

cos

Induction variable

Variable in a loop updated in a regular pattern at each iteration.

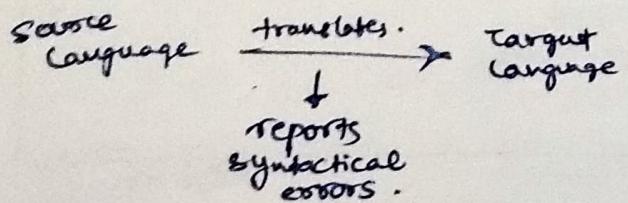
Static Runtime Environment	Stack Based Runtime Environment
Memory allocation	Compile time
Memory lifespan	Entire program execution
Memory organization	Fixed segments
Recursion	Not suitable less flexible.
	Ideal for recursion. More flexible.

CDI

Introduction: The structure of a compiler, Phases of Compilation, The translation process, Major Data Structures in a compiler, Bootstrapping and Postlinking.

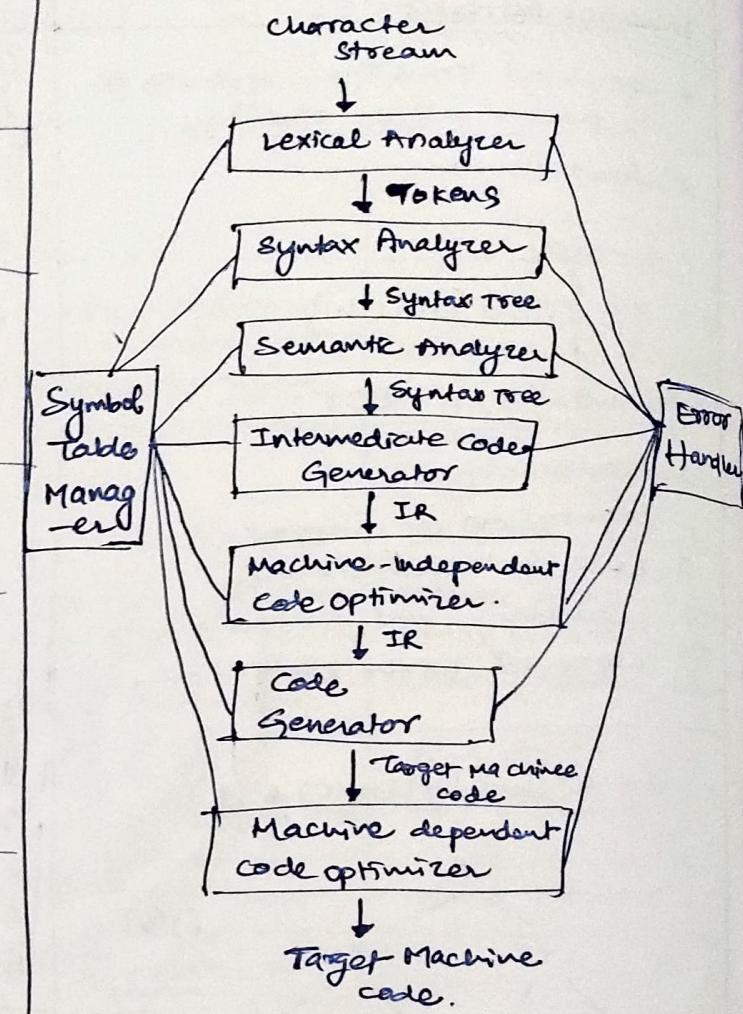
Lexical Analysis (Scanner): The Role of the Lexical Analyzer, Input Buffering, Specification of Tokens, Recognition of Tokens, The lexical Analyzer Generator Lex.

Compiler: A system software which is used to translate the program written in one programming language into machine understandable language to be executed by a computer. It should make the target code efficient and optimized in terms of time & space.



Structure of Compiler

- * Compiler operates as a sequence of phases.
- * each phase transforms one representation of the source program to another.
- * Decomposition of a Compiler.



PHASES OF COMPILER

Symbol Table: Which stores information about the entire source program and is used by all phases of the compiler.

Phases of Compilation:

- | | |
|---------------------------------|---------------------------|
| * Lexical Analysis. | * Symbol Table Management |
| * Syntax Analysis. | * Error Handling. |
| * Semantic Analysis. | |
| * Intermediate code generation. | |
| * Code optimization. | |
| * Code Generation. | |

Expt. No. _____

Date _____

Lexical Analysis (Scanning)

- * First phase of compiler.
- * reads the stream of characters that make up the source program starting from left to right
- * And groups the characters into meaningful sequences called lexemes.
- * For each lexeme, it produces a token of the form:

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

abstract symbol
(used during syntax analysis)
points to an entry in the symbol table.

- * These tokens are passed to syntax analyzer.

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

Ex: $P = i + r * 60$

$P, =, +, *, 60$ are all separate lexemes.

$\langle id, 1 \rangle \leftarrow \langle id, 2 \rangle \leftarrow \langle id, 3 \rangle \leftarrow \langle \dots \rangle$

↓
tokens generated...

Syntax Analysis:

- * Second phase.
- * parser uses the tokens produced by the lexical analyzer to create a tree like intermediate representation that depicts the grammatical structure of the token stream.
- * Syntax tree in which each interior node represents an operation and the children of the node represent the arguments (operands) of the operation.

Ex: $p = i + r * 60$

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

Semantic Analysis:

- * third phase
- * uses syntax tree and the information in the symbol table to check for consistency with the language definition.
- * It performs type conversion of all the data types into real data types.
- * imp part is type checking.
- * Some language specification may permit some type conversions called coercions.

Ex: Int to float.

Kmec

Intermediate Code Generation

* fourth phase

- * In process of translating a source program into target code, a compiler may construct one or more intermediate representations.
- * After semantic analysis of the source program, many compilers generate an explicit low-level or machine-like intermediate representation of the source program.
- * three-address code consists of a sequence of assembly-like instructions with three operands per instruction. Each operand can act like a register.

Ex: $p = i + r * 60$

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

Three address code :

$t1 = \text{intofloat}(60)$
 $t2 = id3 * t1$
 $t3 = id2 + t2$
 $id1 = t3.$

Code Optimization :

- * ~~sixth phase~~.
- * It gets the intermediate code as input and produces optimized intermediate code as output.
- * attempts to improve the intermediate code so that better target code is generated.
- * faster, shorter code, or target code that consumes less power.

Ex :- $P = i + r * 60$

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

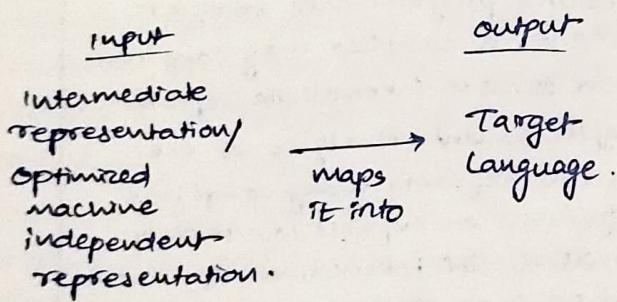
Optimized code is :

$t1 = id3 * 60.0$

$id1 = id2 + t1$

Code Generation :

- * ~~sixth phase~~.



- * If target language is machine code, registers or memory locations are selected for each of the variables used in the program.
- * Then, the intermediate instructions are translated into sequences of machine instructions that perform the same task.

Ex :- $P = i + r * 60$

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

$LDF R2, id3$
 $MULF R2, R2, #60.0$
 $LDF R1, id2$
 $ADD R1, R1, R2$

STF id1, R1.

where R1, R2 are registers.

Symbol Table Manager

- * It is used to store all the information about identifiers used in the program.
- * It is a data structure containing a record for each identifier, with fields for the attributes of the identifier.
- * It allows to find the record for each identifier quickly and to store or retrieve data from that record.

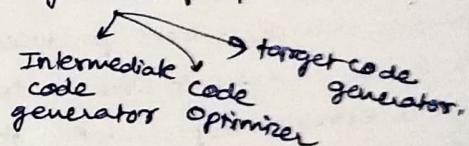
Error Handling:

- * mp function is the detection and reporting of errors in the source program.
- * The error message should allow the programmer to determine exactly where the errors have occurred.
- * Each phase can encounter errors. After detecting an error, a phase must handle the error so that compilation can proceed.

