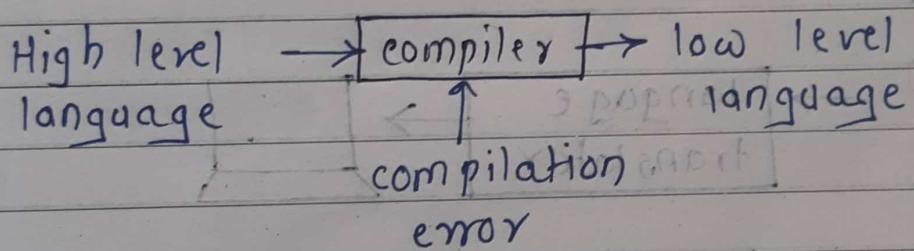


Unit 1:

Compilers : Introduction to compiler phases, introduction to cross compiler, features of mc dependent & independent compilers, types of compilers, Interpreters: compilers vs interpreter, phases, working, Preprocessor: header file and macro expansion.

Compiler :

A compiler is a sw that converts a program written in a high level language (source lang.) to a low level language (object/target mc lang. i.e., assembly or machine code).

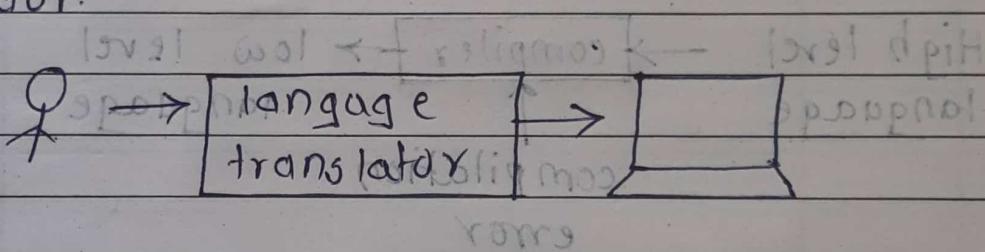


It performs following operations:

1. Performs a pre-processing of source code. Gather all files required for the source code to compile.
2. Parses the entire source code. checks for any syntax errors in the source code.
3. Performs a thorough syntax analysis of the source code. To understand the structure & semantic of the source code.
4. Optionally translates the source code in an intermediate code known as object code to enhance the performance.
5. Translates the object code to binary language known as executable code.

Introduction

- computers are electronic devices so it understand the language of 0's & 1's only.
- In early days of computers, instruction to a computer is given with the help of paper tape & punched cards, so in these the presentation of hole represents a digit if it is 1.
- So writing of 0's & 1's is difficult.
- So with punched cards we are facing a communication barrier to make the communication with the computers, so will require a language translator.



Language Translators

i. Assembler

MOV R1, 02H → Assembler → 0110111011
MOV R2, 03H → 0011011010

Assembly
lang.

mlc
code

So to write the code we require the complete knowledge of computer architecture so its difficult to remember this so the Interpreter & compiler came into picture.

Interpreter eg. -

php, ruby, javascript

compiler eg. - c, c++, erlang.

- compiler is more faster & efficient than interpreter

High level language

```
#include <stdio.h>
```

```
int main()
```

```
{
```

```
int x, a=2, b=3, c=5;
```

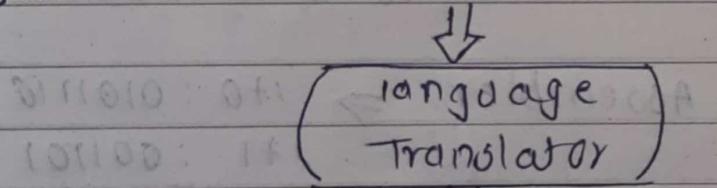
```
x = a+b*c
```

```
printf ("The value of x is %d", x);
```

```
return 0;
```

```
}
```

source code/HL2 code



00110110111
01101101011
01011011100
mlc code .

Language Translator - Internal architecture

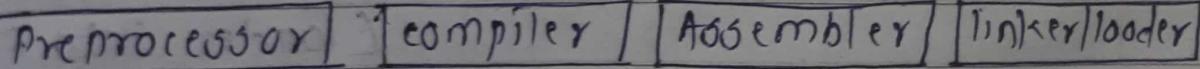


Fig: clang. translator Arch.

Preprocessor -

Reserve all preprocessor dir & remove all comments

#include <stdio.h>

II code → Preprocessor → II code

HLL code

stdio.h

pure HLL

Compiler -

Generate equivalent assembly lang. code

stdio.h

mov \$0, 02h

II code → compiler → mov \$0, 04h

pure HLL

assembly code

Assembler -

Generate relocatable m1c code

assembly → Assembler → iTO : 0101110
code

i+1 : 001101
relocatable
m1c code .

