = \text{FIRST}(+) = \{+\}$$

$$\text{FIRST}(*S) = \text{FIRST}(*) = \{*\}$$

$$\text{FIRST}(R) = \{\text{FIRST}(aS) \cup \text{FIRST}((R)S)\} = \{a, (, \}\}$$

$$\text{FIRST}(S) = \{\text{FIRST}(+RS) \cup \text{FIRST}(aRS) \cup \text{FIRST}(*S) \cup \text{FIRST}(\epsilon)\}$$

$$= \{+, a, *, \epsilon\}.$$

terminal symbol so  
first of that symbol is  
first of same that symbol  
by rule

$$\text{FIRST}((R)S) = \text{FIRST}(C) = \{C\}$$

$$\text{FIRST}(aRS) = \text{FIRST}(a) = \{a\}$$

$$\text{FIRST}(\epsilon) = \text{FIRST}(\epsilon) = \{\epsilon\}$$

opening symbol  
first comes  
and these symbols  
are treated as  
terminal

Example 2: Find FIRST of following grammar symbols,

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

opening brace

Solution:

$$\text{FIRST}(F) = \{\text{FIRST}((E)) \cup \text{FIRST}(id)\} = \{C, id\}$$

$$\text{FIRST}(id) = \{id\}$$

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

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

$$\text{FIRST}(T) = \text{FIRST}(F) = \{C, id\}$$

$$\text{FIRST}(E) = \text{FIRST}(T) = \{C, id\}$$

$$\text{FIRST}(+TE') = \text{FIRST}(+) = \{+\}$$

$$\text{FIRST}(TE') = \text{FIRST}(T) = \{C, id\}$$

$$\text{FIRST}(\epsilon) = \{\epsilon\}$$

$$\text{FIRST}(*FT') = \text{FIRST}(*) = \{*\}$$

$$\text{FIRST}(FT') = \text{FIRST}(F) = \{C, id\}$$

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

Compute FOLLOW: FOLLOW(A) is the set of terminals that can immediately follow non terminal A except  $\epsilon$ .

Rules:

- 1 If A is a starting symbol of given grammar then FOLLOW(A) = { $\$$ }.
- 2 For every production  $B \rightarrow \alpha A \beta$ , where  $\alpha$  and  $\beta$  are any string of grammar symbols and A is non terminal, then everything in FIRST( $\beta$ ) except  $\epsilon$  is FOLLOW(A).
- 3 For every production  $B \rightarrow \alpha A$ , or a production  $B \rightarrow \alpha A \beta$ , FIRST( $\beta$ ) contains  $\epsilon$ , then everything in FOLLOW(B) is FOLLOW(A).

Example 1: Compute FOLLOW of the following grammar.

Follow compute  
जटी left side  
को symbol का मात्र  
जरुर नि चुना

$$R \rightarrow aS | (R)S$$

$$S \rightarrow +RS | aRS | *S | \epsilon$$

Solution:

$$\text{FOLLOW}(R) = \{\text{FOLLOW}(R)S \cup \text{FOLLOW}(+RS) \cup \text{FOLLOW}(aRS)\}$$

$$= \{\text{FIRST}(S) \cup \text{FIRST}(S) \cup \text{FOLLOW}(S)\}$$

according  
to 2nd rule

$= \{\$, ), +, a, *\}$

R starting symbol  
भारकोले  $\$$  घपिएको

$$\text{FIRST}(S) \text{ is } )$$

$\text{FIRST}(S)$  is  
already computed  
in FIRST section  
example 1 so we  
directly use here  
 $\text{FIRST}(S) = \{+, a, *, \epsilon\}$

We have aRS  
but  $S \rightarrow \epsilon$ , so if  
we put  $\epsilon$  in place of S  
we get aR. So, rule  
3rd used

aRS मा  $\epsilon$  नबरको भए  
2nd rule use कुनैयो  $\text{FIRST}(S)$   
जुङ्गी but  $\text{FIRST}(S)$  जस्ता  
आइसेको चिह्नो दो  
repeat कुनै चिह्नन  $\text{FIRST}(S)$   $\cup$   
 $\text{FIRST}(S)$  को कुनै term मध्ये जुङ्गी  
we do not use  $\epsilon$  in FOLLOW

$$\text{FOLLOW}(S) = \{\text{FOLLOW}(R)\}$$

$$= \{\$, ), +, a, *\}$$

Follow(S) same as follow(R) आद्यो दो  
जाँच, घपेने किनालि set मा value repeat  
कुनैनन / फरका आको भए जाँच FOLLOW(S)  
को ठाउँमा ~~value~~ value दें, घपड्यै।

Example 2: Compute FOLLOW of the following grammar.

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

left side symbols E, E', T, T'. र F को  
Follow फिनाले

right side E नाको production R3U मात्र ह  
सो यो मात्र लिएको F  $\rightarrow (E)$

Solution:

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

$$= \{\$, )\}$$

E starting symbol  
पाले हो तो start से  
अप्पा जाए

E लाई follow जी ब्रेस  
which is terminal. Terminal की  
First terminal नि कुनै।

$$\text{FOLLOW}(E') = \{\text{FOLLOW}(E) \cup \text{FOLLOW}(E')\}$$

$E'$  को follow निकालते हैं,  $E' \rightarrow E$  द्वारा  
but  $E'$  में राहमन मिलेंगे  
i.e., निकालने वाले symbol ए राहमन symbol  
same नहीं  $E$  राहमन हैं।

$$E \rightarrow TE' \text{ तो } E' \text{ की कठि दें।}\\ \text{so, FOLLOW}(E) \\ \text{similar for } E'$$

$\text{follow}(E)$  को value लियको मानिए।  
 $\text{FOLLOW}(E')$  निकालते हैं तो  $\text{Follow}(E')$  को  
निकालते हैं तो  $\text{Follow}(E)$  को लिये।

$$\text{FOLLOW}(T) = \{\text{FOLLOW}(TE') \cup \text{FOLLOW}(+TE') \\ \cup \text{FOLLOW}(T) \cup \text{FOLLOW}(E')\}$$

we have following  
productions for this  
 $E \rightarrow TE'$ ,  $E' \rightarrow +TE'$

But we also have  $E' \rightarrow E$ .  
so putting in above two we  
get two new productions as  
 $E \rightarrow T$ ,  $E' \rightarrow +T$

$$E \rightarrow TE' \\ E' \rightarrow +TE' \\ \text{gives FIRST}(E') \\ \text{by 2nd rule}$$

$$= \{\text{FIRST}(E') \cup \text{FOLLOW}(E) \cup \text{FOLLOW}(E')\} \\ = \{+, \$, )\}$$

$\text{FOLLOW}(E)$  is  $\$$  and  
already calculated above.  
 $\text{FOLLOW}(E')$  is  $\$$  already in  
set so not necessary to  
include

$E \rightarrow T$  gives  $\text{FOLLOW}(E)$   
by 3rd rule

by 3rd rule  
 $E' \rightarrow +T$  gives  
 $\text{FOLLOW}(E')$

$$\text{FOLLOW}(T') = \{\text{FOLLOW}(FT') \cup \text{FOLLOW}(*FT')\}$$

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

value of  $\text{FOLLOW}(T)$   
that we already  
calculated above

by 3rd rule.  
but we can  
neglect this, we  
are calculating  $\text{Follow}(T')$   
itself so.

we have two productions  
for this:  $T \rightarrow FT'$  and  $T' \rightarrow *FT'$   
Also we have  $T' \rightarrow E$  but  
 $T'$  के निकालने की राहमन नहीं  
follow तो उसे नहीं करेंगे as we did  
for  $E'$  above

$$\text{FOLLOW}(F) = \{\text{FOLLOW}(FT') \cup \text{FOLLOW}(*FT') \cup \text{FOLLOW}(F) \cup \text{FOLLOW}(*F)\}$$

$$= \{\text{FIRST}(T') \cup \text{FOLLOW}(T')\} \\ = \{+, *, ), \$\}$$

We have two production  
for  $\text{FOLLOW}(F)$  as;

$T \rightarrow FT'$  and  $T' \rightarrow *FT'$   
but we also have  $T' \rightarrow E$  in grammar  
so, we put and get two more:  
 $T \rightarrow F$  and  $T' \rightarrow *F$

Note: LL(1) Parsing table construct को लागि

FIRST ए FOLLOW find जारी At first left  
recursion ए left factorial द्वारा check जारी, अतः  
directly FIRST ए FOLLOW निकालने / भर पहला left  
recursion ए left factorial reduce जारी (that we already read)  
अब FIRST ए FOLLOW की निकाली same process proceed जारी,

# FOLLOW निकालदा कुछ symbol की, ये symbol production को right side से कुछ  
ठाउंगा परनि नज़ेरिये यसको follow start symbol भर दें मात्र हूँड़ और, भर करें परनि हूँड़ empty.

We have calculated FIRST and FOLLOW, now we can construct LL(1) parsing table easily as in the following examples:

Example 1: Construct LL(1) parsing table of following grammar.

$$R \rightarrow aS \mid (R)S$$

$$S \rightarrow +RS \mid aRS \mid *S \mid \epsilon$$

Solution:

Non-terminals	Terminal Symbols					
	a	(	)	+	*	\$
R	$R \rightarrow aS$	$R \rightarrow (R)S$				
S	$S \rightarrow aRS$		$S \rightarrow \epsilon$	$S \rightarrow +RS$ , $S \rightarrow \epsilon$	$S \rightarrow *S$ , $S \rightarrow \epsilon$	$S \rightarrow \epsilon$

मात्रिकी production हैं।  
इन्हें non-terminal कहते हैं। इनके बारे में लेखा जाएगा।

2)   
Matiiki production or grammar हैं।  
terminal \$ है।  
सबे लेखा जाएगा। परन्तु अपने last मा

4)   
हमें LL(1) parsing table दे दिया गया है। लेकिन यहाँ but here इसमें cell में more than 1 production होती है। so in this type of case, the given grammar is not feasible for LL(1) parser. इसमें cell में 2 production होती है। इसको बेला LL(1) parser को लागी feasible होता है।

Now fill एवं सुरु जाएँ। so R ले किसी production हैं मात्रिकी grammar मा।  $R \rightarrow aS$  व  $R \rightarrow (R)S$  हैं। Production मा कुछ case आउन सकदूँ सज्जा है। जोको आको  $\epsilon$  नभएको।  $\epsilon$  नभएको आस (i.e.,  $A \rightarrow \alpha$  form where  $\alpha$  can be terminals and non-terminals 1 or more than one) यातिबेला FIRST( $\alpha$ ) है। For e.g.  $R \rightarrow aS$  मा  $\epsilon$  नभएको case है, so  $FIRST(aS)$  है। जुन हामिल पटिले ने निकोलका चिये जसको value  $\{a\}$  है। so  $R \rightarrow aS$ , 'R' row and 'a' column मा डुने जाया। same for  $R \rightarrow (R)S$ ,  $S \rightarrow +RS$ ,  $S \rightarrow aRS$  and  $S \rightarrow *S$ .

$\epsilon$  नभएको case है जो, (i.e.,  $A \rightarrow \epsilon$ ) यातिबेला FOLLOW(A) है। For e.g.  $S \rightarrow \epsilon$  है last production हामिल so we see FOLLOW(S) which we already calculated and it has value  $\{\$, ), +, *, a\}$  है, so  $S$  row अभि यो सबे symbol जोको बल  $S \rightarrow \epsilon$  production डुने जाया,

Example 2: Construct LL(1) parsing table for following grammar.

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

Solution:

Non-terminals	Terminal Symbols					
	+	*	(	)	id	#
E			$E \rightarrow TE'$		$E \rightarrow TE'$	
E'	$E' \rightarrow +TE'$			$E' \rightarrow E$		$E' \rightarrow E$
T			$T \rightarrow FT'$		$T \rightarrow FT'$	
T'	$T' \rightarrow E$	$T' \rightarrow *FT'$		$T' \rightarrow E$		$T' \rightarrow E$
F			$F \rightarrow (E)$		$F \rightarrow id$	

## 2. Bottom-up Parser:

ये topic काट पनि कमितमा हुई question  
क्षेत्र नहीं कहिले 3 पनि सोड़ / ये topic मात्र  
shift-reduce parsing, SLR, LR(1), LALR(1) प्रदर्शन

Bottom-up Parser is the parser which generates the parse tree for the given input string with the help of grammar productions by starting from input string and ends on the start symbol, by the reduction process. Reduction is the process of replacing a substring by a non-terminal in bottom-up parsing. It is the reverse process of production.

Example: Consider the grammar

$$S \rightarrow aABe$$

$$A \rightarrow Abc \mid b$$

$$B \rightarrow d$$

Production मात्र left side को symbol  
लाए right side को symbol से replace  
जादू, but Reduction is just opposite.  
rightside को symbol अतः left side  
symbol से replace जाए

Now, the sentence abbcde can be reduced to S as follows:  
abbcde

aAbcde (replacing b by using  $A \rightarrow b$ )

aAde (replacing Abc by using  $A \rightarrow Abc$ )

aABe (replacing d by using  $B \rightarrow d$ )

S (replacing aABe by using  $S \rightarrow aABe$ )

Hence, S is the starting symbol of grammar.

### ④ Handle:

A substring that can be replaced by a non-terminal when it matches its right sentential form is called a handle. If the grammar is unambiguous, then every right-sentential form of the grammar has exactly one handle.

Example: Let we have production as:  $E \rightarrow E + T \mid T$ , Now if we reduce  $E + T$  to  $E$  then,  $E + T$  is handle. Similarly let we have production as:  $F \rightarrow (E) \mid d$ , Now if we reduce  $d$  to  $F$  then,  $d$  is handle.

### a) Shift-Reduce Parsing:

A shift-reduce parser tries to reduce the given input string into the starting symbol. At each reduction step, a substring of the input matching to the right side of a production rule is replaced by non-terminal at the left side of that production rule. The process of reducing the given input string into the starting symbol is called shift-reduce parsing.

A string  $\xrightarrow{\text{Reduced to}}$  the starting symbol.

### Stack Implementation of Shift-Reduce Parser:

There are mainly following four basic operations used in shift-reduce parser:

Shift: This involves moving of symbols from input buffer onto the stack.

Reduce: If the handle appears on top of the stack then, its reduction by using appropriate production rule is done i.e., RHS of production rule is popped out of stack and LHS of production rule is pushed onto the stack.

Accept: If only start symbol is present in the stack and the input buffer is empty then, the parsing action is called accept. Accept means successful parsing is done.

Error: This is the situation in which the parser can neither perform shift action nor reduce action and not even accept action.

### Algorithm:

- Initially stack contains only the sentinel \$, and input buffer contains the input string w\$.

2. While stack not equal to  $\$$  do

a. While there is no handle at the top of the stack, do shift input buffer and push the symbol onto the stack.

b. If there is a handle on the top of the stack, then pop the handle, and reduce the handle with its non-terminal and push it onto stack.

3. Done.

Example 1: Use the following grammar

$$S \rightarrow S+S \mid S*S \mid (S) \mid id$$

Note that  $S \rightarrow F$  and  $S \rightarrow (F)$  are different  
not same

production rule को right side  
तिर अंडे हैं यो सबी handle  
इन तीने  $S+S$  is one,  $S*S$  is another,  
 $(S)$  and  $id$  are also handle.

Perform Shift Reduce parsing for input string:  $id+id+id$

Solution:

→ स्टैक अंडे हैं मात्र कुन्हा  
दुन्हा कुन्हा मात्र कुन्हा  
→ Input buffer मा अंडे हैं  
input string एकदृष्टि  
last मा कुन्हा symbol आपर

Handle हैं stack  
मा वा input buffer  
empty भयो जाने  
यो तिक्का ERROR  
हो यो ERROR लेख  
यदि stop जाने,

→ shift/reduce conflict  
→ यो case मा  $S$  handle  
होइन but  $S+S$   
handle, तो so यो  
यो reduce हो जाने  
कुन्हा  $S \rightarrow S+S$ , फल  
पढ़ने यो गुण मा

accept हुने stack  
मा कुन्हा start symbol मात्र जाने पढ़ि व अंडे हैं व input buffer मा कुन्हा symbol मात्र

Stack	Input Buffer	Parsing Action
$\$$	$id+id+id\$$	Shift $id$
$\$id$	$+id+id\$$	Reduce by $S \rightarrow id$
$\$S$	$+id+id\$$	Shift $+$
$\$S+$	$+id+id\$$	Shift $id$
$\$S+id$	$+id\$$	Reduce by $S \rightarrow id$
$\$S+S$	$+id\$$	Shift $+$
$\$S+S+$	$+id\$$	Shift $id$
$\$S+S+id$	$\$$	Reduce by $S \rightarrow id$
$\$S+S+S$	$\$$	Reduce by $S \rightarrow S+S$
$\$S+S$	$\$$	Reduce by $S \rightarrow S+S$
$\$S$	$\$$	Accept

Parsing action मा  
algorithm  
अनुसार mostly  
shift & Reduce  
operation होती है।

stack को top  
मा handle हैं जाने shift जाने  
handle वा या  
Reduce जाने

shift जाने input  
buffer मा भयो  
first symbol  
stack को top मा  
move जाने

Reduce जाने  
terminal symbol  
यो non-terminal  
वा replace जाने  
rule अनुसार

Example 2: Use the grammar, and perform shift reduce parsing.