Classification of Compilers:

- * are classified as single-pass, multipass, load-and-go, debugging, or optimizing.
- ⇒ The processes of translation of program to machine understandable code is divided into 2 parts.

- 1) Lexical Part → lexical analysis, syntax analysis,
- 2) Synthetic part. → semantic analysis



The Translation Process

Ex :- $\text{Position} = \text{initial} + \text{rate} * 60$



Lexical Analyzer

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



Syntax Analyzer



Major Data Structures in a Compiler:

Tokens:

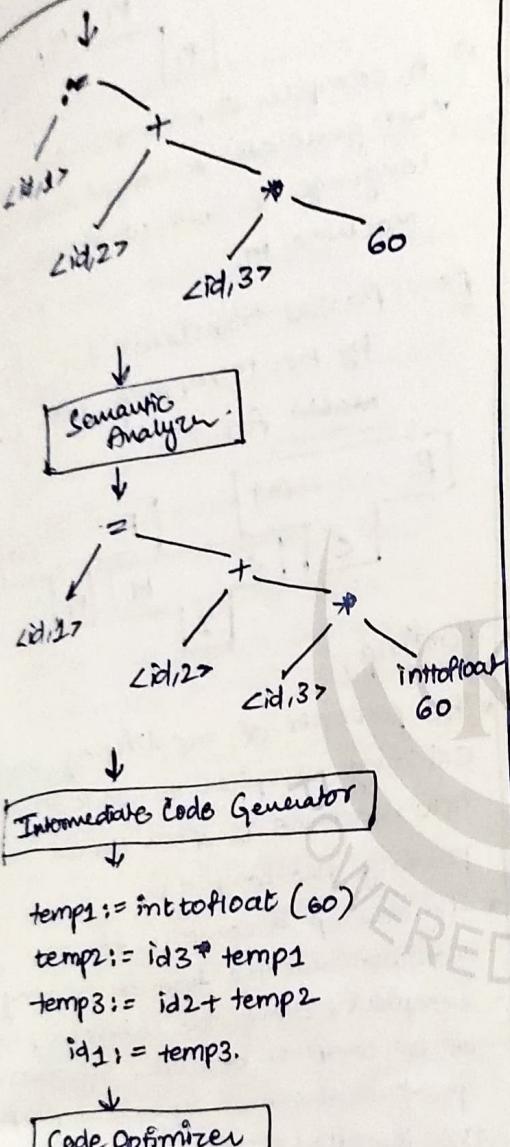
- ⇒ Scanner collects characters/lexeme into a token.
- ⇒ It represents the token symbolically as a value of an enumerated data type representing a set of tokens of the source language.
- ⇒ Sometimes, a single global variable can be used to hold the token information.
- ⇒ In other cases, an array is required to hold the lexeme @ token.

Syntax Tree:

- ⇒ Parser generates Syntax Tree.
- ⇒ It is constructed as a standard pointer-based structure that is dynamically allocated.
- ⇒ Entire tree can be kept as a single variable pointing to the root.
- ⇒ Each node is a record.

Symbol Table:

- ⇒ created and maintained by compiler in order to store information about the occurrences of various entities such as variable names, function names, objects, classes, interfaces etc..
- ⇒ To store the names of all entities in a structured form at one place.
- ⇒ To verify if a variable has been declared.
- ⇒ To implement type checking, by verifying assignments and expressions in the source code are semantically correct.
- ⇒ To determine the scope of a name, (scope resolution).



temp1 := inttofloat (60)
 temp2 := id3 * temp1
 temp3 := id2 + temp2
 id1 := temp3.

↓
Code Optimizer

Temp1 := id3 * 60.0
 id1 := id2 + temp1

↓
Code Generator

MOVF H3, R2
 MULF *60.0, R2
 MOVF id2, R2
 ADDF T2, R2
 MOVF T2, id1

Literal Table:

- Stores constant and strings used in the program.
- One literal table applies globally to the entire program.
- Used by code generator to:
 - Assign addresses for literals.
 - Enter data definitions in the target code file.
- Avoids the replication of constants and strings.
- Quick insertion and lookup are essential.
- Deletion is not allowed.

Temporary files:

- Some memory during compilation is used to store temporary files to hold the products of intermediate steps.

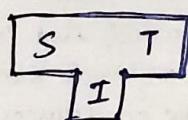
Bootstrapping:

- A compiler is characterized by 3 languages:
 1. Source Language (S).
 2. Target language (T).
 3. Implementation Language (I).

$S \rightarrow$ Translated to $\rightarrow T$
with the help
of (I).

- represented in the form of a

T-diagram $\Leftrightarrow S \begin{smallmatrix} T \\ I \end{smallmatrix}$

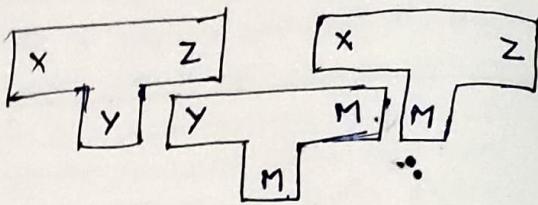


- When a translator is required for an old language on a new machine, or a new language on an old machine, use of existing compilers on either machine is the best choice for developing.

- Developing a compiler for a new language by using an existing compiler is called bootstrapping.

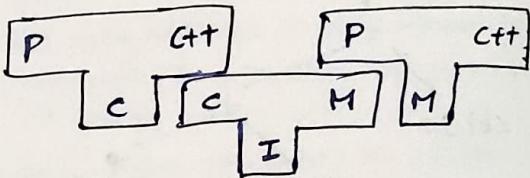
- A process by which a single language is used to translate more complicated programs, which intern may handle an even more complicated program.

- Bootstrapping is used to create compilers and to move them from one machine to another by modifying the backend.



- A compiler for source language X that generates a target code in language Z and which runs on machine M.

Ex:- Pascal translator for C++.
by bootstrapping with C++
machine for M.



Porting:

- The process of modifying an existing compiler to work on a new machine is often known as porting the compiler.

- To develop a compiler for new hardware machine from an existing compiler, change the synthesis part of the compiler because, synthesis part is machine dependent part. This is called Porting.

Quick and Dirty Compiler:

- Native Compiler: Compilers that generate code for the same platform on which it runs.

High language \longrightarrow computer's native language.

- Cross Compiler: Compiler capable of creating executable code for a platform other than the one which the compiler is running.

Lexical Analysis [Scanner]

Tokens:

- Keywords / Reserve words / predefined words [lexemes]
- Identifiers: user defined strings (ID)
- Special Symbols.
- * Any value associated to a token is called as an attribute of the token.
String is numerical.
- * `getNextToken()` command, causes the lexical analyzer to read characters from its input until it can identify the next lexeme and produce for it the next token, which it returns to the parser.

The Role of Lexical Analyzer:

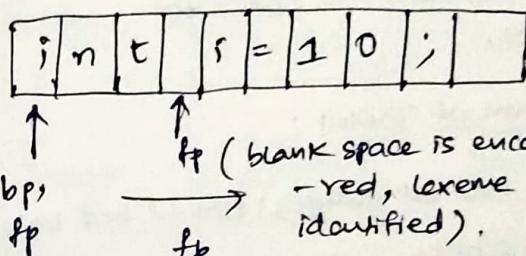
- * To read the source program character by character and form them into logical units called as tokens.
- * These tokens are given as inputs to next phase of compiler.
- * Scanner rarely converts the entire program into tokens at once, but conversion always depends on the parser.
- * Lexical analyzer interacts with the symbol table when it discovers a lexeme constituting an identifier and enters that lexeme into the symbol table.
- * Identification of lexemes.
- * Stripping out comments and white space.
- * Correlating error messages generated by the compiler with the source program.

Input Buffering

Expt. No. _____

Date _____

- * Scanner scans the input from left to right one character at a time.
- ⇒ begin ptr (bp)
- ⇒ forward ptr (fp).