Linker/loader -

iTO : 0100110

i+1 : 0110110 → linker →

relocatable
m1c code

loader → absolute m1c
code

[loader into RAM &
ready for execution]

Compiler

It has 6 different phases these are,

characters \rightarrow HLL

tokens

Lexical Analysis

Syntax Analysis

Semantic Analysis

Intermediate code Gen.

Code Optimization

Target code generation

Analysis phase

Front end

Error Handler

symbol
table
manager

Synthesis Phase

back end

Fig : Compiler Design

1) Lexical Analyzer :

$$x = a + b * c$$



Lexical Analysis



Lexemes	Tokens
x	identifier
=	operator
a	identifier
+	operator
b	identifier
*	operator
c	identifier

- Lexical Analyzer takes lexemes as input & generate tokens
- Lexemes - similar to words, where words have their individual meaning & lexemes contain group of words.
For eg. - analysis - detail analysis of anything
where $x = a+b$ if we consider only x which does not have a meaning but it is meaningful when we say $x = a+b$.
- tokens - meaning of lexemes.

- i. The job of lexical analyzer is to find out the meaning of every lexemes.
- It recognizes tokens using regexs

Eg. - Regex for identifier

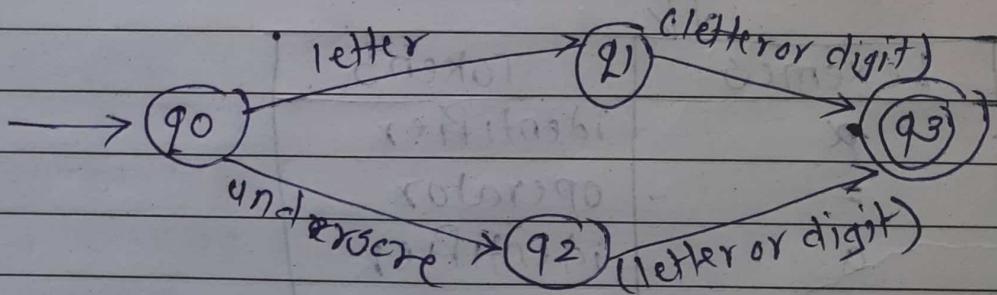
Regular expression $\rightarrow l(l+d)^* l - (l+d)^*$ where,

l - letter

d - digit

$_$ - underscore

Illustrate this using Finite automata.



- So for examining the lexemes the lexical analyzer uses the regular grammars

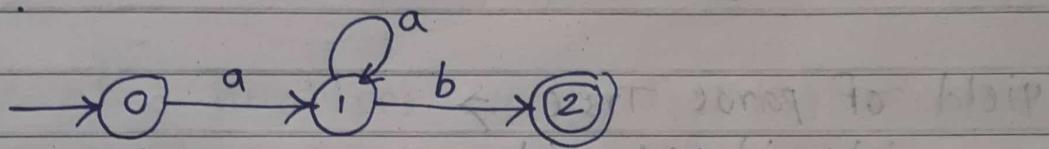
Regular expression :

it is composed of simple characters such as /abc/ or combination of simple & special characters such as /ab*c/

Finite automata :

It is a state machine that takes a string of symbols as input & changes its state accordingly. Finite automata is a recognizer for regular expression. When a regular expression string is fed into finite automata, it changes its state for each literal.

Eg :



Initial state

Final state

state

states	a	b
0	1	2
2		1

Syntax Analyzer (Parser) :

Stream of tokens passes to syntax analysis phase. Syntax analysis depends on type 2 grammar or context free grammar