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

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

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

Input string:  $id+id^*id$

Solution:

Stack	Input Buffer	Parsing Action.
\$	id + id * id \$	Shift id
\$ id	+ id * id \$	Reduce by $F \rightarrow id$
\$ F	+ id * id \$	Reduce by $T \rightarrow F$
\$ T	+ id * id \$	Reduce by $E \rightarrow T$
\$ E	+ id * id \$	Shift +
\$ E+	id * id \$	Shift id
\$ E+id	* id \$	Reduce by $F \rightarrow id$
\$ E+F	* id \$	Reduce by $T \rightarrow F$
\$ E+T	* id \$	Shift * (OR Reduce by $E \rightarrow T$ ) CONFLICT
\$ E+T*	id \$	Shift id
\$ E+T*id	\$	Reduce by $F \rightarrow id$
\$ E+T*T	\$	Reduce by $T \rightarrow T*T$
\$ E+T	\$	Reduce by $E \rightarrow E+T$
\$	\$	Accept

we find possible ways to accept string in this case

T alone is handle while  $E+T$  is also handle so, this is case of reduce/reduce conflict

we have production  $E \rightarrow T$  so T is handle but if we reduce T by E it will be  $E+E$  but later  $E+E$  can't be reduced which is not valid production and may lead to error so we did shift operation in this type of conflict case.

Practice more questions from tec book

### ④ Conflicts in Shift-Reduce Parsing:

There are two kinds of shift-reduce conflicts:

i) Shift/Reduce Conflict: Here, the parser is not able to decide whether to shift or to reduce. (like in 6th step of example 1)

ii) Reduce/Reduce Conflict: Here, the parser cannot decide which sentential form to use for reduction. (like in 9th step of example 2).

मा आर  
exam  
example परि  
ज्ञान

↳ LR Parser: ↳ shift/reduce JTR  
 ↳ shift/reduce conflict 31/3/2021 ↳ JTR method use JTR conflict आठदेश /

LR parsing is one type of bottom up parsing. It is used to parse large class of grammars. In the LR parsing, L stands for left-to-right scanning of the input. R stands for constructing a right most derivation in reverse. We will study SLR, LR(1) and LALR(1) here in this section.

### LR Parsers: General Structure

The LR algorithm requires stack, input, output and parsing table. In all type of LR parsing, input, output, and stack are same but parsing table is different.

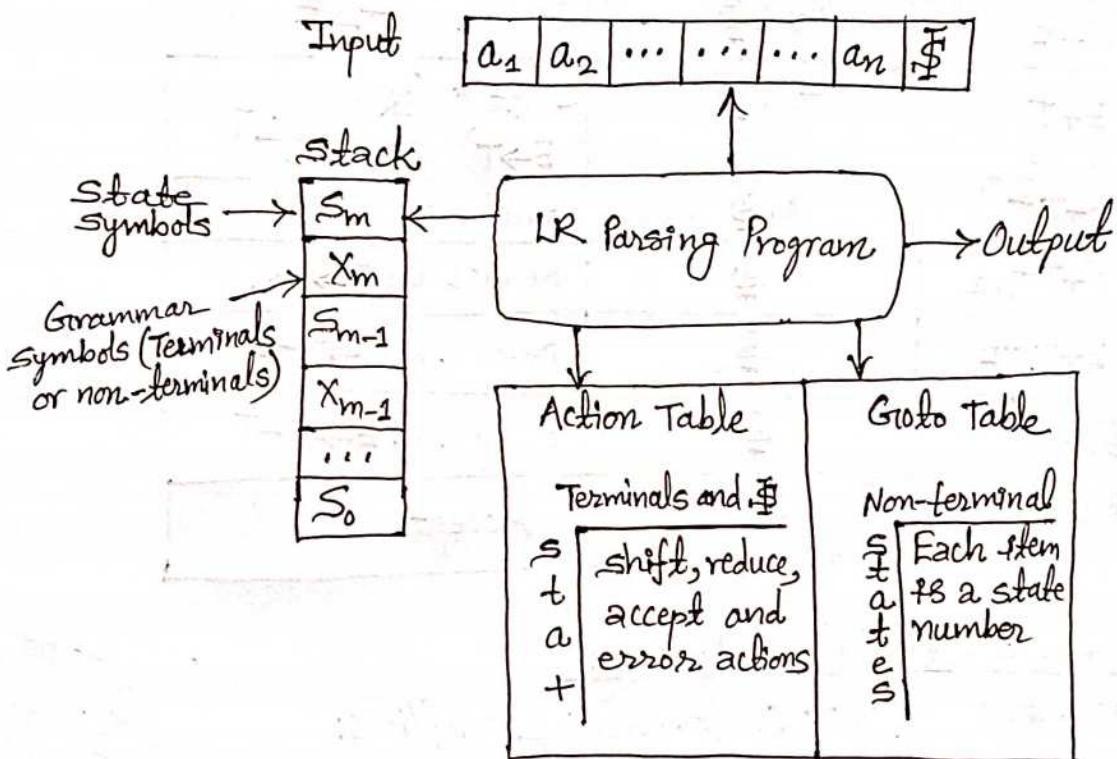


Fig: Block diagram of LR parser.

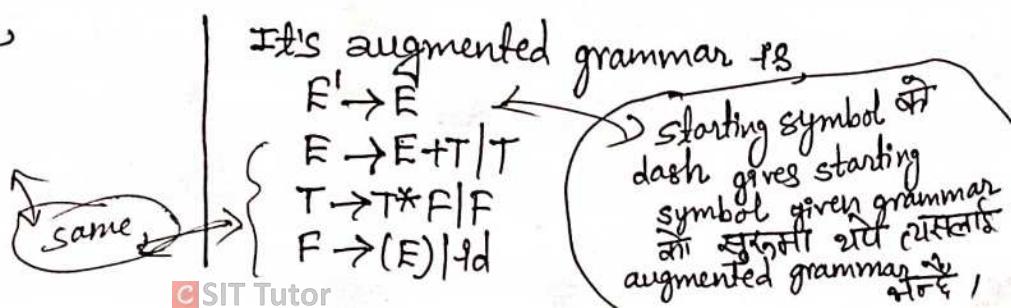
### SLR (Simple LR Parser):

Basic terminologies used for LR parsing table:

Augmented grammar: If  $G_1$  is a grammar with start symbol  $S$ , then the augmented grammar  $G'_1$  of  $G_1$  is a grammar with a new start symbol  $S'$  and production  $S' \rightarrow S$ .

Example: the grammar,

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$



→ LR(0) Item: An 'item' (LR(0) item) is a production rule that contains a dot (•) somewhere in the right side of the production. For example, the production  $A \rightarrow a \cdot AB$  has four items:

$$\begin{aligned} A &\rightarrow \cdot a \cdot AB \\ A &\rightarrow a \cdot \cdot AB \\ A &\rightarrow a \cdot A \cdot B \\ A &\rightarrow a \cdot A \cdot B \end{aligned}$$

↑ production की right side में  
कुन कुन टाइमा • राखने सकिन्दा  
यो सबे possible outcome लाई  
LR(0) item नामिन्दा। • सबे टाइमा  
राखने मिल्दे terminal, non-terminal  
symbol के मतलब भरने

A production  $A \rightarrow \epsilon$ , generates only one item  $A \rightarrow \cdot$ .

→ Closure Operation: If  $I$  is a set of items for a grammar  $G_1$ , then closure( $I$ ) is the set of LR(0) items constructed from  $I$  using the following rules:

1. Initially, every LR(0) item in  $I$  is added to closure( $I$ ). symbol like  $\beta$  which can be terminal or non-terminal
2. If  $A \rightarrow a \cdot B \beta$  is in closure( $I$ ) and  $B \rightarrow Y$  is a production rule of  $G_1$  then add  $B \rightarrow \cdot Y$  in the closure( $I$ ) repeat until no more new LR(0) items added to closure( $I$ ).

Example: Consider the grammar:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T^* F \mid F \\ F &\rightarrow (E) \mid id \end{aligned}$$

It's augmented grammar is;

$$\begin{aligned} E' &\rightarrow E \\ E &\rightarrow E + T \mid T \\ T &\rightarrow T^* F \mid F \\ F &\rightarrow (E) \mid id \end{aligned}$$

If  $I = \{[E' \rightarrow \cdot E]\}$ , then closure  $I$  contains the items,

$$\begin{aligned} E' &\rightarrow \cdot E \\ E &\rightarrow \cdot E + T \\ E &\rightarrow \cdot T \\ T &\rightarrow \cdot T^* F \\ T &\rightarrow \cdot F \\ F &\rightarrow \cdot (E) \\ F &\rightarrow \cdot id \end{aligned}$$

closure निकालदा • को पढ़ाई  
जून symbol द्वायसेल हिन सकने OR  
यो जट दामि कुन कुन production मात्र  
पूजन सकदै, मुखमा • लेखेर यो  
सबे production लेखदै। Foreg we have  
 $E' \rightarrow \cdot E$  initially. We have  $E$  after. so  
we write productions of  $E$  using dot at  
start as  $\cdot E + T$  and  $\cdot T$  again we have  
 $T$  which can produce again and so on.

→ Goto Operation:

If  $I$  is a set of LR(0) items and  $X$  is a grammar symbol (terminal or non-terminal), then goto( $I, X$ ) is defined as follows:

If  $A \rightarrow a \cdot X \beta$  in  $I$  then every item in closure ( $\{A \rightarrow a X \cdot \beta\}$ ) will be in goto( $I, X$ )).

goto निकालदा मुखमा • symbol, 1 step  
पढ़ाई साने then यो पढ़को symbol की  
साथ उसे गरि closure operation जरे यहि हो।

Example:

Set of LR(0) item,

Let,  $I = \{E^1 \rightarrow E^0, E \rightarrow E^0 + T, E \rightarrow E^0 T, T \rightarrow T^0 F, T \rightarrow T^0 * F, F \rightarrow F^0 (E), F \rightarrow F^0 \text{id}\}$

Now,

Now we are finding goto of  $E$ , so  $E$  produce non-terminal term  $E^0$  मात्र लिए रखो set को जो goto जादू! goto means short one step right then do closure

$$\text{goto}(I, E) = \text{closure}(\{[E^1 \rightarrow E^0, E \rightarrow E^0 + T]\}) \\ = \{E^1 \rightarrow E^0, E \rightarrow E^0 + T\}$$

$$\text{goto}(I, T) = \text{closure}(\{[E \rightarrow T^0, T \rightarrow T^0 * F]\}) \\ = \{E \rightarrow T^0, T \rightarrow T^0 * F\}$$

$$\text{goto}(I, F) = \text{closure}(\{[T \rightarrow F^0]\}) \\ = \{T \rightarrow F^0\}$$

$$\text{goto}(I, C) = \text{closure}(\{[F \rightarrow (E)]\}) \\ = \{F \rightarrow (E), E \rightarrow E^0 + T, E \rightarrow E^0 T, T \rightarrow T^0 * F, T \rightarrow T^0 F, F \rightarrow F^0 (E), F \rightarrow F^0 \text{id}\}$$

$E^1 \rightarrow E^0$  मा पढ़िए कुनै term नहीं तो no closure or no other new production occur.  $E \rightarrow E^0 + T$  मा पढ़िए + symbol है तो, terminal symbol को closure करोगा नon-terminal को मात्र तो same form वह closure of that

• पढ़िए non-terminal परिणाम सबको goto मा • पढ़िए opening brace आगे तो वहसको goto

पढ़िए  $E$  इसके closure निकालने मेंलौं means  $E$  वाले जून जून production कुनै भी कोई योग्य स्ट्रक्चर नहीं है।

$$\text{goto}(I, \text{id}) = \text{closure}(\{[F \rightarrow \text{id}^0]\}) \\ = \{F \rightarrow \text{id}^0\}$$

## ↳ Canonical LR(0) collection:

To construct canonical LR(0) collection of grammars we require augmented grammar, closure, and goto functions.

Algorithm:

1. Start

2. Augment the grammar by adding production  $S^1 \rightarrow S$ .

3.  $C = \{\text{closure}(\{S^1 \rightarrow S\})\}$

4. Repeat the followings until no more set of LR(0) items can be added to  $C$ .

for each  $I$  in  $C$  and each grammar symbol  $X$

if  $\text{goto}(I, X)$  is not empty and not in  $C$

add  $\text{goto}(I, X)$  to  $C$ .

an algorithm solving process जिसका मतलब है इनपुट पद्धति

i.e., each LR(0) item की goto function का बहाव देती है।

5. Repeat step 4 until no new sets of items are added to  $C$ .

6. Stop.

Example 1: Find canonical collection of LR(0) items of following grammar.

$$C \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow a$$

Solution:

The augmented grammar of given grammar is;

$$C' \rightarrow C$$

$$C \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow a$$

Next, We obtain the canonical collection of sets of LR(0) items as follows;

$$I_0 = \text{closure}(\{C' \rightarrow \cdot C\}) = \{C' \rightarrow \cdot C, C \rightarrow \cdot AB, A \rightarrow \cdot a\}$$

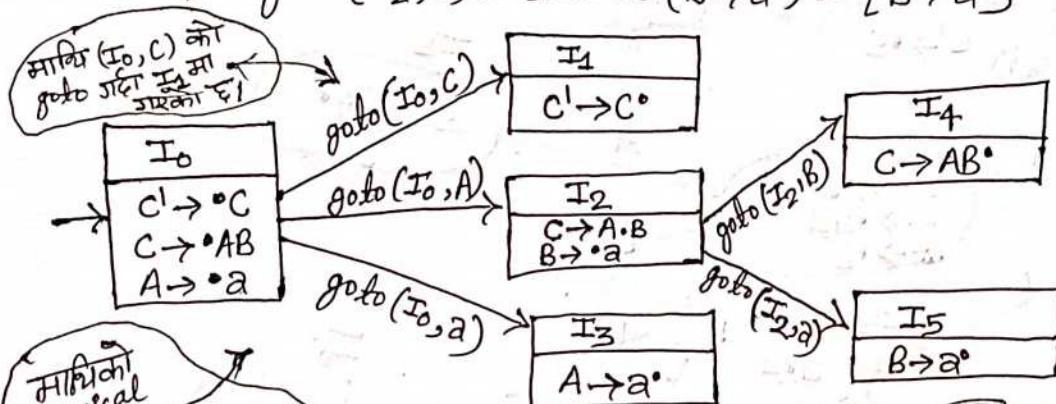
$$I_1 = \text{goto}(I_0, C) = \text{closure}(C' \rightarrow C \cdot) = \{C' \rightarrow C \cdot\}$$

$$I_2 = \text{goto}(I_0, A) = \text{closure}(C \rightarrow A \cdot B) = \{C \rightarrow A \cdot B, B \rightarrow \cdot a\}$$

$$I_3 = \text{goto}(I_0, a) = \text{closure}(A \rightarrow a \cdot) = \{A \rightarrow a \cdot\}$$

$$I_4 = \text{goto}(I_2, B) = \text{closure}(C \rightarrow AB \cdot) = \{C \rightarrow AB \cdot\}$$

$$I_5 = \text{goto}(I_2, a) = \text{closure}(B \rightarrow a \cdot) = \{B \rightarrow a \cdot\}$$



If you see  
canonical  
LR(0)  
collection  
then  
just  
call it  
DFA

सुरक्षा state होई  
इनिले I<sub>0</sub> मानते/  
सुरक्षा state तो  
augmented grammar  
की first production  
की closure निकलते/

