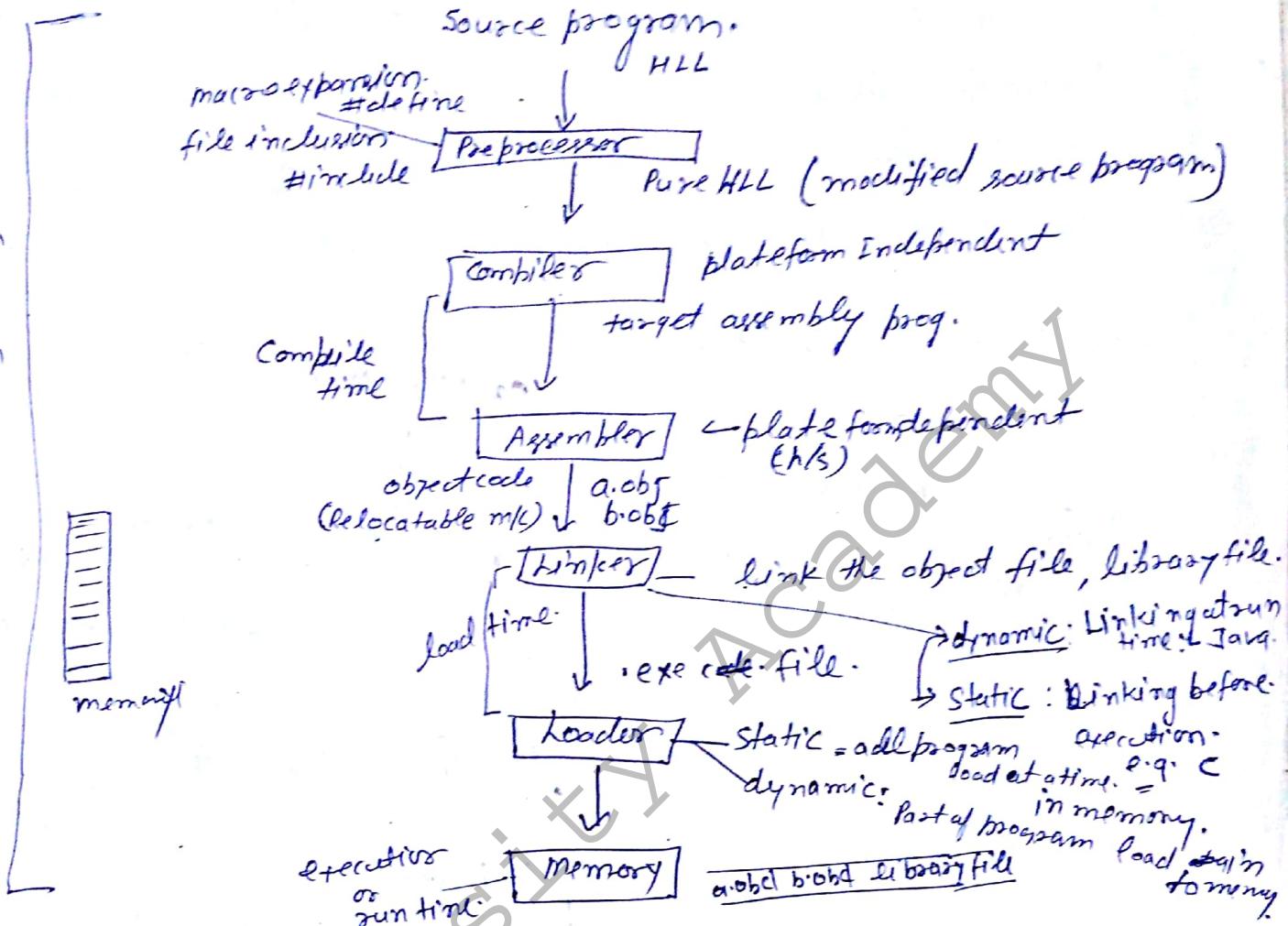


Compiler design

language processing System

translation
and
execution
process



source program

e.g

```

#include <stdio.h>
main()
{
    int a,b;
    int c;
    c = add(a,b);
    printf("%d",c);
}
    
```

a.c

```

int add(int x,int y)
{
    int z;
    z = x+y;
    return z;
}
    
```

b.c

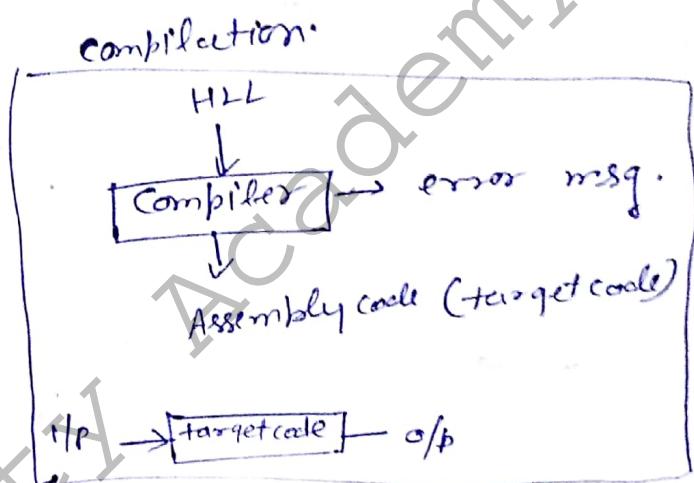
We know that computer understand only machine code i.e. 0 and 1. It is difficult to write computer program directly in machine code. Now days the programs are written in high-level language i.e. BASIC, PASCAL, C, C++ etc. A program written HLL is called **source program**.

the source code cannot be executed directly by the computer.
the source program must be converted into machine code to run
it on computer. Every language has its own translator.

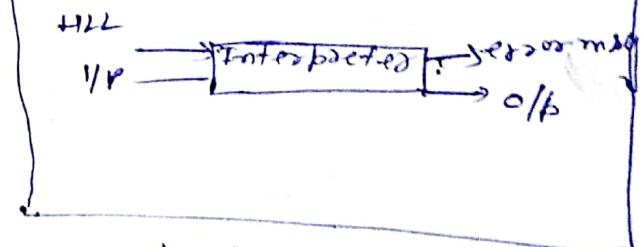
The high-level language is converted into binary language
in various phase. A compiler is a program that convert HLL
to Assembly language, similarly Assembler is a program that
convert the Assembly language to machine level language.

Compiler vs Interpreter

speed: Compiler is fast
error diagnosis: → Interpreter is better



Interpreter



Compiler

1. It translate source code into object code as a whole
2. It create object file
3. execution is fast
4. translator program is not required to translate each time to run the program
5. it does not make easier to correct the mistakes in the source code.
6. most HLL uses compiler

Interpreter

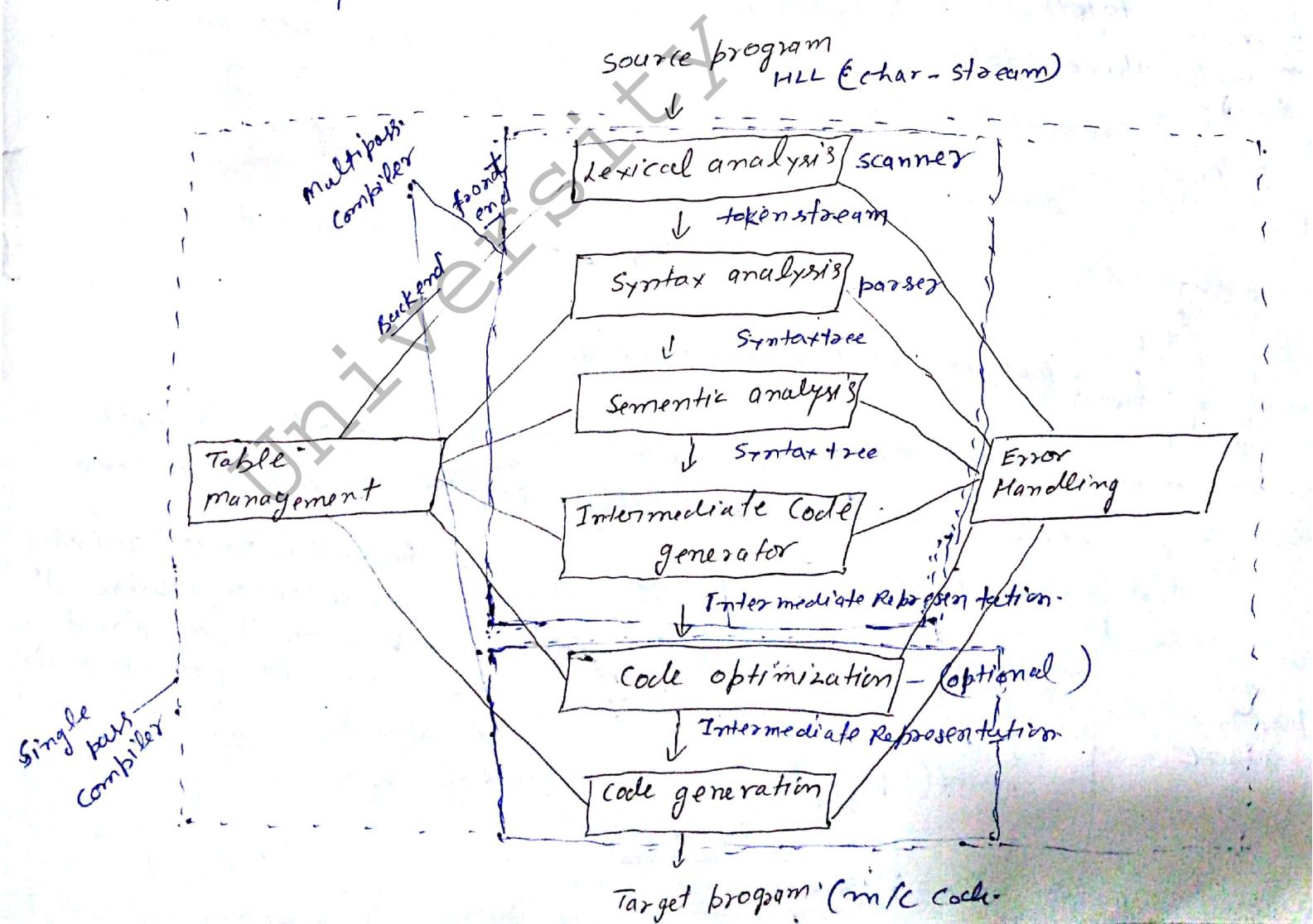
1. It translate the statement of the source code one by one and execute immediately.
2. It slower.
3. execution is slow
4. translator program is required to translate the program each time to run.
5. it makes easier to correct the mistakes in the source code.
6. few HLL uses Interpreter.

1- Translators and Compiler :-

A translator is a program that takes as input a program written in one programming language and produces as output a program in another language.

If the source language is a high level language such as FORTRAN, C, C++ and the object language is low level language such as assembly or machine language then the such translator is called compiler.

The compilation process is so complex that it is not reasonable either logical or implementation point of view to understand this process we partition it into series of subprocess called Phases.



3. Semantic Analysis

Once the syntax is checked in the syntax analyzer phase the next phase semantic analysis determines the meaning of the source string. It also detect the types of data or variable used in statement.

e.g. $\text{total} = \text{count} + \text{rate} * 60$

Note: Semantic analyzer perform the implicit type or explicit typecasting.

Implicit → perform by compiler

Explicit → perform by user (program)

Implicit → $\text{int} + \text{int}$ e.g. $\text{float } a = 3.2;$

Explicit → $\text{int} + \text{float}$ e.g. $\text{float } a = 3.2;$

Implicit: $\text{float} + \text{float}$ e.g. $a = 3.2;$

Explicit: $\text{int} + \text{int}$ e.g. $a = 3;$

4. Intermediate Code Generation

The intermediate code is kind of code which is easy to generate and this code can be easily converted to target code. This code is in form of such three address code, quadruple, triple, etc.

e.g. three address code of the statement

$A/B * C$ is

$$T_1 = A/B$$

$$T_2 = T_1 * C$$

Note - the work till ICG will be platform independent.

T_1 & T_2 are temporary variable.

5. Code optimization: It is an optional phase designed to improve the intermediate code so that the ultimate object program runs faster and take less space.

Its output is another intermediate code program that does the same job as the original.

$$d_1 = id_3 * 60.0$$

$$id_1 = id + t_1$$

6. Code generation: produce the object code by deciding on memory location for data. Designing a code generator that produce truly efficient object program is one of the most difficult part of compiler design.

e.g. machine code for

$$A := B + C \quad 18$$

$R_2 \leftarrow id_3$

$R_2 \leftarrow R_2 * 60.0$

$R_1 \leftarrow id_2$

$R \leftarrow R_1 + R_2$

$id_1 \leftarrow R_1$

LD R_2, id_3
MUL $R_2, R_2, \#60.0$ LOAD B
LD R_1, id_2
ADD R_1, R_1, R_2
ST id_1, R_1 ADD C
STORE A

7. Table Management or Book keeping

table management portion of the compiler keeps record of the name used by the program and essential information about each data type used in program (integer, real etc.) the data structure used to record this information is called symbol table.

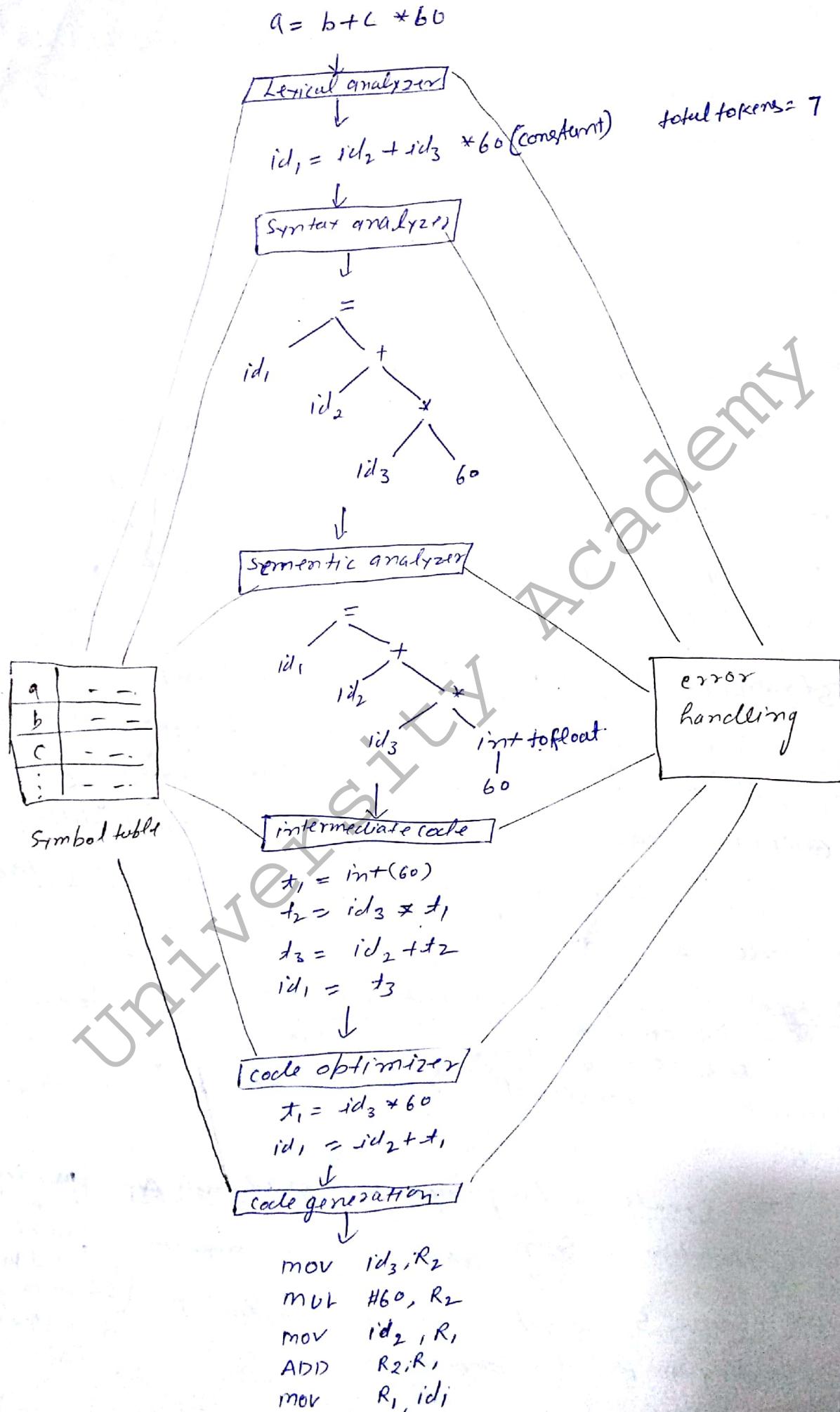
8. Error handling

error handler is invoked when an error in the source program is detected. It must warn the programmer by issuing a message.

both the table management and error handling interact with all phase of the compiler.

Question

Show how an input $a = b + c * 60$ get processed in Compiler

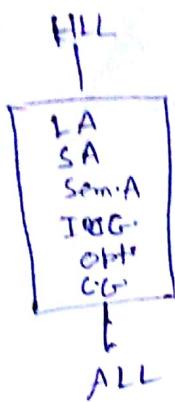


Passes:

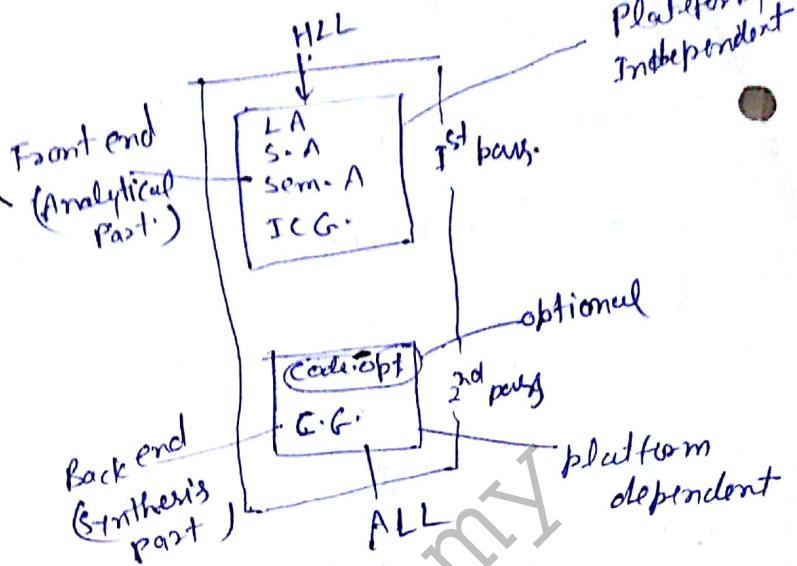
In a compiler portions of one or more phases are combined into a module called pass. A pass read the source program or output of previous pass, makes the transformation specified by its phase and write output into an intermediate file which may be read by subsequent pass.

generally there are two passes in compiler. In the ~~first~~ Front-end phase, Front-end-pass of lexical analysis, syntax analysis, semantic analysis and Intermediate code generation might be grouped into one pass. Code generation optimization is an optional phase. Then Back-end-pass consisting of code generation for a particular target ~~as~~ code.

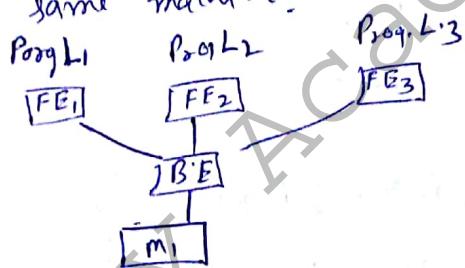
single pass compiler



multipass compiler



if we want to design compiler for three different programming language for same machine:



if want to same programming language to different machine.



Communication gap.

speed.

memory

time.

portability

single pass
No

fast

more

less

No

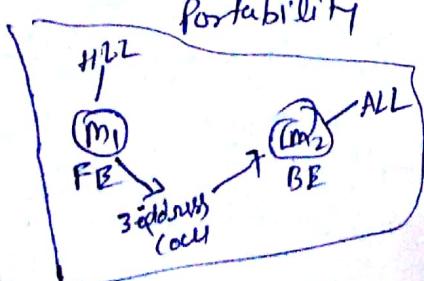
multipass
Yes

slow

less

more

Yes



platform independent

optional

platform dependent

Bootstrapping is the process of writing a compiler in the source programming language that it intends to compile. This technique leads to a self-hosting compiler. Many prog. lang. are bootstrapped including compiler for BASIC, C, JAVA, LISP, PYTHON etc.

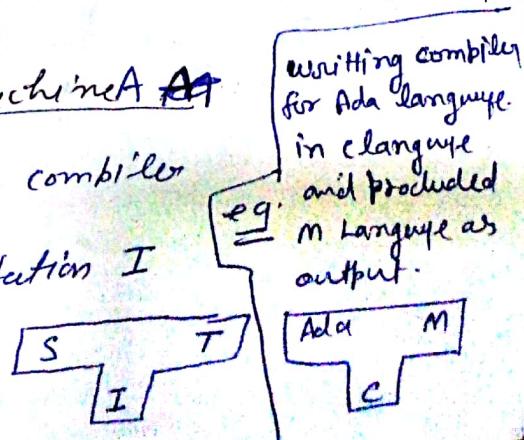
A compiler is characterized by three languages, its source language, its object(^{target}) language and the language in which it is written. These languages may all be quite different.

A compiler may run on one machine produce object code for another machine such a compiler called Cross-compiler. A cross compiler is a compiler capable of creating executable code for a platform other than the one on which the compiler is running. e.g. A compiler that runs on a windows but generates code that runs on Android. Suppose we have a new language L which we want to construct compiler on 2 machines A and B.

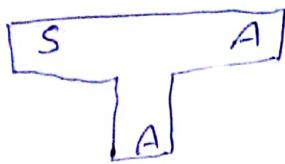
To create a language for Machine A

Notation $S \rightarrow^I T$ represent a compiler for Source S, Target T, Implementation I

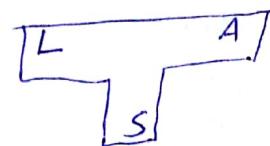
T-diagram for this compiler is



- 1- Create S_{CA}^A a compiler for subset S of the desired language L , using language A which run on machine A (language A may be Assembly language)

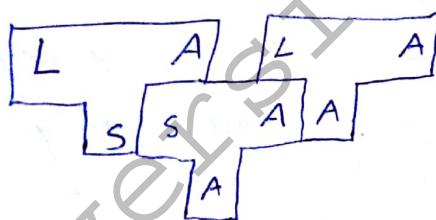


2. create LC_S^A a compiler for language L written in a subset of L



3. combine LC_S^A using S_{CA}^A to obtain LC_A^A :-
a compiler for language L , which runs on machine A
and produce^{code for} machine A

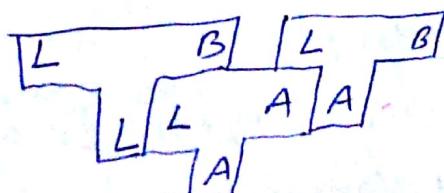
$$LC_S^A \rightarrow S_{CA}^A \rightarrow LC_A^A$$



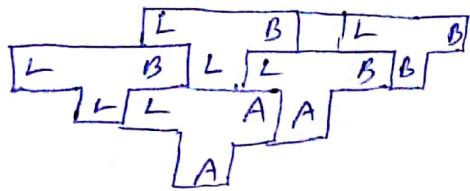
To produce a compiler for different machine:

- 1- Convert LC_S^A into LC_L^B , it is easy if machine A and B both are not different. the language S is a subset of language L .

- 2- combine LC_L^B using LC_A^A to produce LC_A^B



3. Now we cross compiler L_A^B using L_C^B to produce L_C^B .



Regular Expressions.

The regular expression are useful for representing certain sets of string in an algebraic fashion. Actually these describe the language accepted by finite automata.

We call regular expression over alphabet Σ are exactly those expression that can be constructed from following rule.

- 1- ϵ is a regular exp. denoting $\{\epsilon\}$, that is language containing only the empty string.
2. for each ' a ' in Σ ' a ' is a regular expr. denoting $\{a\}$ for language only one string that consisting of single symbol.
3. if R and S are regular expression denoting language L_R and L_S respectively them.
 - (i) $(R)/(S)$ is regular expression denoting $L_R \cup L_S$.
 - (ii) $(R).(S)$ is Regular expression denoting $L_R \cdot L_S$.
 - (iii) $(R)^*$ is regular expression denoting L_R^*

Example! let $\Sigma = \{a, b\}$

- 1- RE a/b denotes the language $\{a, b\}$

2. $(a/b) \cdot (a/b)$ denote the language $\{aa, ab, ba, bb\}$

3. a^* denote the language $\{ \epsilon, a, aa, aaa, \dots \}$

4. $(a/b)^*$ denotes the set of all strings consisting of zero or more instances of a or b that is $\{ \epsilon, a, b, aa, ab, bb, aaa, \dots \}$

5. a/a^*b denotes the language $\{ a, b, ab, aab, aaab, \dots \}$

Finite Automata

Finite Automata are recognizers they simply say yes or no about each possible input string.