For e.g. - $x = a + b * c$

where,

$E \Rightarrow \text{exp.}$

$T \Rightarrow \text{Term}$

$F \Rightarrow \text{Factor}$

$S \rightarrow id = E ;$

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow id$

These are CFG

rule to form a

parse tree

Parse Tree

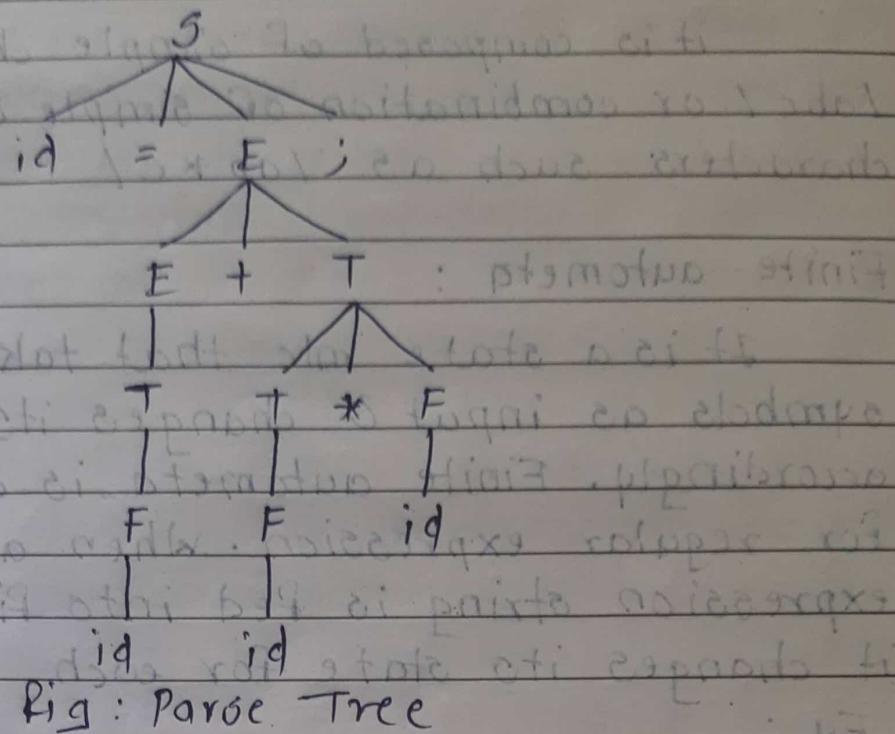


Fig: Parse Tree

yield of parse tree \Rightarrow

$$id = id + id * id;$$

if yield of parse & expression is same
then there is no error.

Context Free Grammer (CFG):

is a set of recursive rules used to generate pattern of strings. A CFG can describe all regular languages.

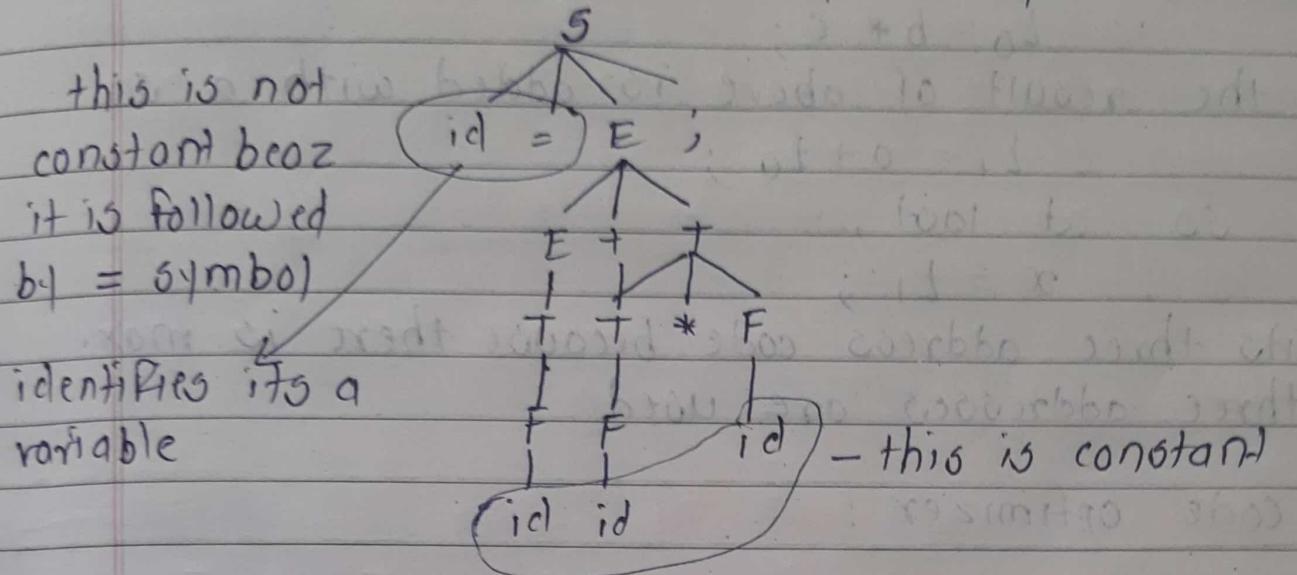
3) Semantic Analyzer:

parse Tree \Rightarrow semantic \Rightarrow semantically
Analyzer verified

parse

Tree

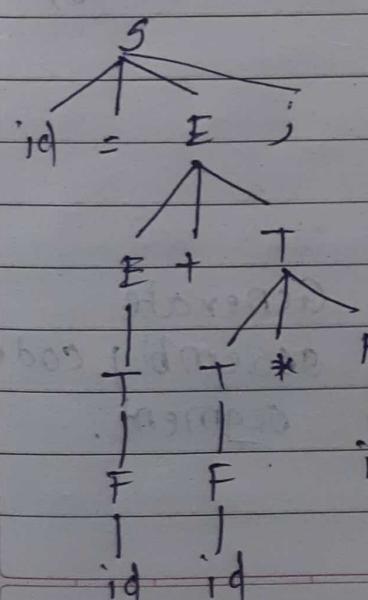
- Semantic Analyzer is responsible for type checking, error bound checking & correctness of scope resolution
- It does the logical analysis of parse tree



- Semantic analyzer looks for the meaningfulness of parse tree

4) Intermediate code Generator

semantically verified parse tree \Rightarrow Intermediate code generation \Rightarrow Intermediate code



Three Address Code (TAC)

$$\begin{aligned} \text{id} &= \text{id} + \text{id} * \text{id}; & \text{shot. } b_0 &= b * c; \\ \Rightarrow & & \Rightarrow 2. t_1 &= a + z_0 \\ & & & \Rightarrow 3. x = t_1; \end{aligned}$$

$$x = a + b * c;$$

- TAC - popular intermediate code
- if we took the parse tree from bottom to up, so at lowest level * operator is there, so * performs Pirot

$$\therefore t_0 = b * c;$$

the result of above is added with a so,

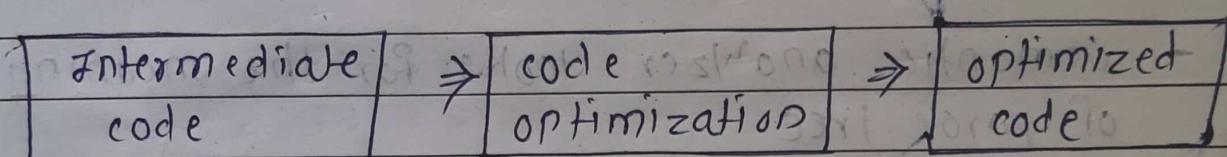
$$t_1 = a + t_0;$$

so at last,

$$x = t_1;$$

its three address code because there is max. three addresses are used

5) Code Optimizer :



- code optimization can be mlc dependent or mlc independent.

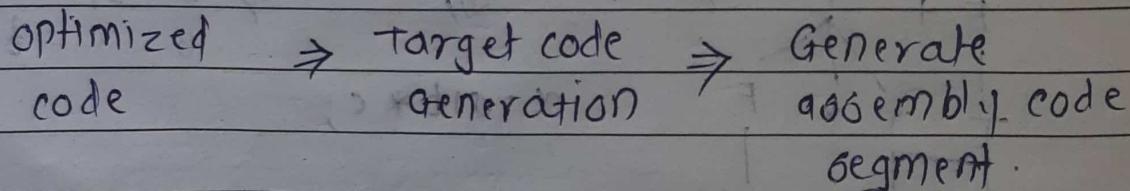
TAC

1. $t_0 = b + c;$
2. $t_1 = a + t_0;$
3. $x = t_1;$

$$1. t_0 = b + c;$$

$$2. x = a + t_0;$$

6) Target code Generation



$$t_0 = b + c$$

$$x = a + t_0$$

extended
version of eax

`mov eax, DWORD PTR [rbp - 8]` → pointing to b.
 |
 accumulator variable
 (32 bits stored)
 signed multiplication

`imul eax, DWORD PTR [rbp - 12]` — value pointing to variable c = eax.

`mov edx, eax`

`mov eax, DWORD PTR [rbp - 4]` ⇒ a in accumulator

`add eax, edx` eax & edx.

`mov DWORD PTR [rbp - 16], eax` add of x.

7) Symbol Table:

- Symbol table is basically a data structure which stores the occurrences of entities
 - variable & function names
 - objects
 - classes
 - Interfaces

Symbol Table Usage by phases

Lexical analysis -

creates entries for identifiers.

Syntax analysis -

Adds information regarding attributes.

Semantic analysis -

using the available info. stored in symbol table checks the semantics of identifier created by lexical analysis phase & update the symbol table.

Intermediate code generation -

available info. in symbol table helps in adding temporary variable info.

Code optimization -

available info. in symbol table is used in mc dependent optimization

Target code generation -

Generates the target code using address info. of identifiers.