Now, I<sub>0</sub> मा • भाको  
symbol •C, •A, •a  
दून्हे, तो अब यो symbol  
हरेको goto निकालो  
र नेहा states I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub>  
मा राखो

अब I<sub>2</sub> मा  
goto calculate, जोरको  
नेहा symbol •B यायो  
तो अब तो set मा भस्का  
सको •B को goto निकालो  
•B र •a को

यो diagram बोलाउना सिफारिश  
information निरापत्ति खाली states  
I<sub>0</sub>, I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub> मात्र राखदा निकलो

Example 2: Find canonical collection of LR(0) items of following grammar

$$\begin{array}{l} S \rightarrow AA \\ A \rightarrow aA \mid b \end{array}$$

Solution:

The augmented grammar of given grammar is;

$$S' \rightarrow S$$

$$S \rightarrow AA$$

$$A \rightarrow aA$$

$$A \rightarrow b$$

Now we obtain the canonical collection of sets of LR(0) items as follows;

$$I_0 = \text{closure}(\{S' \rightarrow \cdot S\}) = \{S' \rightarrow \cdot S, S \rightarrow \cdot AA, A \rightarrow \cdot aA, A \rightarrow \cdot b\}$$

$$I_1 = \text{goto}(\{I_0, S\}) = \text{closure}(\{S' \rightarrow S \cdot\}) = \{S' \rightarrow S \cdot\}$$

$$I_2 = \text{goto}(\{I_0, A\}) = \text{closure}(\{S \rightarrow A \cdot A\}) = \{S \rightarrow A \cdot A, A \rightarrow \cdot aA, A \rightarrow \cdot b\}$$

$$I_3 = \text{goto}(\{I_0, a\}) = \text{closure}(\{A \rightarrow a \cdot A\}) = \{A \rightarrow a \cdot A, A \rightarrow \cdot aA, A \rightarrow \cdot b\}$$

$$I_4 = \text{goto}(\{I_0, b\}) = \text{closure}(\{A \rightarrow b^*\}) = \{A \rightarrow b^*\}$$

$$I_5 = \text{goto}(\{I_2, A\}) = \text{closure}(\{S \rightarrow AA^*\}) = \{S \rightarrow AA^*\}$$

$$I_3 = \text{goto}(\{I_2, a\}) = \text{closure}(\{A \rightarrow a^*A\}) = \{A \rightarrow a^*A, A \rightarrow \cdot aA, A \rightarrow \cdot b\}$$

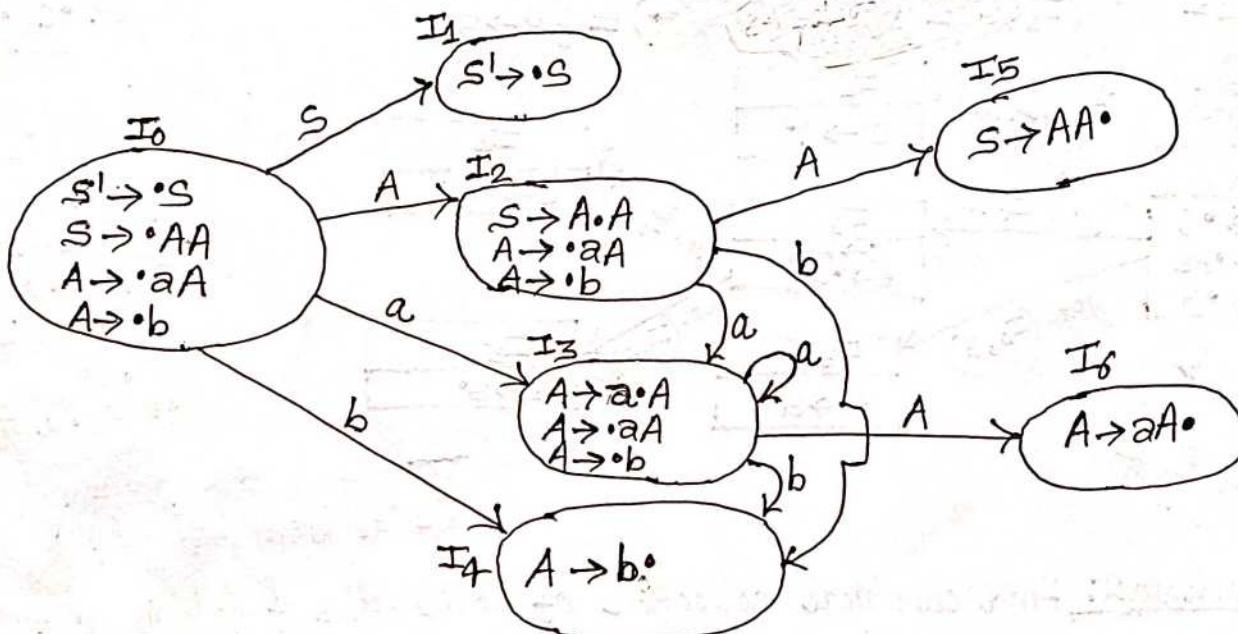
$$I_4 = \text{goto}(\{I_2, b\}) = \text{closure}(\{A \rightarrow b^*\}) = \{A \rightarrow b^*\}$$

$$I_6 = \text{goto}(\{I_3, A\}) = \text{closure}(\{A \rightarrow aA^*\}) = \{A \rightarrow aA^*\}$$

$$I_3 = \text{goto}(\{I_3, a\}) = \text{closure}(\{A \rightarrow a^*A\}) = \{A \rightarrow a^*A, A \rightarrow \cdot aA, A \rightarrow \cdot b\}$$

$$I_4 = \text{goto}(\{I_3, b\}) = \text{closure}(\{A \rightarrow b^*\}) = \{A \rightarrow b^*\}$$

Drawing DFA: For DFA of above canonical LR(0) collection we have 7 states  $I_0$  to  $I_6$  as follows:



### Constructing SLR Parsing Tables:

Algorithm: just to understand

1. Construct the canonical collection of sets of LR(0) items for  $G_1$ .

$$C \leftarrow \{I_0, I_1, \dots, I_n\}$$

2. Create the parsing action table as follows

- If  $A \rightarrow \alpha \cdot a \beta$  is in  $I_i$  and  $\text{goto}(I_i, a) = I_j$  then set action  $[i, a] = \text{shift } j$ .
- If  $A \rightarrow \alpha \cdot$  is in  $I_i$ , then set action  $[i, a]$  to "reduce  $A \rightarrow \alpha$ " for all ' $a$ ' in  $\text{FOLLOW}(A)$ . where,  $A \neq S^*$ .
- If  $S^* \rightarrow S^*$  is in  $I_i$ , then action  $[i, \$] = \text{accept}$ .
- If any conflicting actions generated by these rules, the grammar is not SLR(1).

goto calculate J of first state I1  
it's assign J of because case for SLR

This is same as  $I_3$  state so no need to assign to new state  $I_6$ . Now, the transition  $\text{goto}(I_2, a)$  will go to  $I_3$  as it is same to  $I_3$

Same as  $I_4$

Same as  $I_3$

3. Create the parsing goto table.

• for all non-terminals A, if  $\text{goto}(I_i, A) = I_j$  then  $\text{goto}[i, A] = j$

4. All entries not defined by (2) and (3) are errors.

5. Initial state of the parser contains  $S' \rightarrow \cdot S$ .

Example 1: Construct SLR parsing table of following grammar.

$$C \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow a$$

Solution:

The augmented grammar of given grammar is;

$$C' \rightarrow C$$

$$C \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow a$$

Next, we obtain the canonical collection of sets of LR(0) items as follows,

$$I_0 = \text{closure}(\{C' \rightarrow \cdot C\}) = \{C' \rightarrow \cdot C, C \rightarrow \cdot AB, A \rightarrow \cdot a\}$$

$$I_1 = \text{goto}(I_0, C) = \text{closure}(C' \rightarrow C \cdot) = \{C' \rightarrow C \cdot\}$$

$$I_2 = \text{goto}(I_0, A) = \text{closure}(C \rightarrow A \cdot B) = \{C \rightarrow A \cdot B, B \rightarrow \cdot a\}$$

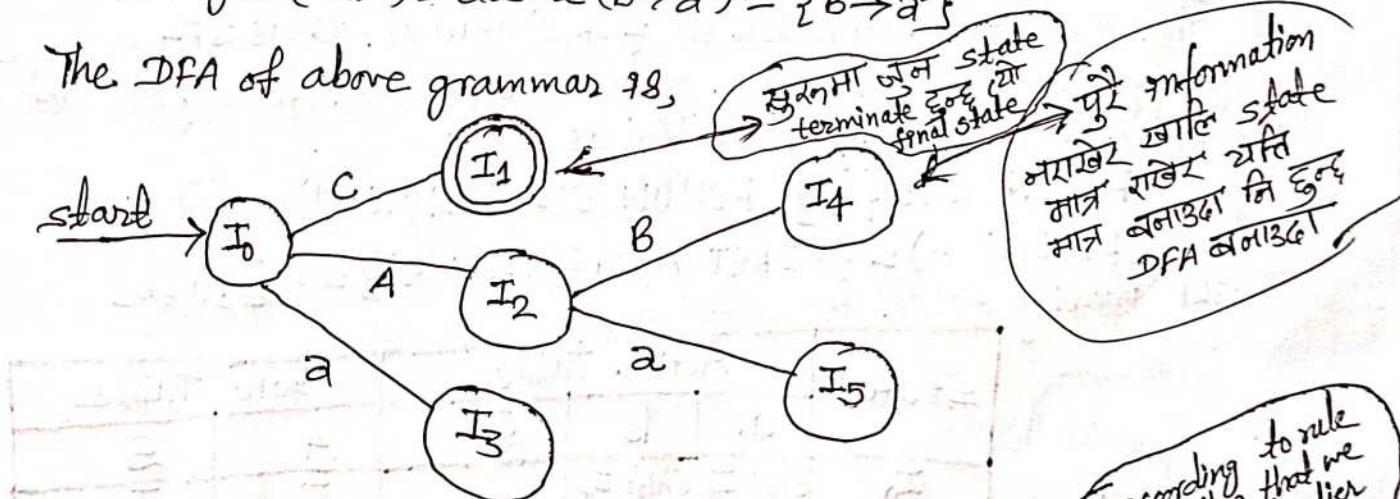
$$I_3 = \text{goto}(I_0, a) = \text{closure}(A \rightarrow a \cdot) = \{A \rightarrow a \cdot\}$$

$$I_4 = \text{goto}(I_2, B) = \text{closure}(C \rightarrow AB \cdot) = \{C \rightarrow AB \cdot\}$$

$$I_5 = \text{goto}(I_2, a) = \text{closure}(B \rightarrow a \cdot) = \{B \rightarrow a \cdot\}$$

These processes are same  
as we already did before  
in example of canonical  
collection of LR(0) items  
all same

The DFA of above grammar is,



Now we calculate FOLLOW function as;

$$\text{FOLLOW}(C') = \{\$\}$$

$$\text{FOLLOW}(C) = \{\text{FOLLOW}(C')\} = \{\$\}$$

$$\text{FOLLOW}(A) = \{\text{FIRST}(B) \cup a\} = \{a\}$$

$$\text{FOLLOW}(B) = \{\text{FOLLOW}(C)\} = \{\$\}$$

grammar मा  
right side C'  
मुख्य हुआ पाइ  
हो तो empty only  
\\$ symbol

according to rule  
of follow that we  
studied earlier  
have a look

Now finally we construct SLR parsing table as below:

States	Action Table			Goto Table	
	a	\$	c	A	B
I <sub>0</sub>	Shift I <sub>3</sub>		I <sub>1</sub>	I <sub>2</sub>	
I <sub>1</sub>		Accept			
I <sub>2</sub>	Shift I <sub>5</sub>				I <sub>4</sub>
I <sub>3</sub>		Reduce R <sub>2</sub>			
I <sub>4</sub>			Reduce R <sub>1</sub>		
I <sub>5</sub>			Reduce R <sub>3</sub>		

(2) I<sub>0</sub>, I<sub>1</sub>... I<sub>n</sub> को set बनाना जैसे की production हैं  
mainly 3 case  
 1) A → a, shift जैसा,  
 2) A → aβ form  
 (i.e., पहली terminal आरे shift जैसा),  
 3) A → aC form (last मा भाँको) र augmented grammar की 1st production आरे reduce A → a जैसा,  
 i.e., FOLLOW(A) है,  
 4) Augmented grammar की first production आरे Accept जैसा,  
 ये 3 case में से कोने पनि लगाए दिए neglect जैसा table मा रखेंगे

(1) Action table मा terminals a र \$, GOTO table मा non-terminals C, A, B  
 state I<sub>5</sub> non-terminal आउट का जाकी है i.e., goto  
 Short मा shift नहीं, Accept नहीं Acc, र Reduce को R ले लिए

(3) सुरक्षा I<sub>0</sub> को set है, C → C र C → AB ले करने परि 3 case satisfy होते हैं तो we neglect them, A → a first case सेट match है। So, I<sub>0</sub> मा a की तरह shift है। DFA मा जो मा a आउट I<sub>3</sub> मा जाने वाला है तो we write shift I<sub>3</sub>. Now we are in I<sub>1</sub> set, we have C → C. which satisfies 3rd case so, \$ की मुदि accept होते हैं according to rule/algorithm, and so on...  
 3rd Reduce वाला case आरे, जसे की produce होते हैं यसको FOLLOW निकाले। यसमें FOLLOW मा जैसे symbol की तरह Reduce ले लेने,

Example 2: Construct the SLR parsing table for the grammar.

$$\begin{array}{l} S \rightarrow AA \\ A \rightarrow aA \mid b \end{array}$$

Solution:

The augment and canonical set of LR(0) item as well as DFA is same as we solved before section in example 2, so we directly construct table here calculating FOLLOW as below:

$$\text{FOLLOW}(S') = \{\$\}$$

$$\text{FOLLOW}(S) = \{\text{FOLLOW}(S')\} = \{\$\}$$

$$\text{FOLLOW}(A) = \{\text{FIRST}(A) \cup \text{FOLLOW}(S)\} = \{a, b\}$$

SLR parsing table is as follows:

States	Action Table			Goto Table	
	a	b	\$	A	S
I <sub>0</sub>	Shift I <sub>3</sub>	Shift I <sub>4</sub>		I <sub>2</sub>	I <sub>1</sub>
I <sub>1</sub>			Accept		
I <sub>2</sub>	Shift I <sub>3</sub>	Shift I <sub>4</sub>			I <sub>5</sub>
I <sub>3</sub>	Shift I <sub>3</sub>	Shift I <sub>4</sub>			I <sub>6</sub>
I <sub>4</sub>	Reduce I <sub>3</sub>	Reduce I <sub>3</sub>	Reduce I <sub>3</sub>		
I <sub>5</sub>			Reduce I <sub>1</sub>		
I <sub>6</sub>	Reduce I <sub>2</sub>	Reduce I <sub>2</sub>	Reduce I <sub>2</sub>		

we can practice more examples from KEC book onwards page no. 46

table बनाओ इसे cell मा more than one action आरे shift, shift I शift, shift & reduce etc. possibility आरे parseable grammar यसलाई ambiguity in SLR मिहाने।

## LR(1) Grammars:

→ 2nd type of LR parser (also called General LR)  
 after shift-reduce parsing. यहां canonical set of LR(1) निकालें  
 instead of LR(0) अब table draw करें, goto निकालें, closure  
 निकालें सभी same as we did in shift-reduce.

LR(1) parsing uses look-ahead to avoid unnecessary conflicts in parsing table.

LR(1) item = LR(0) item + look-ahead.

LR(1) item जैसे को LR(0) पर ही यहां यहां look-ahead जैसे हैं।

LR(0) item  
जौही form को  
हैं।

LR(0) item	LR(1) item
$[A \rightarrow \alpha \cdot \beta]$	$[A \rightarrow \alpha \cdot \beta, a]$

LR(1) item यो form को  
हैं where a is assumed  
look-ahead symbol.