FA comes in two type:

(i) Non-Deterministic Finite Automata:

Have no restriction on the label on their edge. A symbol can have several edge out of the same state, and ϵ , the empty string is a possible string.

A NFA consist of Five tuples (Q, Σ, S, q_0, F) where

Q is finite non-empty set of states.

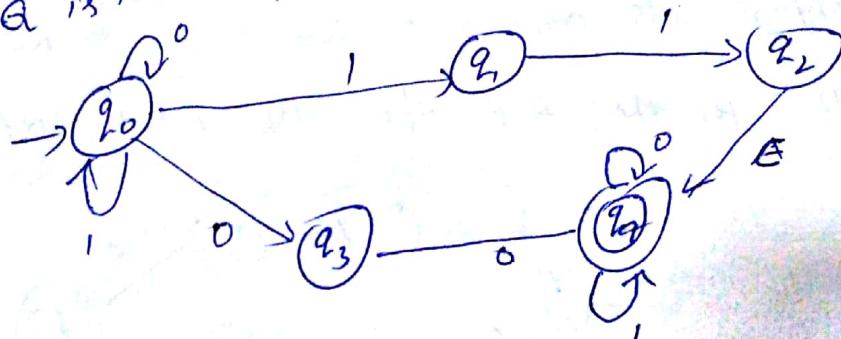
Σ is a finite nonempty set of inputs

S is transition function mapping from $Q \times \Sigma$ into 2^Q which is power set of Q .

$q_0 \in Q$ is the initial state.

$F \subseteq Q$ is the set of final state.

e.g.

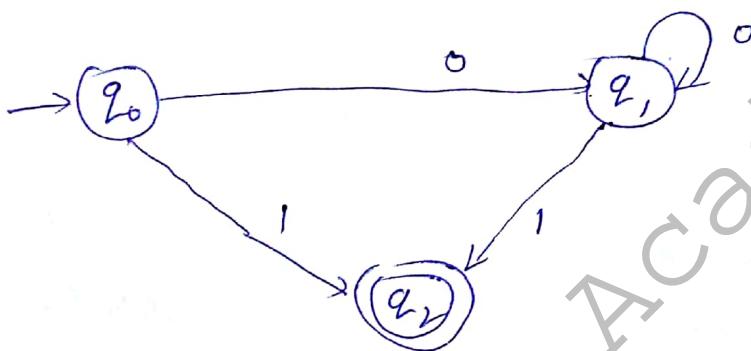


Deterministic Finite Automata (DFA)

DFA have each state and for each symbol of its input alphabet exactly one edge with that symbol leaving that state.

A DFA is a special case of an NFA where -

- 1- there are no moves in input ϵ .
- 2- for each state 's' and input 'a' there is exactly one edge out of 's' labeled 'a'.



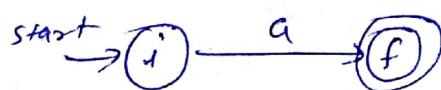
Regular Expression to NFA

McNaughton - Yamada - Thompson algorithm

- 1- For ϵ , construct the NFA

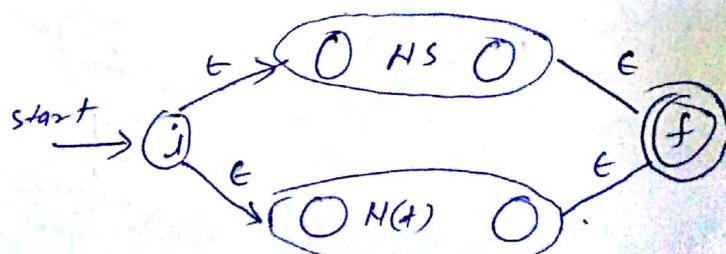


2. For ' a ' in Σ construct NFA



3. Suppose $N(s)$ and $N(t)$ are NFA for RE s and t .

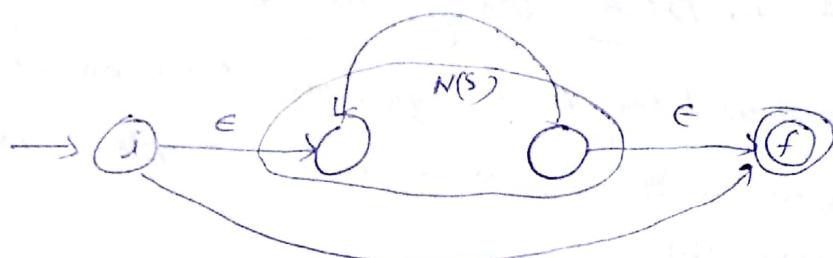
- (a) for the R.E s/t the NFA $N(s/t)$



(b) For the R.E 'st' the NFA $N(st)$



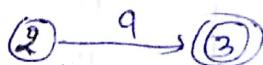
(c) For RE s^* the NFA $N(s^*)$



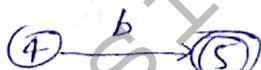
Example: Construct the NFA for $RE = \underline{(a/b)}^* a b b$.

Note
 $a/b = a + b$

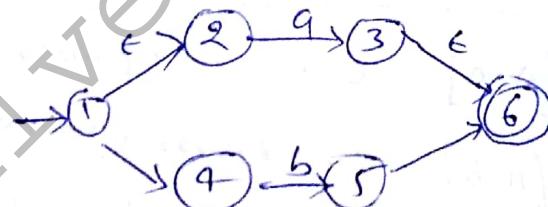
Step-1



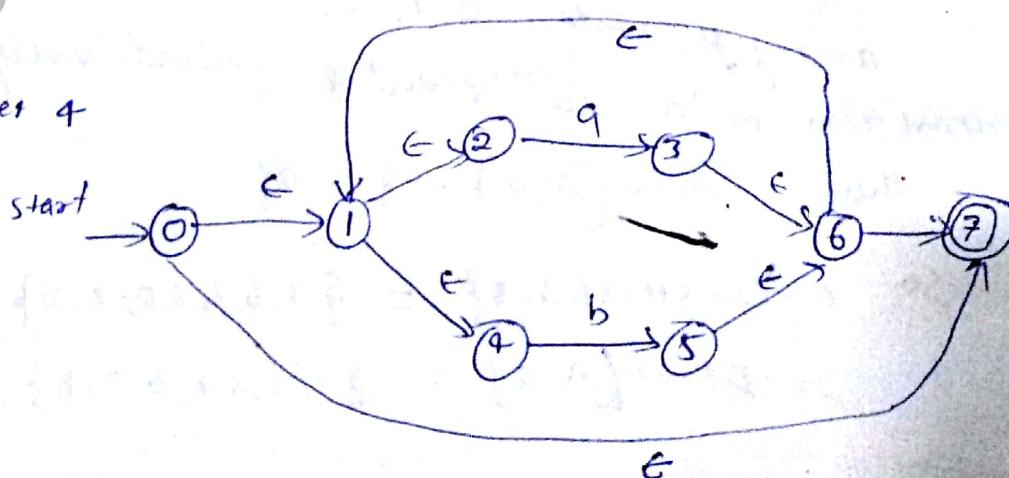
Step-2



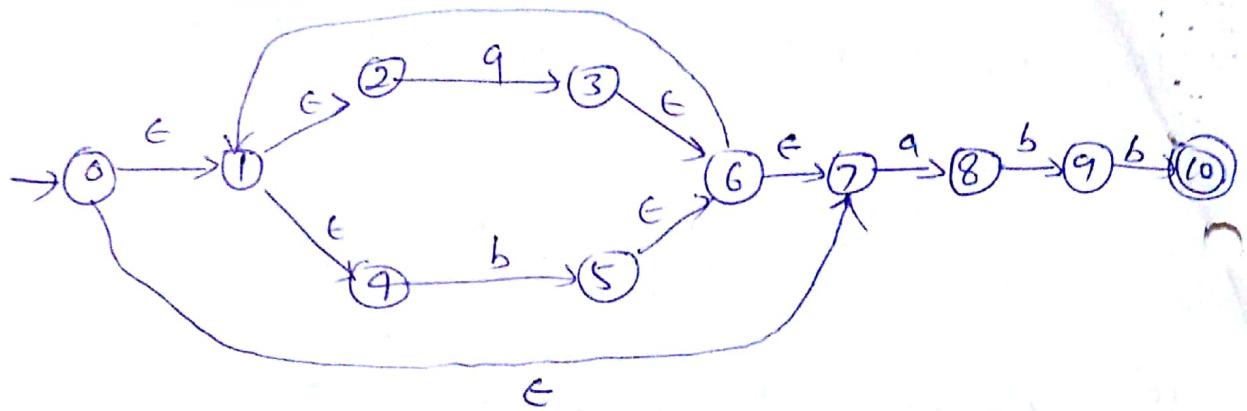
Step 3



Step 4



NFA for RE $(a/b)^*abb$



Conversion NFA to DFA using Subset construction

Example: construct the DFA equivalent to above NFA
we start state A of the equivalent DFA is
 $\epsilon\text{-closure}(0)$

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

| NFA state | DFA state | a | b |
|----------------------------|-----------|---|---|
| $\{0, 1, 2, 4, 7\}$ | A | B | C |
| $\{1, 2, 3, 4, 6, 7, 8\}$ | B | B | D |
| $\{1, 2, 4, 5, 6, 7\}$ | C | B | C |
| $\{5, 6, 1, 2, 4, 7, 9\}$ | D | B | E |
| $\{1, 2, 4, 5, 6, 7, 10\}$ | E | B | C |

Compute $D\text{tran}[A, a] = \epsilon\text{-closure}(\text{move}(A, a))$

$D\text{tran}[A, b] = \epsilon\text{-closure}(\text{move}(A, b))$

among the state 0, 1, 2, 4 and 7 only 2, 7 have transition on 'a' to 3, and 8 respectively.

thus $\text{move}(A, a) = \{3, 8\}$

so $\epsilon\text{-closure}\{3, 8\} = \{3, 6, 1, 2, 4, 7, 8\}$

so $D\text{tran}[A, a] = \{1, 2, 3, 4, 6, 7, 8\} = B$

Now we compute

$$D_{\text{tran}}[A, b] = \epsilon\text{-closure}(\{5\})$$

$$= \{5, 6, 7, 1, 4, 2\}$$

$$= C$$

$$D_{\text{tran}}[B, a] = \epsilon\text{-closure}(\text{move}(B, a))$$

$$\Rightarrow \epsilon\text{-closure}(\{3, 8\})$$

$$= B$$

$$D_{\text{tran}}[B, b] = \epsilon\text{-closure}(\text{move}(B, b))$$

$$= \epsilon\text{-closure}(\{5, 9\})$$

$$= \{5, 6, 1, 2, 4, 7, 9\}$$

$$= D$$

$$D_{\text{tran}}[C, a] = \epsilon\text{-closure}(\text{move}(C, a))$$

$$= \epsilon\text{-closure}(\{3, 8\})$$

$$= B$$

$$D_{\text{tran}}[C, b] = \epsilon\text{-closure}(\text{move}(C, b))$$

$$= \epsilon\text{-closure}(\{5\})$$

$$= C$$

$$D_{\text{tran}}[D, a] = \epsilon\text{-closure}(\text{move}(D, a))$$

$$= \epsilon\text{-closure}(\{3, 8\})$$

$$= B$$

$$D_{\text{tran}}[D, b] = \epsilon\text{-closure}(\text{move}(D, b))$$

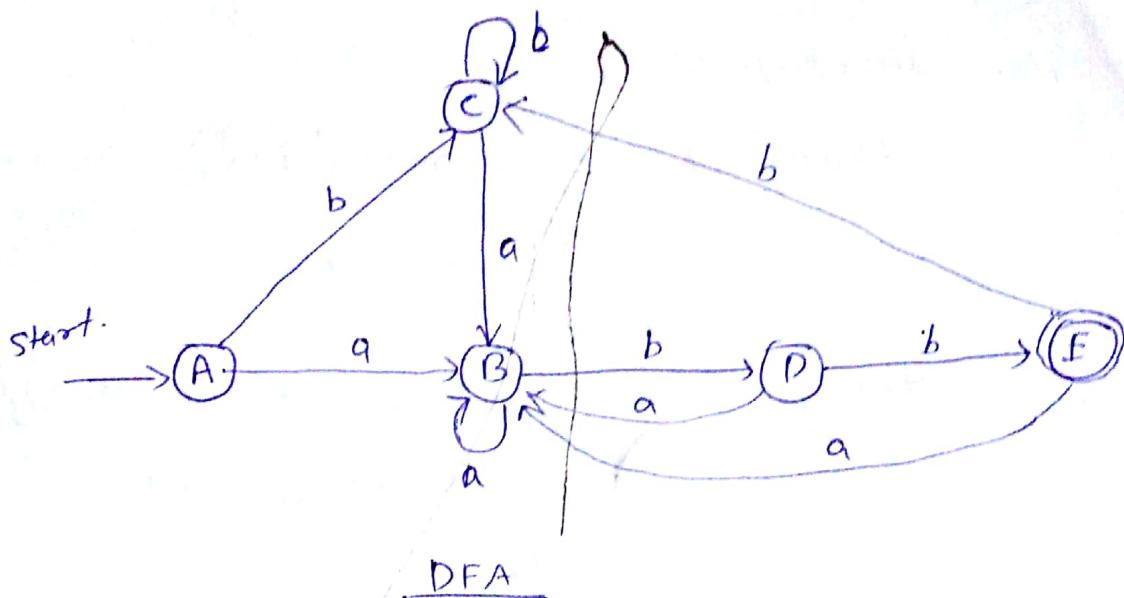
$$= \epsilon\text{-closure}(\{5, 10\})$$

$$= \{5, 6, 7, 12, 4, 10\}$$

$$= E$$

$$D_{\text{tran}}[E, a] = B$$

$$D_{\text{tran}}[E, b] = C$$



DFA

Example: design DFA for $10 + (0+11)0^*1$

Optimization of DFA-Based pattern Matchers:

In this section we present three algorithm that have been used to implement and optimize pattern matchers constructed from regular expression.

1. An algorithm for construct a DFA directly from a RE.
2. The algorithm to minimize the number of state of DFA.
3. The algorithm to produce representation of transition table.

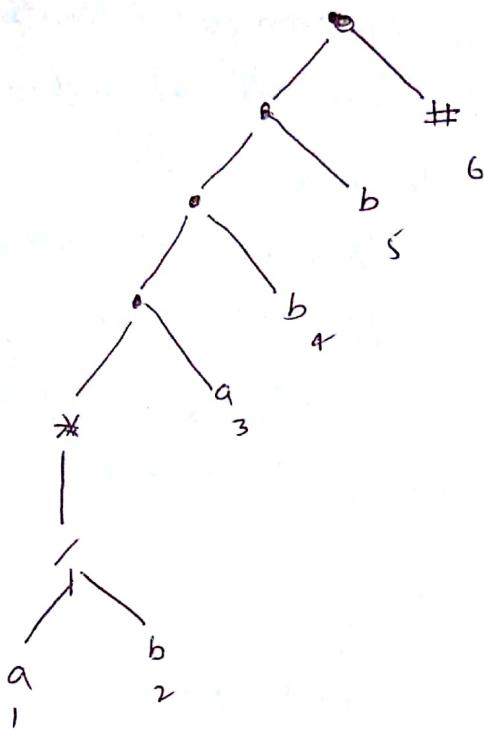
1- Converting RE Directly to DFA :

- (i) construct a Syntax tree T from an augmented RE $(P)\#$.
- (ii) compute nullable, firstpos, lastpos, and followpos for T .
- (iii) construct D states, the set of states of DFA D and D_{tran} , the transition table for D .

example: $r = (a/b)^* abb$

$$(r)\# = (a/b)^* abb\#$$

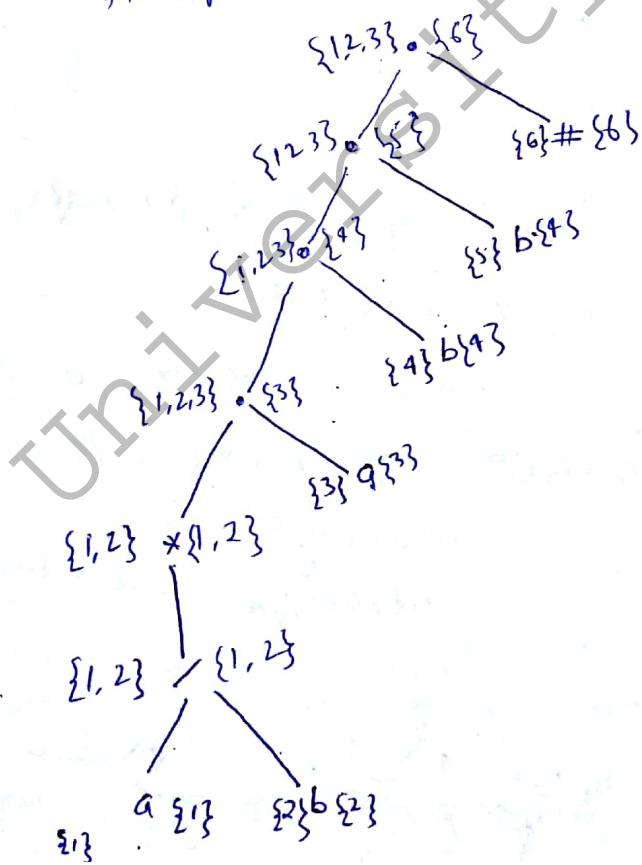
(i)



(8)

(ii) Nullable (η) true for only star node.

first pos and last pos as:



$\text{followpos}(i)$ tell us what position can follow position i in the syntax tree. two rule define all all followpos.

1- if n is a cat-node with child c_1 and c_2 and i

is a position in $\text{lastpos}(c_1)$, then all position in $\text{firstpos}(c_2)$ are in $\text{followpos}(i)$.

2. if n is a star node ~~then all~~ ^{i is a} possible position in $\text{lastpos}(n)$, then all position in $\text{firstpos}(n)$ are in $\text{followpos}(i)$.

| node | followpos |
|------|--------------------|
| 1 | $\{1, 2, 3\}$ |
| 2 | $\{1, 2, 3\}$ |
| 3 | $\{4\}$ |
| 4 | $\{5\}$ |
| 5 | $\{6\}$ |
| 6 | - |

in the syntax tree firstpos of root is $\{1, 2, 3\}$. Let a be

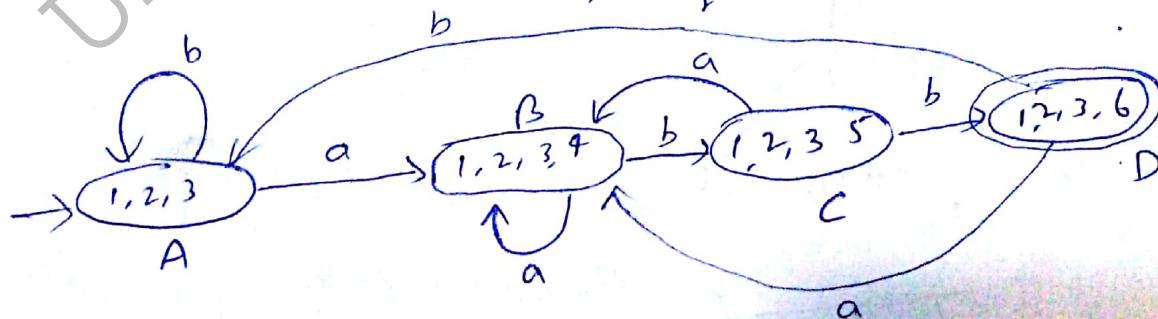
A and consider input symbol a

position 1 and 3 are for a so let $B =$

$$\text{followpos}(1) \cup \text{followpos}(3) = \{1, 2, 3, 4\} = B$$

$$D \xrightarrow{a} A, a = B$$

$$D \xrightarrow{a} B = \text{followpos}(2) = \{1, 2, 3\} = A$$

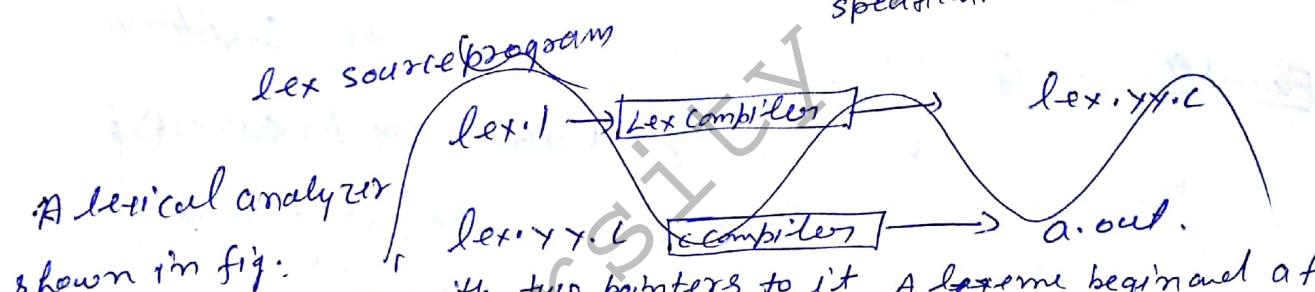


Lexical Analyzer Generator: Lex

LEX is tool that allows one to specify a lexical analyzer by specifying regular expressions to describe pattern for tokens.

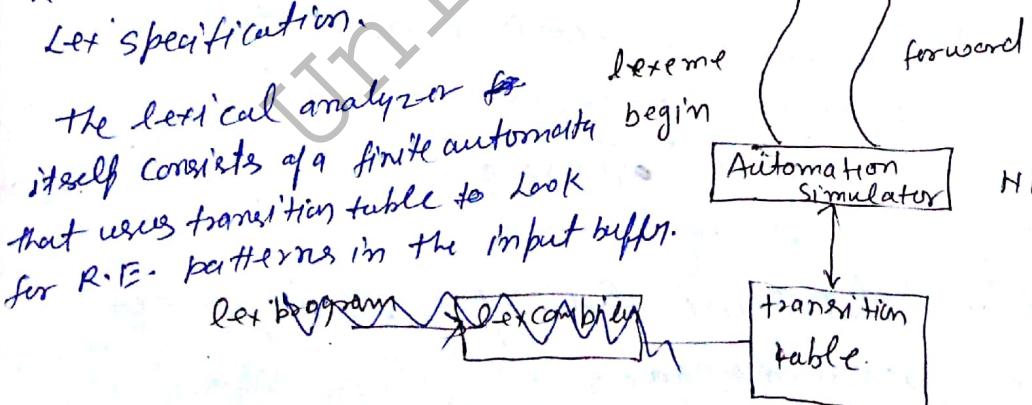
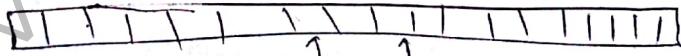
The input notation for the Lex tool is referred as a Lex language and the tool itself is the lex compiler. the lex compiler transforms the input patterns into a transition diagram and generates code.

e.g. An input file which we call lex.1 is written in the lex language and describes the lexical analyzer to be generated. the lex compiler transforms lex.1 to C program and is compiled by C compiler



A lexical analyzer shown in fig.

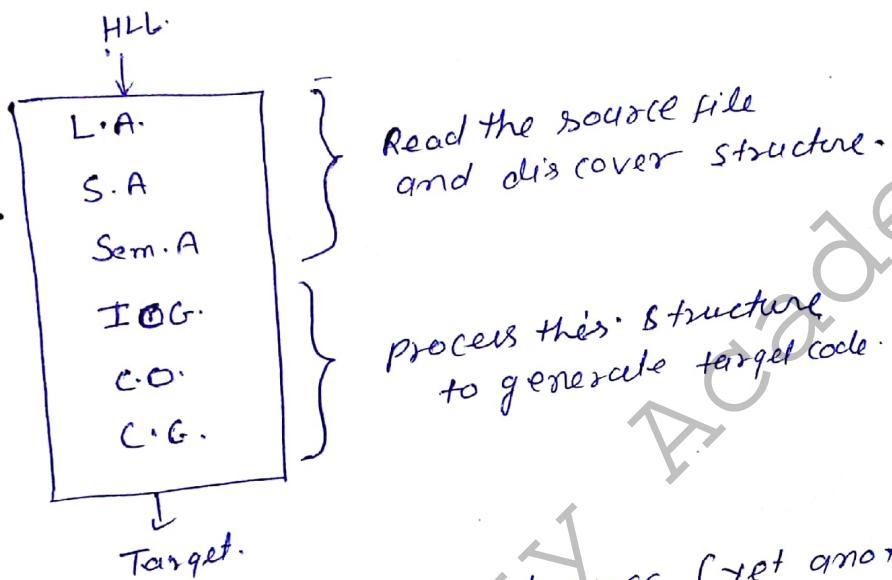
There is an input buffer with two pointers to it. A lexeme begins and a forward pointer. the lex compiler constructs a transition table for finite automata from R.E. patterns in the Lex specification.