Symbol table: entries

Name	Type	size	dimension	line of declaration	line of usage	Address
count	int	2	0	-	-	-
x	char	3	1	-	-	-

int count;
char x[] = "xyz";

- for primitive data type the dimension is 0 & for 1D, 2D its 1 or 2 respectively.

Symbol table: operations

- operations perform on symbol table depends on the language
- Early days of langs we used
 - No block structured language
 - contains single instance of the variable declarations
- so on symbol table we will perform

i) insert() ii) lookup() - access variable.
 Eg. For non structured block lang. is
 FORTRAN lang.

- Block structured lang:
 - variable declaration may happen multiple times
 - operations
 - i) insert() ii) lookup()
 - iii) set() iv) Reset()

* Features of machine dependent & independent compilers

- A compiler operates in phases. A phase is a logically interrelated operation that takes source pgm in one representation & produces output in another representation.
- There are two phases of compilation
 - 1) Analysis (mlc independent lang. dependent)
 - 2) Synthesis (mlc dependent lang. independent)

Machine dependent

- The allocation of the registers as well as their machine instructions rearrangement for the execution of the efficiency improvement process.
- It is the form which is considered for the code optimization technique of the pgm which has to be compiled.

Machine Independent

- attempts to improve the intermediate code to get a better target code.
- The process of intermediate code generation introduces much inefficiency like: using variable instead of constants, extra copies of variable, repeated evaluation of expression. Through code optimization you can remove such inefficiencies & improves code.

code optimization can perform in the following different ways :

1) compile time evaluation -

$$x = 5.7$$

$$y = x/3.6$$

evaluate $x/3.6$ as $5.7/3.6$ at compile time

2) variable propagation -

Before optimization

$$c = a * b$$

$$x = a$$

fill

$$d = x * b + y$$

After optimization

$$c = a * b$$

$$x = a$$

fill

$$d = a * b + y$$

Here, after variable propagation, $a+b$ & $x+b$ identified as common sub expression

3) Dead code elimination -

Before elimination

$$c = a+b$$

$$x = b$$

HII

$$d = a+b+4$$

Here, $x=b$ is a dead state bcoz it will never subsequently used in the pgm. so, we can eliminate this state.

after elimination

$$c = a+b$$

HII

$$d = a+b+4$$

4) code motion -

- It reduces the evaluation freq. of expr
- It brings loop invariant stmts out of the loop.

```

do item=10;
{ value2 = value + item
item=10
value2 = value+item
} while (value<100); } while (value<100);
    
```

5) Induction variable & strength Reduction -

- strength reduction is used to replace the high strength operator by the low strength.

- An induction variable is used in loop for the following kind of assignment like $i = i + \text{constant}$.

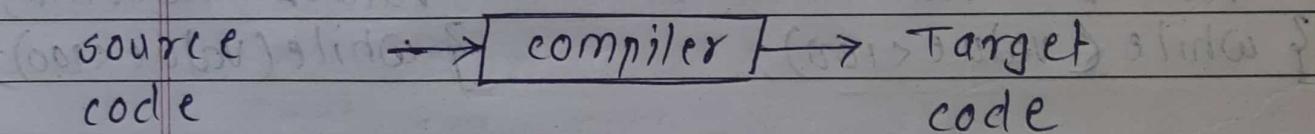
1 /

<u>Before reduction</u>	<u>after reduction</u>
$i = 1;$ $\text{while}(i < 10)$ $\{$ $\quad \text{for } t = 0; t < 4; t++$ $\quad \{$ $\quad \quad d = t * 4;$ $\quad \}$ $\quad \}$ $\quad \text{for } d = 0; d < 6; d++$ $\quad \{$ $\quad \quad \text{if } d == 0$	$i = 1;$ $t = 4$ $\{$ $\quad \text{while}(t < 40)$ $\quad \{$ $\quad \quad 4 = t;$ $\quad \quad t = t + 4;$ $\quad \}$ $\quad \}$ $\quad \text{for } d = 0; d < 6; d++$ $\quad \{$ $\quad \quad \text{if } d == 0$

* Types of compiler

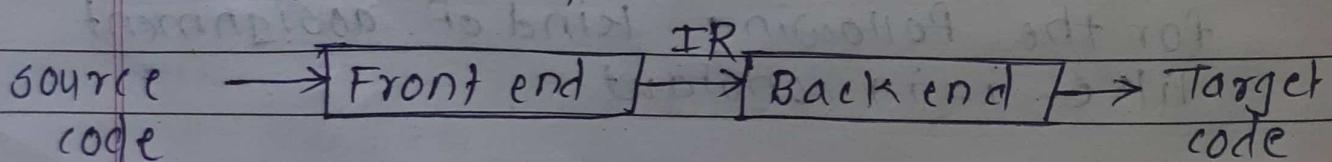
- 1) Single Pass compiler
- 2) Two Pass compiler
- 3) Multipass compiler
- 4) Single Pass compiler \Rightarrow

When we merge all the phases of compiler design in a single module, then it is called a single pass compiler. In a case of single pass compiler, the source code converts into machine code.