Example 1: Construct LR(1) parsing table of following grammar,

$$S \rightarrow AA$$

$$A \rightarrow aA$$

$$A \rightarrow b$$

Solution: The augmented grammar of given grammar is:

$$S' \rightarrow S$$

$$S \rightarrow AA$$

$$A \rightarrow aA$$

$$A \rightarrow b$$

The canonical collection of LR(1) items of augmented grammar is:

$$I_0 = \text{Closure}(S' \rightarrow^* S) = \{ S' \rightarrow^* S, \$ \}$$

same as LR(0) item. comma (.)  
पर इसकी look-ahead मात्र फरक यहां

मुख्य production मा  
look-ahead नहीं रखें अनि  
प्रस्तुत rule अनुसार  $A \rightarrow [\alpha \cdot \beta, a]$   
देंगा compare जैसे 1st पटिकों  
यहां non-terminal symbol होइए  
जौही string को FIRST निकालें, वो value

2nd production को look-ahead value हैं and so on...

For e.g.  $[S' \rightarrow^* S, \$]$  compared with  $[A \rightarrow \alpha \cdot \beta, a]$ , now

$\text{FIRST}(\$) = \{\$\}$  which is look-ahead of  $S \rightarrow^* AA$

पटिकों द्वारा symbol  
द्वारा जौही को FIRST

$$\begin{aligned} &S \rightarrow^* AA, \$ \\ &A \rightarrow^* aA, a/b \\ &A \rightarrow^* b, a/b \end{aligned}$$

Rule for look-ahead symbol  
 If  $A \rightarrow [\alpha \cdot \beta, a] \in \text{closure}(I)$   
 then, add the item  $[B \rightarrow^* \gamma, b]$   
 to I if not already in I.  
 where,  $b \in \text{FIRST}(\beta)$

similarly  $A \rightarrow^* aA$  production is  
 given by  $S \rightarrow^* AA, \$$ , so,  $\text{FIRST}(A\$)$   
 $= \text{FIRST}(A) = \{a, b\}$ . Hence,  
 look-ahead is a/b.

Similarly  $A \rightarrow^* b$  production परन्तु  $S \rightarrow^* AA, \$$   
 के गांठ भासको हैं so, same  $\text{FIRST}(A) = \{a, b\}$ .

$$I_1 = \text{Goto}(I_0, S) = \text{closure}(S' \rightarrow^* S, \$) = \{ S' \rightarrow^* S, \$ \}$$

$$I_2 = \text{Goto}(I_0, A) = \text{closure}(S \rightarrow A \cdot A, \$) = \{ S \rightarrow A \cdot A, \$ \}$$

copy  
 $A \rightarrow^* aA, \$$   
 $A \rightarrow^* b, \$ \}$   
 यो production  $S \rightarrow A \cdot A, \$$  के  
 गांठ आके so  $\text{FIRST}(\$) = \{ \$ \}$   
 is look-ahead by rule.  
 Similarly this.

$$I_3 = \text{Goto}(I_0, a) = \text{closure}(A \rightarrow a \cdot A, a/b) = \{ A \rightarrow a \cdot A, a/b \}$$

$$\begin{aligned} &A \rightarrow^* aA, a/b \\ &A \rightarrow^* b, a/b \end{aligned}$$

$\beta$  absent तो  
 मात्र जौही हैं तो  
 तल परन्तु look-ahead  
 new calculate  
 नहीं पढ़ें

$$I_4 = \text{Goto}(I_0, b) = \text{closure}(A \rightarrow b^*, a/b) = \{A \rightarrow b^*, a/b\}$$

$$I_5 = \text{Goto}(I_2, A) = \text{closure}(S \rightarrow AA^*, \$) = \{S \rightarrow AA^*, \$\}$$

$$I_6 = \text{Goto}(I_2, a) = \text{closure}(A \rightarrow a^*A, \$) = \{A \rightarrow a^*A, \$\} \quad \begin{array}{l} \beta \text{ absent here} \\ \text{so same } \$, \$ \\ \text{in look-ahead} \\ \text{as before} \end{array}$$

$$I_7 = \text{Goto}(I_2, b) = \text{closure}(A \rightarrow b^*, \$) = \{A \rightarrow b^*, \$\}$$

$$I_8 = \text{Goto}(I_3, A) = \text{closure}(A \rightarrow aA^*, a/b) = \{A \rightarrow aA^*, a/b\}$$

$$\text{Goto}(I_3, a) = \text{closure}(A \rightarrow a^*A, a/b) = \{A \rightarrow a^*A, a/b\}$$

$$A \rightarrow a^*A, a/b$$

$$A \rightarrow b, a/b\}$$

can be assigned  
to  $I_3$  or left  
without assigned  
 $\text{Goto}(I_3, a) = I_3$   
same as  $I_3$

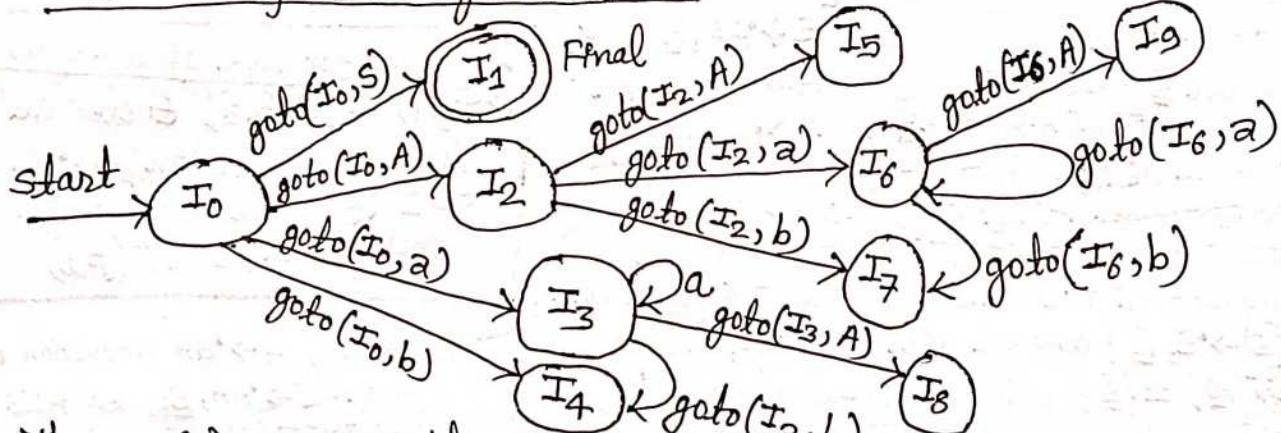
$$\text{Goto}(I_3, b) = \text{closure}(A \rightarrow b^*, a/b) \text{ same as } I_4.$$

$$I_9 = \text{Goto}(I_6, A) = \text{closure}(A \rightarrow aA^*, \$) = \{A \rightarrow aA^*, \$\}$$

$$\text{Goto}(I_6, a) = \text{closure}(A \rightarrow a^*A, \$) \text{ same as } I_6.$$

$$\text{Goto}(I_6, b) = \text{closure}(A \rightarrow b^*, \$) \text{ same as } I_7.$$

The DFA of above grammar is:



The LR(1) parsing table is:

States	Action Table			Goto Table	
	a	b	\$	S	A
0	Shift $I_3$	Shift $I_4$		$I_1$	$I_2$
1			Accept		
2	Shift $I_6$	Shift $I_7$			$I_5$
3	Shift $I_3$	Shift $I_4$			$I_8$
4	Reduce $I_3$	Reduce $I_3$			
5			Reduce $I_1$		
6	Shift $I_6$	Shift $I_7$			$I_9$
7			Reduce $I_3$		
8	Reduce $I_2$	Reduce $I_2$			
9			Reduce $I_2$		

### iii) LALR(1) Grammars:

everything same as LR(1) only difference here is we combine states having same productions, but can have different look-ahead symbol.

Example 1: Construct LALR parsing table for following grammar,

$$S \rightarrow L = R$$

$$S \rightarrow R$$

$$L \rightarrow * R$$

$$L \rightarrow id$$

$$R \rightarrow L$$

= is terminal symbol and  $L = R$  is one single string. Only next string when separated by |

Solution:

The augmented grammar of above grammar is,

$$S' \rightarrow S$$

$$S \rightarrow L = R$$

$$S \rightarrow R$$

$$L \rightarrow * R$$

$$L \rightarrow id$$

$$R \rightarrow L$$

At first we find the canonical collection of LR(1) items of the given augmented grammar as;

$$I_0 = \text{closure}(S' \rightarrow \cdot S, \$)$$

$$= \{ S' \rightarrow \cdot S, \$ \\ S \rightarrow \cdot L = R, \$ \\ S \rightarrow \cdot R, \$ \\ L \rightarrow \cdot * R, \$ \\ L \rightarrow \cdot id, \$ \\ R \rightarrow \cdot L, \$ \}$$

$$I_1 = \text{closure}(\text{goto}(I_0, S))$$

$$= \text{closure}(S' \rightarrow S^*, \$)$$

$$= \{ S' \rightarrow S^*, \$ \\ R \rightarrow L^*, \$ \}$$

$$I_2 = \text{closure}(\text{goto}(I_0, L))$$

$$= \text{closure}(S' \rightarrow S, L)$$

$$= \{ S' \rightarrow S, L \\ (R \rightarrow L^*, \$) \}$$

$$I_3 = \text{closure}(\text{goto}(I_0, R))$$

$$= \text{closure}(S \rightarrow R^*, \$)$$

$$= \{ S \rightarrow R^*, \$ \}$$

$$I_4 = \text{closure}(\text{goto}(I_0, *))$$

$$= \text{closure}(L \rightarrow * \cdot R, =)$$

$$= \{ (L \rightarrow * \cdot R, =), (R \rightarrow \cdot L, =), (L \rightarrow \cdot * R, =), (L \rightarrow \cdot id, =) \}$$

$$I_5 = \text{closure}(\text{goto}(I_0, id))$$

$$= \text{closure}(L \rightarrow id^*, =)$$

$$= \{ L \rightarrow id^*, = \}$$

$$I_6 = \text{closure}(\text{goto}(I_2, =))$$

$$= \text{closure}(S \rightarrow L = R^*, \$)$$

$$= \{ S \rightarrow L = R^*, \$ \\ R \rightarrow \cdot L, \$ \\ L \rightarrow \cdot * R, \$ \\ L \rightarrow \cdot id, \$ \}$$

$$I_7 = \text{closure}(\text{goto}(I_4, R))$$

$$= \text{closure}(L \rightarrow * R^*, =)$$

$$= \{ L \rightarrow * R^*, = \}$$

$$I_8 = \text{closure}(\text{goto}(I_4, L))$$

$$= \text{closure}(R \rightarrow L^*, =)$$

$$= \{ R \rightarrow L^*, = \}$$

$$I_{12} = \text{closure}(\text{goto}(I_6, id))$$

$$= \text{closure}(L \rightarrow id^*, \$)$$

$$= \{ L \rightarrow id^*, \$ \}$$

$$I_{13} = \text{closure}(\text{goto}(I_{11}, R))$$

$$= \{ L \rightarrow * R^*, \$ \}$$

प्रो 12 र 13, 11 पद्धति का इसे मात्रा देने वाली मात्रा जानकारी दिये गए हैं।

$$I_9 = \text{closure}(\text{goto}(I_8, R)) \\ = \text{closure}(S \rightarrow L = R^0, \#) \\ = \{ S \rightarrow L = R^0, \# \}$$

$$I_{10} = \text{closure}(\text{goto}(I_8, L)) \\ = \{ R \rightarrow L^0, \# \}$$

$$I_{11} = \text{closure}(\text{goto}(I_8, *)) \\ = \text{closure}(L \rightarrow * \cdot R, \#) \\ = \{ L \rightarrow * \cdot R, \# \\ R \rightarrow \cdot L, \$ \\ L \rightarrow \cdot * R, \$ \\ L \rightarrow \cdot id, \$ \}$$

यह सभी सब same होंगे as  
before अब जो combine होने  
मात्र different

Now we combine states as follows:

Combine state 4 and 11 as:

$$I_4: \{ (L \rightarrow * \cdot R, =), (R \rightarrow \cdot L, =), (L \rightarrow \cdot * R, =), (L \rightarrow \cdot id, =) \} \\ I_{11}: \{ (L \rightarrow * \cdot R, \$), (R \rightarrow \cdot L, \$), (L \rightarrow \cdot * R, \$), (L \rightarrow \cdot id, \$) \}$$

$$I_{4,11}: \{ (L \rightarrow * \cdot R, =/\$), (R \rightarrow \cdot L, =/\$), (L \rightarrow \cdot * R, =/\$), (L \rightarrow \cdot id, =/\$) \}$$

Combine state 5 and 12 as;

$$I_5: \{ L \rightarrow id^0, = \}$$

$$I_{12}: \{ L \rightarrow id^0, \$ \}$$

$$I_{5,12}: \{ L \rightarrow id^0, =/\$ \}$$

Combine state 7 and 13 as;

$$I_7: \{ L \rightarrow * R^0, = \}$$

$$I_{13}: \{ L \rightarrow * R^0, \$ \}$$

$$I_{7,13}: \{ L \rightarrow * R^0, =/\$ \}$$

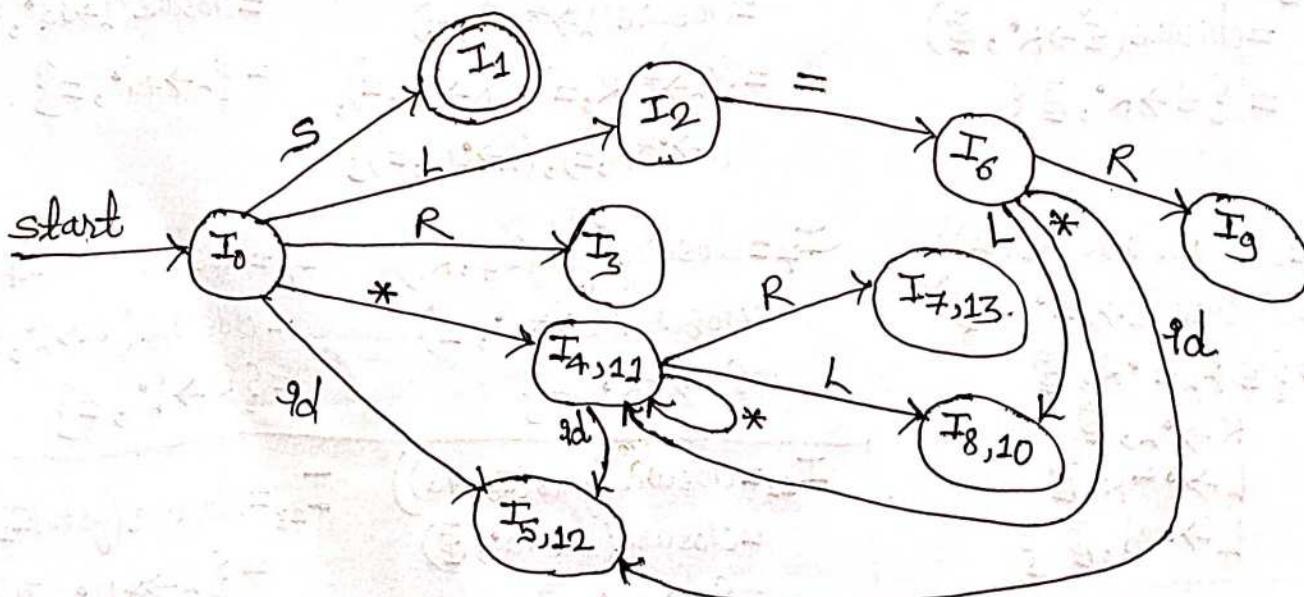
Combine state 8 and 10 as;

$$I_8: \{ R \rightarrow L^0, = \}$$

$$I_{10}: \{ R \rightarrow L^0, \$ \}$$

$$I_{8,10}: \{ R \rightarrow L^0, =/\$ \}$$

Now the DFA of LALR parsing is:



## The LALR parsing table :-

states	Action Table				Goto Table		
	Id	*	=	\$	S	L	R
0	Shift I <sub>5</sub>	Shift I <sub>4</sub>			I <sub>1</sub>	I <sub>2</sub>	I <sub>3</sub>
1				Accept			
2			Shift I <sub>6</sub>	Reduce I <sub>5</sub>			
3				Reduce I <sub>2</sub>			
4	Shift I <sub>5</sub>	Shift I <sub>4</sub>				I <sub>8</sub>	I <sub>9</sub>
5			Reduce I <sub>4</sub>	Reduce I <sub>4</sub>			
6	Shift I <sub>12</sub>	Shift I <sub>11</sub>				I <sub>10</sub>	I <sub>9</sub>
7			Reduce I <sub>3</sub>	Reduce I <sub>3</sub>			
8			Reduce I <sub>5</sub>	Reduce I <sub>5</sub>			
9				Reduce I <sub>1</sub>			

table also same method as before  
 दो इनके अनेको अब I<sub>4,11</sub> हों  
 combine करके भारत similarly others which are combined.