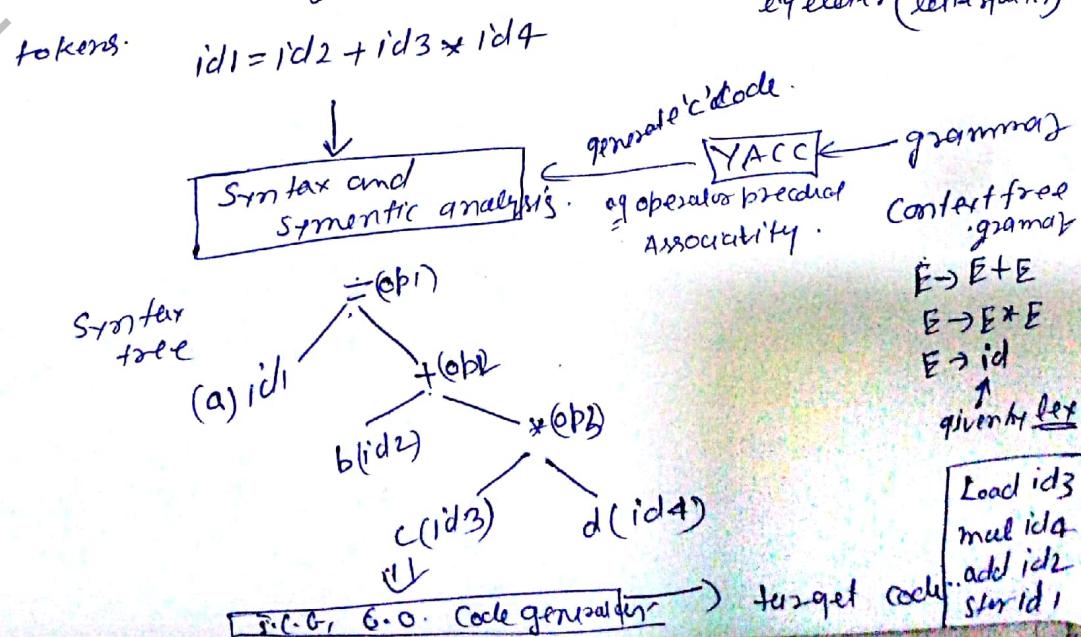
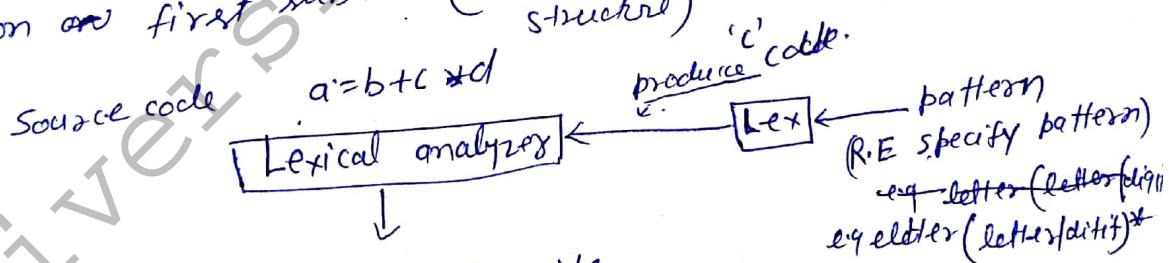
Design of lexical analyzer generator.

Lex and YACC

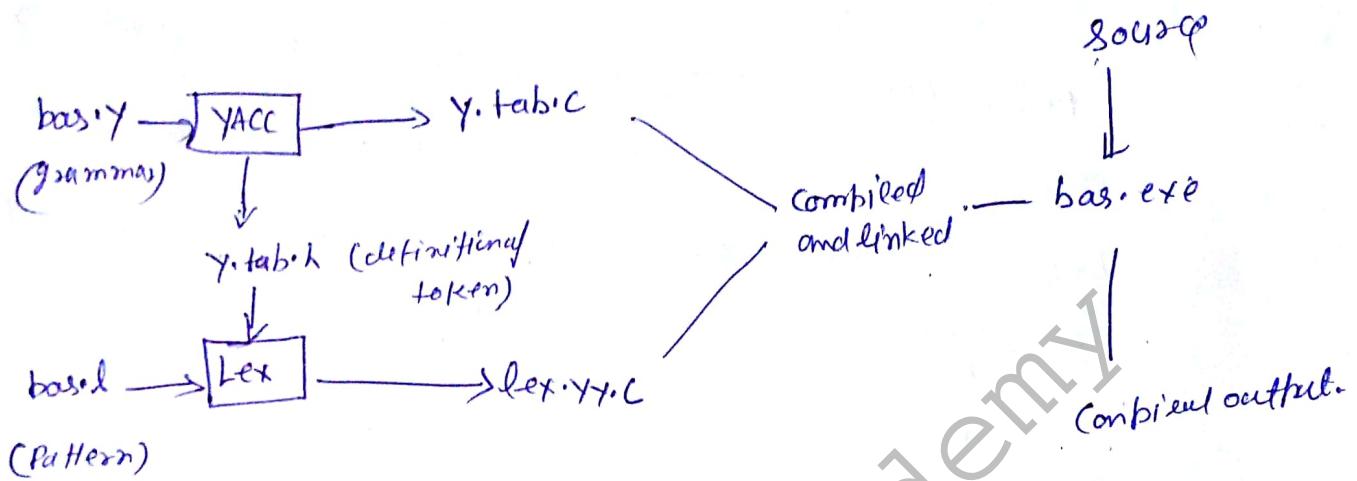
Before 1975 writing a compiler was very time consuming process. In 1975 Lesk and Johnson published papers on Lex and YACC, these utilities greatly simplify compiler writing.



the role of Lex and YACC (yet another compiler combiner) is on ~~and~~ first subtask (Read the source file and discover structure)



Writing a BASIC Compiler.



let is program

Syntax:

... definition ...
% %.
... rule ...
% %.
... subroutine ...

Yacc is program

... definition :
% %.
-- salary --
% %.
--- subroutine ---

Formal Grammar

A formal grammar is a set of rules for rewriting strings, along with a "start symbol" from which rewriting starts. Therefore a grammar is usually act as language generator.

A grammars have four table (V_N, Σ, P, S)

V_N is a set of non-terminal sometime called syntactic variable.

Σ is a set of terminal symbol sometime called tokens.

P is a set of production rule

S is a start symbol.

Example: $G = (V_N, \Sigma, P, S)$ is a grammar consisting,

$V_N = \{\text{Sentence}, \text{Noun}, \text{Verb}, \text{Adverb}\}$

$\Sigma = \{\text{Ram}, \text{Sam}, \text{ate}, \text{sang}, \text{well}\}$

$S = \{\text{Sentence}\}$

P consists of following production.

$\text{Sentence} \rightarrow \langle \text{Noun} \rangle \langle \text{Verb} \rangle$

$\text{Sentence} \rightarrow \langle \text{Noun} \rangle \langle \text{Verb} \rangle \langle \text{Adverb} \rangle$

$\text{Noun} \rightarrow \text{Rame} / \text{Sam}$

$\text{Verb} \rightarrow \text{age} / \text{sang}$

$\text{well} \rightarrow \text{well}$.

there are two important type of Formal grammar.

1. Context-Free-Grammars

2. Context-Sensitive-Grammar.

BNF Notation:

BNF (Backus-Naur Form) is main notation techniques for Grammar. It is developed by John Backus and Peter Naur to describe the syntax of a given language. It is formally define the grammar of a language.

A BNF Notation is written as:

$\langle \text{Symbol} \rangle ::= - \text{expression}$

where $\langle \text{Symbol} \rangle$ is a non-terminals and
- expression consist one or more sequence of symbol
separated by vertical bar "/" indicating the choice
symbol that never appear on a left side are terminals
and symbol that appearing left side are non-terminals
the non-terminal always enclosed between the pair

→ the $::=$ means the symbol on left side must be
replace with the expression on the right.

example: $\langle \text{number} \rangle ::= \langle \text{digit} \rangle / \langle \text{number} \rangle \langle \text{digit} \rangle$
 $\langle \text{digit} \rangle ::= 0 / 1 / 2 / 3 / 4 / 5 / 6 / 7 / 8 / 9$

Context-Free Grammars:

A CFG consist of terminals, non-terminals, a start symbol and production rule. terminals are basic symbols from which strings are formed. Non-terminal are syntactic variable that denote set of string.

the production of a grammar, specify the manner in which the terminals and non-terminals can be combined to form string.

Terminals: the following symbols are terminals.

- (i) lower case letter a, b, ... z
- (ii) operator symbol +, -, /, *
- (iii) punctuation symbol (,), ;, etc.
- (iv) the digit 0, ..., 9
- (v) Bold ~~for~~ string such as id.

Non-terminals:

- (i) uppercase letters A, B, C ... Z
- (ii) S - start symbol
- (iii) lower case italic name such as expr or stmt.
~~(id)~~

example:

$$E \rightarrow EAE / CE / -E / id$$
$$A \rightarrow + / - / * / / / \uparrow$$

E and A are non-terminal and E is a start symbol.
the remaining symbols are terminals.

Derivations:

The central idea here is that a production is treated as a rewriting rule in which the non-terminals on the left is replaced by the string on the right hand side of the production.

For example: the derivation for string -(id+id) is a sentence of grammar $E \rightarrow E+E / E*E / (E) / -E / id$

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

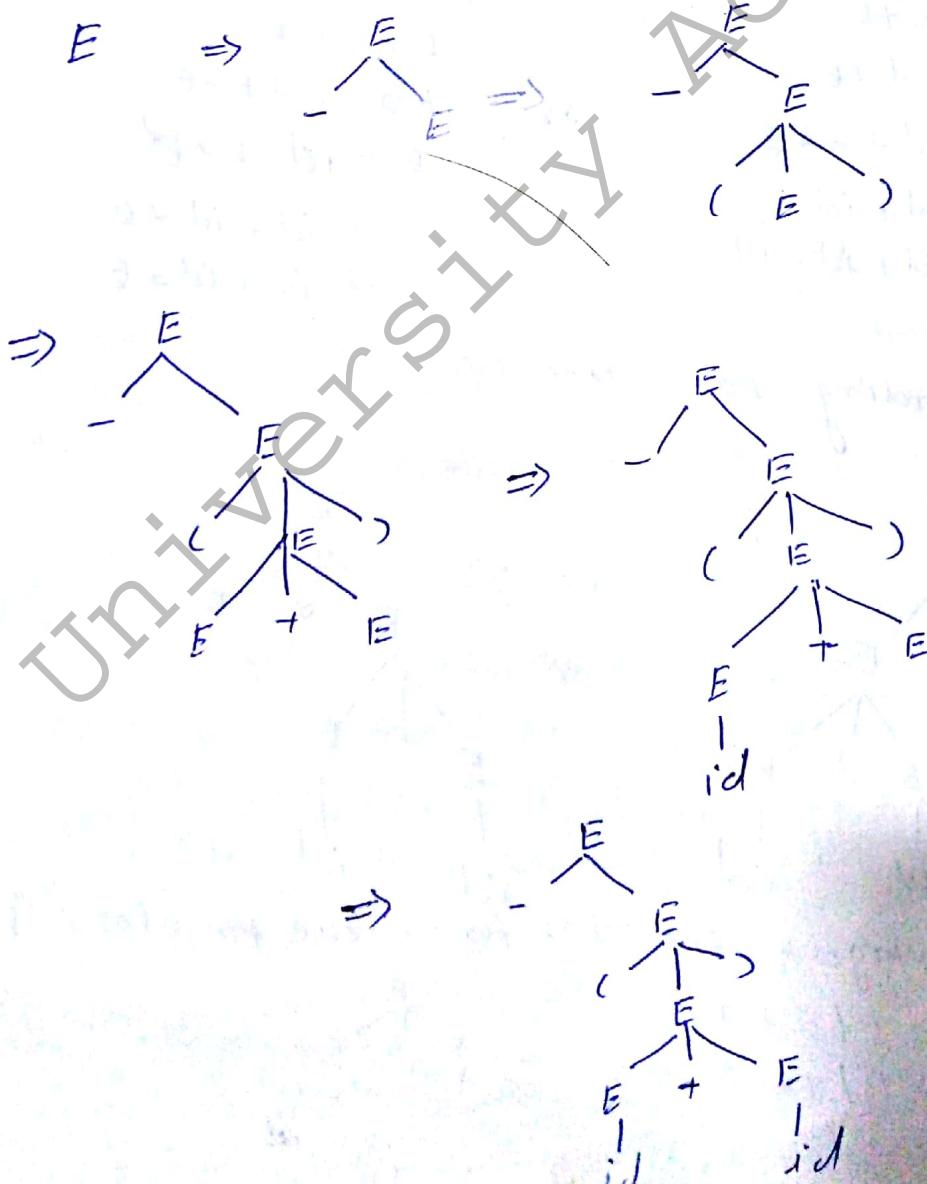
Parse trees:

A parse tree is a graphical representation of a derivation that filters out the choice regarding replacement order. Each interior node of a parse tree is labeled by some nonterminal and the leaves of the parse tree are labeled by nonterminal or terminals and read from left to right.

for example: the parse tree for $-(id+id)$.

the derivation for the string is

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



Ambiguity:

A Grammar that produce more than one parse tree for some language is said to be ambiguous. An ambiguous grammar is one that produce more than one left most or more than one right most derivation for the same sentences.

Example: the grammar is

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

and sentence $id + id * id$

Solution: it permit two distinct left most derivation for sentence

$$E \Rightarrow E+E$$

$$E \Rightarrow id+E$$

$$\Rightarrow id+E * E$$

$$\Rightarrow id+id * E$$

$$\Rightarrow id+id * id$$

$$E \Rightarrow E * E$$

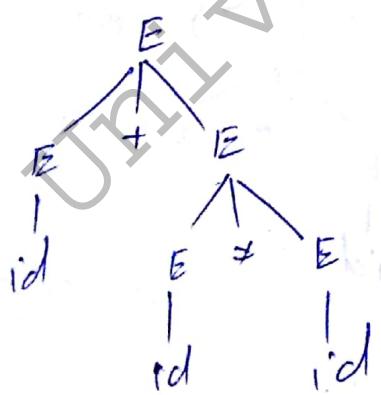
$$E \Rightarrow E+E * E$$

$$E \Rightarrow id+E * E$$

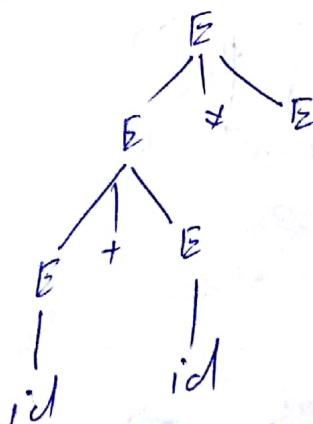
$$\Rightarrow id+id * E$$

$$\Rightarrow id+id * id$$

corresponding parse tree are



or

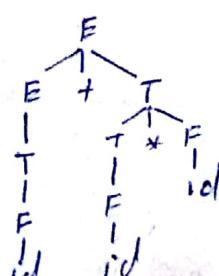


The unambiguous grammar for above for above is

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

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

$$F \rightarrow id$$

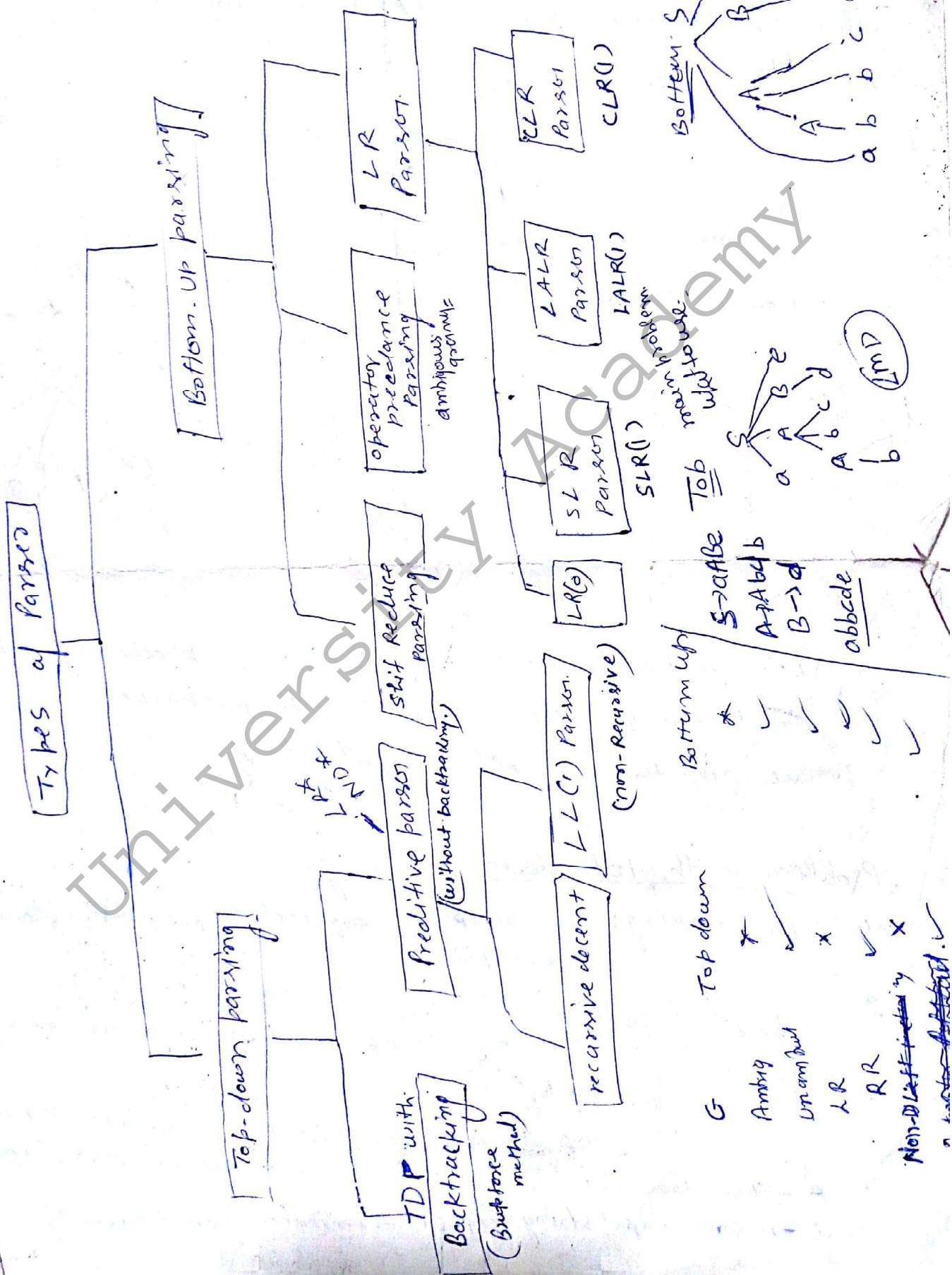


UNIT

Basic Parsing

Techniques

2



Top-down parser

In the top-down parsing the parse tree is generated from top to bottom (from root to leaves). The derivation terminates when the required input string terminates.

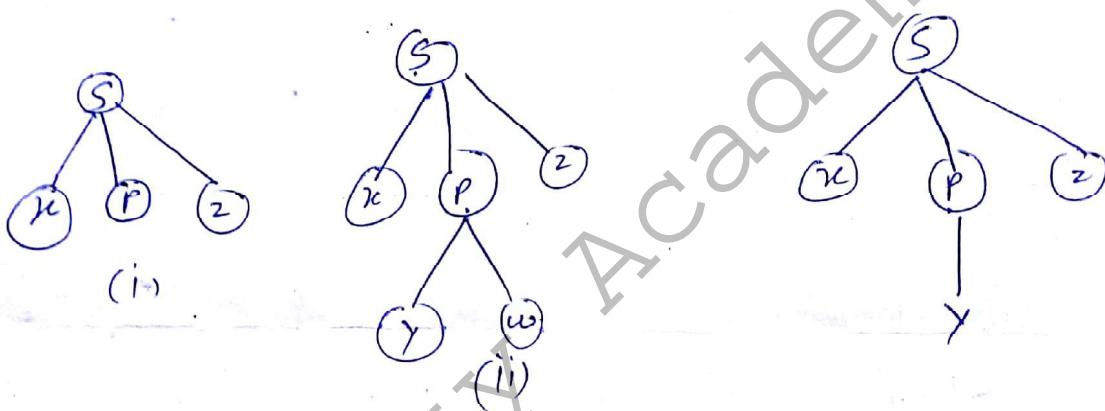
Consider the grammar

$$S \rightarrow xPz$$

$$P \rightarrow yw/y$$

and input string xyz

Now we construct the parse tree in top-down approach.



Note: if the production rule is not producing the correct input string then we need to Backtrack and then we get the correct input string.

Problem with top down parsing:

1. Backtracking: the concept of Backtracking problem shown in above example:

Left Recursion:

The left recursive grammar is a grammar which is as given below

$$A \xrightarrow{*} A\alpha$$

here $\xrightarrow{*}$ means the input in one or more steps.

A - Non-terminal

α - some input string (either terminal or non-terminal)

(2)

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

consider the example

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

then we eliminate left recursion by following rule:

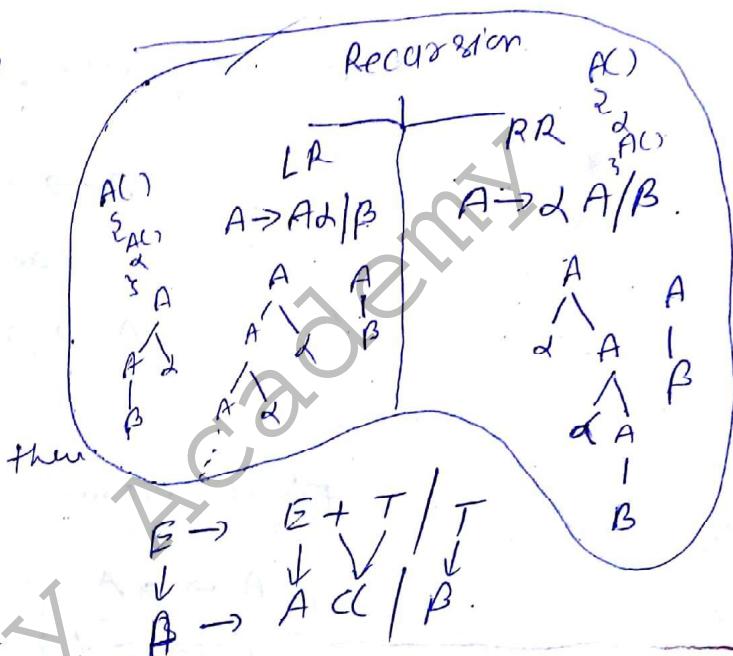
$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' \\ A' &\rightarrow \epsilon \end{aligned}$$

consider

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

we make $A \rightarrow A\alpha\beta / \beta$

E



then

$$E \rightarrow TE'$$

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

similarly

$$T \rightarrow TF / F$$

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

then

$$T \rightarrow FT'$$

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

Example: Consider the following grammar.

$$A \rightarrow ABd / Aa / a$$

$$B \rightarrow Be / b$$

Consider the rule

$$A \rightarrow ABd / Aa / a$$

we map this gram with rule $A \rightarrow Aa / B$

$$A \rightarrow ABd / a \Rightarrow A \rightarrow aA' \\ \text{and} \quad A' \rightarrow BdA' / \epsilon$$

$$A \rightarrow Aa / a \Rightarrow A \rightarrow aA'$$

$$A' \rightarrow aA' / \epsilon$$

Final grammar without Left recursive

$$A \rightarrow aA' \\ A' \rightarrow BdA' / aA' / \epsilon$$

the rule $B \rightarrow Be / b$ map with $A \rightarrow Aa / B$ grammar

becomes

$$B \rightarrow bB' \\ B' \rightarrow eB' / \epsilon$$

3-Left Factoring if the grammar is left factored then it becomes suitable for the use. Basically left factoring is used when it is not clear that which of two alternative is used to expand the non-terminal.

in general form

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

then left factored are

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

(3)

for example: (1)

$$S \rightarrow iEts \mid iEdses \mid a$$

$E \rightarrow b$

the left factored grammar becomes.

$$S \rightarrow iEts \mid iEdses \mid a$$

$$A \rightarrow aB_1 \mid aB_2 \mid \cancel{aB_3}$$

$$A \rightarrow 2B_1 \mid 2B_2 \mid 2B_3$$

 $\boxed{\alpha B_3}$

$$\begin{array}{c} A \swarrow \\ 2B_1 \end{array} \quad \begin{array}{c} A \swarrow \\ 2B_2 \end{array} \quad \begin{array}{c} A \swarrow \\ 2B_3 \end{array}$$

$$A \rightarrow AA'$$

$$A' \rightarrow B_1 \mid B_2 \mid B_3$$

$$a \rightarrow iEts, B_1 \rightarrow e$$

$$B_2 \rightarrow es$$

$$S \rightarrow iEtss' \mid a$$

$$s' \rightarrow es \mid e$$

Example 2 do left factoring in the grammar

$$A \rightarrow aAB \mid aA \mid a$$

$$A \rightarrow c(B_1) \mid cc(B_2) \mid c(B_3)$$

$$a \rightarrow a$$

$$B_1 \rightarrow AB$$

$$B_2 \rightarrow A$$

$$B_3 \rightarrow e$$

$$A \rightarrow aA'$$

$$A' \rightarrow AB \mid A \mid e$$

(4)

Ambiguity: Already discussed in previous lecture

Predictive parsing

there are two type of predictive parsers

- 1- Recursive descent
- 2- LL(1) Parser

1- Recursive descent parser :

A parser that uses collection of recursive procedure for parsing the given input string is called Recursive descent parser. In this type of parser the CFG is used to build recursive routines.