(blank space is ignored and both pointers are placed at the next character).

- * Input characters are always read from secondary storage, but this reading is costly.
- * To speed up the scanning process, buffering technique is used.
- * A block of data is first read into a buffer and lexical analysis process is continued on buffer.

Buffering techniques:

1. Buffer (one/two buffers are used).
2. Sentinels.

Buffering:

One buffer scheme:

- * Only one buffer is used to store the input string.
- * Problem: 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: overcomes one buffer Scheme problem.

* two buffers are used to store the input string.

* first and second buffer are scanned alternately.

Sentinel:

* Special character introduced at the end of the buffer is called as Sentinel which is not a part of program.

* 'eof' is the natural choice for sentinel.

Specification of Tokens:

1. Strings and Languages: ASCII and Unicode

* Alphabets

* Strings.

* Special Symbols

* Language.

* Longest match rule.

* Operations. (union, concat, Kleene closure)

* Notations '(L(σ) U L(δ))

* Regular Expression. (x^* , x^+)

* Finite automata. ($(Q, \Sigma, q_0, \delta, q_f)$).

Recognition of Tokens:

* Patterns are created by using regular expression.

* These patterns are used to build a piece of code that examines the input strings to find a prefix that matches the required lexemes.

Ex:- Stmt \rightarrow If expr then Stmt
| If expr then else Stmt
| E

expr \rightarrow term relop term
| term

term \rightarrow id
| number.

digit \rightarrow [0-9]

digits \rightarrow digit

number \rightarrow digits(. digits)?

(E [+ -] ? digits) ?

letter \rightarrow [A-Za-z]

id \rightarrow letter (letter | digit)*

if \rightarrow if

then \rightarrow then

else \rightarrow else

relop \rightarrow < | > | <= | >= | = | <>

) with the help of these definitions, lexical analyzer will recognize tokens.

for white space:

ws \rightarrow (blank / tab / newline) *

The Lexical Analyzer Generator Lex

* we can also produce a lexical analyzer automatically by specifying the lexeme patterns to a lexical analyzer generator and compiling those patterns into code that functions as a lexical analyzer.

* Also called Lex (flex).

Lex input file consists of three parts:

1. A collection of definitions.

2. A collection of rules.

3. A collection of auxiliary routines
(or) user routines.

* All the sections are separated by double percent signs. Default layout of a lex file is :

{ definitions }

%/%

{ rules }

%/%

{ auxiliary routines }.

Expt. No. _____

Date _____

declaration section :

- * declarations of variables, identifiers and regular definitions.

obj declarations

obj

Rules Section :

- * The translation rules each have the form [Rule_i & Action_i].
- * Each rule is a regular expression.
- * The actions are fragments of code, typically written in C.

obj

Rule₁ { Action₁ }

Rule₂ { Action₂ }

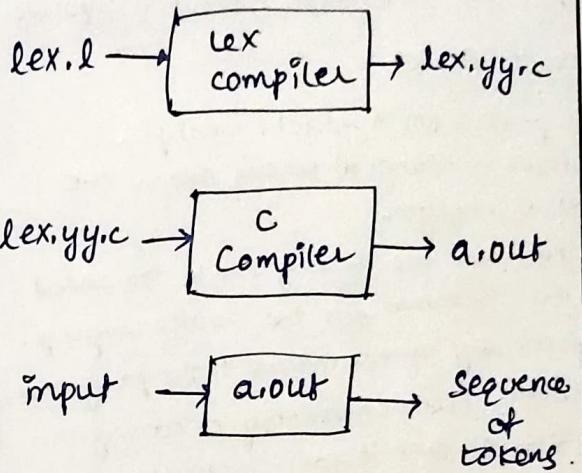
obj

Auxiliary Routines :

- * contain additional functions that are required.
- * These functions may be compiled separately and loaded with the lexical analyzer.

yyflex() } are predefined
yywrap() } procedures of LEX.

- * Lex program (or) Lex files are saved with '.l' extension (dot l).



lex predefined functions and variable

int yyflex() → call to invoke lexer,
returning token.

char *yytext → pointer to matched string.

yylen → length of matched string.

ECHO → write matched string.

CD2:

Syntax Analysis (Parser): The role of the parser, Syntax Error handling and Recovery, Top-Down Parsing, Bottom-up parsing, Simple LR parsing, More powerful LR parsing, using Ambiguous Grammars, Parser Generator Yacc.

Syntax Analysis (Parsers)

- * Second phase.
 - * It gets the input from the tokens and generates a syntax tree or parse tree.
- Parsing: Determining the syntax or the structure of a program. Parse tree specifies the statements execution sequence.

The Role of Parser:

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

Functions of the parser:

1. Verifies the structure generated by the tokens.
2. It constructs the parse tree.
3. It reports the errors.
4. It performs error recovery.

Issues:

- Parser cannot detect errors such as:
1. Variable re-declaration.
 2. Variable initialization before use.
 3. Data type mismatch for an operation.
- These can be handled by Semantic Analysis phase.

- * There are 3 general types of parsers for grammars:

[I. Universal] [II. top-down] [III. bottom-up]

Universal parsing methods can parse any grammar.

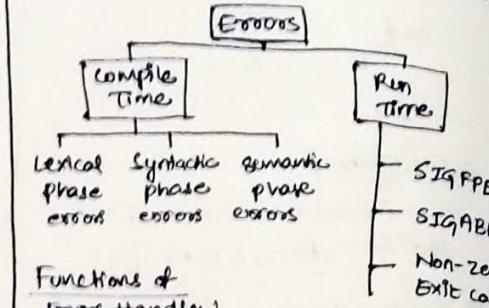
i. the Cocke-Younger-Kasami algorithm.
ii. Earley's algorithm.

Commonly used methods in compilers are either top-down or bottom-up.

Syntax Error Handling and Recovery:

Errors at many different levels. Ex.

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



Functions of error handler:

- * Should report the presence of errors clearly and accurately.
- * Should recover from each error quickly.
- * Should not significantly slow down the processing of correct programs.

Error Recovery Strategies:

1. Panic mode
 2. Phrase level
 3. Error productions
 4. Global correction.
- ↳ Synchronizing token, delimiter.
- ↳ augmented grammar.
- ↳ local correction.
- ↳ no. of insertions, deletions and changes of tokens is as small as possible.

PARSING:

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

Parse Tree:

- Graphical representation of a derivation or deduction is called a parse tree.
- interior node → non-terminal.
child → terminals (or) non-terminals.

KESHAV MEMORIAL E

Top-Down Parsing

- * A parser can start with the start symbol and try to transform it to input string.

Ex: LL Parser.

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

Two Types:

1. Recursive Descent Parsing.
2. Predictive parsing.

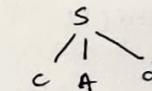
1. Recursive Descent Parsing:

- * uses a set of recursive procedures to scan its input.
- * may involve backtracking, making repeated scans of the input.

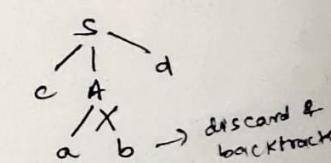
Ex: G: $S \rightarrow cAd$.
 $A \rightarrow ab/a$.

Input string is w = cad.

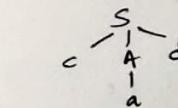
①



②



③





If input is non-terminal

(1) call corresponding function.

If input is terminal

(1) compare terminal with corresponding input symbol.

(2) if same \rightarrow input++.

If a non-terminal produces ≥ 2 prod. they we have to write corresponding non-terminal.

There is no need to declare main function and variables

Ex 1: $E \rightarrow TE'$

$E' \rightarrow +TE' | E$

non-terminal.