### ④ Kernel and Non-Kernel Items:

Kernel item includes the initial items,  $S' \rightarrow S$  and all items whose dot are not at the left end. Similarly non-kernel items are those items which have their dots at the left end except  $S' \rightarrow S$ .

Example: Find the kernel and non-kernel items of following grammar,

$$C \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow a$$

Solution:

The augmented grammar of given grammar is,

$$C' \rightarrow C$$

$$C \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow a$$

Now, we obtain the canonical collection of sets of LR(0) items as follows,

$$I_0 = \text{closure}(\{C' \rightarrow \cdot C\}) = \{C' \rightarrow \cdot C, C \rightarrow \cdot AB, A \rightarrow \cdot a\}$$

$$I_1 = \text{goto}(I_0, C) = \text{closure}(C' \rightarrow C \cdot) = \{C' \rightarrow C \cdot\}$$

$$I_2 = \text{goto}(I_0, A) = \text{closure}(C \rightarrow A \cdot B) = \{C \rightarrow A \cdot B, B \rightarrow \cdot a\}$$

$$I_3 = \text{goto}(I_0, a) = \text{closure}(A \rightarrow a \cdot) = \{A \rightarrow a \cdot\}$$

$$I_4 = \text{goto}(I_2, B) = \text{closure}(C \rightarrow AB^*) = \{C \rightarrow AB^*\}$$

$$I_5 = \text{goto}(I_2, a) = \text{closure}(B \rightarrow a^*) = \{B \rightarrow a^*\}$$

List of kernel and non-kernel items are listed below:

States	Kernel Items	Non-kernel items
$I_0$	$C^* \rightarrow C$	$C \rightarrow A^*B$ $A \rightarrow a^*$
$I_1$	$C^* \rightarrow C$	
$I_2$	$C \rightarrow A^*B$	$B \rightarrow a^*$
$I_3$	$A \rightarrow a^*$	
$I_4$	$C \rightarrow AB^*$	
$I_5$	$B \rightarrow a^*$	

Augmented grammar  
का सहायता  
production &  
सुरक्षा dot (\*)  
निम्नलिखी सबै  
kernel item  
सुरक्षा production  
 $C^* \rightarrow C$  एक  
सुरक्षा dot (\*) भएको  
सबै non-kernel item

### Top down Parsing vs. Bottom up Parsing:

S.N.	Top down parsing	Bottom up parsing
1.	It is a parsing strategy that first looks at the highest level of the parse tree and works down the parse tree by using the rules of grammar.	It is a parsing strategy that first looks at the lowest level of the parse tree and works up the parse tree by using the rules of grammar.
2.	This parsing technique uses left most derivation.	This parsing technique uses right most derivation.
3.	Its main decision is to select what production rule to use in order to construct the string.	Its main decision is to select when to use a production rule to reduce the string to get the starting symbol.
4.	Error detection is easy	Error detection is difficult.
5.	Parsing table size is small.	Parsing table size is bigger.
6.	Less Power	High Power.

## Semantic Analysis

→ 2 अट 1 short question को 5 marks  
की ओर

④ Introduction: Semantic Analysis is the third phase of compiler which provides meaning to its constructs, like tokens and syntax structure.

CFG + semantic rules = Syntax Directed Definitions.

### Examples of Semantic Errors:

Example 1: Use of a non-initialized variable.

```
int i;  
i++; //the variable i is not initialized.
```

Example 2: Type incompatibility

```
int a = "hello"; //the types String and int are not compatible.
```

Example 3: Errors in expressions

```
char s = 'A';  
int a = 15 - s; //the '-' operator does not support type char.
```

Example 4: Array index out of range.

```
int a[10];  
a[12] = 25; //12 is not legal index, we have index 0 to 9  
for array of size 10.
```

Example 5: Unknown references

```
int a;  
printf("%d", a); //a is not initialized.
```

④ Type Checking: Compiler must check that the source program follows both the syntactic and semantic conventions of the source language. Type checking is the process of checking the data type of different variables. The design of a type checker for a language is based on information about the syntactic construct in the language, the notation of type, and the rules for assigning types to the language constructs.

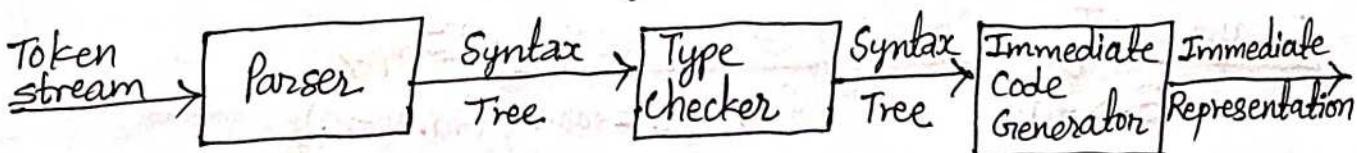


Fig: Position of Type Checker.

④ Type Systems: The collection of different data types and their associated rules to assign types to programming language constructs is known as type systems. It is an informal type system. Rules, for example "if both operands of addition are of type integer, then the result is of type integer." A type checker implements type system.

⑤ Type Expressions: The type of a language construct will be denoted by a "type expression". A type expression is either a basic type or is formed by applying an operator called a type constructor to other type expressions. The sets of basic types and constructors depend on the language to be checked.

Constructors include: *just understand only*

Arrays: If  $T$  is a type expression then  $\text{array}(I, T)$  is a type expression denoting the type of an array with elements of type  $T$  and index set  $I$ . Example:  $\text{array}(0\dots 99, \text{int})$ .

Products: If  $T_1$  and  $T_2$  are type expressions, then their Cartesian product  $T_1 \times T_2$  is a type expression. Example:  $\text{int} \times \text{int}$ .

Records: The record type constructor will be applied to a tuple formed from field names and field types.

Pointers: If  $T$  is a type expression, then  $\text{pointer}(T)$  is a type expression denoting the type "pointer to an object of type  $T$ ." For example,  $\text{var } p: \uparrow$  now declares variable  $p$  have type  $\text{pointer}(\text{row})$ .

Functions: A function in programming languages maps a domain type  $D$  to a range type  $R$ . The type of such function is denoted by the type expression  $D \rightarrow R$ . Example:  $\text{int} \rightarrow \text{int}$  represents the type of a function which takes an  $\text{int}$  value as parameter, and its return type is also  $\text{int}$ .

Example: Type Checking of Expressions: [Imp]

$E \rightarrow \text{id}$

{ $E.\text{type} = \text{lookup}(\text{id}.\text{entry})$ }

$E \rightarrow \text{charliteral}$

{ $E.\text{type} = \text{char}$ }

$E \rightarrow \text{intliteral}$

{ $E.\text{type} = \text{int}$ }

$E \rightarrow E_1 + E_2$	$\{E.type = (E_1.type == E_2.type) ? E_1.type : type\_error\}$
$E \rightarrow E_1 \uparrow$	$\{E.type = (E_1.type == pointer(t)) ? t : type\_error\}$
$E \rightarrow E_1 [E_2]$ <small>This bracket denotes index of array</small>	$\{E.type = (E_2.type == int \text{ and } E_1.type == array(s, t)) ? t : type\_error\}$
$S \rightarrow id = E$	$\{S.type = (id.type == E.type) ? void : type\_error\}$
$S \rightarrow \text{if } E \text{ then } S_1$	$\{S.type = (E.type == boolean) ? S_1.type : type\_error\}$
$S \rightarrow \text{while } E \text{ do } S_1$	$\{S.type = (E.type == boolean) ? S_1.type : type\_error\}$
$S \rightarrow S_1; S_2$	$\{S.type = (S_1.type == void \text{ and } S_2.type == void) ? void : type\_error\}$

### Static vs. Dynamic Type Checking: [Imp]

Static Type Checking	Dynamic Type Checking
The type checking at compilation time is known as static type checking.	The type checking at runtime is known as dynamic type checking.
A language is statically-typed if the type of variable is known at compile time.	A language is strongly-typed, if every program it's compiler accepts will execute without type errors.
It includes languages such as C, C++, C#, Java, FORTRAN, Pascal etc.	It includes languages such as JavaScript, LISP, PHP, Python, Ruby etc.
Typical examples of static checking are: <ul style="list-style-type: none"> <li>→ Type checks</li> <li>→ Flow-of-control checks</li> <li>→ Uniqueness checks</li> <li>→ Name-related checks</li> </ul>	Example: array out of bounds as below; <pre>int a[10]; for(int i=0; i&lt;20; i++)     a[i]=i;</pre>
It is the process of reducing possibilities for bugs in programs and then checking that parts of program have been connected in a consistent way.	It is the process of verifying the type-safety of a program at runtime.

## ② Type Casting vs. Type Conversion:

Type Casting	Type Conversion (Coercion)
In type casting, a data type is converted into another data type by a programmer using casting operators.	In type conversion, a data type is converted into another data type by a compiler.
Type casting can be applied to compatible data types as well as incompatible data types.	Type conversion can only be applied to compatible data types.
In type casting, casting operator is needed in order to cast the data type to another data type.	In type conversion there is no need for a casting operator.
In type casting, the destination data type may be smaller than the source data type, when converting the data type to another data type.	In type conversion, the destination data type can't be smaller than source data type.
Type casting takes place during the program design by programmers.	Type conversion is done at the compile time.

## i) Syntax-Directed Translation (STD):

26.  
STD जैसे describe करा  
आए उसका defn I STD,  
SDTS short AT  
describe  
STD

Syntax-Directed Translation (STD) refers to a method of compiler implementation where the source language translation is completely driven by the parser. Almost all modern compilers are syntax-directed. The general approach to Syntax-Directed Translation is to construct a parse tree or syntax tree and compute the values of attributes at the nodes of the tree by visiting them in same order. There are two ways to represent semantic rules associated with grammar symbols.

→ Syntax-Directed Definitions (SDD).

→ Syntax-Directed Translation Schemes (SDTS).

### ii) Syntax-Directed Definitions (SDD): [Imp]

A syntax-directed definition (SDD) is a context-free grammar together with attributes and rules. Attributes are associated with grammar symbols and rules are associated with productions. A SDD is a generalization of a context-free grammar in which each grammar symbol is associated with a set of attributes.

Example: The syntax directed definition for a simple desk calculator.

Production	Semantic Rules
$L \rightarrow E \text{ return}$	Print (E.val)
$E \rightarrow E_1 + T$	$E.\text{val} = E_1.\text{val} + T.\text{val}$
$E \rightarrow T$	$E.\text{val} = T.\text{val}$
$T \rightarrow T_1 * F$	$T.\text{val} = T_1.\text{val} * F.\text{val}$
$T \rightarrow F$	$T.\text{val} = F.\text{val}$
$F \rightarrow \text{digit}$	$F.\text{val} = \text{digit.lexval}$

any variable  
that does not  
produce further

### ii) Syntax-Directed Translation Schemes (SDTS): [Imp]

SDTS embeds program fragments called semantic actions within production bodies. The position of semantic action in a production body determines the order in which the action is executed. It is used to evaluate the order of semantic rules.

## Syntax

$$A \rightarrow \{ \dots \} \times \{ \dots \} \vee \{ \dots \}$$

where, within the curly brackets we define semantic actions.

Example: A simple translation scheme that converts infix expressions to the corresponding postfix expressions.

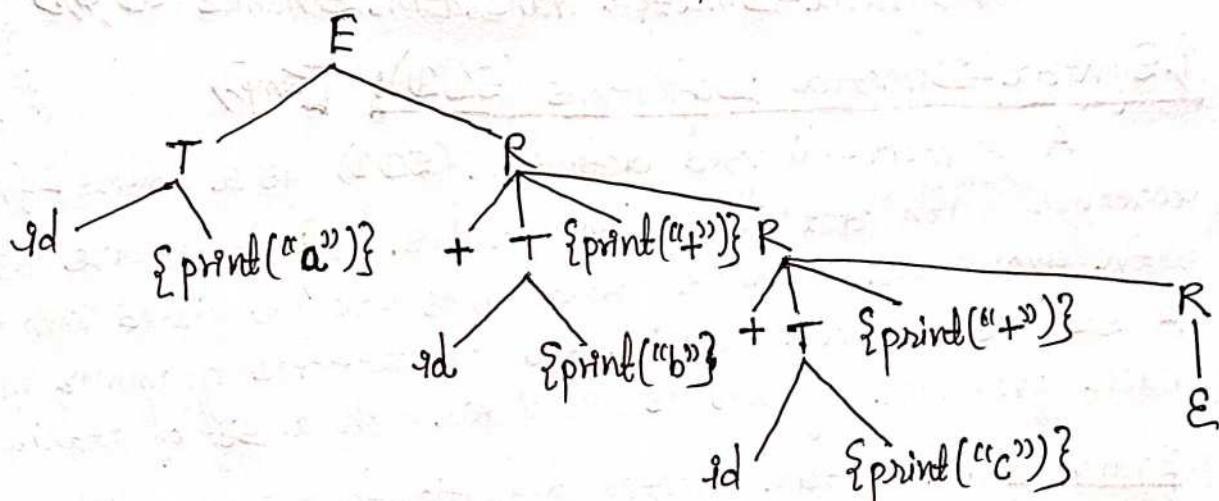
$$E \rightarrow TR$$

$$R \rightarrow +T \{ \text{print} ("+) \} R_1$$

$$R \rightarrow E$$

$$T \rightarrow id \{ \text{print} (\text{id.name}) \}$$

Infix expression  $(a+b+c) \rightarrow$  postfix expression  $(ab+c+)$ .



## ④ ATTRIBUTE TYPES:

1) Synthesized Attribute ( $\uparrow$ ): If the value of the attribute only depends upon its children then it is synthesized attribute. Simply we can also say that the attributes of a node that are derived from its children nodes are called synthesized attributes. Let we have  $S \rightarrow ABC$  as our production, now if  $S$  is taking values from its child nodes ( $A, B, C$ ) then it is said to be a synthesized attribute.

Example for Synthesized Attributes:

Production	Semantic Rules
$L \rightarrow E \text{return}$	$\text{print}(E.\text{val})$
$E \rightarrow E_1 + T$	$E.\text{val} = E_1.\text{val} + T.\text{val}$
$E \rightarrow T$	$E.\text{val} = T.\text{val}$
$T \rightarrow T_1 * F$	$T.\text{val} = T_1.\text{val} * F.\text{val}$
$F \rightarrow \text{digit}$	$F.\text{val} = \text{digit.lexval}$

Same example as SDD

2) Inherited Attribute ( $\rightarrow, \downarrow$ ): A node in which attributes are derived from the parent or siblings of node is called inherited attribute of that node. If the value of the attribute depends upon its parent or siblings then it is inherited attribute. Let we have production as  $S \rightarrow ABC$ , now if 'A' gets values from S, B and C, similarly if B can take values from S, A, C, likewise C can take values from S, A, and B.