~~steps for construction of RD parser:~~

the R.H.S of the production rule is directly converted to a program. for each Non-terminal a separate procedure is written and body of the procedure is R.H.S. of the corresponding non-terminal.

example

consider the grammar

$$E \rightarrow \text{num} T$$

$$T \rightarrow * \text{num} T / \epsilon$$

$$\text{num} \rightarrow 0/1/2/\dots/9$$

and consider the input string $3 * 4$

| | | | | |
|---|---|---|--|----|
| 3 | * | 4 | | \$ |
|---|---|---|--|----|

$$E \rightarrow \text{num} T$$

| | | | | |
|---|---|---|--|----|
| 3 | * | 4 | | \$ |
| 3 | * | 4 | | \$ |

$$T \rightarrow * \text{num} T$$

| | | | | |
|---|---|---|--|----|
| 3 | * | 4 | | \$ |
| 3 | * | 4 | | \$ |

$$T \rightarrow * \text{num} T$$

| | | | | |
|---|---|---|--|----|
| 3 | * | 4 | | \$ |
| 3 | * | 4 | | \$ |

declare success =

Recursive descent parser

$$E \rightarrow E + E' i$$

$$E \rightarrow i E'$$

$$E' \rightarrow + i E' / e$$

$E()$

1. if ($l == '+'$)

2. { match ('+');

3. $E()$;

4. }

$E'()$

1. if ($l == '+'$)

2. { match ('+');

3. $E()$;

4. $E'()$;

5. else { return; }

match(char, it)

if ($l == t$)

$l = getchar();$
else printf("Error")

}

main()

1. $E()$;

2. if ($l == '$'$)

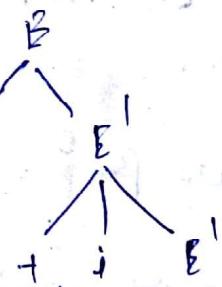
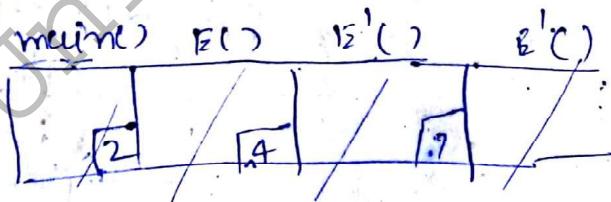
3. printf("Parsing successful")

}

Q.9.



Stack
Record



Predictive LL(1) parser.

This top-down parsing algorithm is non-recursive type parsing. In this type of parsing a table is built.

for LL(1) - the first L means the input is scanned from left to right. the second L means it uses leftmost derivation for input string and (1) means the input symbol uses only one symbol.

Construction of predictive LL(1) parser:

- 1- computation of FIRST and FOLLOW functions.
- 2- construct table using FIRST and FOLLOW function.
- 3- Parse the input string with help of parsing table.

FIRST function:

- 1- if the terminal symbol 'a' then $\text{FIRST}(a) = \{a\}$
- 2- if there is a rule $X \rightarrow \epsilon$ then $\text{FIRST}(X) = \{\epsilon\}$
- 3- for the rule $A \rightarrow X_1 X_2 X_3$ then

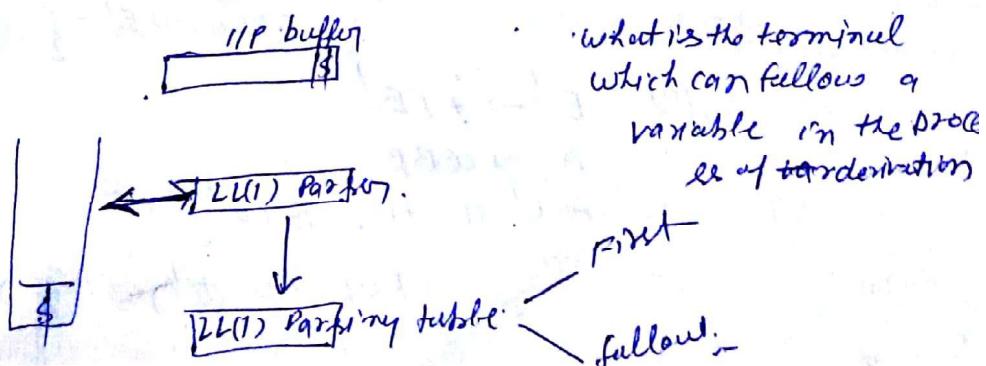
$$\text{FIRST}(A) = \text{FIRST}(X_1)$$

FOLLOW function:

- 1- For the start symbol S place \$ in follow(S)
- 2- if there is a production $A \rightarrow \alpha B \beta$, $B \neq \epsilon$ then every thing in $\text{FIRST}(B)$ without ϵ is placed in $\text{FOLLOW}(B)$.
- 3- if there is a production $A \rightarrow \alpha B$ or $A \rightarrow \alpha B \beta$ where $\text{FIRST}(B) \text{ contains } (\beta \Rightarrow \epsilon)$ then everything in $\text{FOLLOW}(A)$ is in $\text{FOLLOW}(B)$.

e.g $S \rightarrow aABCD$ $\text{First}(S) = a$

| | |
|--------------------------|------------------------------|
| $A \rightarrow b$ | $\text{first}(A) = b$ |
| $B \rightarrow C$ | $\text{first}(B) = C$ |
| $C \rightarrow d$ | $\text{first}(C) = d$ |
| $D \rightarrow \epsilon$ | $\text{first}(D) = \epsilon$ |



Architecture of LL(1)

Example:

Consider the grammar

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

$$\begin{aligned} S &\rightarrow E+T/T \\ T &\rightarrow FT/F \\ F &\rightarrow (E)/id \end{aligned}$$

First function:

$$FIRST(E) = FIRST(T) = FIRST(F) = \{\epsilon, id\}$$

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

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

FOLLOW function:

$$FOLLOW(E) \rightarrow$$

the production $F \rightarrow (E)$ shows that E followed by ')'. Hence ')' is placed in $FOLLOW(E)$. Since E is a start symbol add \$ to $FOLLOW(E)$.

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

$$FOLLOW(E') \div$$

$$(i) \quad E \rightarrow TE' \quad \text{by rule 3}$$

$$A \rightarrow CCB\beta$$

$$A \Rightarrow E, C = T, B = E', \beta = \epsilon$$

then everything in $fOLLOW(A)$ is in $fOLLOW(B)$ since $fOLLOW(A)$ is in $fOLLOW(E')$.

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

$$(ii) \quad E' \rightarrow +TE'$$

$$A \rightarrow CCB\beta$$

$$A \Rightarrow E', C = +T, B = E', \beta = \epsilon$$

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

(5)

$\text{FOLLOW}(T)$

(i) $E \rightarrow TE'$

$A \rightarrow \alpha B \beta$

$A = E, \alpha = \epsilon, B = T, \beta = E'$ by rule 2

every thing $\text{FIRST}(\beta) - \epsilon$ is in $\text{FOLLOW}(B)$

since

$$\text{FIRST}(E') - \epsilon = \{+\}$$

$$\text{FOLLOW}(\beta) = \{+\}$$

(ii) $E' \rightarrow +TE'$

$A \rightarrow \alpha B \beta$ by rule 2.

$$\text{FIRST}(E') - \epsilon = \{+\}$$

since in both production T following E'
all follow E' is also in $\text{FOLLOW}(T)$.

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

$\text{FOLLOW}(T')$

$T \rightarrow FT'$ by rule 3

$A \rightarrow \alpha B$

$\text{Follow}(A)$ is in $\text{FOLLOW}(B)$ since

$\text{FOLLOW}(T)$ is in $\text{FOLLOW}(T')$

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

$\text{FOLLOW}(F)$

$T \rightarrow FT'$ or $T' \rightarrow *FT'$ by rule 3

$A \rightarrow \alpha B \beta$

$A = T, \alpha = \epsilon, B = F, \beta = T'$

$$\begin{aligned}\text{FOLLOW}(F) &= \{\text{FIRST}(T') - \epsilon\} \\ &= \{* \}\end{aligned}$$

$A \rightarrow \alpha B \beta$

$A = T', \alpha = *, B = F, \beta = T'$

$$\begin{aligned}\text{FOLLOW}(T') &= \text{FOLLOW}(F) \\ &= \{+,), \$\} \\ \text{FOLLOW}(F) &= \{+, *,), \$\}\end{aligned}$$

Summarize of FIRST and FOLLOW

| Symbol | FIRST | FOLLOW |
|--------|---------|---------------|
| E | {c, id} | {), \$} |
| E' | {+} | {), \$} |
| T | {c, id} | {+,), \$} |
| T' | {*, id} | {+,), \$} |
| F | {(, id} | {+, *,), \$} |

Parsing table for grammar

| | id | + | * | (|) | \$ |
|----|---------------------|-----------------------|-----------------------|---------------------|--------------------|--------------------|
| E | $E \rightarrow TE'$ | | | $E \rightarrow TE'$ | | |
| E' | | $E' \rightarrow +TE'$ | | | $E' \rightarrow G$ | $E' \rightarrow E$ |
| T | $T \rightarrow FT'$ | | | $T \rightarrow FT'$ | | |
| T' | | $T' \rightarrow *$ | $T' \rightarrow *FT'$ | | $T' \rightarrow G$ | $T' \rightarrow G$ |
| F | $F \rightarrow id$ | | | $F \rightarrow (E)$ | | |

Now if the input string $id + id * id \neq id \neq$ can be parsed.

stack

- \$E
- \$E'T
- \$E'T'F
- \$E'T'id
- \$E'T'id
- \$E'T+
- \$E'T+
- \$E'T
- \$E'T'F
- \$E'T'id
- \$E'T'

| Input | Output |
|---------------------|-----------------------|
| $id + id * id \neq$ | $E \rightarrow TE'$ |
| $id + id * id \neq$ | $T \rightarrow FT'$ |
| $id + id * id \neq$ | $F \rightarrow id$ |
| $+ id * id \neq$ | $T' \rightarrow E$ |
| $+ id * id \neq$ | $E' \rightarrow +TE'$ |
| $id * id \neq$ | $T \rightarrow FT'$ |
| $id * id \neq$ | $F \rightarrow id$ |
| $id * id \neq$ | |

Bottom Up Parsing:

The parse tree is constructed from bottom to up that is from leaves to root is called Bottom up parsing. In this process, the input symbols are placed at the leaf nodes after successful parsing. The parse tree is created starting from leaves, the leaf nodes reduced further to internal node and these internal node further reduced root and root node obtained.

example consider the grammar

$$S \rightarrow aABe$$

$$A \rightarrow Abc/b$$

$$B \rightarrow d$$

the input string abbcde can be reduced to S by following steps.

abbcde

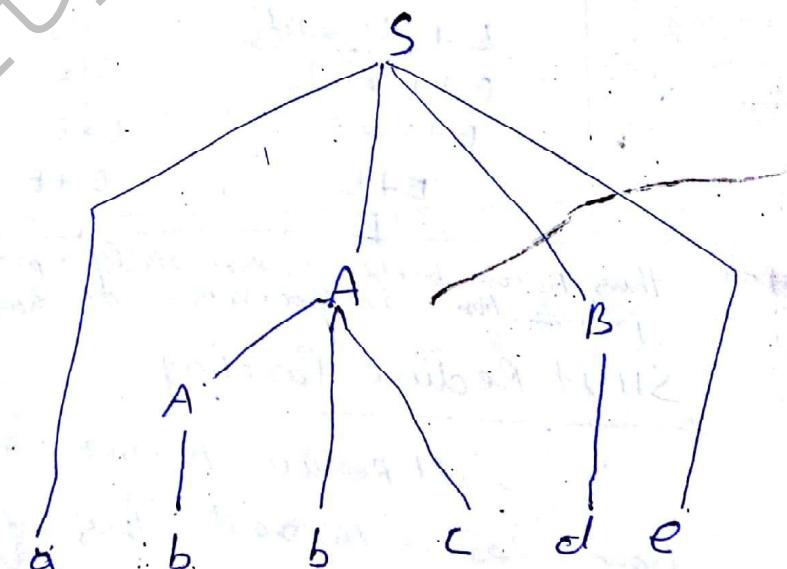
aAbcde

aAde

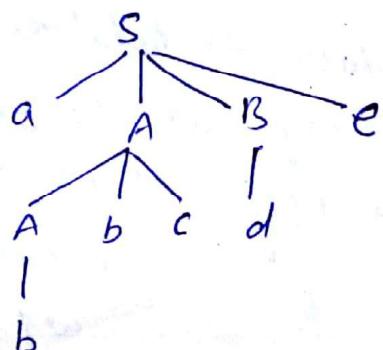
aABe

S

bottom up procedure



or



Handles and handle pruning:

In bottom-up parsing is to find substring that that could be reduced by appropriate non-terminal. such a substring is called handle.

Handle is a string of substring that matches the right side of production and we can reduce such string by a non-terminal on left hand side production.

example: Consider the grammar

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow id$$

and input string $id_1 + id_2 * id_3$

| Right sentential form | Handle | Reducing production. |
|-----------------------|---------|-----------------------|
| $id_1 + id_2 * id_3$ | id_1 | $E \rightarrow id$ |
| $E + id_2 * id_3$ | id_2 | $E \rightarrow id$ |
| $E + E * id_3$ | id_3 | $E \rightarrow id$ |
| $E + E * E$ | $E * E$ | $E \rightarrow E * E$ |
| $E + E$ | $E + E$ | $E \rightarrow E + E$ |
| E | | |

thus Bottom parser is essentially a process of detecting handles and using them in reduction. This process is called handle pruning.

Shift Reduce Parsing:

Shift Reduce parser attempts to construct parse tree from leaves to root. thus it works on Bottom up parser principle. A convenient way to implement a shift-Reduce parser is to use a stack to hold the grammar symbols and an buffer to hold the string w to be parsed. initially the stack is empty and the string w is on input as follows:

stack
\$

Input
 $w \#$

there are four possible actions a shift-reduce parser can make:

1. shift: In a shift action, the next input symbol is shifted onto the top of the stack.
2. Reduce: In Reduce Action, the parser knows the right end of the handle is at the top of the stack. Then reduction of it by appropriate rule.
3. Accept: In accept action the parser announces successful completion of parsing.
4. Error: A situation in which parser can not either shift or reduce the symbols, it cannot even perform accept actions is called error.

Example:

Consider the following grammar

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E - id$$

and input string $id_1 + id_2 * id_3$.

| Stack | Input | Action |
|---------------|-------------------------|---------------------------------|
| \$ | $id_1 + id_2 * id_3 \$$ | shift |
| \$ id, | $+ id_2 * id_3 \$$ | reduce by $E \rightarrow id$ |
| \$ E | $+ id_2 * id_3 \$$ | shift |
| \$ E + | $id_2 * id_3 \$$ | shift |
| \$ E + id | $* id_3 \$$ | reduce by $E - id$. |
| \$ E + E | $id \$$ | shift |
| \$ E + E * | $\$$ | shift |
| \$ E + E * id | $\$$ | Reduce by $E \rightarrow id$ |
| \$ E + E + E | $\$$ | Reduce by $E \rightarrow E * E$ |
| \$ E + E | $\$$ | Reduce by $E \rightarrow E + E$ |
| \$ E | $\$$ | accept. |

Question:

Consider the grammar

$$S \rightarrow TL;$$

$$T \rightarrow \text{int} \mid \text{float}$$

$$L \rightarrow L, id \mid id$$

Parse the input string `int id, id;` using shift-reduce parser.

Operator Precedence Parser

A grammar is said to be operator precedence if there are no production with ϵ in right side and no two adjacent non-terminal in any production.

Consider the grammar for arithmetic expression

$$E \rightarrow EAE \mid (E) \mid -E \mid id$$

$$A \rightarrow + \mid - \mid * \mid / \mid ^$$

This grammar is not an operator precedence grammar. If we substitute for A each of its alternatives, we obtain the following operator grammar

$$E \rightarrow E+E \mid E-E \mid E*E \mid E/E \mid E^E \mid (E) \mid -E \mid id$$

In operator precedence parsing, we define three precedence relations, \prec, \equiv, \succ between certain pairs of terminals.

These precedence relations have following meaning

$a \prec b$ a "yields precedence to" b

$a \equiv b$ a "has same precedence as" b

$a \succ b$ a "takes precedence over" b

(b)

the operator precedence relation for id, +, *, & shown
in table

| | id | + | * | \$ |
|----|----|---|---|----|
| id | > | > | > | |
| + | < | > | < | > |
| * | < | > | > | > |
| \$ | < | < | < | |

Left associativity

Now consider the string $id + id * id$
we will insert \$ symbol at the start and end of the
input string. we will also insert operator precedence
by referring table.

$\$ < id > + < id > * < id > \$$

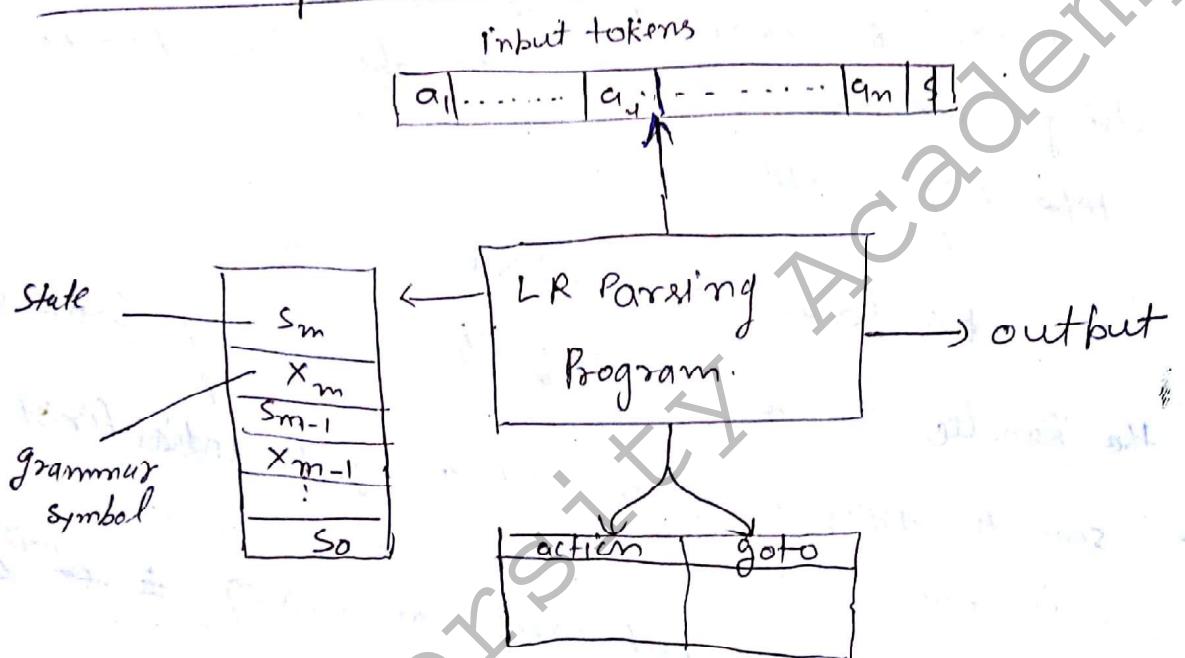
- the handle can be found by the following process.
- 1- scan the string from the left to right until first $>$ is encountered.
 - 2- scan then scan backwards over any $=$ or $<$'s encountered
 - 3- the handle is a string between $<$ and $>$
- the parsing can be done as follows:

| obtain handle and Reduction. | |
|----------------------------------|---|
| $\$ < id > + < id > * < id > \$$ | handle id is obtained $E \rightarrow id$ |
| $E + < id > * < id > \$$ | handle id |
| $E + E * < id > \$$ | " " " " |
| $E + E * E$ | Remove all non-terminal. |
| $+ *$ | insert \$ at begin and end, insert precedence operators |
| $+ < + > \$$ | + becomes handle. $E \rightarrow E * E$ |
| $\$ \$$ | + becomes handle. $E \rightarrow E + E$ |
| | Parsing done. |

LR Parsers

This is most efficient method of bottom-up parsing which can be used to parse the large class of CFG. This method is also called LR(K) Parsing. Here the LR parser detects syntactical errors very efficiently. Here L stands for left to right scanning R stands for right most derivation in reverse K is no. of input symbol. When K is omitted K is assumed to be 1.

Structure of LR Parser



LR Parser consists of an input, an output, a stack, and a driver program and a parsing table that has two part (action and goto).

- A driver program is the same for all LR parsers any the parsing table changes from one parser to another.
- The parsing table consist two part action and goto. it determines s_m the state currently on top of the stack and a_i the current input symbol. it then consult action $[s_m, a_i]$, the parsing action table entry for state s_m and input a_i , which can have one of the four value.

- 1 - Shift s, where s is state
- 2 - Reduce by a grammar
- 3 - Accept
- 4 - error

→ the goto takes a state and grammar symbol as arguments and produce a state

Types of LR Parser:

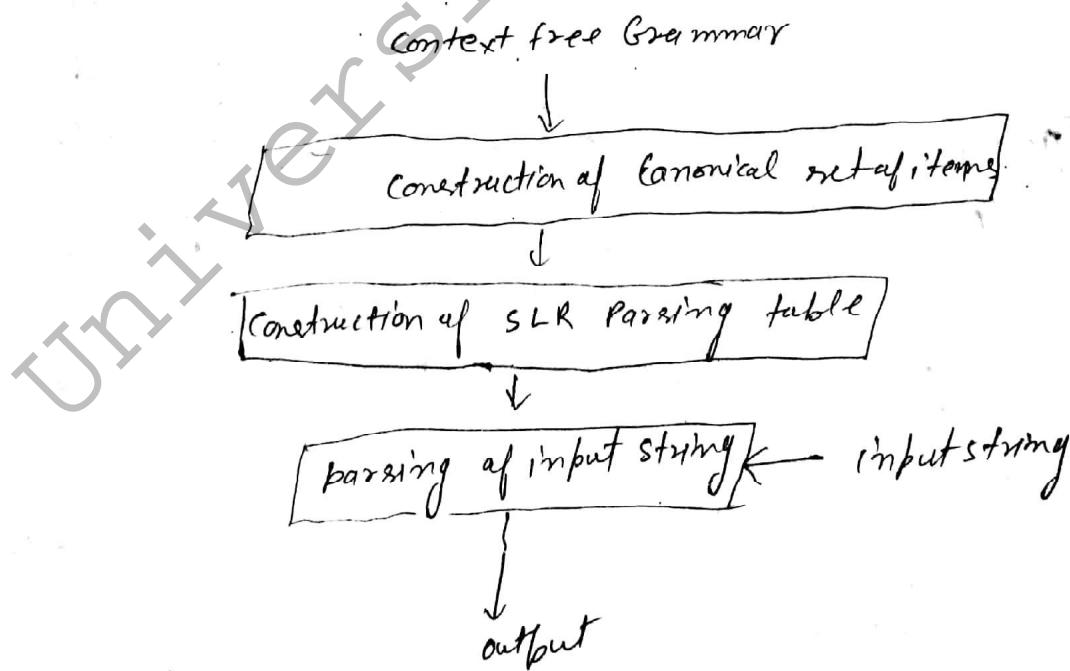
1- LR(0) SLR (Simple LR parser)

2- LALR (Lookahead LR parser)

3- CLR (canonical LR parser)

Simple LR parser. (LR(0))

It is the weakest of the three method but it is easiest to implement. the parsing can be done as follows.