- 2) Two pass compiler \Rightarrow

A processor that runs through the program to be translated twice is considered two-pass compiler



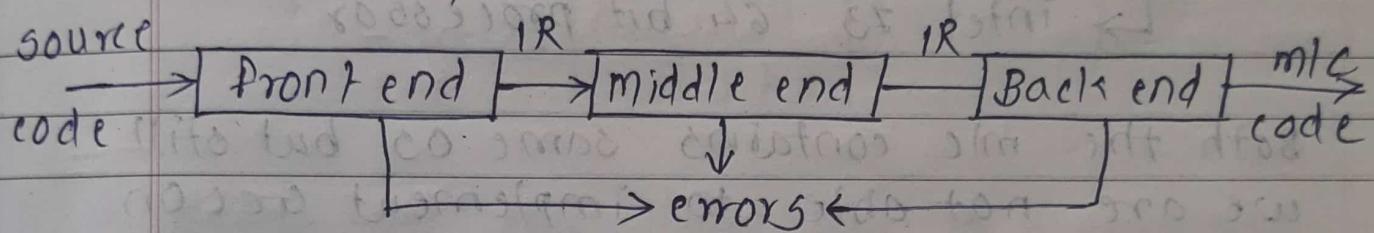
Incremental - executes the recompilation of
only a changed source instead of
compiling the complete code

3) Multipass compiler

A pgm source code or syntax tree is processed many times by the multipass compiler. It breaks down a huge pgm into numerous smaller pgms & runs them all at the same time. It creates n. of intermediate codes.

All of these multipasses use the previous phases output as an input. As a result, it necessitates less memory.

'wide computer' is another name for it.

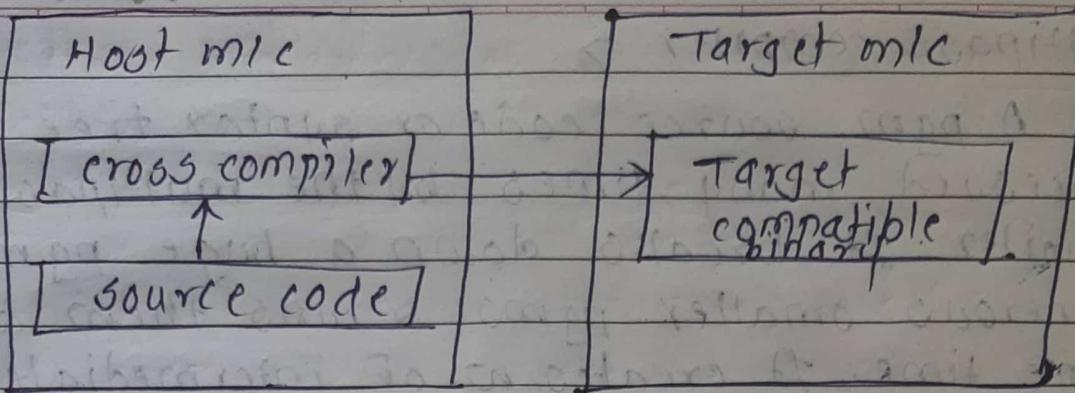


* Cross compiler

- A cross compiler is a type of compiler, that generates mc code targeted to run on a system different than one generating it. For eg. - A compiler that runs on windows platform also generates a code that runs on linux platform is a cross compiler.
- The process of creating executable code in different machines is called "retargeting".
- The cross compiler is also known as retargetable compiler.
- Eg. of cross compiler \Rightarrow GNU gcc

windows

linux



$\text{GCC} \rightarrow$ old mic with 8 bit processor
 \rightarrow intel Z3 64 bit processor

Both the mic contains same OS but still we are not able to implement GCC on second mic.