Example for Inherited Attributes:

Production	Semantic Rules
$D \rightarrow TL$	$L.in = T.type$
$T \rightarrow int$	$T.type = integer$
$T \rightarrow real$	$T.type = real$
$L \rightarrow L_1, id$	$L_1.in = L.in;$ $addtype(id.entry, L.in)$
$L \rightarrow id$	$addtype(id.entry, L.in)$

④ Synthesized attributes vs. Inherited attributes:

Synthesized attribute	Inherited attributes.
Attributes of a node that are derived from its children nodes are called synthesized attributes.	Attributes which are derived from the parent or siblings of node are called inherited attribute.
The production must have non-terminal as its head.	The production must have non-terminal as a symbol in its body.
It can be evaluated during a single bottom-up traversal of parse tree.	It can be evaluated during a single top-down and sideways traversal of parse tree.
Synthesized attributes can be contained by both the terminals and non-terminals.	Inherited attributes can be only contained by non-terminals.
<u>Example:</u> $E \rightarrow F$ $E.val = F.val$	<u>Example:</u> $E \rightarrow F$ $E.val = F.val$
$E.val$ $\uparrow$ $F.val$	$E.val$ $\downarrow$ $F.val$

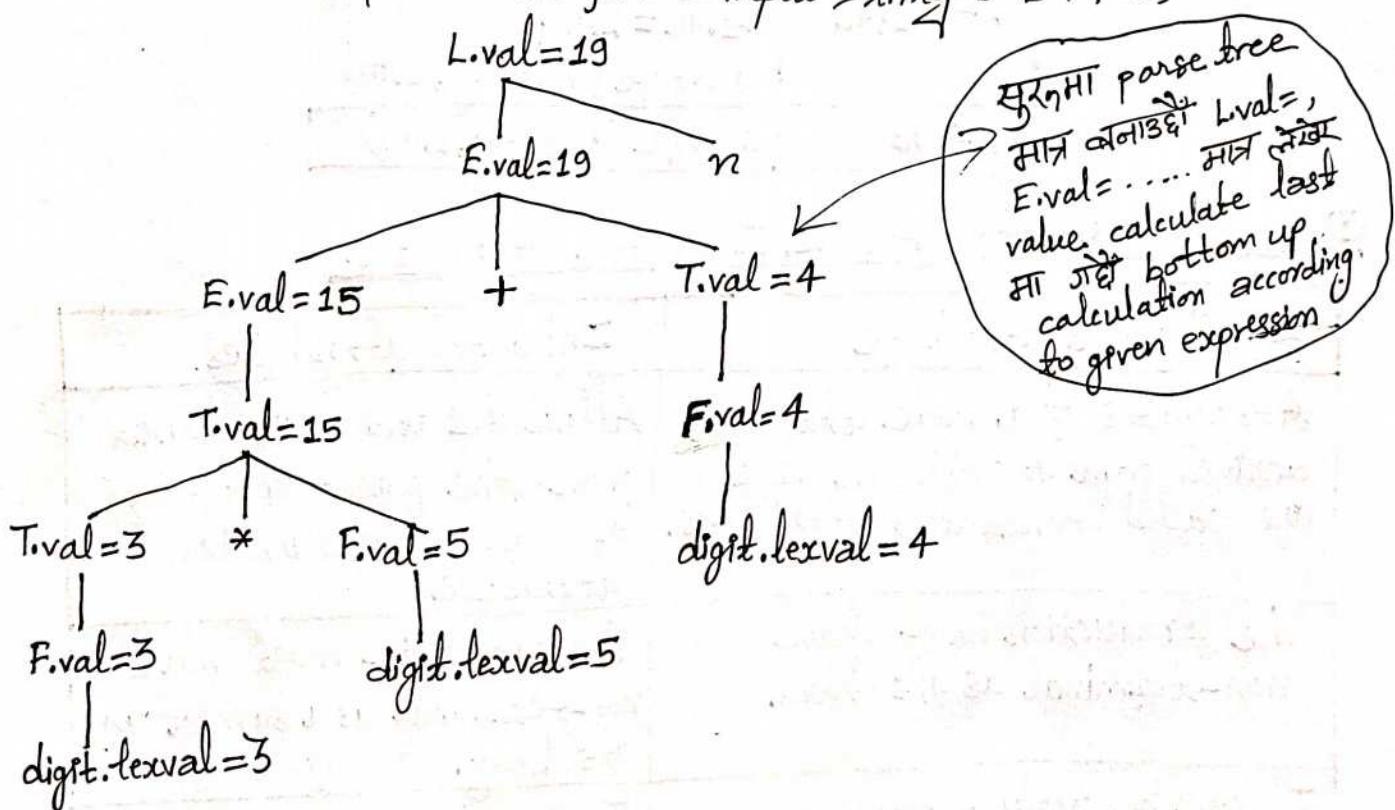
## ④ Annotated Parse Tree:

A parse tree constructing for a given input string in which each node showing the values of attributes is called an annotated parse tree. It is a parse tree showing the values of the attributes at each node. The process of computing the attribute values at the nodes is called annotating or decorating the parse tree.

Example 1: Let's take a grammar,

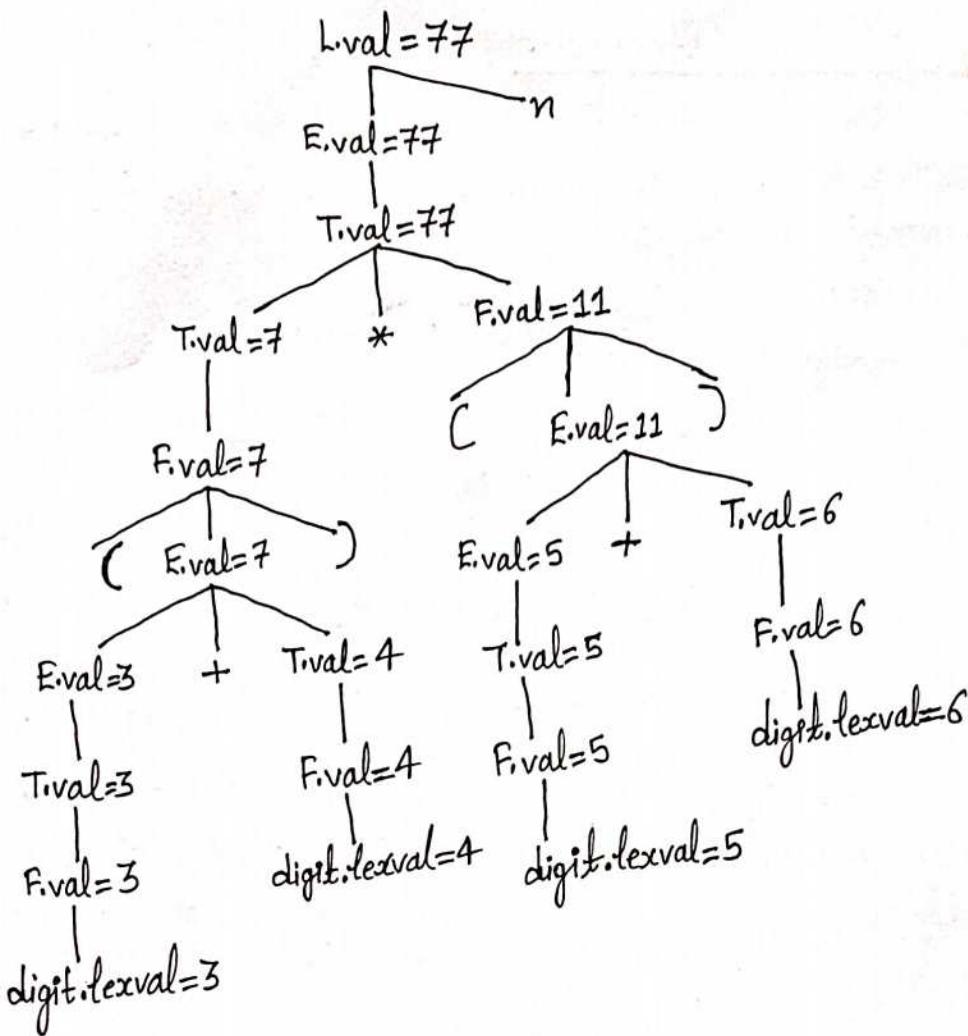
$$\begin{aligned} L &\rightarrow E \cdot n \\ E &\rightarrow E_1 + T \\ E &\rightarrow T \\ T &\rightarrow T_1 * F \\ T &\rightarrow F \\ F &\rightarrow (E) \\ F &\rightarrow \text{digit.} \end{aligned}$$

Now the annotated parse tree for the input string  $3 * 5 + 4 \cdot 98$ ,



Example 2: For the Syntax-Directed Definitions (SDD) of the grammar,  $(3+4)* (5+6)n$ , give annotated parse tree.

Solution:



### ④ S-Attributed Definitions:

If an SDT uses only synthesized attributes, it is called as S-attributed SDT. Attribute values for the non-terminal at the head is computed from the attribute values of the symbols at the body of the production. The attributes of an S-attributed SDT can be evaluated in bottom up order of nodes of the parse tree.

Example for S-attributed definitions:

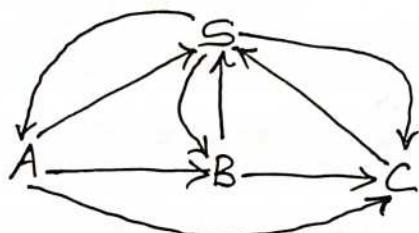
Production	Semantic Rules
$L \rightarrow E \text{ return}$	$\text{print}(E.\text{val})$
$E \rightarrow E_1 + T$	$E.\text{val} = E_1.\text{val} + T.\text{val}$
$E \rightarrow T$	$E.\text{val} = T.\text{val}$
$T \rightarrow F$	$T.\text{val} = F.\text{val}$
$F \rightarrow \text{digit}$	$F.\text{val} = \text{digit}.lexval$

Same example as before

## ④ L-Attributed Definitions:

This form of SDT uses both synthesized and inherited attributes with restriction of not taking values from right siblings. In L-attributed SDTs, a non-terminal can get values from its parent, child, and sibling nodes, as in the following production:

$$S \rightarrow ABC$$



S can take values from A, B, and C (synthesized). A can take values from S only. B can take values from S and A. C can get values from S, A, and B. No non-terminal can get values from the sibling to its right like A does not take value from B or C. Attributes in L-attributed SDTs are evaluated by depth-first and left-to-right parsing manner.

Example:

Production	Semantic Rules
$T \rightarrow FT'$	$T'.inh = F.val$
$T' \rightarrow *FT_1'$	$T_1'.inh = T'.inh * F.val$

UNIT-3

1 question is asked  
for 5 marks from this  
unit

SYMBOL TABLE DESIGN AND RUNTIME STORAGE MANAGEMENT④ Symbol Table Design: Impl upto functions

Symbol Table is an important data structure created and maintained by the compiler in order to keep track of semantics of variable. It is built in lexical and syntax analysis phase. It is used by various phases of compiler as follows:

Lexical Analysis: Creates new table entries in the table, example like entries about token.

Syntax Analysis: Adds information regarding attribute type, scope, dimension, line of reference, use etc. in the table.

Semantic Analysis: Uses available information in the table to check for semantics.

Intermediate Code generation: Helps in adding temporary variable information.

Code Optimization: Uses information present in symbol table for machine dependent optimization.

Target Code generation: Generates code by using address information of identifier present in the table.

Functions of a symbol table:

- To store the names of all entities in a structured form at one place.
- To verify if a variable has been declared or not.
- To determine the scope of a name (scope resolution).
- To access information associated with a given name.
- To add new information with a given name.
- To delete a name or group of names from the table.

## ⊗. Information Used by Compiler from Symbol Table:

Name:  
→ Name of identifier  
→ May be stored directly or as a pointer to another character string in an associated string table names can be arbitrarily long.

Type:  
→ Type of identifier: variable, label, procedure name etc.  
→ For variable, its type: basic types, derived types etc.

Location:  
→ Offset within the program where the current definition is valid.

Other attributes: array limits, fields of records, patterns parameters, return values etc.

Scope:  
→ Region of the program where the current definition is valid.

## ⊗ Basic Operations on a symbol table:

allocate → to allocate a new empty symbol table.

free → to remove all entries and free the storage for symbol table.

insert → to insert a name in a symbol table and return a pointer to its entry.

lookup → to search for a name and return a pointer to its entry.

set\_attribute → to associate an attribute with a given entry.

get\_attribute → to get an attribute associated with a given entry.

## ⊗ Storing Names in Symbol Table:

There are two types of name representation:

1. Fixed Length Name: A fixed space for each name is allocated in symbol table. In this type of storage if name is too small then there is wastage of space.

2. Variable Length Name: The amount of space required by string is used to store the names. The name can be stored with the name of starting index and length of each name.

## ④ Data Structures for Symbol Table:

Following are the commonly used data structures for implementing symbol table:

1) List Data Structure: In this method, an array is used to store names and associated information. A pointer "available" is maintained at end of all stored records and new names are added in the order as they arrive.

Name 1	Info 1
Name 2	Info 2
Name 3	Info 3
:	:
Name n	Info n

Available  
(start of empty slot) →

Fig: list data structure for symbol table.

2) Self Organizing List Data Structure: This implementation is using linked list. A link field is added to each record. A pointer "First" is maintained to point of first record of symbol table. Searching of names is done in order pointed by link of link field.

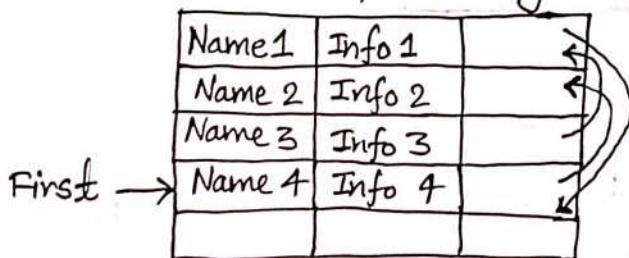


Fig: Self organizing list data structure for symbol table.

3) Binary Trees: When the organization of symbol table is by means of binary tree, the node structure will be as follows:

Left child	Symbol	Information	Right child
------------	--------	-------------	-------------

The left child field stores the address of previous symbol and right child field stores the address of next symbol. The symbol field is used to store the name of symbol, and information field is used to give information about symbol.

4) Hash Tables: In hashing scheme two tables are maintained: A Hash Table and Symbol Table. The hash table consist of  $k$  entries from 0,1, to  $k-1$ . These entries are basically pointers to symbol table pointing to the names of symbol table. To determine whether the 'Name' is in symbol table, we use a hash function ' $h$ ' such that  $h(\text{name})$  will result any integer between 0 to  $k-1$ . We can search any name by position =  $h(\text{name})$

Using this position we can obtain the exact location of name in symbol table.

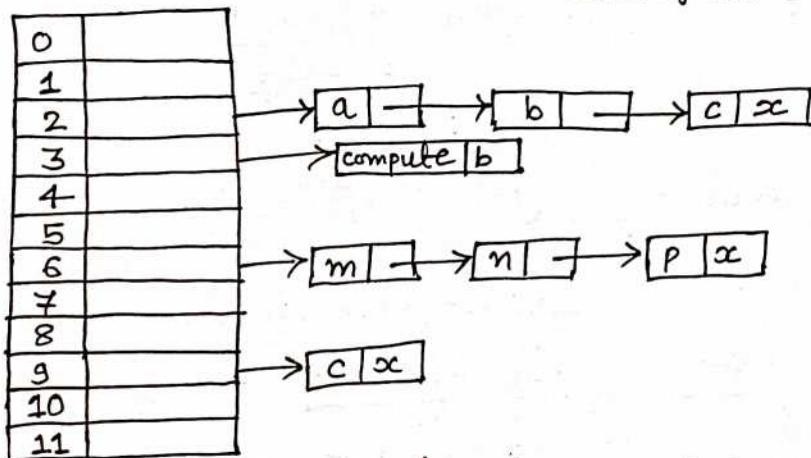


Fig: Hash table data structure for symbol table.

## ② Run-time Storage Management:

Activation record: Activation record is a block of memory used for managing information needed by a single execution of a procedure. The following three-address statements are associated with the run-time allocation and de-allocation of activation records:

1. Call
2. Return
3. Halt
4. Action, a placeholder for other statements.

## Activities performed by caller and callee during procedure call and return:

### On caller's side:

- Save registers for later use.
- Arrange for the procedure arguments to be found by the procedure.
- Call the procedure.
- Copy the return value.
- Throw away procedure arguments.
- restore saved registers
- continue

### On the callee's side:

- Control is transferred from the caller to the starting address.
- Save registers for later use.
- allocate storage for automatic local variables.
- allocate storage for temporaries.
- put return value in appropriate place.
- Throw away local variables and temporaries.
- restore saved registers
- work

Different storage allocation strategies are:-

1) Static Storage Allocation: In static allocation, if memory is created at compile time, memory will be created in the static area and only once. It does not support dynamic data structure i.e., memory is created at compile time and de-allocated after program completion.

The drawback with static storage allocation is recursion is not supported and size of data should be known at compile time.

Its advantage is, it is faster than stack storage allocation.

2) Stack Storage Allocation: Stack allocation is based on the idea of control stack.

→ A stack is Last In First Out storage device where new storage is allocated and de-allocated at only one end, called the top of the stack.

→ Storage is organized as a stack and activation records are pushed and popped as activations begin and end, respectively.

→ Storage for the locals in each call of procedure is contained in the activation record for that call.

→ The values of local are deleted when the activation ends.

→ At run time, an activation record can be allocated and de-allocated by incrementing and decrementing top of the stack respectively.