$E()$

```
{ if (inp == 'T') input++;
    EPRIME(); }
```

}

EPRIME()

```
{
    if (inp == '+')
        input++;
    if (inp == 'E')
        input++;
    E();
}
else
    return;
```

Ex 2: b

$E \rightarrow TE'$

$E' \rightarrow +TE' | E$

$T \rightarrow FT'$

$T' \rightarrow *FT' | E$

$F \rightarrow (E) | id$

KESHAV MEMORIAL ENGINEERING COLLEGE

Expt. No. _____

Date _____

$E()$

```
{
    T();
    EPRIME();
}
```

EPRIME()

```
{
    if (inp == '+')
        input++;
```

T();

```
    EPRIME();
}
```

else if

y return;

}

T()

```
{
    F();
    TPRIME();
}
```

TPRIME()

```
{
    if (inp == 'E')
        input++;
```

F();

```
    TPRIME();
}
```

else if

y return;

}

F()

```
{
    if (inp == '(')
        input++;
```

E();

id+id * id.

if (inp == ')')

input++;

y

} else if (inp == 'id')

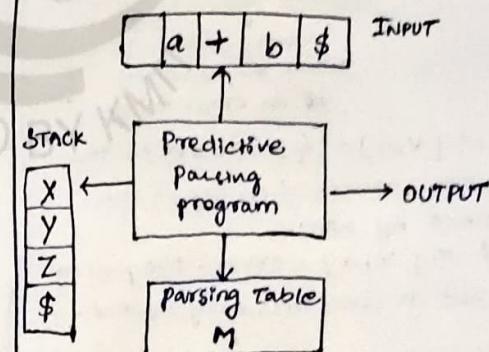
input++;

y

2) Predictive Parsing :

* It is a special case of recursive descent parsing where no back tracking is required.

* The key problem of predictive parsing is to determine the production to be applied for a non-terminal in case of alternatives.



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

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

Parsing table: It is 2-D Array
 $M[x, a]$, $X \rightarrow$ non-terminal,
 $a \rightarrow$ terminal.

Predictive parsing program:

- * The parser is controlled by a program that considers
 - { $X \rightarrow$ the symbol on top of stack
 - $a \rightarrow$ current input symbol.
- () these two determine the parser action.

There are three possibilities:

1. If $X = a = \$$, the parser halts and announces successful completion of parsing.
2. If $X = a \neq \$$, the parser pops off the stack and advances the input pointer to the next input symbol.
3. If X is a non-terminal, the program consults entry $M[x, a]$ of the parsing table M .

\overline{J}
 entry will
 either X -production
 or an error entry.

- \Rightarrow If $M[x, a] = \{X \rightarrow UVW\}$, the parser replaces X on top of the stack by WVU .
- \Rightarrow If $M[x, a] = \text{error}$, the parser calls an error recovery routine.

FIRST:

rules:

- * $X \rightarrow$ terminal ? $\text{FIRST}(X)$ is $\{X\}$.
- * $X \rightarrow E$? E is a production ? then add E to $\text{FIRST}(X)$.
- * $X \rightarrow$ non-terminal & $X \rightarrow a \alpha$? $a \alpha$ is a production ? then add a to $\text{FIRST}(X)$.

* $X \rightarrow$ non-terminal & $X \rightarrow Y_1 Y_2 \dots Y_n$ is a production ? then place a^* in $\text{FIRST}(X)$ if for some i , a is in $\text{FIRST}(Y_i)$ and ϵ is in $\text{FIRST}(Y_{i-1}) \dots$ all of $\text{FIRST}(Y_1) \dots$

$$\text{Ex: } A \rightarrow BC \mid EFGH \mid \epsilon$$

$$\begin{aligned} B &\rightarrow b \\ C &\rightarrow c \mid \epsilon & \epsilon \rightarrow h \mid \epsilon \\ E &\rightarrow e \mid \epsilon \\ F &\rightarrow CE \\ G &\rightarrow g \end{aligned}$$

$$\text{first}(H) = \text{first}(h) \cup \text{first}(e) \\ = \{h, \epsilon\}$$

$$\text{first}(G) = \{g\}$$

$$\text{first}(C) = \text{first}(c) \cup \text{first}(e) \\ = \{c, \epsilon\}.$$

$$\text{first}(E) = \{e, \epsilon\}.$$

$$\begin{aligned} \text{first}(F) &= \text{first}(CE) \\ &= (\text{first}(C) - \{\epsilon\}) \cup \text{first}(E) \\ &= (\{c\} - \{\epsilon\}) \cup (e, \epsilon) \\ &= \{c, e, \epsilon\}. \end{aligned}$$

$$\text{first}(B) = \{b\}.$$

$$\text{first}(A) = \text{first}(BC) \cup \text{first}(EFGH) \\ \cup \text{first}(H).$$

$$= \text{first}(B) \cup (\text{first}(E) - \{\epsilon\}) \cup \text{first}(FGH) \cup \{h, \epsilon\}.$$

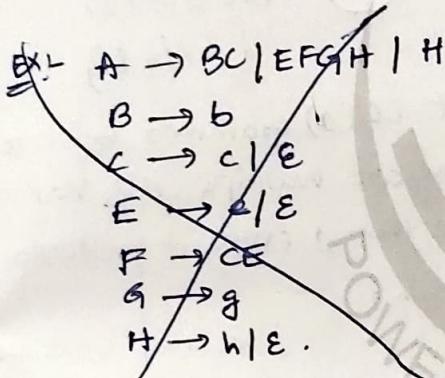
$$= \{b\} \cup \{e\} \cup (\text{first}(F) - \{\epsilon\}) \\ \cup \text{first}(GH) \cup \{h, \epsilon\}$$

$$= \{b, c, e, g, h, \epsilon\}.$$

FOLLOW:

- * If S is a start symbol, then $\text{FOLLOW}(S)$ contains $\$$.
- * $A \rightarrow \alpha B \beta$, then everything in $\text{FIRST}(\beta)$ except ϵ is placed in $\text{FOLLOW}(B)$.

- * ~~$A \rightarrow \alpha B, A \rightarrow \alpha B \beta$~~ where $\text{FIRST}(\beta)$ contains ϵ , then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$.



$$\text{Follow}(A) = \{ \$ \}.$$

$$\text{Follow}(B) =$$

~~FIRST:
Rules:-
Notes:-~~

$\Rightarrow \text{FOLLOW}(A)$ contains set of all terminals present immediate in right of 'A'

$$FO(A) = \{ \$ \}$$

$$\begin{array}{l} * S \rightarrow ACD \\ C \rightarrow a/b \end{array}$$

$$\begin{array}{l} IP \\ S \rightarrow AaCD \\ FO(A) = \{ a \} . \end{array}$$

$$FO(A) \rightarrow FI(C) = \{ a, b \}$$

$$FO(D) \rightarrow FO(S) = \{ \$ \} .$$

$$* S \xrightarrow{\substack{aSbS \\ \downarrow \\ aSbS}} / \xrightarrow{\substack{bSaS \\ \downarrow \\ bSaS}} / \epsilon$$

$$FO(S) = \{ \$, b, a \} .$$

Start symbol

'Follow never contain ϵ '

$$* S \rightarrow AaAb / BbBa .$$

$$A \rightarrow E$$

$$B \rightarrow E$$

$$FO(A) = \{ a, b \}$$

$$FO(B) = \{ b, a \} .$$

$$* S \xrightarrow{\substack{ABC \\ \downarrow \\ ABC}} / \xrightarrow{\substack{DEF \\ \downarrow \\ DEF}} / \epsilon$$

$$FO(A) = \text{FIRST}(B) .$$

$$B \rightarrow E$$

$$C \rightarrow E$$

$$D \rightarrow E$$

$$E \rightarrow E$$

$$F \rightarrow E$$

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

$$(\epsilon)$$

$$= \{ \$ \} .$$

$$(\epsilon)$$

$$\{ \$ \} .$$

$$\begin{array}{l} Q) \quad S \rightarrow (L) / a \\ L \rightarrow SL' \\ L' \rightarrow , SL' / e \end{array}$$

Find
First &
Follow.

$$\begin{aligned} \rightarrow \text{FIRST}(S) &= \{c, a\} \\ \text{FIRST}(L) &= \{c, a\} \\ \text{FIRST}(L') &= \{\}, \{\epsilon\}. \end{aligned}$$

$$\begin{aligned} \rightarrow \text{FOLLOW}(S) &= \{\$y \cup \text{FIRST}(L') \\ &\cup \text{FOLLOW}(L)\} \\ &= \{\$, , , \} \end{aligned}$$

$$\begin{aligned} \text{FOLLOW}(L) &= \{)y\} \\ \text{FOLLOW}(L') &= \text{FOLLOW}(L) \\ &= \{)y\}. \end{aligned}$$

(LL(1)) Non-recursive Predictive Parser:

$$\begin{array}{l} S \rightarrow (L) / a \\ L \rightarrow SL' \\ L' \rightarrow \epsilon / , SL' \end{array}$$

Step 1: Find first and follow of above.