Construction of Canonical set of items:

- i- for grammar G initially add $S' \rightarrow S$ in set of grammar production and construct the augmented grammar for G.

2. for each set of item I_i in G and for each grammar symbol α (may be terminal or non-terminal), add closure (I_i, α) . This process should be repeated by applying $\text{goto}(I_i, \alpha)$ for each α in I_i such that $\text{goto}(I_i, \alpha)$ is not empty and not in G .

Example:

consider the grammar

1. $E \rightarrow E + T$
2. $E \rightarrow T$
3. $T \rightarrow T * F$
4. $T \rightarrow F$
5. $F \rightarrow (E)$
6. $F \rightarrow \text{id}$

in this grammar
 $E' \rightarrow E$ and we will add
make augmented
grammar such as

| $I_0:$ |
|---------------------------------|
| $E \rightarrow \cdot E$ |
| $E \rightarrow \cdot E + T$ |
| $E \rightarrow \cdot T$ |
| $T \rightarrow \cdot T * F$ |
| $T \rightarrow \cdot F$ |
| $F \rightarrow \cdot (E)$ |
| $F \rightarrow \cdot \text{id}$ |

| $\text{goto}(I_0, E)$ |
|---|
| $I_1: E' \xrightarrow{E} E \cdot$ $E \xrightarrow{} E \cdot + T$ |

| $\text{goto}(I_0, T)$ |
|--|
| $I_2: E \xrightarrow{T} T \cdot$ $T \xrightarrow{} T \cdot * F$ |

| $\text{goto}(I_0, F)$ |
|--|
| $I_3: T \xrightarrow{F} F \cdot$ $F \xrightarrow{} F \cdot E$ |

| $\text{goto}(I_0, F)$ |
|--|
| $I_4: F \xrightarrow{F} (\cdot E)$ $E \xrightarrow{} \cdot E + T$ $E \xrightarrow{} \cdot T$ $T \xrightarrow{} \cdot T * F$ $T \xrightarrow{} \cdot F$ $F \xrightarrow{} \cdot (E)$ $F \xrightarrow{} \cdot \text{id}$ |

| $\text{goto}(I_0, \text{id})$ |
|--|
| $I_5: F \xrightarrow{F} \text{id} \cdot$ |

| $\text{goto}(I_1, +)$ |
|--|
| $I_6: E \xrightarrow{E} E + \cdot T$ $T \xrightarrow{} \cdot T * E$ $T \xrightarrow{} \cdot F$ $F \xrightarrow{} \cdot (E)$ $F \xrightarrow{} \cdot \text{id}$ |

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

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

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

goto($I_2, *$)

I_2 :

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

goto(I_4, E)

I_8

$$\begin{aligned} F &\rightarrow (E) \\ E &\rightarrow E \cdot + T \end{aligned}$$

goto(I_6, T)

I_9 :

$$\begin{aligned} E &\rightarrow E + T \\ T &\rightarrow T \cdot * E \end{aligned}$$

goto(I_7, F)

$$I_{10} \quad T \rightarrow T * F.$$

goto($I_8, ()$)

$$I_{11} \quad F \rightarrow (E).$$

$$\text{goto}(I_4, T) = I_2$$

~~$$\text{goto}(I_4, ()) = I_4$$~~

~~$$\text{goto}(I_4, id) = I_5$$~~

~~$$\text{goto}(I_4, F) = I_3$$~~

$$\text{goto}(I_8, +) = I_6$$

$$\text{goto}(I_6, F) = I_3$$

$$\text{goto}(I_6, ()) = I_4$$

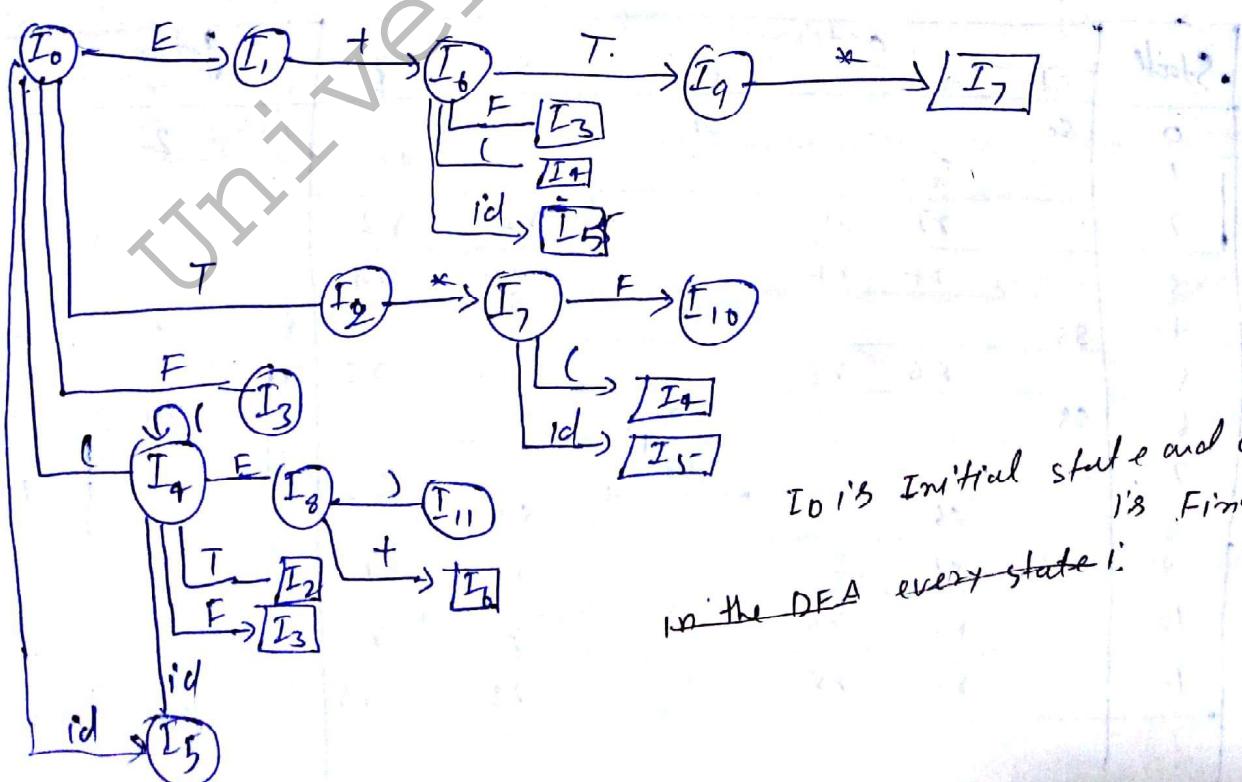
$$\text{goto}(id, id) = I_5$$

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

$$\text{goto}(I_7, id) = I_5$$

$$\text{goto}(I_9, *) = I_7$$

Now we can derive the DFA for above set of item!



Construction of SLR Parsing Table:

- ① Initially construct set of item $C = \{I_0, I_1, \dots, I_n\}$ for input grammar.
- ② The parsing action are based on each item I_i , as given below.
 - (a) if $A \rightarrow a \cdot aB$ is in I_i and $\text{goto}(I_i, a) = I_j$ then set action $[i, a]$ as "shift J" or "ST"
 - (b) if there is a rule $A \rightarrow a \cdot$ is in I_i then set action $[i, a]$ to reduce $A \rightarrow a\ell$ for all symbols a , where $a \in \text{FOLLOW}(A)$. ℓ must not be an augmented grammar S' .
 - (c) if $S' \rightarrow S_0$ is in I_i then the entry in action table action $[i, \$] = \text{accept}$.
- ③ the goto part of SLR table can be filled as the goto transition for state i is considered for non terminable only if $\text{goto}(I_i, A) = I_j$ then $\text{goto}[I_i, A] = j$.
4. all the entries not defined considered as error.

| Stack | action | | | | | | E | T | F |
|-------|--------|-------|-------|---|-------|----------|---|-------|----|
| | id | * | (|) | \$ | acc. | | | |
| 0 | S_5 | | | | S_7 | | 1 | | |
| 1 | | S_6 | | | | | | 2 | |
| 2 | | r_2 | S_7 | | | r_2 | | r_2 | |
| 3 | | r_4 | r_4 | | | r_4 | | r_4 | |
| 4 | S_5 | | | | S_8 | | 6 | | |
| 5 | | r_6 | r_6 | | | r_6 | | r_6 | |
| 6 | S_5 | | | | S_9 | | | 2 | 3 |
| 7 | S_5 | | | | S_9 | | | | 3 |
| 8 | | S_6 | | | | S_{11} | | | 10 |
| 9 | | r_1 | S_7 | | | r_1 | | r_1 | |
| 10 | | r_3 | r_3 | | | r_3 | | r_3 | |
| 11 | | r_5 | r_5 | | | r_5 | | r_5 | |

Parsing of Input string:

let us take the example input: $id * id + id$

| Stack | Input Buffer | Action Table | Goto Table | Parsing action. |
|------------|-------------------|----------------|---------------|--------------------------------------|
| \$0 | $id * id + id \$$ | $[0, id] = S5$ | | shift |
| \$0ids | $* id + id \$$ | $[5, *] = R6$ | $[0, F] = 3$ | reduce $F \rightarrow id$ |
| \$0F3 | $* id + id \$$ | $[3, *] = R4$ | $[0, T] = 2$ | reduce $T \rightarrow F$ |
| \$0T2 | $* id + id \$$ | $[2, *] = S7$ | | shift |
| \$0T2*7 | $id + id \$$ | $[7, id] = S5$ | | shift |
| \$0T2*7ids | $+ id \$$ | $[5, +] = R6$ | $[7, F] = 10$ | Reduce $F \rightarrow id$ |
| \$0T2*7F10 | $+ id \$$ | $[10, +] = R3$ | $(0, T) = 2$ | Reduce by $T \rightarrow T \times F$ |
| \$0T2 | $+ id \$$ | $[2, +] = R2$ | $(0, E) = 1$ | Reduce by $E \rightarrow T$ |
| \$0E1 | $+ id \$$ | $[1, +] = S6$ | | shift |
| \$0E1+6 | $id \$$ | $[6, id] = S5$ | | shift |
| \$0E1+6*7 | $\$$ | $[5, \$] = R6$ | $[6, F] = 3$ | $F \rightarrow id$ |
| \$0E1+6F3 | $\$$ | $[3, \$] = R4$ | $[6, T] = 9$ | Reduce $T \rightarrow F$ |
| \$0E1+6T9 | $\$$ | $[9, \$] = S1$ | $(0, E) = 1$ | $E \rightarrow E + T$ |
| 012F10 | $\$$ | accept | | accept. |

Example 2. Construct the SLR Parsing table for following grammar.

$$S \rightarrow SS$$

$$S \rightarrow a$$

$$S \rightarrow c$$

Ans: It is not SLR(0) parsing.

Bimul 2

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

$$T \rightarrow TF / F$$

$$F \rightarrow F \times / a / b$$

Constructed SLR Parsing table

Tables

Ans: Follow(B) = {+, \$}

$$D = \{+, \$, a, b\}$$

$$R = \{+, \$, a, b\}$$

CLR (Canonical LR parser)

The canonical set of item is the parsing technique in which a lookahead symbol is generated while constructing set of item. Hence collection of set of items is referred as LR(1).

We follow the following steps in CLR:

1. Construction of Canonical set of items along with lookahead.
2. Building Canonical LR Parsing table.
3. Parsing the put string using Canonical LR table.

Construction of canonical set of items along with lookahead

1. For the grammar G initially add $s \rightarrow \cdot s$ in set of item.
2. The closure function can be computed as follows.
for each item $A \rightarrow C \cdot B \beta$ [$A \rightarrow C \cdot B \beta, a$] and rule $\beta \rightarrow Y$ and $b \in \text{FIRST}(\beta, a)$
such that $[B \rightarrow \cdot Y, b]$ is not in I then add $[B \rightarrow \cdot Y, b] \rightarrow I$.
3. Similarly the goto function can be computed as follows.
for each item $[A \rightarrow C \cdot B \beta, q]$ is in I and rule $[A \rightarrow C(B \cdot \beta, a)]$ is not in goto items then add $[A \rightarrow C(B \cdot \beta, a)]$ to goto items.

(12) Consider the grammar

$$S \rightarrow CC$$

$$C \rightarrow cC/d \Rightarrow$$

$$S' \rightarrow S$$

$$S \rightarrow .EC$$

$$C \rightarrow .cC/d$$

Augmented grammar

initially we add $S' \rightarrow S, \$$ as the first rule in I_0
now match

$$S' \rightarrow S, \$ \text{ with}$$

$$A \rightarrow C \cdot BB, a$$

~~so~~ $A = S'$, $cL = \epsilon$, $B = S$, $\beta = \epsilon$, $a = \$$

if there is a production $B \rightarrow r$

$\therefore S \rightarrow .CC$
and

$$\begin{aligned} b &\in \text{FIRST}(\beta a) \\ \epsilon &\in \text{FIRST}(a) \\ \epsilon &\in \text{FIRST}(\$) \\ \epsilon & (\$) \end{aligned}$$

$$\boxed{\begin{array}{l} I_0 \\ S' \rightarrow S, \$ \\ S \rightarrow .CC, \$ \\ C \rightarrow .cC, c/d \\ C \rightarrow .d, c/d \end{array}}$$

so $S \rightarrow .CC, \$$ will be added to I_0 .

Now $S \rightarrow .CC, \$$ match with

$$A \rightarrow C \cdot BB, a \quad A = S, cL = \epsilon, B = C, \beta = C, a = \$$$

if there is a production $B \rightarrow r$

$$\Rightarrow C \rightarrow .cC \quad \text{and } b \in \text{First}(\beta a)$$

$$\epsilon \in \text{First}(ca)$$

$$\epsilon \in \text{First}(c)$$

$$\epsilon \in \{c, d\}$$

them

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

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

will added to I_0

Now we compute goto function.

goto(I_0, S)
 $I_1 \quad S' \rightarrow S, \$$

$S \rightarrow C \cdot C, \$$

$A \rightarrow C \cdot B \beta, a$

$A \rightarrow S, CC = C, B = C, \beta = \epsilon, a = \$$

the production
 $B \rightarrow r$

$\Rightarrow C \rightarrow \cdot cC$ and $b \in \text{FIRST}(Ba)$
 $C \rightarrow \cdot d$
 $\in \text{FIRST}(a)$
 $\in \text{First}(\$)$
 $\in \{\$, \}\}$

since $C \rightarrow \cdot CC, \$$
 $C \rightarrow \cdot d, \$$ } added to I_2

goto(I_0, a)
 $I_3:$
 $C \rightarrow c \cdot C, c/d$
 $C \rightarrow \cdot CC, c/d$
 $C \rightarrow \cdot d, c/d$

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

$A \rightarrow C \cdot B \beta, a$

$A = C, CC = C, B = C, \beta = \epsilon, a = \$$

the production

$B \rightarrow r \Rightarrow C \rightarrow \cdot cC$

$C \rightarrow \cdot d$

and $b \in \text{first}(Ba)$

$\in \text{First}(a)$

$\in \text{First}(c, d)$

$\in \{c, d\}\}$

so

$C \rightarrow \cdot CC, c/d$
 ~~$C \rightarrow \cdot d, c/d$~~ added to I_3

goto(I_0, d)
I_d: $C \rightarrow d \cdot a/d$

goto(I_2, C)
I_c: $S \rightarrow CC \cdot \$$

I_6 :
 goto(I_2, c)
 $c \rightarrow c.C, \$$
 $c \rightarrow .CC, \$$
 $c \rightarrow .d, \$$

Now: $C \rightarrow c.C, \$$
 $A \rightarrow CC, BB, \$$
 $B \rightarrow r \Rightarrow C \rightarrow .CC$ and $b \in \text{First}(B)$
 $c \rightarrow .d$
 $\epsilon, \$$

goto(I_2, d)
 $I_7: c \rightarrow d, \$$

goto(I_3, C)
 $I_8: c \rightarrow aC, a/c$

goto(I_6, C)
 $I_9: C \rightarrow a, \$$

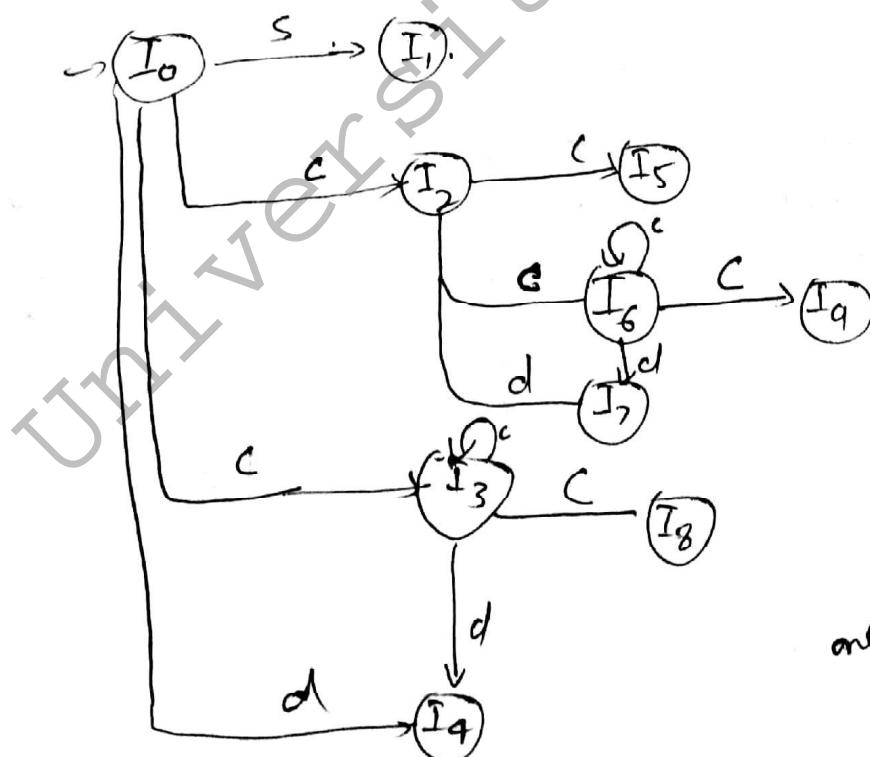
goto(I_3, c) = I_3

goto(I_3, d) = I_4

goto(I_6, c) = I_6

goto(I_6, d) = I_7

DFA for above canonical set.



only I_0 is initial all other state is final state.

Construction of CLR Parsing table:

The algorithm for CLR parsing table is similar as SLR parsing table.

| state | action | | | goto |
|-------|--------|----|--------|------|
| | a | d | f | |
| 0 | s3 | s4 | f | s |
| 1 | s2 | s4 | accept | c |
| 2 | s6 | s7 | accept | 1 |
| 3 | s3 | s7 | | 2 |
| 4 | r3 | r5 | | s5 |
| 5 | r3 | r3 | r1 | 6 |
| 6 | s6 | s7 | | 8 |
| 7 | . | . | r3 | 9 |
| 8 | r2 | r2 | | |
| 9 | . | . | r2 | |

Parsing the input using CLR or LR(1) table.
Input string is "add"

| state | input Buffer | action | goto | parsing action |
|-----------|--------------|---------------------------------|--------------|----------------|
| \$0 | add \$ | $(\epsilon, \epsilon) = s_3$ | | shift |
| \$0c3 | cdd \$ | $(3, \epsilon) = s_3$ | | shift |
| \$0c3c3 | dd \$ | $(3, d) = s_4$ | | shift |
| \$0c3c3d4 | d \$ | $(4, d) = r_3$ | $(3, c) = 8$ | C → d |
| \$0c3c3C8 | d \$ | $(8, d) = r_2$ | $(3, c) = 8$ | C → CC |
| \$0c3C8 | d \$ | $(8, d) = r_2$ | $(0, c) = 9$ | C → CC |
| \$0C2 | d \$ | $(2, d) = s_7$ | | shift |
| \$0C2d7 | \$ | $[7, \epsilon] = r_3$ | $(2, c) = 5$ | C → d |
| \$0C2C5 | \$ | $[5, \epsilon] = r_1$ | | S - CC |
| \$0S1 | \$ | $[1, \epsilon] = \text{accept}$ | | |

LALR Parser (Lookahead-LR Parser)

This method is often used in practice because the table obtained by it are considerably smaller than CLR table. In fact the states of SLR and LALR parsing are always same. Most of the programming languages use LALR parser.

We follow the same steps as discussed in SLR and CLR parsing techniques and those are

1. Construction of canonical set of items along with lookahead.
2. Building LALR Parsing table.
3. Parsing the input string using LALR Parsing table.

1. Construction of canonical set of items along with lookahead.
Same procedure as CLR used in LALR canonical set of items along with lookahead symbol.

For Example consider the previous example of CLR.

$$\begin{array}{l} S \rightarrow CC \\ C \rightarrow cC \\ C \rightarrow d \end{array}$$

$$\begin{array}{l} I_0: \begin{array}{l} S' \rightarrow \cdot S, \$ \\ S \rightarrow \cdot CC, \$ \\ C \rightarrow \cdot cC, c/d \\ C \rightarrow \cdot d, c/d \end{array} \end{array}$$

$$\begin{array}{l} I_3: \text{goto}(I_0, C) \\ \quad C \rightarrow \cdot cC, a/d \\ \quad C \rightarrow \cdot aC, a/d \\ \quad C \rightarrow \cdot d, a/d \end{array}$$

$$\begin{array}{l} I_7: \text{goto}(I_2, d) \\ \quad C \rightarrow d, \$ \\ I_8: \text{goto}(I_3, C) \\ \quad C \rightarrow \cdot CC, c/d \\ I_9: \text{goto}(I_6, C) \\ \quad C \rightarrow \cdot cC, \$ \end{array}$$

$$\begin{array}{l} I_1: \text{goto}(I_0, S) \\ S' \rightarrow S, \$ \end{array}$$

$$\begin{array}{l} I_4: \text{goto}(I_0, d) \\ \quad C \rightarrow d, a/d \end{array}$$

$$\begin{array}{l} I_2: \text{goto}(I_0, C) \\ S \rightarrow C \cdot C, \$ \\ C \rightarrow \cdot CC, \$ \\ C \rightarrow \cdot d, \$ \end{array}$$

$$\begin{array}{l} I_5: \text{goto}(I_2, C) \\ S \rightarrow CC, \$ \\ I_6: \text{goto}(I_2, C) \\ \quad C' \rightarrow C \cdot C, \$ \\ \quad C \rightarrow \cdot CC, \$ \\ \quad C \rightarrow \cdot d, \$ \end{array}$$

Construction of LALR Parsing table:

- Step 1. may construct $C = \{I_0, I_1, \dots, I_n\}$ the collection of set of LR(1) items.

2. merge the two states I_i and I_j if the first component (the production rule with dots) are matching and create a new state replacing one of the older state such as

$$I_{ij} = I_i \cup I_j$$

3. the parsing action are given below.

 - (a) if $[A \rightarrow \alpha \cdot aB, b]$ is in I_i and $\text{goto}(I_i, a) = I_j$ then create an entry in the action table $\text{action}[I_i, a] = \text{shift } j$
 - (b) if there is a production $[A \rightarrow \alpha \cdot, a]$ in I_i then in the action table $\text{action}[I_i, a] = \text{reduce by } A \rightarrow \alpha - A \text{ should not be } S'$.
 - (c) if there is a production $S' \rightarrow S, \$$ in I_i then action
 $(\cancel{I_i}, A) - \text{action}[i, \$] = \text{accept}$.

4. the goto part of LALR table can be filled as:

if $\text{goto}(I_i, A) = I_j$ then $\text{goto}(I_i, A) = j$

5. if the parsing action conflict then the algorithm fails to produce LALR parser. and grammar is not LALR. all the entries not defined by rules 3 and 4 are considered to be error.

In the previous example after the step 2 of LALR¹ parsing tableau algorithm, the canonical item set will be.