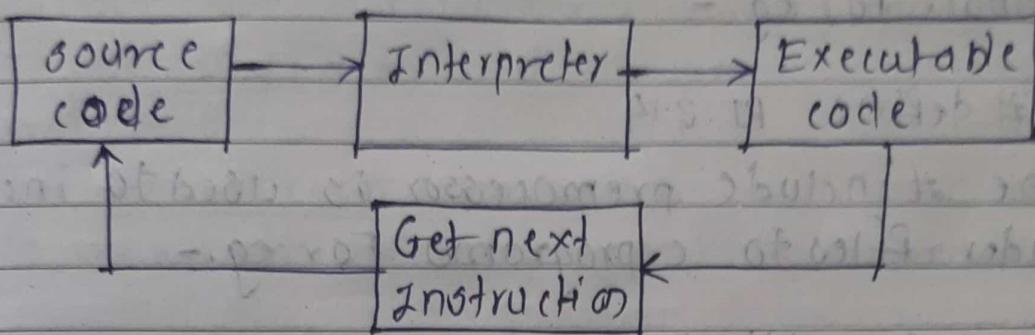
so, GCC will get upgrade to implement on intel Z3, so this is known as cross compiler.

* Interpreter

- The translation of a single statement of the source program into machine code is done by a language processor & executes immediate before moving on to the next line is called an interpreter.
- If there is an error in the statement, the interpreter terminates its translating process at that stmt & displays an error message.
- The interpreter moves onto the next line for

execution only after the removal of the error.

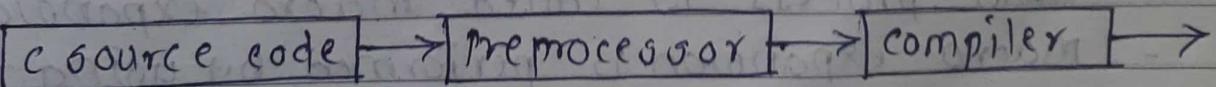
- An interpreter directly executes instructions written in a programming or scripting language without previously converting them to an object code or machine code.
eg. - perl, python, matlab



* Interpreter Phases & Working

- Interpreter works like the front end of a compiler. That is, an interpreter does lexing, parsing & semantic analysis.
- The process of interpretation can be carried out in following phases
 - lexical analysis
 - syntax analysis
 - semantic analysis
 - direct execution

* Preprocessor; header file & macro expansion



Working of C Preprocessor

- The c preprocessor is a macro preprocessor (allows you to define macros) that transforms pgm before it is compiled.
- These transformation can be the inclusion of header files, macro expansion etc.
- All preprocessing directives begin with a # symbol. for eg. -

#define PI 3.14

- The #include preprocessor is used to include header files to c programs. For eg. -

#include <stdio.h>

- Here, stdio.h is a header file. The #include preprocessor directive replaces the above line with the contents of stdio.h header file.
- Macro provide the feature for defining new operations & data in programming lang.
- A macro definition in a pgm defines either new operation or new method of declaring data.
- A macro call in the pgm is an innovation of the new operation or the new method of declaring data defined in the macro.
- It leads to a pgm generation activity during which the macro call is replaced by a sequence of stmts in the programming lang. This process is called macro expansion.
- Macro expansion is performed by using two kinds of actions -
 - 1) lexical substitution
 - 2) Semantic expansion

- A macro assembler performs expansion of each macro call into a sequence of assembly stmts & also assembles the resulting assembly pgm.
- A macro preprocessor merely performs expansion of macro calls & produces an assembly program.

Macros are useful for following purposes

- 1) To simplify & reduce the amount of repetitive coding
- 2) To reduce errors caused by repetitive coding
- 3) To make an assembly pgm more readable.

A macro consist of name, set of formal parameters & body of code. The use of macro name with set of actual parameters is replaced by some code generated by its body. This is called macro expansion.

Eg. -

```
#define max(a,b) a>b ? a : b
main()
{
    int x,y;
    x=4, y=6;
    z=max(x,y); }
```

After preprocessing

becomes $z = x > y ? x : y;$

After macro expansion, the whole code would look like :

```
#define max(a,b) a>b ? a : b
main()
{
    int x,y;
```

$x=4; y=6; z = x > y ? x:y; \}$

* Two pass macro processor

General Design steps

1. Specification of problem
2. Specification of databases
3. Specification of database Formats
4. Algorithm

There are four basic tasks that a macro processor must perform

1) Recognize macro definition -

A macro inst. processor must recognize macro def. identified by MACRO & MEND pseudo operations.

2) Save definition -

The processor must store macro definition which will be needed at the time of expansion of calls.

3) Recognize macro call -

The processor must recognize macro calls that appear as operation mnemonics.

4) Expand calls & substitute arguments -

The macro call is replaced by the macro definition. The dummy arguments are replaced by the actual data.