Parse Table:

S	()	a	,	\$
S	1		2		
L	3		3		
L'	L'	4	5		

- ① $S \rightarrow (L) \quad S \rightarrow a$
 \downarrow
First
- ② $L \rightarrow SL'$
- ③ $L' \rightarrow \epsilon \Rightarrow \text{Follow}(L')$
- ④ $L' \rightarrow , SL'$

* If it is LL(1) Grammar because in every cell there is only one entry of production.

$$P) S \rightarrow aSbs / bSas / \epsilon$$

$$\text{first}(S) = \{a, b, \epsilon\}$$

$$\text{follow}(S) = \{b, a, \$\}$$

	a	b	\$
S	1/3	2/3	3

$$① S \rightarrow aSbs$$

$$② S \rightarrow bSas$$

$$③ S \rightarrow \epsilon \xrightarrow{\text{Follow}(S)} \{a, b, \$\}$$

\Rightarrow Not LL(1) grammar because there are multiple cells having more than 1 entry of production.

Algorithm:

repeat

$x \rightarrow \text{top stack symbol}$
 $a \rightarrow \text{symbol pointed to by ip}$.

if $x = \text{terminal or } \epsilon$ then
 if $x = a$ then
 pop x and advance ip
 else error()
else
 if $M[x, a] = x \rightarrow y_1 y_2 \dots y_k$ then
 pop x
 push y_k, y_{k-1}, \dots, y_1 , \top
 output $x \rightarrow y_1 y_2 \dots y_k$
 else error().

until $x = \epsilon$

problem :-

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

Elimination of Left Recursion :-

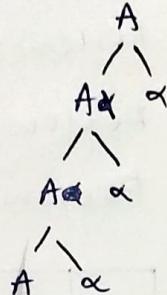
$$i) A \rightarrow A\alpha \mid \beta$$

① write all the productions with non-left recursive

② introduce new non-terminal related to A'

$$A \rightarrow \beta A'$$

$$A' \rightarrow \underbrace{\alpha A' \mid \epsilon}_{\text{except first non-terminal symbol}} \quad \text{a write remaining.}$$



Solution of problem :-

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

$$\Rightarrow E \rightarrow TE'$$

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

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

$$\Rightarrow T \rightarrow FT'$$

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

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

After removing left recursion

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

$$\Rightarrow FI(E) \rightarrow FI(T) \rightarrow FI(F)$$

$$\rightarrow \{c, id\}.$$

$$FI(E') \rightarrow \{+, \epsilon\}.$$

$$FI(T) \rightarrow FI(F) \rightarrow \{c, id\}$$

$$FI(T') \rightarrow \{*, \epsilon\}.$$

$$FI(F) \rightarrow \{c, id\}.$$

$$\Rightarrow FO(E) \rightarrow \{$,)$\}.$$

$$FO(E') \rightarrow \{$, $\}.$$

$$FO(T) \rightarrow \{+, \$,)\}.$$

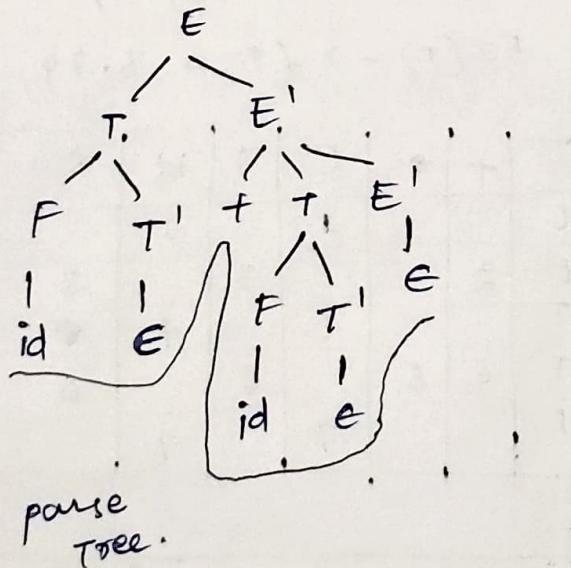
$$FO(T') \rightarrow \{+, \$,), \}\}$$

$$FO(F) \rightarrow \{*, +, \$,)\}$$

	+	*	()	id	\$
E			1		1	
E'	2			3		3
T			4		4	0
T'	6	5		6		6
F			7		7	

$\Rightarrow id + id\$$ Stack.

<u>Stack</u>	<u>Input string</u>	<u>Action</u>
\$ E	id + id \$	$E \rightarrow TE'$
\$ E' T	id + id \$	$T \rightarrow FT'$
\$ E' F	id + id \$	$F \rightarrow id$
\$ E' T' id	id + id \$	POP
\$ E' T'	id + id \$	$T' \rightarrow E$
\$ E'	id + id \$	$E' \rightarrow +TE'$
\$ E' T ±	id + id \$	POP
\$ E' T	id + id \$	$T \rightarrow FT'$
\$ E' T F	id + id \$	$F \rightarrow id$
\$ E' T \$ id	id + id \$	POP
\$ E' T'	id + id \$	$T' \rightarrow E$
\$ E'	id + id \$	$E' \rightarrow E$
\$	id + id \$	Accepted.



$$\text{ii) } S \rightarrow (L)/a \\ L \rightarrow L, S/S.$$

$$S \rightarrow (\zeta)/a$$

$$L \rightarrow S L'$$

\hookrightarrow , $s \in$

$$\text{FIRST}(S) \rightarrow \{c, a\}$$

$\text{FIRST}(L) \rightarrow \{c, a\}$

$\text{FIRST}(L') \rightarrow \{\}, \in\}$

`FOLLOW(S) → { , , , }`

$\text{Follow}(L) \rightarrow \{\}\}$

Follow(L) $\rightarrow \{ \} \}$

	()	,	a	\$
S	$S \rightarrow (L)$			$S \rightarrow a$	
L	$L \rightarrow SL'$			$L \rightarrow SL'$	
L'	$L' \rightarrow S$	$L' \rightarrow \epsilon$			LL'

$$(iii) S \rightarrow aAB/bA/e$$

$$A \rightarrow aAb/e$$

$$B \rightarrow b\bar{B} / \epsilon.$$

$\text{FIRST}(S) \Rightarrow \{a, b, e\}$

$$\text{FIRST}(A) \Rightarrow \{a, e\}$$

FIRST(B) = {b, e}

$\text{Follow}(S) \Rightarrow \{ \$ \}$.

$\text{Follow}(A) \Rightarrow \{b, \$\}$.

$\text{Follow}(B) \Rightarrow \{\$\}$.

	a	b	
S	$s \rightarrow aAB$	$s \rightarrow bAf$	$s \rightarrow E$
A	$A \rightarrow aAb$	$A \rightarrow E$	ASSOC AX
B	aAb	$B \rightarrow bB$	$B \rightarrow E$

Expt. No. _____

Date _____

BOTTOM UP PARSING

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

Shift Reduce Parser

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

$$\text{Ex: } S \rightarrow aABe$$

$$A \rightarrow Abc | b$$

$$B \rightarrow d$$

input: abbcde.

$$S \rightarrow aABe$$

$$S \rightarrow aAbcBe \quad (A \rightarrow Abc)$$

$$S \rightarrow abbcBe \quad (A \rightarrow b)$$

$$S \rightarrow abbcde \quad (B \rightarrow d)$$

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

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

$$E \rightarrow E^*E$$

$$E \rightarrow (E)$$

$$E \rightarrow id$$

Input: $id_1 + id_2 * id_3$.