#### Advantages:

- It supports recursion as memory is always allocated on block entry.
- It allows creating data structure dynamically.

#### Disadvantage:

- Memory addressing can be done using pointers and index registers. Hence, this type of allocation is slower than static allocation.

3) Heap Storage Allocation: The de-allocation of activation records need not occur in a last-in first-out fashion, so storage cannot be organized as a stack. Heap allocation parcels out pieces of contiguous storage, as needed for activation records or other objects. Pieces may be de-allocated in any order. So over time the heap consists of alternate areas that are free and in use. Heap is an alternate for stack.

UNIT-44.1: INTERMEDIATE CODE GENERATOR

There are three parts in this unit, this first part asks 10 marks in exam

⊗ Intermediate Code Generation:

The front end translates the source program into intermediate representation from which the backend generates target code. Intermediate codes are machine independent codes, but they are close to machine instructions.

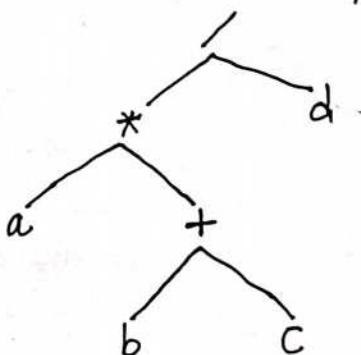
Advantages of using intermediate code representation:

- If a compiler translates the source language to its target machine language without having the option for generating intermediate code, then for each new machine, a full native compiler is required.
- Intermediate code eliminates the need of a new full compiler for every unique machine by keeping the analysis portion same for all compilers.
- The second part of compiler, synthesis, is changed according to the target machine.
- It becomes easier to apply the source code modifications to improve code performance by applying code optimization techniques on the intermediate code.

⊗ Intermediate Representations:1. Graphical Representations:

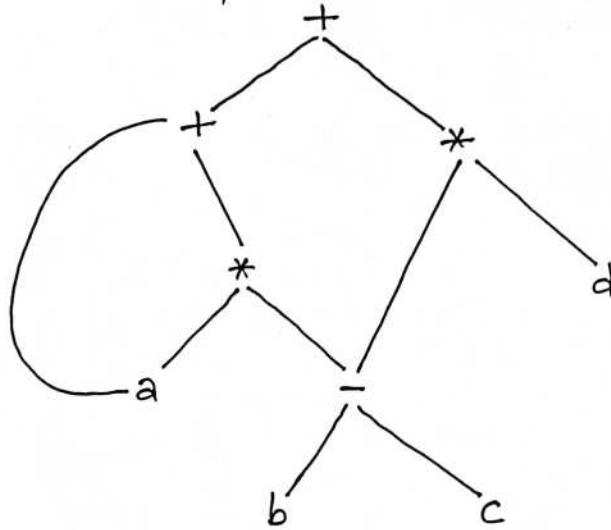
**i) Syntax tree:** Syntax tree is a graphical representation of given source program and it is also called variant of parse tree. A tree in which each leaf represents an operand and each interior node represents an operator is called syntax tree.

Example: Syntax tree for the expression  $a^*(b+c)/d$



1) Directed acyclic graph (DAG): A DAG for an expression identifies the common sub expressions in the expression. It is similar to syntax tree, only difference is that a node in a DAG representing a common sub expression has more than one parent, but in syntax tree the common sub expression would be represented as a duplicate subtree.

Example: DAG for the expression  $a + a * (b - c) + (b - c) * d$ .



2) Postfix notation: The representation of an expression in operators followed by operands is called postfix notation of that expression. In general if  $x$  and  $y$  be any two postfix expressions and  $OP$  is a binary operator then the result of applying  $OP$  to the  $x$  and  $y$  in postfix notation by " $x\ y\ OP$ ".

Examples:

1.  $(a+b)*c$  in postfix notation is:  $ab+c*$
2.  $a*(b+c)$  in postfix notation is:  $abc+*$

3) Three Address Code: The address code that uses at most three addresses, two for operands and one for result is called three address code. Each instruction in three address code can be described as a 4-tuple: (operator, operand<sub>1</sub>, operand<sub>2</sub>, result). A quadruple (three address code) is of the form:

$x = y \text{ op } z$  where  $x, y$  and  $z$  are names, constants or compiler-generated temporaries and  $\text{op}$  is any operator.

We use the term "three-address code" because each statement usually contains three addresses (two for operands, one for the result).

Example 1: Three address code for expression:  $(B+A) * (Y - (B+A))$

Solution:

$$t1 = B+1$$

$$t2 = Y - t1$$

$$t3 = t1 * t2$$

per line JTT maximum 3  
variable हुन्ह परी 1, 2, 3  
वरी नह तो हो

Maths मा जस्तै विधि precedence

जस्तै होसलाई विधि  
convert JTT here  
bracket first.

Example 2: Three address code for expression:

$$g = 2 * n + k$$

while  $g \neq 0$

$$g = g - k$$

Solution:

$$t1 = 2$$

$$t2 = t1 * n$$

$$t3 = t2 + k$$

$$g = t3$$

This first four lines  
for  $g = 2 * n + k$

जस्तै होने 3  
variable होने  
L1, L2, T1

L1: if  $g=0$  goto L2

$$t4 = g - k$$

$$g = t4$$

goto L1

L2: .....

We use two variables  
L1 and L2 when loop comes.  
L1 for true condition and  
L2 for false condition. Here  
instead of  $g=0$  we can also  
write  $g=false$

These three lines are for  
 $g=g-k$  which is body of loop  
goes out of body i.e. L2 when  $i=0$

## ②. Implementation of Three-Address Statements:

Three-address statements are implemented or represented by Quadruples and Triples.

Quadruples: A quadruple is a record structure with four fields, which are op, arg1, arg2 and result. The op field contains an internal code for the operator. The three-address statement  $x = y \text{ op } z$  is represented by placing y in arg1, z in arg2 and x in result.

	OP	arg1	arg2	result
(0)	-	c		t1
(1)	*	b	t1	t2
(2)	-	c		t3
(3)	*	b	t3	t4
(4)	+	t2	t4	t5
(5)	=	t5		a

(0), (1), (2), ...  
numbering nothing  
else

first do three-  
address in rough  
as below then it is  
easy:

$$\begin{aligned} t1 &= -c \\ t2 &= b * t1 \\ t3 &= -c \\ t4 &= b * t3 \\ t5 &= t2 + t4 \\ a &= t5 \end{aligned}$$

Fig: Quadruple representation of three-address statement  $a = b * -c + b * -c$

Triples: A triple is a record structure with three fields, which are [Imp] op, arg1 and arg2. The result field of quadruples is absent here all other things are same, as quadruples. The fields arg1 and arg2, for the arguments of op, are either pointers to the symbol table or pointers into the triple structure.

	op	arg1	arg2
(0)	-	c	
(1)	*	b	(0)
(2)	-	c	
(3)	*	b	(2)
(4)	+	(1)	(3)
(5)	=	a	(4)

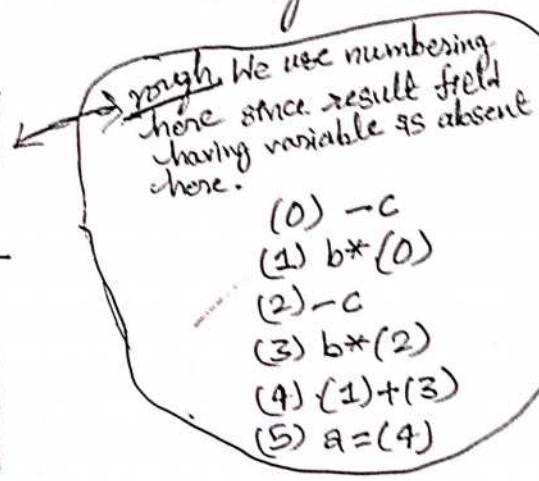


Fig: Triple representation of three-address statement  $a = b * -c + b * -c$ .

Example: Translate the expression  $x = (a+b)* (c+d) - (a+b+c)$

- a. Quadruples
- b. Triples.

Solution: Three address code is:

$$\begin{aligned}
 t1 &= a+b \\
 t2 &= c+d \\
 t3 &= t1 * t2 \\
 t4 &= t1 + c \\
 t5 &= t3 - t4 \\
 x &= t5
 \end{aligned}$$

Quadruples:

	op	arg1	arg2	result
(0)	+	a	b	t1
(1)	+	c	d	t2
(2)	*	t1	t2	t3
(3)	+	t1	c	t4
(4)	-	t3	t4	t5
(5)	=	t5		x

Triples:

	op	arg1	arg2
(0)	+	a	b
(1)	+	c	d
(2)	*	(0)	(1)
(3)	+	(0)	c
(4)	-	(2)	(3)
(5)	=	x	(4)

## Naming conventions for three address code:

S.code: three-address code for evaluating S.

S.begin: label to start of S.

S.after: label to end of S.

E.code: three-address code for evaluating E

E.place: a name that holds the value of E.

To represent three-address statement.

Gen (E.place = E1.place + E2.place)

code generation

$$t3 = t1 + t2$$

जहाँ से एक एकत्रित तौर पर  
expression का three-address  
code बनाया जाता है। इसकी उपर्युक्त  
प्रोडक्शन्स को 3-address code  
कहा जाता है।

## Syntax-Directed Translation into Three-Address Code

### Assignment statements:

#### Productions

$$S \rightarrow id = E$$

// is concatenation operator here

$$E \rightarrow E1 + E2$$

$$E \rightarrow E1 * E2$$

$$E \rightarrow -E1$$

$$E \rightarrow (E1)$$

$$E \rightarrow id$$

three address-code  
variable (i.e., terminal  
तो इसका दाला further  
produce जरूर सकता non-terminal  
को मात्र निभाता)

#### Semantic rules

$$S.place = newtemp();$$

$$S.code = E.code || gen(id.place = E.place);$$

$$E.place = newtemp();$$

$$E.code = E1.code || E2.code || gen(E.place = E1.place + E2.place)$$

$$E.place = newtemp();$$

$$E.code = E1.code || E2.code || gen(E.place = E1.place * E2.place)$$

$$E.place = newtemp();$$

$$E.code = E1.code || gen(E.place = 'minus' E1.place)$$

$$E.place = newtemp();$$

$$E.code = E1.code$$

$$E.place = id.name$$

$$E.code = null$$

gen एक expression की  
3-address code बनाने की use  
होती है।

right side three address  
फलोंका जिक्र non-terminal  
के दृष्टि, so newtemp() जरूर पड़े।

id का तो name आउट पड़े।  
E की place जो बदला  
code फलोंका तो so E.code  
ग्रे null.

## 2) Boolean Expressions:

[Imp]

### Production

$$B \rightarrow B_1 \parallel B_2$$

production मा  $\parallel$  symbol  
आप logical OR होता है।  
Semantic rule मा आप  
मात्र concatenation

### Semantic Rules

$$B_1.\text{true} = B.\text{true}$$

$$B_1.\text{false} = \text{newlabel}()$$

$$B_2.\text{true} = B.\text{true}$$

$$B_2.\text{false} = B.\text{false}$$

$$B.\text{code} = B_1.\text{code} \parallel \text{label}(B_1.\text{false}) \parallel B_2.\text{code}$$

logical OR सा कुनै रक्ती साले  
true नहीं तो true होता है।  
false होता है।  
So, मुख्या  $B_1.\text{true}$  होता है।  
But if  $B_1.\text{false}$  होता है।  
अब  $B$  true हो तो  $B_2$  भी depend  
गई। we do not know  $B_2$  so  
we give it newlabel(). Now  
 $B_2$  is true  $B$  is true  
otherwise false.  
finally we  
write  $\parallel$   
address code.

$$B \rightarrow B_1 \& B_2$$

logical AND.  
 $B$  true होने के  
 $B_1$  र  $B_2$  true होने पर्याप्त।

$$B_1.\text{true} = \text{newlabel}()$$

$$B_1.\text{false} = B.\text{false}$$

$$B_2.\text{true} = B.\text{true}$$

$$B_2.\text{false} = B.\text{false}$$

$$B.\text{code} = B_1.\text{code} \parallel \text{label}(B_1.\text{true}) \parallel B_2.\text{code}$$

$B_1.\text{true}$  होता है  $B$  true होता है।  
जिसे decide जैसे सवेदने, logical  
AND मा  $B$  true होने  $B_1$  र  $B_2$  को  
true होने पर्याप्त। So,  $B$  to be  
true now it depends on value  
of  $B_2$  so newlabel() is given.

$$B \rightarrow !B_1$$

$$B_1.\text{true} = B.\text{false}$$

$$B_1.\text{false} = B.\text{true}$$

$$B.\text{code} = B_1.\text{code}$$

$$E \rightarrow \text{id}_1 \text{ relop } \text{id}_2$$

relation operator (can be  
 $<$ ,  $>$ ,  $=$ ,  $\neq$ ,  $\geq$ ,  $\leq$  etc.)

$$E.\text{code} = \text{gen}(\text{'if' } \text{id}_1.\text{place} \text{ relop } \text{id}_2.\text{place} \\ \text{'goto' } E.\text{true}) \parallel \text{gen}(\text{'goto' } E.\text{false})$$

$$B \rightarrow \text{true}$$

$$B.\text{code} = \text{gen}(\text{goto } B.\text{true})$$

$$B \rightarrow \text{false}$$

$$B.\text{code} = \text{gen}(\text{goto } B.\text{false})$$

## 3) Flow of control statements:

### Semantic Rules

$$B.\text{true} = \text{newlabel}()$$

$$B.\text{false} = S.\text{next}$$

$$S_1.\text{next} = S.\text{next}$$

$$S.\text{code} = B.\text{code} \parallel \text{label}(B.\text{true}) \parallel S_1.\text{code}$$

$$S \rightarrow \text{if}(B) S_1 \text{ else } S_2$$

$$B.\text{true} = \text{newlabel}()$$

$$B.\text{false} = \text{newlabel}()$$

$$S_1.\text{next} = S.\text{next}$$

$$S_2.\text{next} = S.\text{next}$$

$$S.\text{code} = B.\text{code} \parallel \text{label}(B.\text{true}) \parallel S_1.\text{code} \parallel$$

$$\text{gen}(\text{goto } S.\text{next}) \parallel \text{label}(B.\text{false}) \parallel S_2.\text{code}.$$

$S \rightarrow \text{while } (B) S_1$

$\text{begin} = \text{newlabel}()$   
 $B.\text{true} = \text{newlabel}()$   
 $B.\text{false} = S.\text{next}$   
 $S_1.\text{next} = \text{begin}$

$S.\text{code} = \text{label}(\text{begin}) // B.\text{code} // \text{label}(B.\text{true}) //$   
 $S_1.\text{code} // \text{gen(goto begin)}$ .

Example : Generate three address code for the expression

$\text{if } (x < 5 \text{ || } (x > 10 \text{ & } x == y)) x = 3;$

Solution:

L1: if  $x < 5$  goto L2  
else  
  goto L3

L2:  $x = 3$  goto L5

L3: if  $x > 10$   
  goto L4  
else  
  goto L5

L4: if  $x == y$   
  goto L2  
else  
  goto L5

L5: .....

think of if condition in C program  
 This is equivalent to as:  
 $\text{if } (x < 5 \text{ || } (x > 10 \text{ & } x == y))$   
 $\{ x = 3$   
 Logical OR  
 $\};$ .....  
 i.e., if  $x < 5$  it goes to  $x = 3$   
 OR if  $x > 10$  and  $x == y$  then  
 also it goes to  $x = 3$   
 otherwise goes out of if block .....

#### 4) Switch/case statements: [Impl]

Example: Convert the following switch statement into three address code:

Switch ( $i+j$ )

{ Case1:  $x = y + z$

Case2:  $u = v + w$

Case3:  $p = q * w$

Default:  $s = u / v$

}

Solution:

$t = i + j$

→ 3 address code द्वारा  
use करना चाहिए  
कैसे करें

L1: if  $t == 1$  goto L2

else goto L3

L2:  $x = y + z$  goto L8

L3: if  $t == 2$  goto L4

else goto L5

L4:  $u = v + w$  goto L8

L5: if  $t == 3$  goto L6

else goto L7

L6:  $p = q * w$  goto L8

L7:  $s = u / v$  goto L8

L8: .....