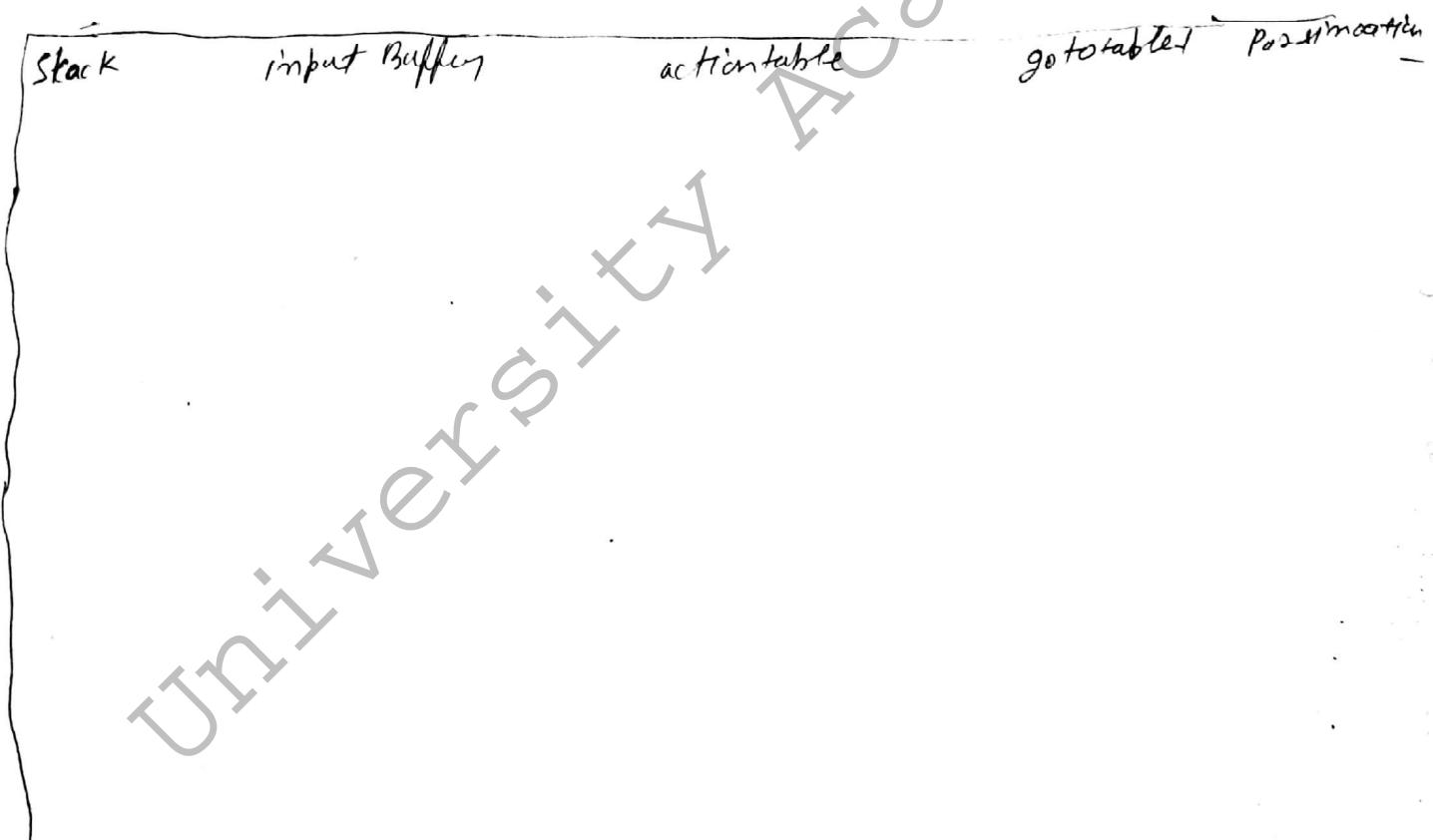
| | | |
|-----------------------------|------------------------------------|---|
| I_0 | $I_2 \text{ goto}(I_0, C)$ | $I_4 \rightarrow \text{ goto}(I_0, d)$ $C \rightarrow d, a/d / \$$ |
| $s \rightarrow . s, \$$ | $s \rightarrow . C, \$$ | |
| $s \rightarrow . CC, \$$ | $C \rightarrow . CC, \$$ | $I_5 : \text{ goto}(I_2, C)$ |
| $C \rightarrow . CC, a/d$ | $C \rightarrow . d, \$$ | $s \rightarrow CC, \$$ |
| $C \rightarrow . d, a/d$ | | |
| $I_1, \text{ goto}(I_0, S)$ | $I_{36} : \text{ goto}(I_0, C)$ | $I_{89} : \text{ goto}(I_3, C)$ |
| $S' \rightarrow S, . \$$ | $C \rightarrow @. C, a/d / \$$ | $C \rightarrow @C, a/d / \$$ |
| | $C \rightarrow . @C, @/d / \$$ | |
| | $C \rightarrow . d, @/d / \$$ | |
| | $\text{ goto}(I_{36}, C) = I_{89}$ | $\text{ goto}(I_0, C) = I_{36}$ |
| | $\text{ goto}(I_{36}, d) = I_{89}$ | $\text{ goto}(I_2, d) = I_9$ |

Construction of LALR Parsing table:

(15)

| State | action | | | goto | |
|-------|--------|-----|-----|------|-----|
| | c | d | d | s | c |
| 0 | s36 | | | | |
| 1 | | s47 | | 1 | 2 |
| 2 | s36 | | s47 | | |
| 36 | s36 | | s47 | | -55 |
| 47 | r3 | | r3 | r3 | 89 |
| 5 | | | | r1 | |
| 89 | r2 | | r2 | r2 | |

parse the input string "aadd".



E Show the grammar $S \rightarrow Aa \mid bAc \mid Bc \mid bBc$

$$A \rightarrow d$$

$B \rightarrow d$

is LR(1) but not LALR(1) grammar.

Syntax Directed Translation

$SDT = \text{Grammar} + \text{Semantics rule}$.

Syntax Directed translation are augmented rule to the grammar that facilitate semantic analysis.

There are two notation for associating semantic rule.

(1) Syntax Directed Definition

(2) Syntax Directed translation scheme.

1- SyntaxDirected Definition: In Syntax Directed translation definition. Each grammar production $A \rightarrow c_1 c_2 \dots c_k$ has associated with it a set of semantic rules ~~form~~ of form $b_i = f(b_1, b_2, \dots, b_k)$ where f 's function and either.

i. b_i is synthesized attribute of A and c_1, c_2, \dots, c_k are attributes belonging to the grammar symbol of production

ii. b_i is inherited attribute of one of the grammar symbol on R.H.S Production and c_1, c_2, \dots, c_k are attributes belonging to the grammar symbol of the production.

2. Syntax Directed translation scheme: It indicate the order in which ~~translation rule~~ semantic rule are to be evaluated.

input string — parser tree — dependency graph — translation order.

Arithmetical evaluation

$$E \rightarrow E + T \quad \{ E.value = E.value + T.value \}$$

$$E \rightarrow T \quad \{ E.value = T.value \}$$

$$T \rightarrow T * F \quad \{ T.value = T.value * F.value \}$$

$$T \rightarrow F \quad \{ T.value = F.value \}$$

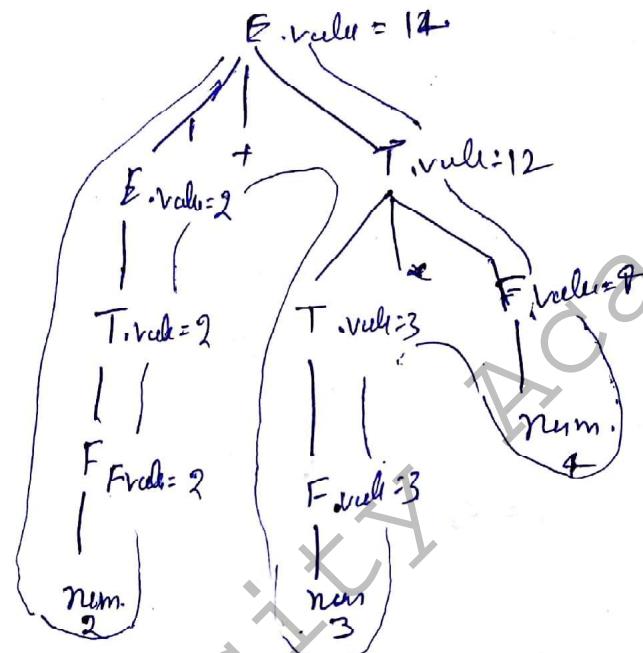
$$F \rightarrow \text{num} \quad \{ F.value = \text{num.value} \}$$

Attribute: Synthesized Attribute:

$$\begin{array}{rcl} 2+3*4 & = & 14 \\ | & | & | \\ \text{num} & \text{num} & \text{num} \end{array}$$

$$\text{left child value} =$$

Parse tree:



fork down left
fork right

reduction
action

Convert infix to Postfix

$$2 + 3 * 4 .$$

$$E \rightarrow E + T \quad \{ \text{print}(+) \} \quad ①$$

$$E \rightarrow T \quad \{ 3 \} \quad ②$$

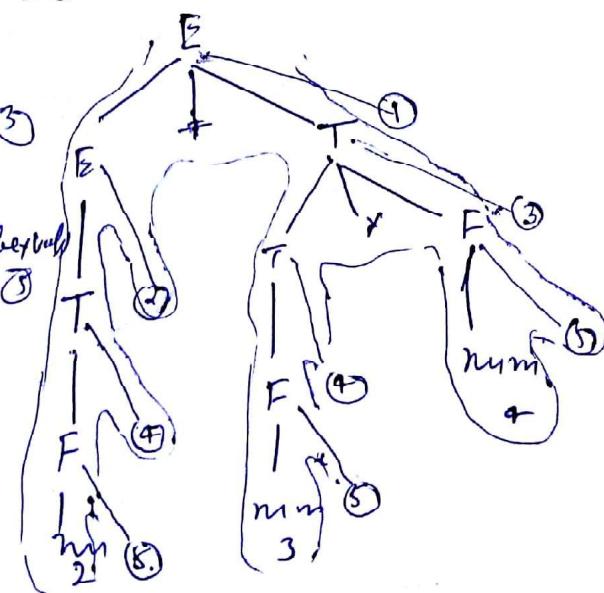
$$T \rightarrow T * F \quad \{ \text{print}(*3) \} \quad ③$$

$$T \rightarrow F \quad \{ 3 \} \quad ④$$

$$F \rightarrow \text{num.} \quad \text{print}(\text{num.value}) \quad ⑤$$

Top down
Parse tree

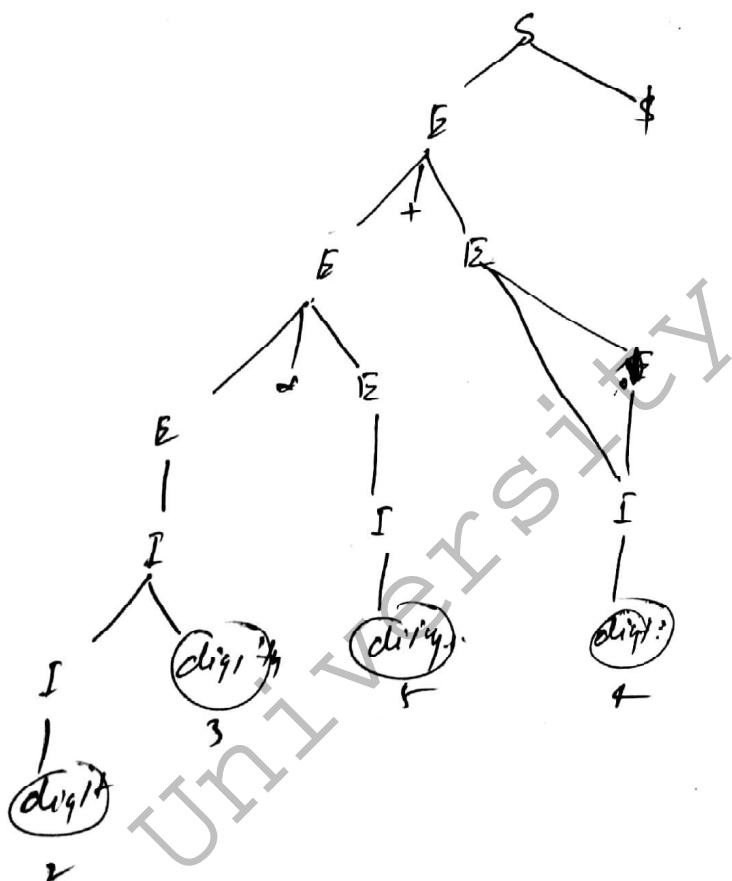
$$2 3 5 * +$$



Q. Give the SDT for desk calculator.

| | |
|--------------------------------|--|
| $S \rightarrow E \$$ | $\{ \text{print } E.\text{val} \}$ |
| $E \rightarrow E' + E^2$ | $\{ E.\text{val} = E'.\text{val} + E^2.\text{val} \}$ |
| $E \rightarrow E' * E^2$ | $\{ E.\text{val} = E'.\text{val} * E^2.\text{val} \}$ |
| $E \rightarrow (E')$ | $\{ E.\text{val} = E'.\text{val} \}$ |
| $I^2 \rightarrow I$ | $\{ I^2.\text{val} = I.\text{val} \}$ |
| $I \rightarrow I' digit^4$ | $\{ I.\text{val} = 10 * I'.\text{val} + \text{digit}^4.\text{lexval} \}$ |
| $I \rightarrow \text{digit}^4$ | $\{ I.\text{val} = \text{digit}^4.\text{lexval} \}$ |

93 * 5 + 4



SDT one two types

S- attributed SBT

- 1- We only synthesize attributes.
 2. Semantic actions are placed at right end of production

$A \rightarrow B \in \Sigma^*$ $\{ \dots \}$
Postfix S.D.T.

3. attributes are evaluated directly B.R.



Bear $A \rightarrow L^m \quad \{ \quad L_{i,j} = f(A, i) ; \quad m_{\cdot, i} = f(L, s) ; \quad A \cdot s = f(m, s); \}$

$$A \rightarrow Q.R \quad \{ \quad R.i = f(A.i); \quad Q.i = f(R.i); \quad A.i = f(Q.s); \}$$

a) S-enantiomer b) L-enantiomer c) both enantiomers

A → BCD

$$A \circ S = f(B, S, C, S, D, S)$$

C.i = A.i child taken from Parent

$$C_1 = \beta \cdot i$$

or

$\mathbf{G}_1 = \mathbf{D}_{\mathbf{M}_1}$

1

Introducing
the Validation ability

$$A \rightarrow BC \quad \{ B:S \rightarrow A:S \}$$

a) scattered b) Lanthanides c) boron d) none

L-Attributed SDT

1. needs both inherited and synthesized.
each inherited is restricted to inherit either from parent or left sibling.
 2. Semantic actions are performed only when in RHS.

$A \rightarrow \{3BC\}$
 $1DSS3E$
 $1FG\{S\}$

3. Attributes are evaluated by traversing parent to one depth first left to right.



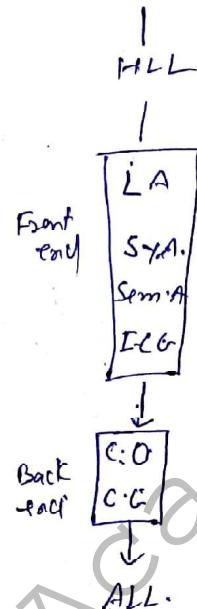
Depth first
left two right

Intermediate code.

"Three Address Code!"

Single pass compiler.

multiblock/two pass compiler.



Benefit of multiblock compiler.

1. Beneficial for design compilers for different machine programming language for same machine.
2. Beneficial for design compilers for same programming language for different machine. [Lecture-6]

There are three forms of Intermediate Code.

1. Postfix Notation.

2. Syntax tree or Parse tree

3. Three-Address Code.

→ quadruples representation
→ Triples representation

Postfix Notation.

infix notation, Prefix Notation, Postfix Notation,

+ab

abt

@bt c*

abc+x

a+b
(a+b)*c

1. No Parentheses are needed

a*(b+c)

2. No operator precedence needed

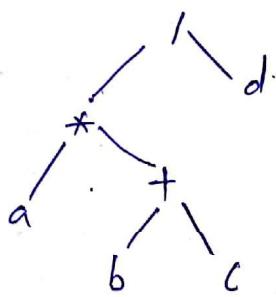
3. No associativity needed

4. Easy Implementation and scanning left to right.

5. memory efficient (stack)

2. Syntax tree or Parse tree: The parse tree is useful intermediate language representation for a source code. Syntax tree is a tree in which each leaf represents operand and interior nodes are operators.

$$a * (b + c) / d$$



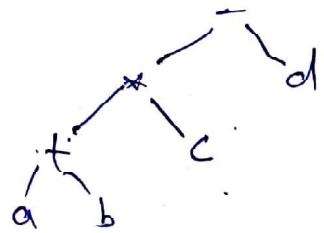
Infix: $a * (b + c) / d$

Prefix: $* a + b / c$

Postfix: $a b c + * d /$

$\text{infix} = \text{left}, \text{root}, \text{right}$
 $\text{prefix} = \text{root}, \text{left}, \text{right}$
 $\text{postfix} = \text{left}, \text{right}, \text{root}$

$$((a+b)*c)-d$$



Prefix: $- * + a b c d$

Postfix: $a b + c * d -$

3. Three Address Code:

Three address code is a sequence of statement of the general form $x := y \text{ op } z$

$x = y + z$ is translated into:

$$t_1 := y * z$$

$$t_2 := x + t_1$$

x, y, z are variable, constant, or compiler generated temporary variables.

ob is operator

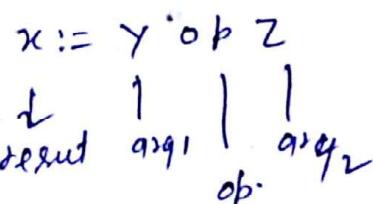
t_1, t_2 are compiler generated temporary variables.

Type of Three address statement

1. Assignment Statement of form $x := y \text{ op } z$ 'op' is binary or logical operation.
2. Unary operation $x := \text{op } y$
3. Copy statement $x := y$
4. Conditional jump if $x \text{ zelob } y$ goto lavel.
zelob: $=, \neq, >, \geq, \leq, \leq$
5. Procedure call:
Param x_1 :
Param x_2 :
6. Index Assignment: $x := y[i]$
 $x[i] = y$
7. Address and pointers $x := \&y$
 $x = *y$

Implementation of three address codes

Quadruples - A quadruple is a record structure with four fields called op, arg1, arg2, result.



$$a = b * -c + b * -c$$

Three address codes

$$\begin{aligned} t_1 &= -c \\ t_2 &= b * t_1 \\ t_3 &= -c \\ t_4 &= b * t_3 \\ t_5 &= t_2 + t_4 \\ a &= t_5 \end{aligned}$$

| | op | arg1 | arg2 | result |
|-----|------------|-------|-------|--------|
| (0) | - Unary | c | | t_1 |
| (1) | * | b | | t_2 |
| (2) | - Unary | c | | t_3 |
| (3) | * | b | | t_4 |
| (4) | + | t_2 | t_4 | t_5 |
| (5) | Assignment | t_5 | | a |

Triples.

A triple is a record with three fields. op, arg₁, arg₂ are op, arg₁, and arg₂ contain variable names or pointers to symbol table.

| | op | arg ₁ | arg ₂ |
|-----|-----------|------------------|------------------|
| (0) | - (unary) | c | |
| (1) | * | b | (0) |
| (2) | - (unary) | c | |
| (3) | * | b | (2) |
| (4) | + | (1) | (3) |
| (5) | assign | a | (4) |

More about translation.

1. Array Reference in arithmetic
2. Procedure Call; obfuscation
3. Declaration and case statement.

1- Array Reference in Arithmetic

1- One Dimensional array.

Base address (B) width (w)
4 byte.

| Actual address in the memory | 1100 | 1104 | 1108 | 1112 | 1120 |
|--------------------------------|------|------|------|------|------|
| Element | 15 | 10 | 11 | 44 | 39 |
| Address with respect to Array. | 0 | 1 | 2 | 3 | 4 |
| lower Bound LB | | | | | |

$$\text{Address of } A[i] = B + w \times (i - LB)$$

e.g. given value B=1020, LB, 1300, w=2 ...

$$\begin{aligned}
 A[1700] &=? \\
 I & 1020 + 2 \times (1700 - 1300) \\
 &= 1020 + 2 \times 400 \\
 &= 1820 \text{ Ans}
 \end{aligned}$$

Translate the following code in three Address code.

~~Notes~~

```

int i
int a[10];
i=1
while(i<=10)
{
    if(a[i]==0)
        i=i+1
}

```

1. $i = 1$
2. if $i \leq 10$ goto (4)
3. goto (B)
4. $a_i = 4i$
5. $a_{[i]} = 0$
6. $i = i + 1$
7. goto (2)
8. Stop

$$a[i] = B + w \times (i - LB) = B + 4 \times (i - 0) = B + 4i$$

Address calculation in two-dimensional array:

Column Index.

- (1) row major
- (2) column major

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| 0 | 8 | 6 | 5 | 4 |
| 1 | 3 | 2 | 9 | 6 |
| 2 | 7 | 3 | 2 | 1 |

$A[2][4]$

(1) row-major :

| | | | | | | | | | | | |
|-------|---|---|---|-------|---|---|---|-------|---|---|---|
| 18 | 6 | 5 | 4 | 3 | 2 | 9 | 6 | 7 | 3 | 2 | 1 |
| row-0 | | | | row-1 | | | | row-2 | | | |

$$\text{Add } A[i][j] = B + w \times [N \times (i - L_r) + (j - L_c)]$$

(2) Column major :

| | | | | | | | | | | | |
|-------|---|---|-------|---|---|-------|---|---|--------|---|---|
| 8 | 3 | 7 | 6 | 2 | 3 | 5 | 9 | 2 | 4 | 6 | 1 |
| Col-0 | | | Col-2 | | | Col-3 | | | Col-3. | | |

$$\text{Add } A[i][j] = B + w \times [(i - L_r) + (j - L_c) \times m]$$

(1) row wise: $B = 1500$, $w = 1$ byte $L_r = 0$ $L_c = 0$ $N = 4$, $M = 3$

$$\begin{aligned}
 A[2][3] &= 1500 + 1 \times [4 \times 2 + 3] \\
 &= 1500 + 1 \times 11 \\
 &= 1511 \checkmark
 \end{aligned}$$

(2) column wise: $A[2][3] : 1500 + [2 + 3 \times 3]$

$$\begin{aligned}
 &= 1500 + 11 \\
 &= 1511 \cdot \checkmark
 \end{aligned}$$

Translate following fragment into three address code.

```
int i  
int a[10][10]  
i=0  
while (i<10)  
{  
    a[i][i] = 1;  
    i++;  
}
```

- | | |
|---|--------------------------|
| 1 | i = 0 |
| 2 | if i < 10 goto(4) |
| 3 | goto(8) |
| 4 | st _i = 44 * i |
| 5 | a[st _i] = 1 |
| 6 | i = i + 1 |
| 7 | goto(2) |
| 8 | <u>Stop</u> |

We assume 32 bit word form and 4 byte allocation per word.

$$\begin{aligned}a[i][i] &= \text{base } B + 4 \times [10 * (i) + i] \\&= B + 4 \times [11 * i] \\&= B + 44 * i\end{aligned}$$

Three Address code

Procedure call

A procedure call like $P(A_1, A_2, A_3, \dots, A_n)$ may have many addresses for one statement in three address code. so it's shown as sequence of $n+1$ statement.

Param A₁

Param A₂

:

Param A_n

call P, n

P is name of procedure.

n is integer indicating number of actual parameters.