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

$$E \rightarrow E+\underline{E^*E} \quad (E \rightarrow E^*E)$$

$$E \rightarrow E+E*\underline{id_3} \quad (E \rightarrow id)$$

$$E \rightarrow E+\underline{id_2 * id_3} \quad (E \rightarrow id)$$

$$E \rightarrow \underline{id_1 + id_2 * id_3} \quad (E \rightarrow id)$$

The underlined substrings are called handles.

Actions in Shift-reduce parser:

* shift: next input symbol is shifted onto the top of the stack. (push)

* reduce: parser replaces the handle within a stack with a non-terminal.

* accept: announces successful completion of parsing.

* error: discovers that a syntax error has occurred and calls an error recovery routine.

Conflicts in shift-reduce parsing:

1. Shift-reduce conflict: The parser cannot decide whether to shift or to reduce.

2. Reduce-reduce conflict: The parser cannot decide which of several reductions to make.

Ex:

$$E \rightarrow E + T / T$$

$$T \rightarrow T * F / F$$

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

id * id

<u>Stack</u>	<u>Input Buffer</u>	<u>Action</u>
\$	id * id \$	shift
\$ id	* id \$	Reduce by $F \rightarrow id$
\$ F	* id \$	Reduce $T \rightarrow F$
\$ T	* id \$	shift (SR)
\$ T*	id \$	shift
\$ T* id	\$	reduce $F \rightarrow id$
\$ T* F	\$	reduce $T \rightarrow T * F$ (RR)
\$ T	\$	reduce $E \rightarrow T$
\$ E	\$	ACCEPTED

Start symbol.

LR Parsers

* used to parse a large class of CFGs is called LR(K) parsing.

* 'L' \Rightarrow left to right scanning of input.

'R' \rightarrow for constructing rightmost derivation in reverse.

'K' \rightarrow number of input symbols.

(K=1)
default
if K is omitted.

Advantages:

* efficient non-backtracking shift-reduce parsing method.

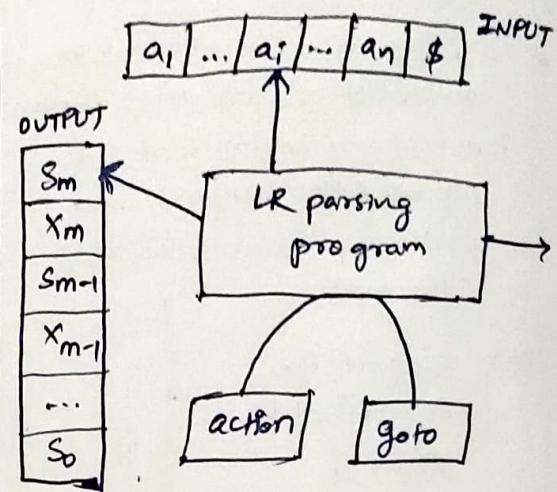
* A grammar that can be parsed

using LR method is a proper superset of a grammar that can be parsed with predictive parser.
* It detects a syntactic error as soon as possible.

Drawbacks

- * too much of work to construct.
- * LR parser generator is needed.
Ex :- YACC.

LR Parsing Algorithm:



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

* The driver program is the same for all LR parser.

* The parsing program reads characters from an input buffer one at a time.

* program uses a stack to store a string

" $S_0 X_1 S_1 X_2 S_2 \dots X_m S_m$ "

where,
 S_m is on top,
 X_i is grammar symbol
 S_i is a state.

Action :- shift, reduce, accept & error.

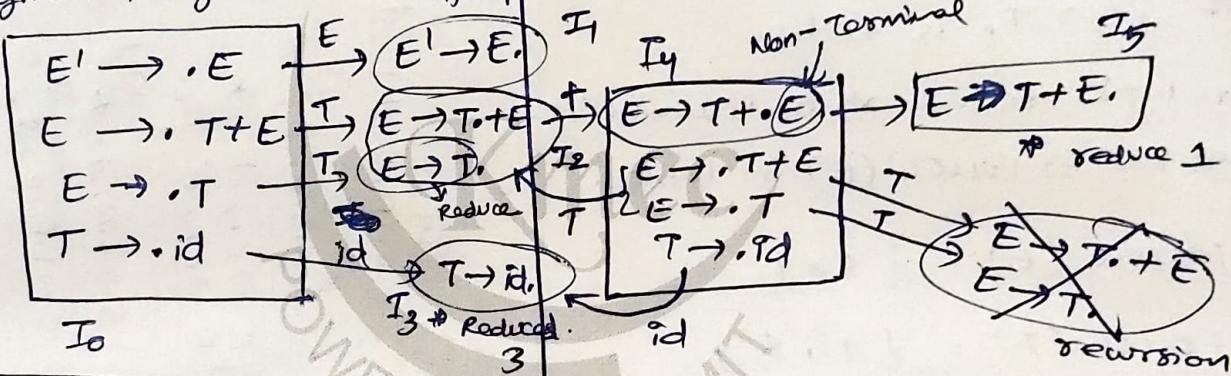
Goto: The function goto takes a state and grammar symbol as arguments and produces a state.

LR(0):

To construct table we use LR(0) canonical items

Ex: $E \rightarrow T+E/T$
 $T \rightarrow id$.

① Augment the grammar. Accept



State	id + \$	E	T
0	S_3	1	2
1	- - Accept		
2	τ_2 S_2		
3	$\tau_3 \tau_3 \tau_3$		
4	S_3	5	2
5	$\tau_1 \tau_1 \tau_1$		
	Action	Goto.	

If there is more than one entry in the cell i.e., Shift-reduce
 (or) reduce-reduce

SLR(1):

- * From the LR(0) example.
- * Difference is from parsing table.

State	Id	+	\$	E	T
0	S_0			1	2
1					
2					
3					
4					
5	S_3			5	2

* If $E \rightarrow T.$ is reduce
then write reduce ~~*~~ in $\text{Follow}(\text{left Hand side}) \rightarrow \text{Follow}(E)$ in table.

$\text{Follow}(E) = \{ \$ \}$

$\text{Follow}(T) = \{ +, \$ \}$

$T \rightarrow \text{id.} \Rightarrow r_3 \Rightarrow (+, \$)$.

* There is no shift-reduce (S)
reduce-reduce conflict in this grammar so TE is SLR(1).

CLR (LR1 Canonical Items).

$$S \rightarrow aAd^0 / bBd^1 / aBe^2 / bAe^3$$

$$A \rightarrow c^5$$

$$B \rightarrow C^6$$

① Augment \downarrow look ahead.

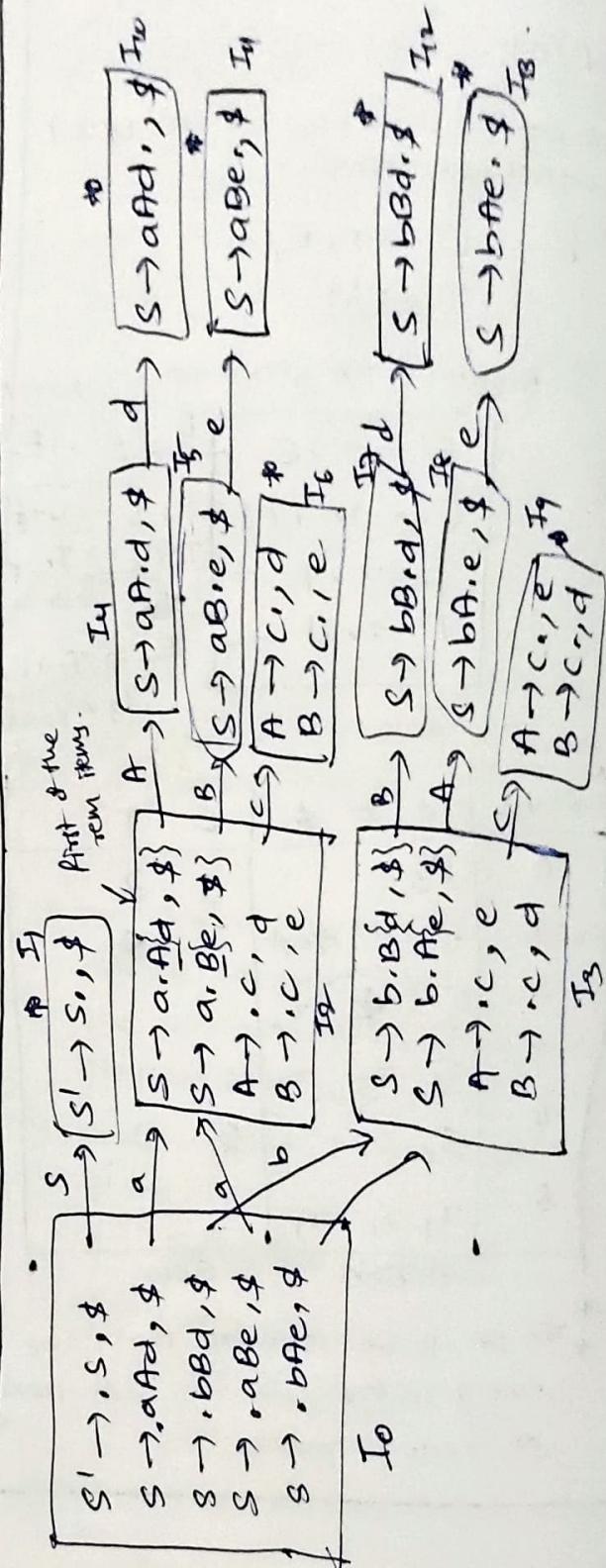
$S' \rightarrow .S\$$ $S \rightarrow .bAe, \$$