## 5) Addressing array elements:

If the width of each array element is  $w$ , then the  $i$ th element of array 'A' begins in location,

$$A[i] = \text{base} + (i - \text{low}) * w$$

where  $\text{low}$  is the lower bound on the subscript and  $\text{base}$  is the relative address of the storage allocated for the array. That is  $\text{base}$  is the relative address of  $A[\text{low}]$ .

Example: Address of 15th element of array is calculated as below,

Suppose base address of array is 100 and type of array is integer of size 4 bytes and lower bound of an array is 10 then,

$$A[i] = \text{base} + (i - \text{low}) * w$$

$$\text{or } A[15] = 100 + (15 - 10) * 4 \\ = 100 + 20 \\ = 120.$$

⇒ Similarly for two dimensional array we have formula,

$$A[i_1, i_2] = \text{base}_A + ((i_1 - \text{low}_1) * n_2 + i_2 - \text{low}_2) * w$$

where  $\text{low}_1, \text{low}_2$  are the lower bounds on values  $i_1$  and  $i_2$ ,  $n_2$  is the number of values that  $i_2$  can take.

Example: Let A be a  $10 \times 20$  array, there are 4 bytes per word, assume  $\text{low}_1 = \text{low}_2 = 1$ . Find addressing array elements  $A[i_1, i_2]$  and three-address code, for 2D array addressing.

Solution:

$$\begin{aligned} A[i_1, i_2] &= \text{base}_A + ((i_1 - 1) * 20 + i_2 - 1) * 4 \\ &= \text{base}_A + (20i_1 - 20 + i_2 - 1) * 4 \\ &= \text{base}_A + (20i_1 - 21 + i_2) * 4 \\ &= \text{base}_A + 80i_1 - 84 + 4i_2 \end{aligned}$$

base<sub>A</sub>,  $i_1$  and  $i_2$  not given so we write them as  $t_1$  &  $t_2$

Now converting into three-address code as follows:

$$t_1 = 80$$

$$t_2 = t_1 * i_1$$

$$t_3 = 4$$

$$t_4 = t_3 * i_2$$

$$t_5 = \text{base}_A + t_2$$

$$t_6 = t_5 - 84$$

$$t_7 = t_6 + t_4$$

$$A[i_1, i_2] = t_7$$

## ④ Procedure Calls:

param  $x$  call  $p$  return  $y$

Here,  $p$  is a function which takes  $x$  as a parameter and returns  $y$ .

Procedure calls have the form:

param  $x_1$

param  $x_2$

... ... ..

param  $x_n$

call  $p, n$

corresponding to the procedure call  $p(x_1, x_2, \dots, x_n)$ .

return statement:

They have the form  $\text{return } y$ , representing a returned value  
is optional.

⑤ Back patching: The easiest way to implement the syntax-directed definitions for boolean expressions is to use two passes. While generating code for boolean expressions and flow-of-control statements during one single pass we may not know the labels that control must go at the time the jump statements are generated. Hence a series of branching statements with the targets of the jumps left unspecified is generated. Each statement will be put on a list of goto statements whose labels will be filled in when the proper label can be determined. We call this subsequent filling in of labels backpatching.

## 4.2 Code Generator

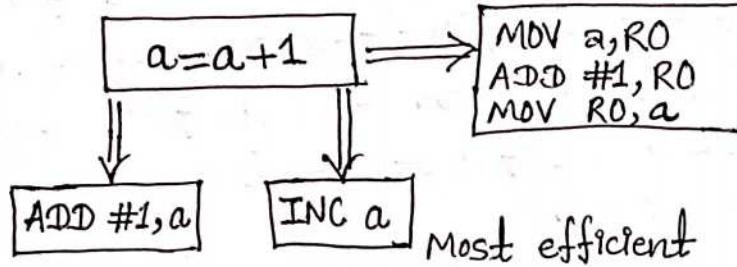
asks for 5 marks in exam  
normally.

### Factors Affecting code generator / code generator design issues: [Impl]

The code generator mainly concern with:

- 1) Input to the code generator: The input to the code generator is intermediate representation together with the information in the symbol table.
- 2) The Target Program: The output of the code generator is target code. Typically, the target code comes in three forms such as: absolute machine language, relocatable machine language and assembly language.  
The advantage of producing target code in absolute machine form is that it can be placed directly at the fixed memory location and then can be executed immediately. The benefit of such target code is that small programs can be quickly compiled.
- 3) The target machine: Implementing code generation requires understanding of the target machine architecture and its instruction set.

- 4) Instruction Selection: Instruction selection is important to obtain efficient code. Suppose we translate three-address code,



Most efficient

- 5) Register Allocation: Since registers are the fastest memory in the computer, the ideal solution is to store all values in registers. We must choose which values are in the registers at any given time. Actually this problem has two parts:

i) Which values should be stored in registers?

ii) Which register should each select to store value?

The reason for the second problem is that often there are register requirements, e.g., floating-point values in floating-point registers and certain requirements for even-odd register pairs for multiplication/division.

6) Evaluation order: The order of the target code can be affected by the order in which the computations are performed.

### ⊗ Basic blocks:

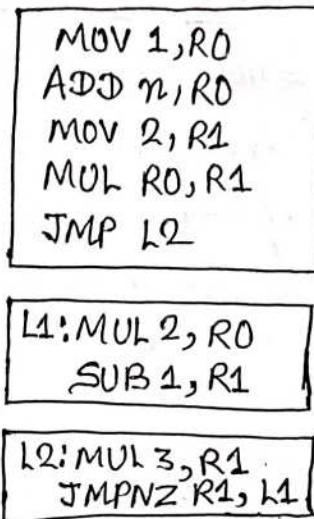
A basic block is a sequence of consecutive instructions in which flow of control enters by one entry point and exit to another point without halt or branching except at the end.

Example:

```
MOV 1, R0  
ADD n, R0  
MOV 2, R1  
MUL R0, R1  
JMP L2
```

=====

```
L1: MUL 2, R0  
SUB 1, R1  
L2: MUL 3, R1  
JMPNZ R1, L1
```



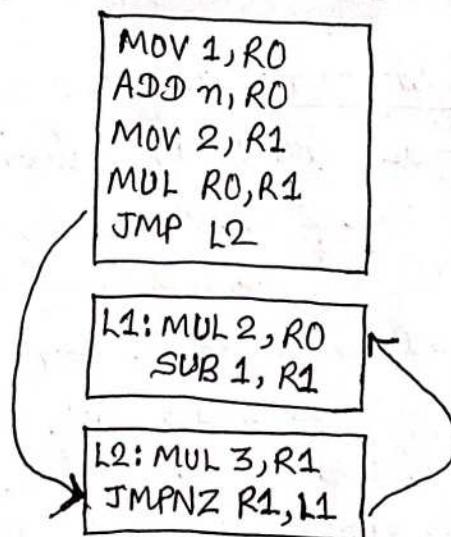
⊗ Flow Graphs: A flow graph is a graphical representation of a sequence of instructions with control flow edges. A flow graph can be defined at the intermediate code level or target code level. The nodes of flow graphs are the basic blocks ~~of given sequence of instructions~~ ~~then such group of blocks~~ and flow-of-control to immediately follow node connected by directed arrow. Simply, if flow of control occurs in basic blocks of given sequence of instructions then such group of blocks is known as flow graphs.

Example:

```
MOV 1, R0  
ADD n, R0  
MOV 2, R1  
MUL R0, R1  
JMP L2
```

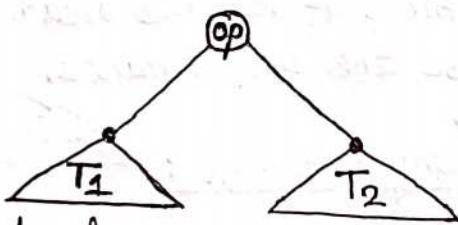
=====

```
L1: MUL 2, R0  
SUB 1, R1  
L2: MUL 3, R1  
JMPNZ R1, L1
```



## \*Dynamic programming code-generation algorithm:

The dynamic programming algorithm can be used to generate code for any machine with  $r$  interchangeable registers  $R_0, R_1, \dots, R_{r-1}$  and load, store, and add instructions. For simplicity, we assume every instruction costs one unit, although the dynamic programming algorithm can easily be modified to work even if each instruction has its own cost.



Contiguous evaluation: Compute the evaluation of  $T_1, T_2$ , then evaluate root.  
Non-contiguous evaluation: First evaluate part of  $T_1$  leaving the value in a register, next evaluate  $T_2$ , then return to evaluate the rest of  $T_1$ .

The dynamic programming algorithm uses contiguous evaluation and proceeds in three phases:

- 1). Compute bottom-up for each node  $n$  of the expression tree  $T$ , an array  $C$  of costs, in which the  $i$ th component  $C[i]$  is the optimal cost of computing the subtree  $S$  rooted at  $n$  into a register assuming  $i$  registers are available for the computation, for  $1 \leq i \leq r$ .
- 2). Traverse  $T$ , using the cost vectors to determine which subtrees of  $T$  must be computed in memory.
- 3). Traverse each tree using the cost vectors and associated instructions to generate the final target code. The code for the subtrees computed into memory locations is generated first.

## 4.3 Code Optimization

The code optimization in the synthesis phase is a program transformation technique, which tries to improve the intermediate code by making it consume fewer resources (i.e., CPU, Memory) so that faster-running machine code will result.

Optimization of the code is often performed at the end of the development stage since it reduces readability and adds code that is used to increase the performance.

### ④ Need of code optimization / why optimize? :

Code optimization def  
and need or importance  
is Impl.

Code optimization is needed because optimization helps to:

- Reduce the space consumed and increases the speed of compilation.
- Manually analyzing datasets involves a lot of time. Hence we make use of software like Tableau for data analysis. Similarly manually performing the optimization is also tedious and is better done using a code optimizer.
- An optimized code often promotes re-usability.

### ⑤ Criteria of code optimization:

- The optimization must be correct, it must not, in any way, change the meaning of program.
- Optimization should increase the speed and performance of the program.
- The compilation time must be kept reasonable.
- The optimization process should not delay the overall compiling process.

## ★ Basic optimization techniques: [Imp]

Optimizations techniques are classified into two categories:

- Machine Independent optimization techniques
- Machine dependent optimization techniques.

### 1) Machine Independent optimization techniques:

Machine independent optimization techniques are program transformations that improve the target code without taking into consideration any properties of the target machine.

a) Constant Folding: As it's name suggests, it involves folding the constants. The expressions that contain the operands having constant values at compile time are evaluated. Those expressions are then replaced with their respective results.

Example: Circumference of Circle =  $(22/7) \times \text{Diameter}$

Here,

- This technique evaluates the expression  $22/7$  at compile time.
- The expression is then replaced with its result  $3.14$ .
- This saves time at run time.

b) Constant Propagation: In this technique, if some variable has been assigned some constant value, then it replaces that variable with its constant value in the further program during compilation. The condition is that the value of variable must not get altered in between.

Example:  $\pi = 3.14$ , Radius = 10, Area of circle =  $\pi \times \text{Radius} \times \text{Radius}$ .

- Here,
- This technique substitutes the variables ' $\pi$ ' and ' $\text{radius}$ ' at compile time.
- It then evaluates the expression  $3.14 \times 10 \times 10$
- The expression is then replaced with its ~~result~~ result  $314$ .
- This saves the time at run time.

c) Redundant Code Elimination: In computer programming, redundant code is source code or compiled code in a computer program that is unnecessary such as; code that is never executed (unreachable code).

Example:  $x = y + z$

$$a = 2 + y + b + z$$

That is reduced by,

$$a = 2 + x + b$$

d) Variable Propagation: Variable propagation means use of one variable instead of another.

Example:  $x = r$

$$A = \pi x^2$$

That is reduced by,  $A = \pi r^2$

e) Strength reduction: It involves reducing the strength of expressions. This technique replaces the expensive and costly operators with the simple and cheaper ones.

Example:

Code before optimization	Code after optimization
$B = A \times 2$	$B = A + A$

→ This is because the cost of multiplication operator is higher than that of addition operator.

f) Loop Optimization:

Defn: Loop optimization is the process of increasing execution speed and reducing the overheads associated with loops. Following are loop optimization techniques.

g) Code Motion/Frequency Reduction: It is a technique which moves the code outside the loop.

Example: `while (i < 5000)`

$$x = i * \sin(A) * \cos(A);$$

This can be optimized as;  
 $t = \sin(A) * \cos(A);$   
`while (i < 5000)`  
 $x = i * t;$

h) Loop jamming/Loop Fusion: In loop fusion method several loops are merged to one loop.

Example: `for i=1 to n do`

`for j=1 to m do`

$$a[i, j] = 10.$$

It can be written as;  
`for i=1 to n*m do`  
 $a[i] = 10$

i) Loop Unrolling: The number of jumps and tests can be reduced by writing code two times.

Example: `int i=1;`

`while (i <= 100)`

$$\{ a[i] = b[i];$$

`i++;`

`}`

It can be written as;

`int i=1;`

`while (i <= 100)`

$$\{ a[i] = b[i];$$

`i++;`

$$a[i] = b[i];$$

`i++;`

`}`

## 2) Machine Dependent Optimization Techniques:

Machine dependent optimization techniques are based on register allocation and utilization of special machine-instruction sequences. It is done after the target code has been generated and when the code is transformed according to the target machine architecture.

Peephole Optimization: Peephole optimization is a simple and effective technique for locally improving target code. This technique is applied to improve performance of the target program by examining the short sequence of target instructions (called the peephole) and replace these instructions replacing by shorter or faster sequence whenever possible. Peephole is a small, moving window on the target program. Peephole optimization techniques are as follows:

a) Redundant Load and Store Elimination: In this technique the redundancy is eliminated.

Example:

Initial code:

$$\begin{aligned}y &= x+5; \\q &= y; \\z &= i; \\w &= z*3;\end{aligned}$$

Optimized code

$$\begin{aligned}y &= x+5; \\q &= y; \\w &= y*3;\end{aligned}$$

b) The Flow of Control Optimization:

Dead Code Elimination: It involves eliminating the dead code i.e., statements of code which never executes or are unreachable.

Example:

Code before optimization

$$\begin{aligned}q &= 0; \\ \text{if } (q = 1) \\ \{ \\ \quad a &= x+5; \\ \}\end{aligned}$$

Code after optimization

$$q = 0;$$

iii) Avoid jump on jump: The unnecessary jumps can eliminate in either the intermediate code or the target code.

We can replace the jump sequence:

Goto L1

L1: goto L2

By the sequence:

Goto L2

L1: goto L2

c) Algebraic Simplification: Peephole optimization is an effective technique for algebraic simplification. The statements such as  $x = x + 0$  or  $x = x * 1$  can be eliminated by peephole optimization.

d) Machine idioms: The target instructions have equivalent machine instructions for performing some operations. Hence we can replace these target instructions by equivalent machine instructions in order to improve the efficiency.

Example: Some machines have auto-increment or auto-decrement addressing modes.

e) Strength reduction: It involves reducing the strength of expressions. This technique replaces the expensive and costly operators with the simple and cheaper ones.

Example: Code before optimization:

$$B = A \times 2$$

Code after optimization:

$$B = A + A$$

# # Compiler Design and Construction (Marking Scheme) With Imp Topics:

Unit 1: (3 hrs): 5 marks (Analysis phase, Synthesis phase, compiler vs interpreter, one-pass vs. Multipass compiler)

Unit 2: (22 hrs):

2.1 Lexical Analysis: 10 marks (lexemes, patterns, tokens, thomsons construction, subset construction)

2.2 Syntax Analysis: 20 marks

2.3 Semantic Analysis: 10 marks

Unit 3: (4 hrs):

Symbol Table Design: 5 marks

Run-time storage management: 5 marks

Unit 4: (16 hrs):

very short unit  
everything imp  
read all

role of syntax analyzer, left recursion, left factoring, LL(1) parsing table, handle, LR(0) item, canonical LR(0) collection, SLR parsing tables, canonical LR(1) collection, LR(1), LALR(1), Top-down vs bottom-up parsing

→ Type checking of expressions, Static vs. dynamic type checking, STD, STDS

4.1 Immediate Code Generator: 10 marks

Role of immediate code generator, Syntax tree, DAC, Postfix notation, finding 3 address code of expressions and statements

4.2 Code Generator: 5 marks

code generator design issues

4.3 Code Optimization: 5 marks

→ defn of code optimization with its need, basic optimization techniques



**If my notes really helped  
you, then you can support  
me on esewa for my  
hardwork.**

**Esewa ID: 9806470952**