Example swapping of two numbers using function.

three Address code

Int main()

{

int x, y

swap(&x, &y)

}

void swap(int *a, *b)

{ int t;

t = *b

*b = *a

*a = t;

}

1. Param &x

2. Param &y

3. call swap, 2

4. t = *b

5. *b = *a

6. *a = t

7

8

Case Statement

switch expression:

{

case value : statement

case value : statement

:

case value : statement

}.

e.g. generate the three address code for
switch (ch)

{

case 1 : c = a + b;
break;

case 2 : c = a - b;

break;

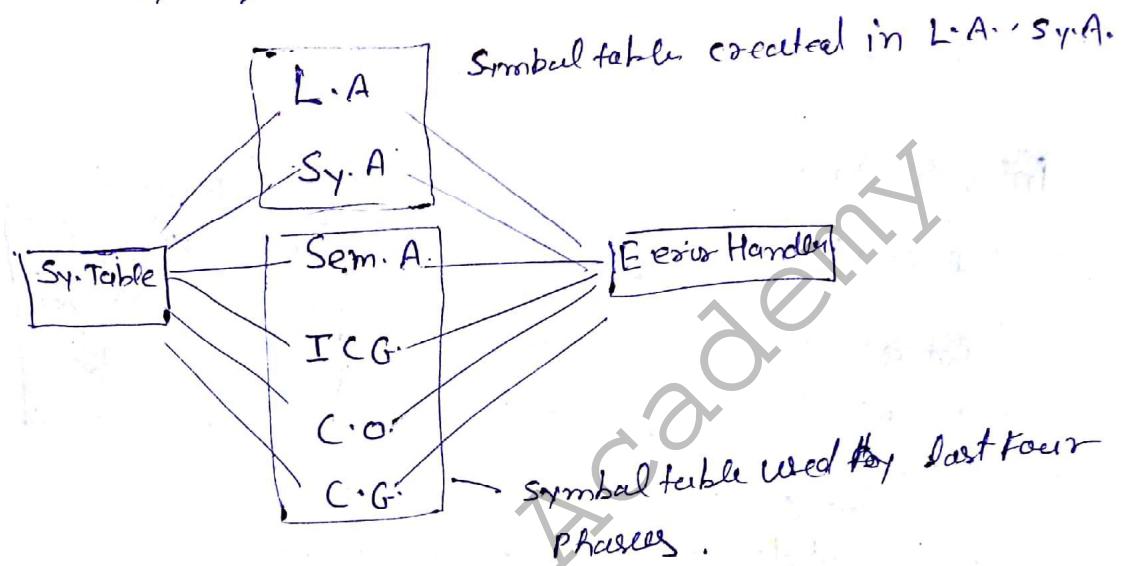
}

| | | |
|---|------------------------|----------|
| 1 | if ch = 1 | goto (3) |
| 2 | if ch = 2 | goto (6) |
| 3 | t ₁ = a + b | |
| 4 | c = t ₁ | |
| 5 | goto 9 | |
| 6 | t ₁ = a - b | |
| 7 | c = t ₁ | |
| 8 | goto 10 - 9 | |
| 9 | stop. | |

UNIT-IV

Symbol table

Symbol table is an data structure created and maintained by compiler in order to store information about. variable name, function name, object, class, compiler generated temp. var. etc.



Format of symbol table

Compiler uses following type of information in symbol table,

- 1) Data type, 2) Name, 3) scope, 4) Address
- 5) other Attributes

e.g. static int a;

| S.R.I. | Name. | Type | Attribute |
|--------|-------|------|-----------|
| 1 | a | int | static |

Symbol table Representation : example- Int calculate; sum, a, b;

1. Fixed length.

| Name. | Type. |
|-----------|-------|
| calculate | int |
| sum | int |
| a | |
| b | |

2. Variable length.

| Name. | Starting index | Length | Type. |
|-------|----------------|--------|-------|
| 0 | 10 | | |
| 10 | 4 | | |
| 12 | 2 | | |
| 16 | 2 | | |

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17
c a l c u l a t e \$ s u m \$ a \$ b \$

- Operations on Symbol table.
1. Insert(): ~~insert eq~~ int a \Rightarrow insert(a, int)
 2. Lookup(): Lookup(Symbol)
 3. Scope mgmt: \Rightarrow global, ornd local: scope:
 4. Delete(): delete(Symbol)

Example

```
int value=10;
```

```
void one()
```

```
{
```

```
int a
```

```
int b;
```

```
{
```

```
int c
```

```
int d;
```

```
}
```

```
int e;
```

```
{
```

```
int f;
```

```
int g;
```

```
}
```

```
void two();
```

```
{
```

```
int x;
```

```
int y;
```

```
{
```

```
int p;
```

```
int q;
```

```
{
```

```
int r;
```

```
{
```

Symbol table (Global)

| value | var | int |
|-------|-----------|------|
| one | Procedure | void |
| two | Procedure | void |

Symbol table (one)

| a | var | int |
|---|-----|-----|
| b | var | int |
| e | var | int |

Symbol table (two)

| x | var | int |
|---|-----|-----|
| y | var | int |
| r | var | int |

inner scope L.S.T.

| c | var | int |
|---|-----|-----|
| d | var | int |

outer scope L.S.T..

| f | var | int |
|---|-----|-----|
| g | var | int |

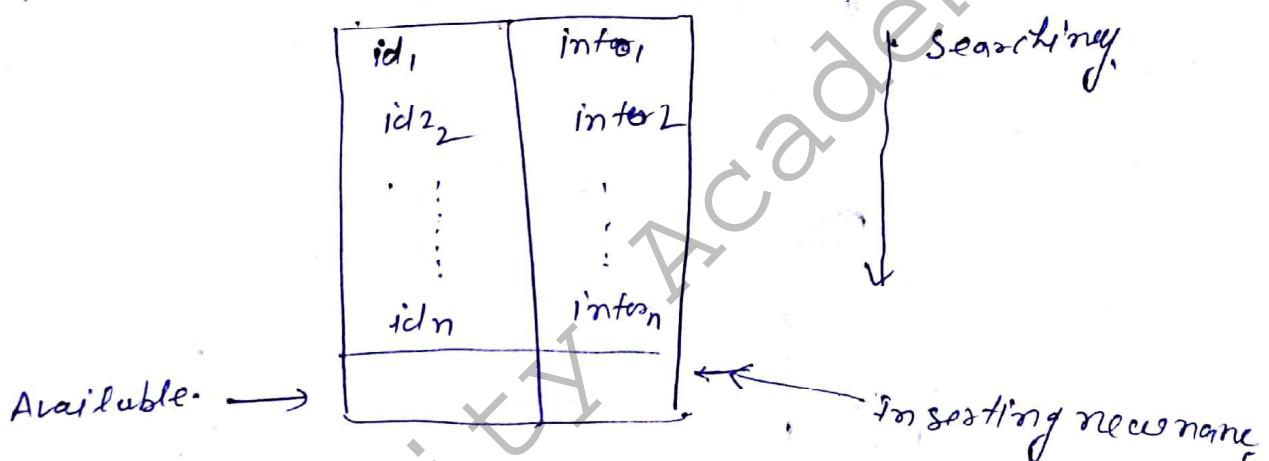
| p | var | int |
|---|-----|-----|
| q | var | int |

Implementation of Symbol table / Data structure for Symbol Table

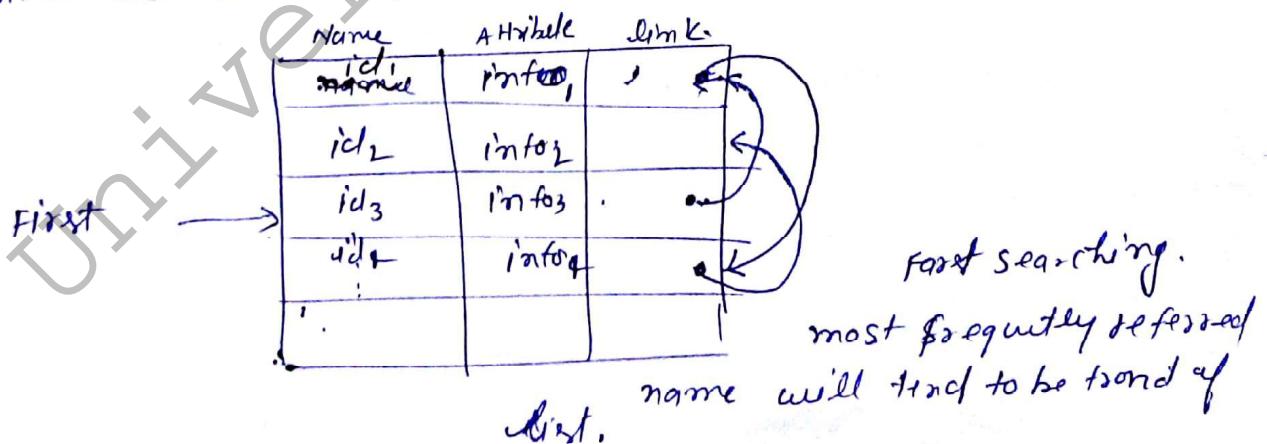
1. Linear list
2. Self organized list
3. Binary search Tree.
4. Hash tables.

1. Linear list:

Linear list is simplest form of mechanism to implement S.T.
An array is used to store name and associated information.

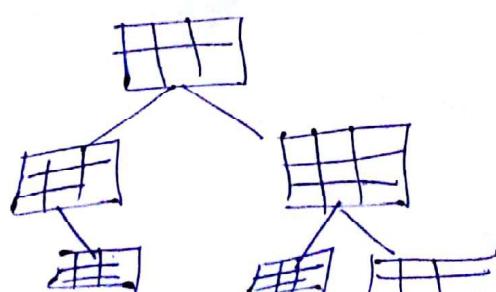


2. Self organizing list: This symbol table implementation is using linked list. A link field added to each record.



(3) Binary search Tree:

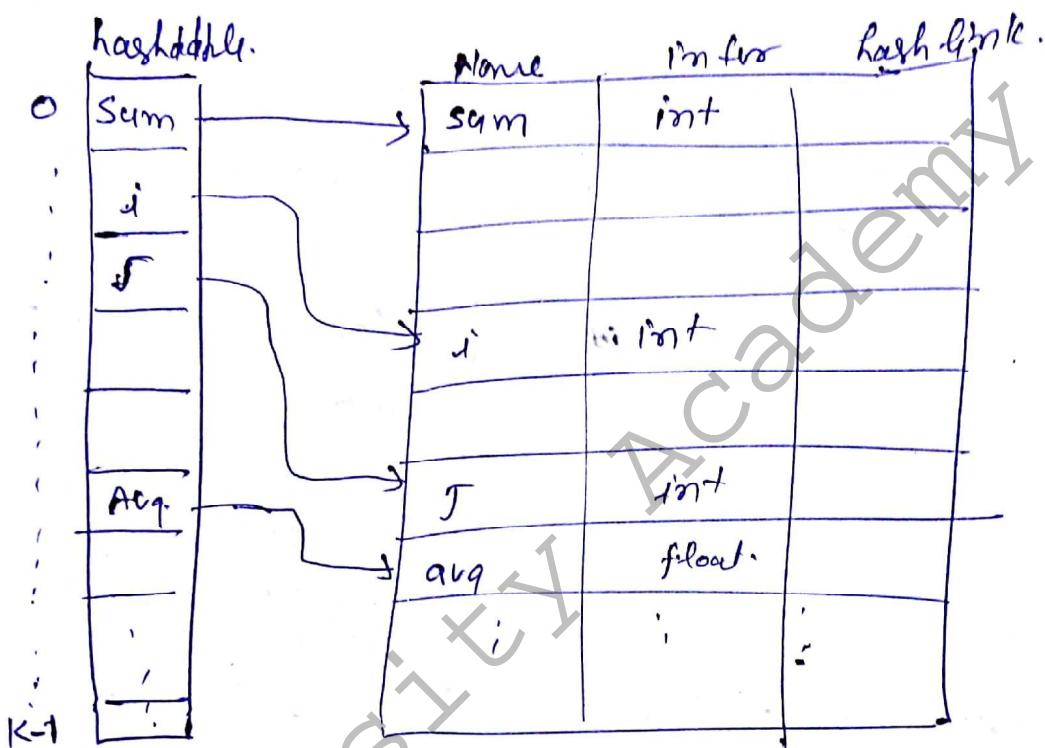
Binary
Tree



Hash table

Having is an important technique used to search the record of symbol table.. In hashing scheme two table are maintained.

1. Hash table
2. Symbol table

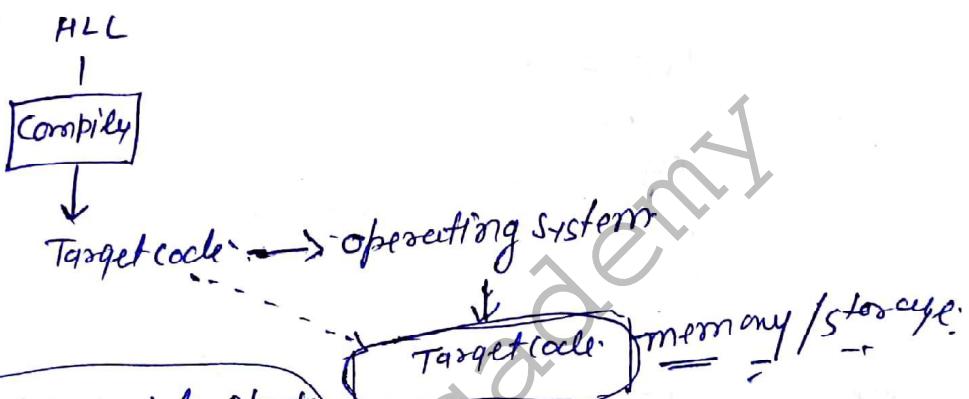


search by has function ' h ' as $h(\text{name})$ that will return integer bet $0 - k-n$.

~~Ques~~ searching of name is so fast int hashing,

Run Time Storage Administration

Compiler Demands for a block of memory to operating system. the Compiler Utilizes this block of memory for running the compiled program. this Block of memory called run time storage.

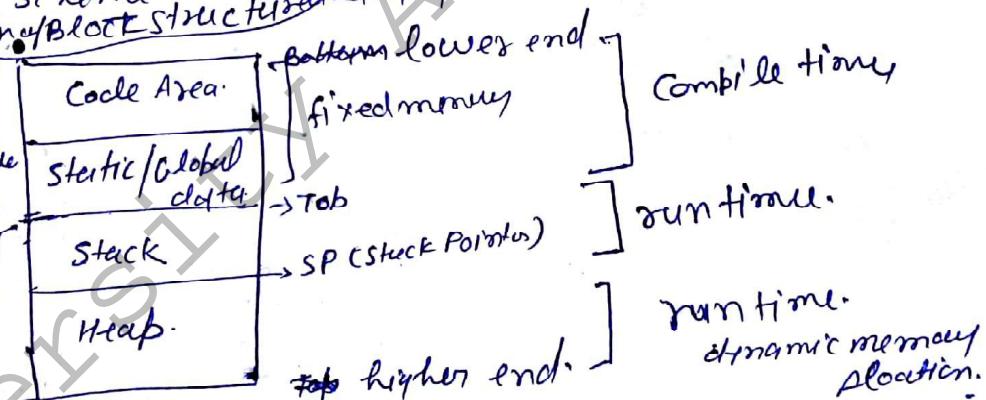


1- Implementation of ASmble Stack

Allocation Scheme.

2- Block - Implementation of BLOCK STRUCTURE Language.

Code Area has fixed size and hold target code.
Conform, Global Variable etc fixed memory.



e.g. main()

{

one();

}

one();

{
two();

}

two();

{
--:
}

two() AR

↓
one() AR.

|
main() AR.

Activation Record

Star Result →

Refer Value.

Actual parameters →

Actual Parameters.

optional, it points to AR of calling procedure.

control link (Dynamic)

optional, it refers to non local data in another AR.

Access link (Static Link)

local variable →

Local Variable

Temp variables →

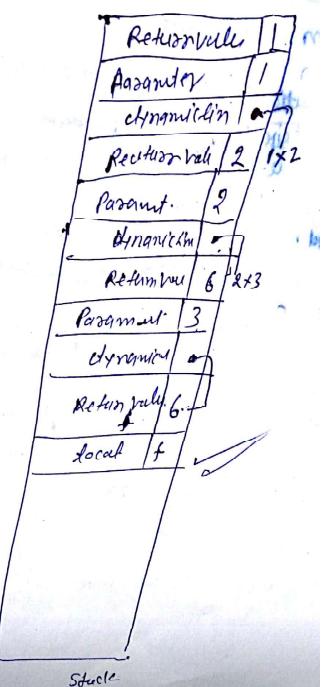
Temp Variables

$\text{Eq1} \quad \text{main}()$

```
{  
    int f;  
    f = fact(3);  
}  
int fact(int n)
```

```
{  
    if(n == 1)  
        return 1;  
    else  
        return (n * fact(n - 1));  
}
```

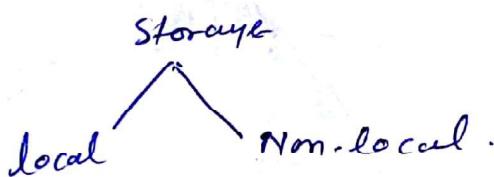
AR - main
AR - fact(3)
AR = fact(2)
AR = fact(1)



Implementation of Block-structured Languages: and Non Block.

Languages with block structure such as ALGOL and PL/I present certain complexities not found in C.

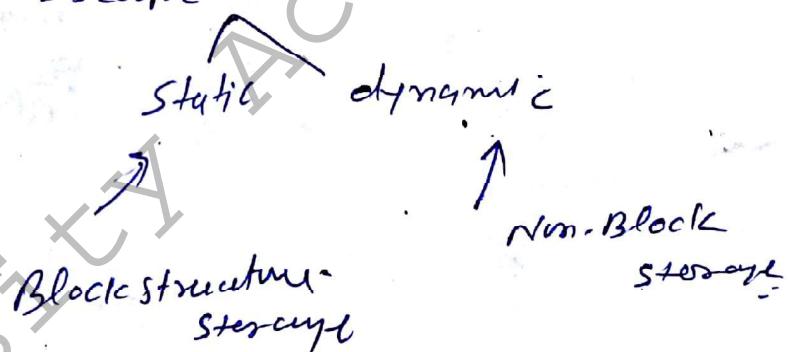
In block structured languages not only procedure but block as well may define their own data.



There are two types of storage allocation:

1- local homogeneous AR

2- Non-local heterogeneous scope information

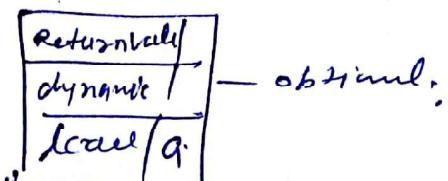


(1) Local data:

local data can be accessed with help of AR e.g.,

```

AC)
{
  int a;
}
  
```



(2) Non-local data:

A procedure may sometime refer to variables which are not local to it; such variable is called Non-local variable.

there are two types of scope static and dynamic.

Static scope rules

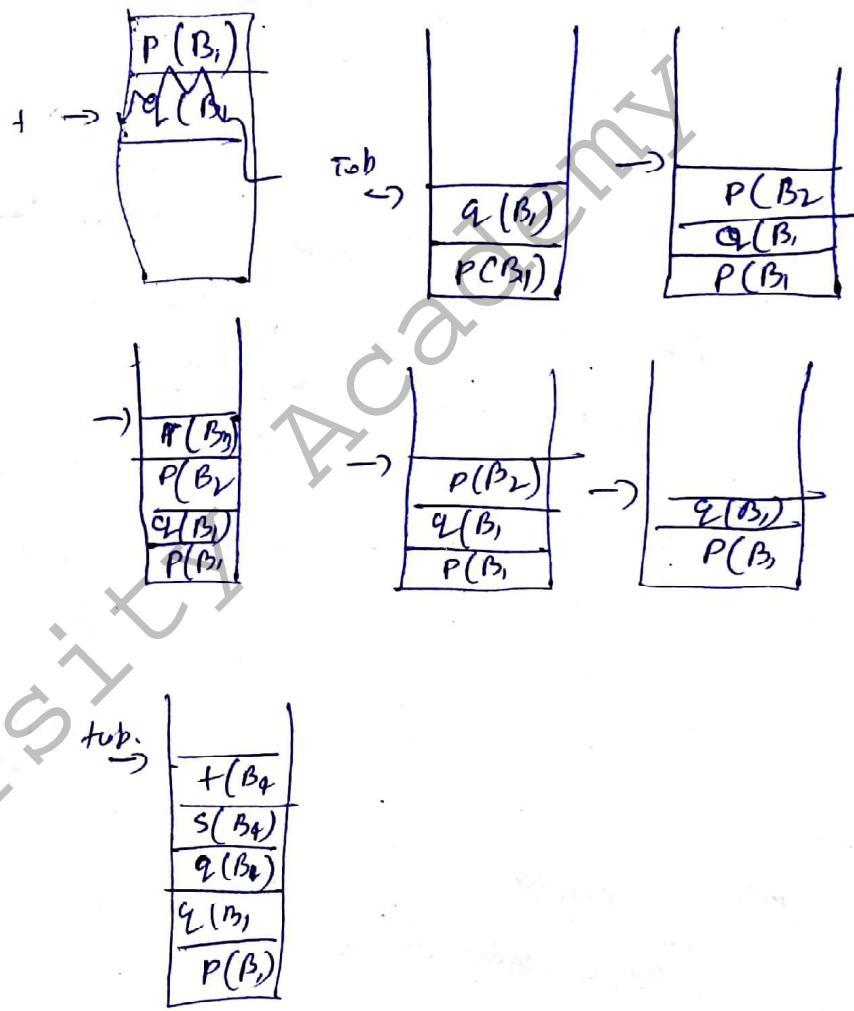
The static scope rule is also called lexical scope. the languages PASCAL, C, ADA are used static scope. the language are also called Block structured language.

e.g.

Scop-test()

```

  {
    int p, q;
    {
      int p;
      {
        B2:
        B3:
        {
          int r;
        }
      }
    }
  }
  
```

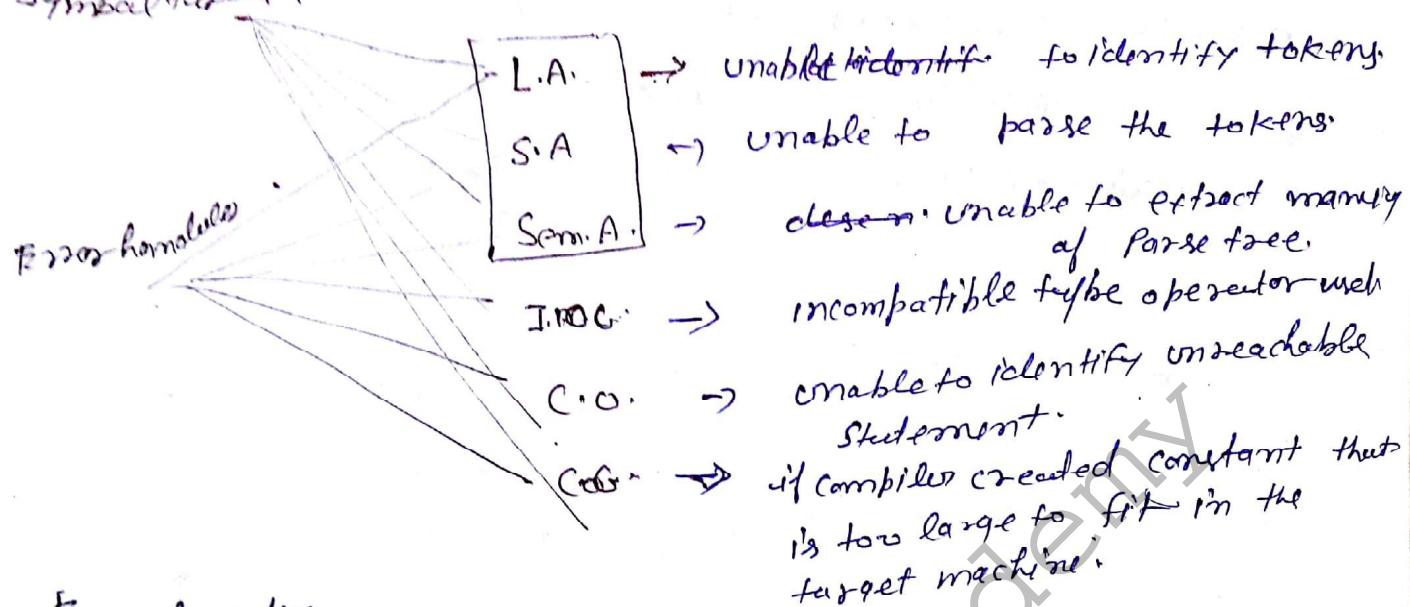


Dynamic scope Rule

for non block structured language this dynamic's scope allocation rule are used. the dynamic's scope rule determines the scope of declaration of the names at such time by considering the current activation. LISP and SPARSH are the language with the dynamic's scope rule.

Error detection and Recovery.

Symbol detector.

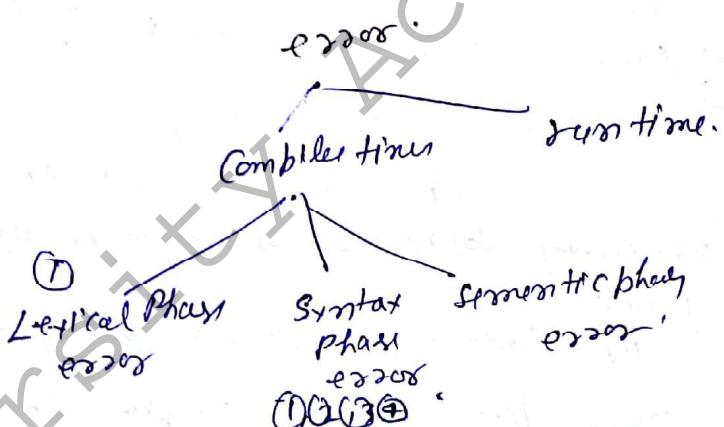


Error handler

detection
reporting

Recovery

- ① Panic mode.
- ② Phase level Recovery
- ③ Error production & Global production.