$S \rightarrow .aAd, \$$ $A \rightarrow .c$

$S \rightarrow .bBd, \$$ $B \rightarrow .c$

$S \rightarrow .aBe, \$$

$$\text{LR}(1) = \text{LR}(0) + \text{look ahead}.$$



STATE	a	b	c	d	e	f	g	S	A	B
0	S_2	S_3						1		
1			S_6						4	5
2			S_6	S_7					8	7
3			S_9							
4										
5										
6										
7										
8										
9										
10										
11										
12										
13										

CLR (canonical LR)

$S \rightarrow CC$
 $C \rightarrow cC/d$

① Augmented grammar.

$S' \rightarrow S$ ⑥

$S \rightarrow CC$ ①

$C \rightarrow cC$ ②

$C \rightarrow d$ ③

② calculating LR(0) items.

$I_0 : S' \rightarrow .S, \$$

$I_0 : S \rightarrow .CC, \$$

$I_0 : C \rightarrow .CC, C/d$

$I_0 : C \rightarrow .d, C/d$

goto (I_0, S) :-

$I_1 : S' \rightarrow S., \$$

goto (I_0, C) :-

$I_2 : S \rightarrow C.C, \$$

$C \rightarrow .CC, \cancel{C/d} \$$

$C \rightarrow .d, \$$

goto (I_0, d) :-

$I_3 : C \rightarrow C.C, C/d$

$C \rightarrow .CC, C/d$

$C \rightarrow .d, C/d$

goto (I_0, d) :-

$I_4 : C \rightarrow d., C/d$

goto (I_2, C) :-

$I_5 : S \rightarrow CC., \$$

goto (I_2, C) :-

$I_6 : C \rightarrow c.C, \$$

$C \rightarrow .CC, \$$

$C \rightarrow .d, \$$

goto (I_2, d) :-

$I_7 : C \rightarrow d., \$$

goto (I_3, C) :-

$I_8 : C \rightarrow CC., C/d$

goto (I_3, C) :-

$C \rightarrow c.C, C/d$

$C \rightarrow .CC, C/d$

$C \rightarrow .d, C/d$

goto (I_3, d):

$C \rightarrow d, \text{id} \} I_4$

goto (I_6, C):

$I_9: C \rightarrow CC, \$ *$

goto (I_6, C):

$C \rightarrow C.C, \$$
 $C \rightarrow .CC, \$$
 $C \rightarrow rd, \$$

goto (I_6, d):

$C \rightarrow d, \$ \} I_7.$

Table:-

State	Action			Goto,
	c	d	\$	
0	S_3	S_4		1
1			ACCEPT	2
2	S_6	S_7		5
3	S_3	S_4		8
4	T_3	T_3		
5			T_1	
6	S_6	S_7		9
7			T_3	
8	T_2	T_2		
9			T_2	

Stack	Input string	Action
$\$ 0$	$dd \$$	shift 4
$\$ 0 d 4$	$d \$$	reduce 3 $C \rightarrow d$
$\$ 0 C 2$	$d \$$	shift 4
$\\$ 0 C 2$	$C \$$	
$\$ 0 C 2 d 7$	$\$$	reduce 3 $C \rightarrow d$
$\$ 0 C 2 C 5$ ↑ RC	$\$$	reduce 1 $S \rightarrow CC$ ← pop 4
$\$ 0 S 1$	$\$$	ACCEPT.

IP file RHS contains 1 symbol
then we have to pop 2 symbols
from the stack

LALR (Lookahead LR)

$$S \rightarrow L = R / R$$

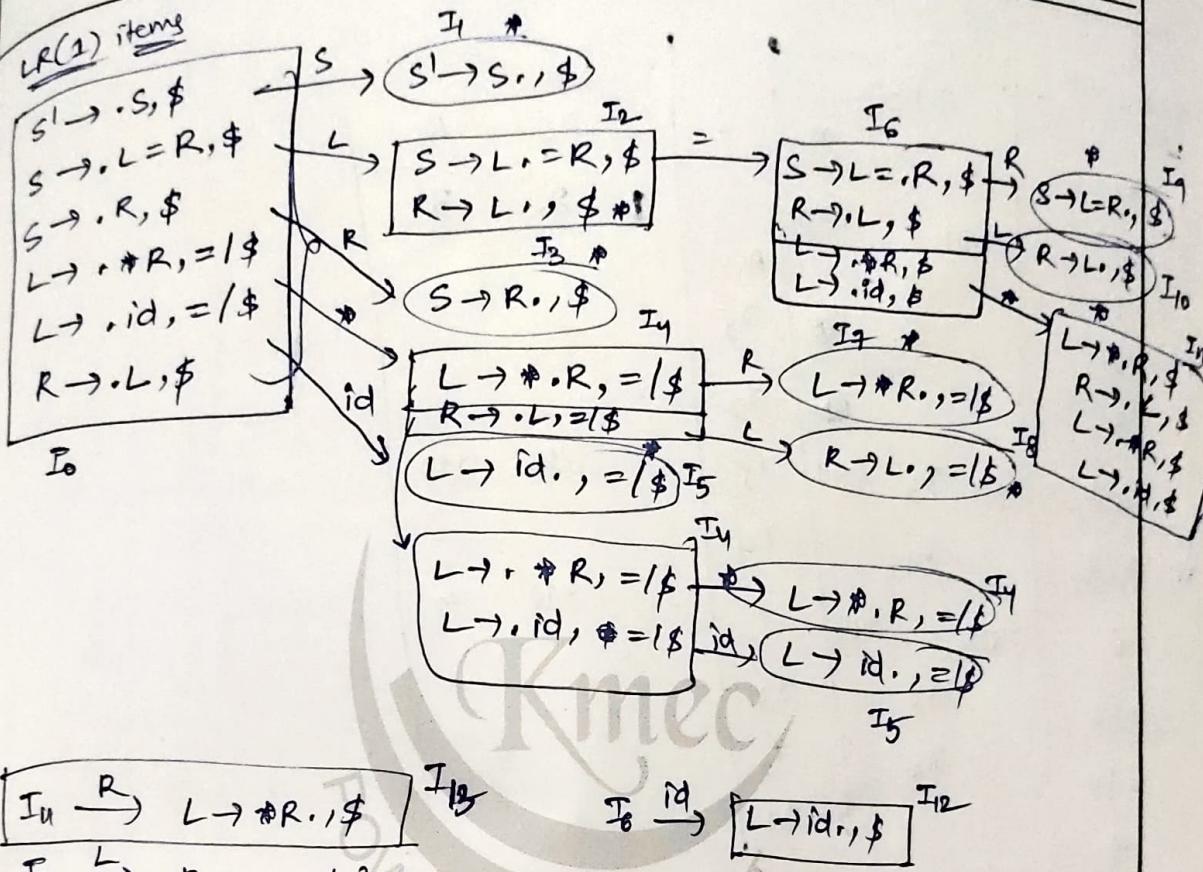
$$L \rightarrow * R / id$$

$$R \rightarrow L$$

also parse $id = id$.

Augmented Grammars:

- ① $S' \rightarrow S$
- ② $S \rightarrow L = R$
- ③ $L \rightarrow * R$
- ④ $R \rightarrow L$
- ⑤ $S \rightarrow R$
- ⑥ $L \rightarrow id$



$$I_0 \xrightarrow{R} L \rightarrow *R., \$ \quad I_{13}$$

$$I_0 \xrightarrow{L} R \rightarrow L., \$ \quad I_{10}$$

$$I_0 \xrightarrow{*} L \rightarrow *., R, \$ \quad I_4$$

$$I_0 \xrightarrow{id} L \rightarrow id., \$ \quad I_{12}$$

Combine I_4 & I_0 into single state

where $I_4 = L \rightarrow *., R, =/\$$ $I_{11} = L \rightarrow *., R, \$$
 $R \rightarrow .L, =/\$$ $R \rightarrow .L, \$$ $\Rightarrow I_{411}$
 $L \rightarrow .*R, =/\$$ $L \rightarrow .*R, \$$
 $L \rightarrow .id, =/\$$ $L \rightarrow .id, \$$

I_5 & I_{12}

$$I_5 = L \rightarrow id., =/\$ \quad I_{12} = L \rightarrow id., \$ \quad \Rightarrow I_{512}$$

$$I_7 = L \rightarrow *R., =/\$ \quad I_{13} = L \rightarrow *R., \$ \quad \Rightarrow I_{713}$$

Similarly :- T810

In LALR, if two states
contains same productions but
different look ahead symbols
then we combine the states.

Table:-

States	Action			goto		
	id	=	*	S	L	R
0	S_{512}	S_{411}	\$	1	2	3
1			Acc			
2	S_6	*	τ_5			
3		*	τ_2			
411	S_{512}	τ_5	S_{411}	τ_5	810	713
512		τ_4		τ_4		
6	S_{512}	S_{411}			810	9
713		τ_3		τ_3		
810		τ_5		τ_5		
9				τ_1		

$id = id$.

Stack	Input String	Action
\$0	<u>id</u> = <u>id</u> \$	S_{512}
\$0 <u>id</u> <u>512</u> <u>pop</u>	= <u>id</u> \$	reduce 4 $\underline{L} \rightarrow id$
\$0 L 2	= <u>id</u> \$	S_6
\$0 L 2 = <u>6</u>	<u>id</u> \$	S_{512}
\$0 L 2 = <u>6</u> <u>id</u> <u>512</u> <u>pop</u>	\$	reduce 4 $\underline{L} \rightarrow id$
\$0 L 2 = <u>6</u> <u>L</u> <u>810</u> <u>pop</u>	\$	reduce 5 $R \rightarrow L$
\$0 L 2 = <u>6</u> <u>R</u> <u>9</u> <u>pop</u>	\$	reduce $S \rightarrow L = R$
\$0 S 1	\$	Accepted.

KESHAV MEMORIAL ENGINEERING COLLEGE

Type	Direction	LALR	Mechanism	Power	Advantages	Limitations	Expt. No.
LL(0)	Top-Down	1	<ul style="list-style-type: none"> Predictive Parsing using FIRST and FOLLOW sets requires left-factorial, non-left-recursive grammars. 	Limited-only works for a subset of CGFs.	<ul style="list-style-type: none"> Simple Implementation, easy to create standard. Implementation grammar - de-facto - only. 	<ul style="list-style-type: none"> cannot handle left-recursion (in ambiguous grammars, responses grammar - de-facto - only. 	Date _____
LR(0)	Bottom-up	0	<ul style="list-style-type: none"> uses dot position no look ahead used for derivation. 	Very limited	<ul style="list-style-type: none"> Simple state construction 	<ul style="list-style-type: none"> High chance of conflicts insufficient LR(0). 	
SLR	Bottom-up	+	<ul style="list-style-type: none"> Enhances LR(0) items by incorporating Follow Sets for reduce derivations. 	More powerful than LR(0), but still limited sizes	<ul style="list-style-type: none"> Simpler than full LR(0) builds smaller table sizes 	<ul style="list-style-type: none"> May produce conflicts due to conflicts in LR(0). 	
CLR	Bottom-up	1	(per item)	<ul style="list-style-type: none"> constructs LR(0) + explicit LALR sets for precise parsing decisions 	<ul style="list-style-type: none"> maximum power nearly handles all determinate CFGs. 	<ul style="list-style-type: none"> parsing tables can be extremely large (state explosion) 	
LALR	Bottom-up	1	(merged)	<ul style="list-style-type: none"> Merges LR(0) states with identical LR(0) entries to reduce table size 	<ul style="list-style-type: none"> as powerful as CLR 	<ul style="list-style-type: none"> compact table size widely used in yacc / bison 	<ul style="list-style-type: none"> State merging can introduce conflicts now seen in parser

Differences between LL(0), LR(0), SLR, CLR, LALR.

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

(*)

More than one LMD or more than one RMD for some string is also called ambiguous grammars.

$$\text{Ex:- } E \rightarrow E+E / E^*E / (E) / -E / \text{id}$$

for $\text{id} + \text{id} * \text{id}$

$$E \rightarrow E+E$$

$$E \rightarrow \text{id}+E$$

$$E \rightarrow \text{id}+E^*E$$

$$E \rightarrow \text{id}+\text{id}^*E$$

$$E \rightarrow \text{id}+\text{id}^*\text{id}$$

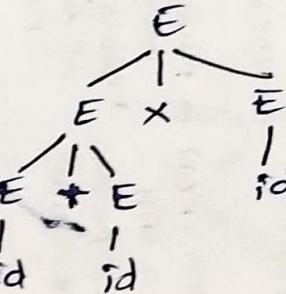
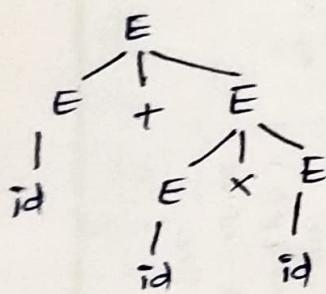
$$E \rightarrow E^*E$$

$$E \rightarrow E+E^*E$$

$$E \rightarrow \text{id}+E^*E$$

$$E \rightarrow \text{id}+\text{id}^*E$$

$$E \rightarrow \text{id}+\text{id}^*\text{id}$$



* Rewriting the grammar can eliminate the ambiguity in grammars.

Eliminating Left-Recursion:

* A grammar is said to be left recursive if it has a non-terminal A such that there is a derivation $A \Rightarrow A\alpha$ for some string α .

* Top-down grammars cannot handle left-recursive grammars.

$$A \rightarrow A\alpha / \beta$$

✓ After removing left recursion.

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' / \epsilon$$

Left factoring:

$$A \rightarrow \alpha\beta_1 / \alpha\beta_2$$



$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 / \beta_2$$

Handling Ambiguous Grammar:

1. Shift-Reduce conflicts.
2. Reduce-Reduce conflicts.
3. Use of Ambiguity-Resolving technique.
4. Preprocessing
5. Grammar modification.

Parser Generator YACC:

* To automate the process of parsing an input string by parser.

- * YACC (Yet Another Compiler Compiler) is one such automatic tool for parser generation.
- * It is an UNIX based utility tool for LR(0) parser generator.
- * LEX and YACC work together to analyse the program syntactically, can report conflicts (in ambiguities) in the form of error messages.

How YACC works :

1. Grammar Specification. ($\text{imp} \rightarrow \text{BNF}$)
2. Parsing Table Generation.
3. Code Generation.
4. Integration with LEX
5. Error Handling.

YACC Specification

① declarations ②
③ token

④
translation rules

⑤
supporting routines, by LEX()