① Lexical phase Errors

1. Exceeding length of identifier or numerical constant C-31, FORTRAN -700.
2. Appearance of illegal character → \$, @.
3. unmatched string. or comments. ≠ ... " or ----- #/

example

void main()

```
{
    int a, @; variable declaration #/
    a = 10;
    printf("%d", a); #
```

void main()

{ int x=10, y=20;

char *a;

a = &x;

x = *ab;

 \ neither variable
 or number

Dia: skipped input without checking my additional or Action type; never go in infinite loop

Panic Recovery: in this method successive characters from the input are removed one at a time until designated set of synchronizing token is found.

Dia: skip input without checking my additional or Action type; never go in infinite loop

Syntactic Phase errors:

- (1) missing parenthesis. e.g. prints("Hello");
- (2) missing operator e.g. axbc
- (3) misspelled keyword. e.g. switch(c)
 {
 --.
 }
 }.
4. Colon in place of semicolon.
 a = 1;
 x = 2;
5. Extra blank. /* comment */

Recovery:

- (1) Panic mode:- same as lexical.
- (2) Statement mode Recovery:
 - (i) when a parser encounters an error, it performs necessary correction on remaining input and parse rest of input.
 - (ii) take care of not going in infinite loop.
- (3) Error production: if user has knowledge of common error that can be encountered then, these errors can be incorporated by augmenting the grammar with error production, difficult to maintain.
4. Global correction. - the parser examines the whole program and tries to find out closest match for it which is error free.
the closest match program has less number of insertion, deletion, and changes of tokens to recover from erroneous input.
due to high time and space complexity it is not implemented practically.

Semantic errors.

- 1- Incompatible type of operands.
- 2- Undeclared variables.
- 3- not matching actual argument with formal one.

e.g. int a[10], b;
 - - -
 a = b;

Recovery. if error "Undeclared Identifier" is encountered
then to recover from this ~~is~~ a symbol table entry
for corresponding identifiers is made.

if data type of two operand are incompatible then
automatic type conversion done by compiler.

UNIT - VCode Optimization & Code generation.Code optimization:

The code optimization is required to produce efficient target code. The improvement over intermediate code by program transformation called optimizations. The code optimization are classified into two categories machine-independent optimization and machine dependent optimization.

Machine-independent optimization need to improve the target code without taking into consideration any properties of the target machine. And machine dependent optimization need for register allocation and utilization of machine instruction sequence. This type of optimization used in code generation.

Three criteria that we have been applied in optimization must preserve the meaning of program.

1. Transformation must preserve the meaning of program.
2. Transformation must, on the average, speed up programs by a measurable amount.
3. A transformation must be worth the effort.

Consider the code for quicksort:

```
i = m - 1; j = n; v = a[n];
while(i) {
    do i = i + 1; while(a[i] < v);
    do j = j - 1; while(a[j] > v)
    if (i >= j) break;
    { x = a[i]; a[i] = a[j]; a[j] = x;
      a[i] = a[j]; a[j] = x;
    }
}
```

three address code for above code.

- 1 $i = m - 1$
- 2 $j = n$
- 3 $t_1 = 4 \times n$
- 4 $v = a[t_1]$
- 5 $i = i + 1$
- 6 $t_2 = 4 \times i$
- 7 $t_3 = a[t_2]$
- 8 if $t_3 < v$ goto (5)
- 9 $j = j - 1$
- 10 $t_4 : 4 \times j$
- 11 $t_5 = a[t_4]$
- 12 if $t_5 > v$ goto (9)
- 13 if $i >= j$ goto (23)
- 14 $t_6 = 4 \times i$
- 15 $x = a[t_6]$

- 16 $d_7 = 4 \times j$
- 17 $t_8 = 4 \times j$
- 18 $t_9 = a[t_8]$
- 19 $a[t_7] = t_9$
- 20 $t_{10} = 4 \times j$
- 21 $a[t_{10}] = x$
- 22 goto (5)
- 23 $d_{11} = 4 \times j$
- 24 $x = a[d_{11}]$
- 25 $t_{12} = 4 \times i$
- 26 $t_{13} = 4 \times j$
- 27 $t_{14} = a[t_{13}]$
- 28 $a[t_{12}] = t_{14}$
- 29 $t_{15} = 4 \times n$
- 30 $a[t_{15}] = x$

Basic Block

our first step to break code into basic block
that is sequence of consecutive statement which may
be entered only at the beginning and executed the
sequence of statement without halt.

to identify determine the basic block we
have to determine leaders. by following rule.

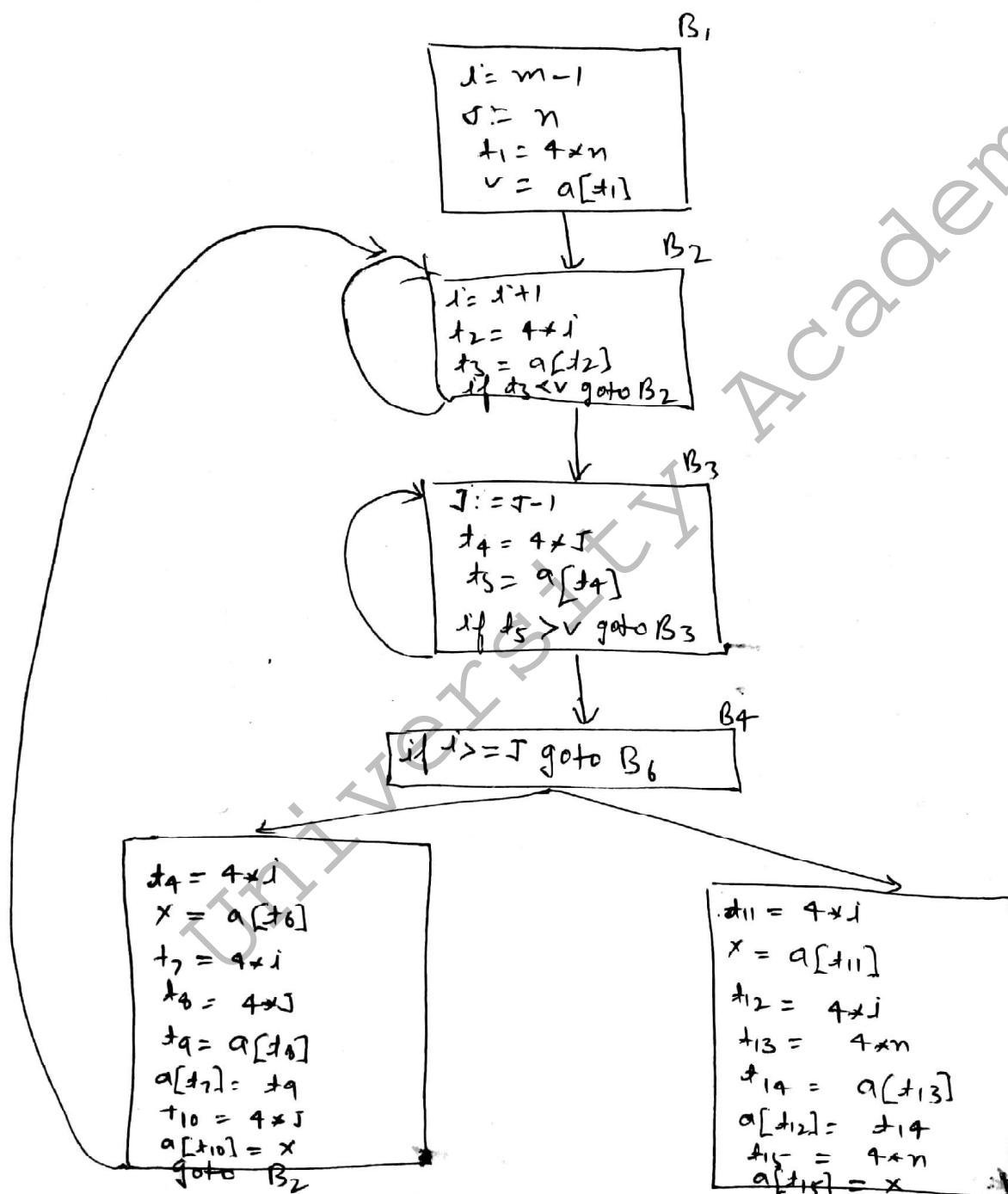
1. the first statement is leaders.
2. conditional or unconditional goto is a leaders
3. A statement immediately follow the conditional
goto is a leaders.

thus we have obtain 6 basic block for above code.

Flow graph: It is useful

to show the relationship between basic block by a directed graph called a flow graph.
the nodes of the flow graph are the basic blocks.

Basic block and flow graph for above code:



flow graph

The Principle source of optimization.

The optimization of program is called local if it can be performed by looking only at the statement in a basic block; otherwise it is called global.

1- Function-Preserving transformation:

there are number of ways in which a compiler can improve a program without changing the function it compute. the function preserving transformation include following optimization:

- (a) common subexpression elimination
- (b) copy propagation.
- (c) dead code elimination
- (d) constant folding.

(a) Common Subexpression elimination:

we observe local common subexpressions in B_5 -

B_5

$$\begin{aligned} t_6 &= 4 \times j \\ x &= a[t_6] \\ t_7 &= 4 \times i \\ t_8 &= 4 \times j \\ t_9 &= a[t_8] \\ a[t_7] &= t_9 \\ t_{10} &= 4 \times i \\ a[t_{10}] &= x \end{aligned}$$

goto B_2

B_5

$$\begin{aligned} t_6 &= 4 \times i \\ x &= a[t_6] \\ t_8 &= 4 \times j \\ t_9 &= a[t_8] \\ a[t_6] &= t_9 \\ a[t_8] &= x \\ \text{goto } B_2 \end{aligned}$$

Now we observe a global common subexpression in B_5 and B_2, B_3 $t_6 = t_2$ and $t_8 = t_4$

B_5

$$\begin{aligned} x &= t_3 \\ a[t_2] &= t_5 \\ a[t_4] &= x \\ \text{goto } B_2 \end{aligned}$$

B_6

$$\begin{aligned} x &= t_3 \\ t_4 &= a[t] \\ a[t_2] &= t_4 \\ a[t_4] &= x \end{aligned}$$

after localizing global common subexpression

(b) Copy propagation

copy propagation in the form of $f := g$:

we observe in B_5

B_5

$x = t_3$
 $a[t_2] = t_5$
 ~~$a[t_4] = t_5$~~
goto B_2

after copy propagation

B_5

$x = t_3$
 ~~$a[t_2] = t_5$~~
 $a[t_4] = t_3$
goto B_2

(c) Dead-Code Elimination

A variable is live at a point in a program if its value can be used subsequently otherwise it is dead at that point.

Now again consider B_5

$x = t_3$
 $a[t_2] = t_5$
 ~~$a[t_4] = t_3$~~
goto B_2

after dead code elimination

$a[t_2] = t_5$
 $a[t_4] = t_3$
goto B_2

that value of x is removed.

(d) constant folding

At compile time that the value of an expression is a constant and using the constant instead it's known as constant folding.

Loop Optimization.

Loop optimization is very important optimization.
there are three techniques are important for Loop optimization.

a- code motion - which move code outside.

b- Induction variable Elimination

c- reduction in strength.

(a) Code motion:

code motion decreases the amount of code in a loop
is code motion. this transformation takes an expression that yield the
same result independent of the number of time a Loop is executed.

examt: $\text{while } (i < \text{limit} - 2)$

can be written as

$t = \text{limit} - 2$

$\text{while } (i < t)$

(b) Induction Variable and Reduction in strength.

A variable n is called an induction variable of a Loop if the value of variable gets changed every time. It is either decrementation or increment by some constant.

Consider the loop in B_3 , note that value of J and t_4 every time the value of J decreases by 1 and t_4 decreases by 4 is assigned to t_4 such identifiers are called induction variable.

Strength Reduction:

(4)

The strength of certain operators is higher than others such as strength of $*$ is higher than $+$. In the strength reduction technique the higher strength operator can be replaced by lower strength operators.

e.g. for ($i=1; i<50; i++$)

{

count = $i * 7$;

}

it can be replaced by following code.

temp = ?
for ($i=1; i<50; i++$)

{

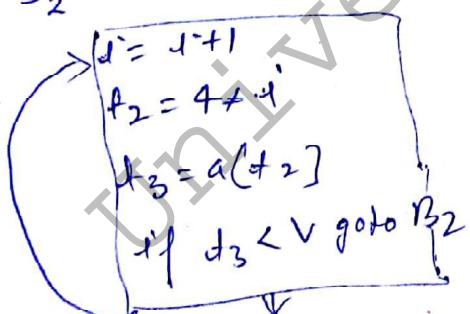
count = temp;

temp = temp + 7

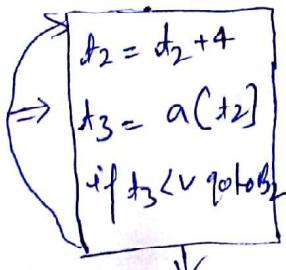
}

Now we can apply the reduction in strength in B_2 and B_3 .

B_2

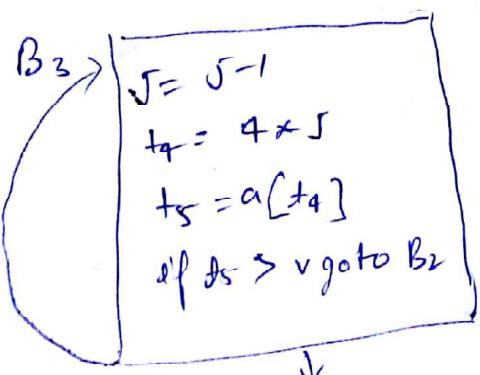


$t_2 = 4 * i$

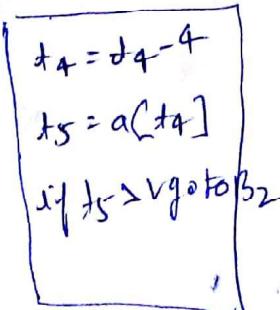


$t_4 = 4 * j$

B_3



\Rightarrow



DAG Representation of Basic Block.

Directed acyclic graph are useful data structure for implementing transformation of Basic Block. Constructing a DAG from three Address code is a good way to determining common subexpression within block, determining which name are used inside and the block and but evaluated outside the block; and determining which statement of the block could have there computed value used outside the block.

A DAG consists Directed Acyclic graph with the following labels on nodes.

1. leaves are labeled by unique identifiers ; either variable or constant.
2. Interior node are labeled by an operators.
3. nodes are also given a sequence of identifiers for labels.

e.g. $\frac{t_1}{t_2} = t_3$

$$1 \quad t_1 = 4 + i$$

$$2 \quad t_2 = \text{a}[t_1]$$

$$3 \quad t_3 = 4 * i$$

$$4 \quad t_4 = b[t_3]$$

$$5 \quad t_5 = t_2 * t_4$$

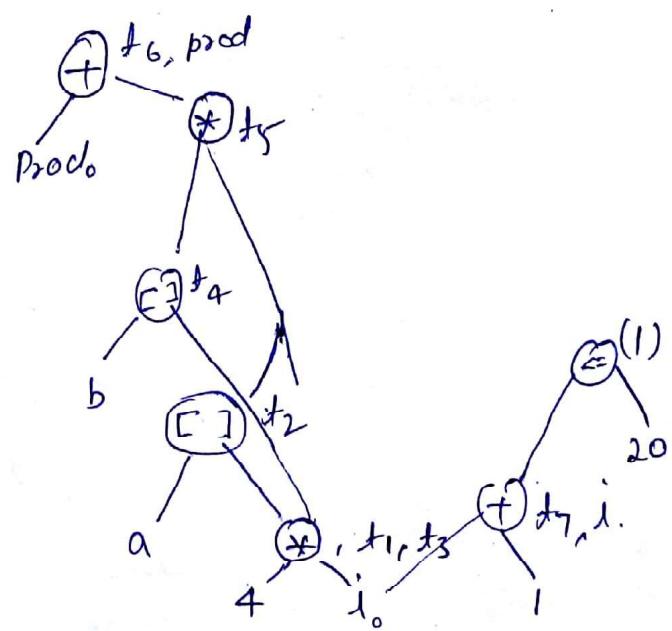
$$6 \quad t_6 = \text{prod} + t_5$$

$$7 \quad \text{prod} = t_6$$

$$8 \quad t_7 = 1 + i$$

$$9 \quad i = t_7$$

$$10 \quad \text{if } t_2 = 20 \text{ goto } l$$



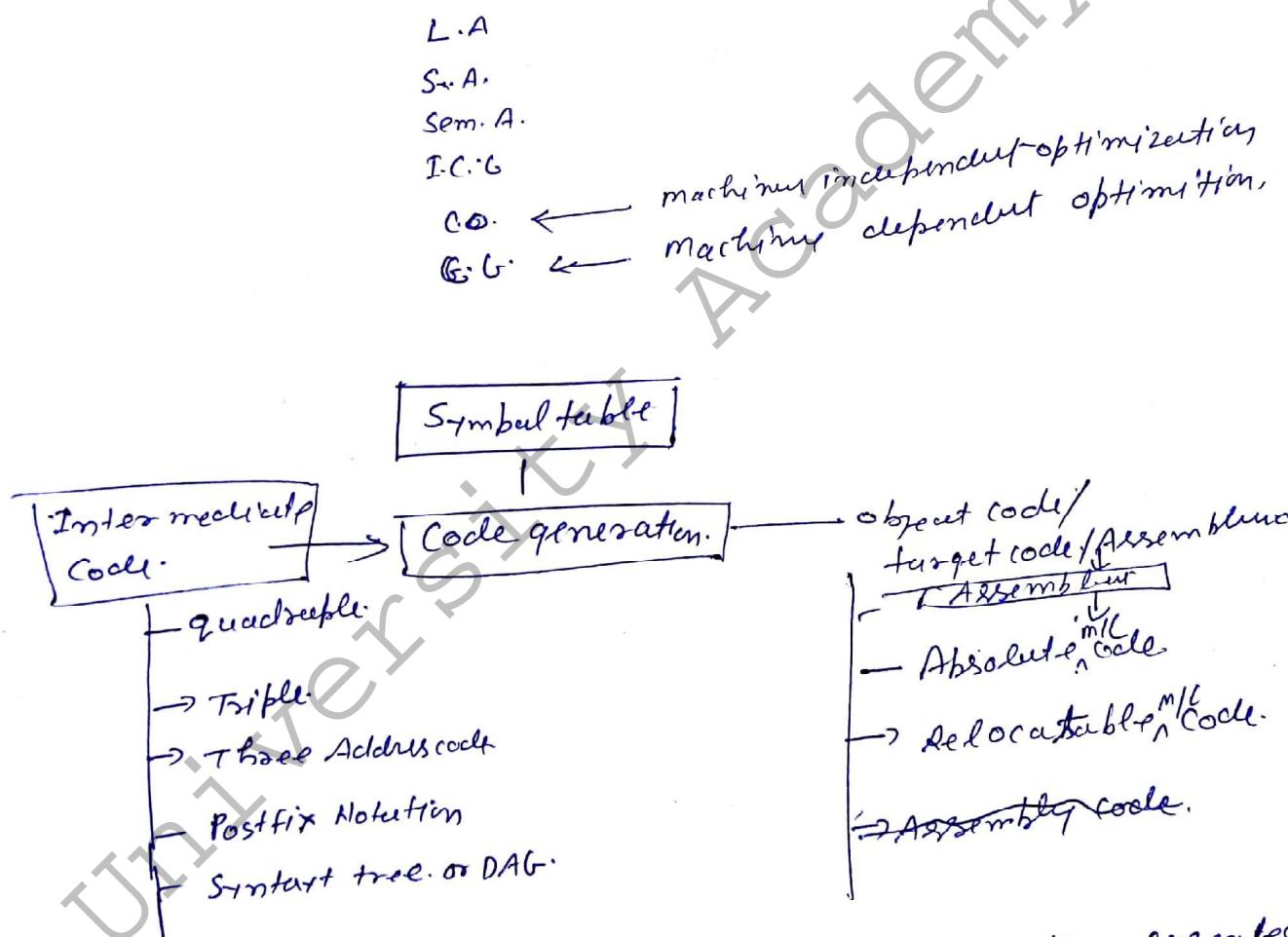
Application of DAG

- 1 - Determining common subexpression and apply DAG based local optimization.
- 2 . Determining which statement of block could have their computed values outside the block.
- 3 determining which name are used inside the block and computed outside the block .

Code Generation.

Code generation is the final activity of compiler. Basically code generation is a process of creating Assembly / machine language. There are some properties of code generation phase.

- (i) Correctness.
- (ii) High quality
- (iii) Efficient use of resources of target machine.
- (iv) Quick code generation.



Absolute code: fixed location in memory and directly executable.
"usefull for small program" exec.

Relocatable code: Not a fixed location in memory, code can be placed wherever wherever the linker finds space. In RAM. Careful for commercial compiler. ^{ex: next file formats}

Assembly code: easy to generate. uses symbolic code, especially for machines with small memory.

Design Issue

1. Input to the Code generator:
 2. Target program:
 3. memory Management:
 4. Instruction Selection:
 5. Register Allocation.
 6. -

$$x = y + z$$

```
MOV Y, R0  
ADD Z, R0  
MOV R0, X.
```

$$a = b + c$$

$$d = \alpha + e$$

mov b, R0

ADD C, RD
MOV RO, DL decrement

~~MOV~~ ~~BB~~ ~~R0~~

~~mov RD, d~~

| | |
|-----|--------|
| MOV | b, RO |
| ADD | c, RO |
| ADD | e, RO |
| MOV | RO, d. |

5. Consider Three-Address-Code -
 $t = a + b$

$$t = a + b$$

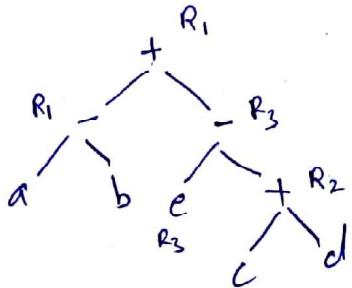
$$t = t + c$$

$$t = it/d$$

MOV a, RO
ADD b, RO
~~ADD~~
MUL 'c, RO
DIV d, RO
MOV RO, t

Quinton

$$(a-b) + (e-(c+d))$$



$R_1 \leftarrow a$
 $R_2 \leftarrow b$
 $R_1 \leftarrow R_1 - R_2$
 ~~$R_2 \leftarrow c$~~
 $R_3 \leftarrow d$
 $R_2 \leftarrow R_2 + R_3$
 $R_3 \leftarrow e$
 $R_3 \leftarrow R_3 - R_2$
 $R_1 \leftarrow R_1 + R_3$

Peephole optimization:

A statement by statement code generation strategy generates target code that contains redundant instruction. The quality of such code is poor. To optimize such target code certain transformations need to be applied on the target code.

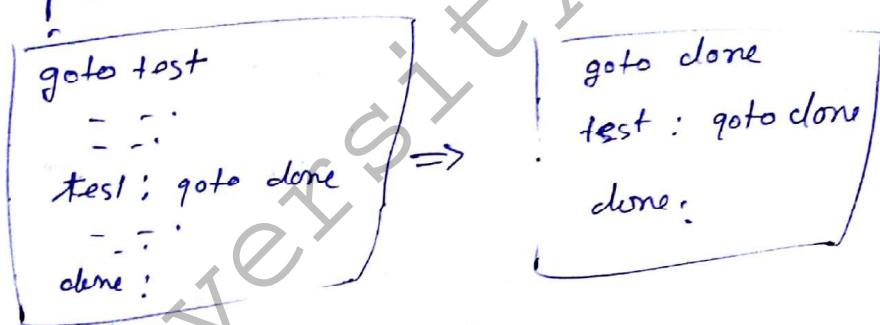
The peephole optimization is simple and effective technique for locally improving target code. The peephole optimization can be applied on the target code using following characteristic.

1. Redundant instruction:

(i) $\text{mov } R0 \ a$
(ii) $\text{mov } a \ R0$

then eliminate induction(2).

2. Flow of control optimization.



3. Algebraic Simplification.

$$\begin{aligned} x &= x + 0 \\ &\text{or} \\ x &= x * 1 \end{aligned}$$

} eliminated such type of instruction.

4. Algebraic Simplification. Reduction strength!

Addition is cheaper than multiplication.
Subtraction is cheaper than division.

5. Machine idioms: Auto increment and Auto decrement.