Specification of databases

Pass 1 databases -

- The input macro source pgm.
- The output macro source pgm to be used by Pass 2.
- Macro definition Table (MDT), to store the body of macros.
- Macro definition table counter (MDTC), to mark next available entry in MDT.
- Macro Name Table (MNT), used to store names of macros.
- Macro Name Table counter, used to indicate the next available entry in MNT.
- Argument List Array (CALA), used to substitute index.

Pass 2 Data bases -

- The copy of the input from Pass 1 (Intermediate File).
- The output expanded source to be given to assembler (Target file).
- MDT, created by pass 1.
- MNT, created by pass 1.
- Macro definition table pointer (MDTP), used to indicate the next line of text to be used during macro expansion.
- Argument list array (CALA), used to substitute macro call arguments for the index markers in the stored macro.

Pass 1 : Algorithm

1. Pass 1 of macro processor makes a line by line scan over its input.
2. Set MDT C = 1 as well as MNT C = 1.
3. Read next line from input program.
4. If it is a MACRO pseudo-op the entire macro definition except this (MACRO) line is stored in MDT.
5. The name is entered into Macro Name Table along with a pointer to the first location of MDT entry of the definition.
6. When the END pseudo-op is encountered all the macro have been processed, so control is transferred to pass 2.

Pass 2 : Algorithm

1. This algo. reads one line of IP pgm at a time.
2. For each line it checks if op-code of that line matches any of the MNT entry.
3. When match is found (ie when call is pointer called MDTF to corresponding macro stored in MNT entry).
4. The initial value of MDT P is obtained from index field of MNT entry.
5. The macro expander prepares the ALA consisting of a table of dummy argument indices & corresponding arguments to the call.
6. Reading proceeds from the MDT, as each successive line is read, the values from the argument list are submitted for the dummy arguments indices in the macro.
7. Reading MEND line in MDT terminates expansion.

of macro & scanning continues from 'IP File'.

8. When END pseudo-op encountered, the expanded source pgm is given to the assembler.

Databases

Macro Name Table (MNT)			Macro defn Table (MDT)		
Index	Macro Name	MDT index	Index	card	
1	xhai	1	1	1	

Argument list array (ALA)	
Index	Argument

card-means inst.
(body of man)

→ positional or actual argument

Example (2 Pass Macro Processor)

SOURCE PGM	109AB SDRA8	20 INSTRUCTION CODE
MACRO		EO
&LAB INCR &ARG1, &ARG2, &ARG3		\$
&LAB ADD AREG, &ARG1		
ADD AREG, &ARG2		
ADD AREG, &ARG3		
MEND		
START		START.
LOOP INCR A,B,C		LOOP INCR A,B,C
LABEL INCR DATA1,DATA2,DATA3		LABEL INCR DATA1,DATA2, DATA3
A DC 2		A DC 2
B DC 2		B DC 2
C DS 2		C DS 2
DATA1 DC 3		DATA1 DC 3
DATA2 DS 2		DATA2 DS 2
DATA3 DC 4		DATA3 DC 4
END		END

MDT (Macro Def. Table)

Index	Code	
01	&LAB INCR &ARG1, &ARG2, &ARG3	
02	#1 ADD AREG #2	#1, #2 ... \Rightarrow defining
03	ADD AREG #3	index pos. of MDT
04	ADD AREG #4	
05	MENID	

MNT (Macro Name Table)

Index	Macro Name	MDT index
01	INCR	01

ALA (Argument list array)

Index	Argument
01	&LAB
02	&ARG1
03	&ARG2
04	&ARG3

So, the output of pass 2 of 2 passes Macro processor is : Intermediate code, MNT, MDT, ALA Table.

Pass 2 of 2Pass Macro Processor [Target code]

Intermediate code

OTART

LOOP INCR A B, C

LABEL1 INCR DATA1, DATA2, DATA3

A DC 2

B DC 2

C DC 2

DATA1 DC 3

DATA2 DS 2

DATA3 DC 4

END

Expanded code

START

LOOP ADD AREG, A

ADD AREG, B

ADD AREG, C

LABEL1 ADD, AREG, DATA1

ADD AREG, DATA2

ADC2 ADD AREG, DATA3

B DC 2

C DC 2

DATA1 DC 3

DATA2 DS 2

DATA3 DC 4

END

Prepare new ALA

ALA		ALA	
Index	Argument	Index	Argu.
01	LOOP	01	LABEL1
02	A	02	DATA
03	B	03	DATA2
04	C	04	DATA3

- Here macro is call . so we can

prepare a new ALA with macro call

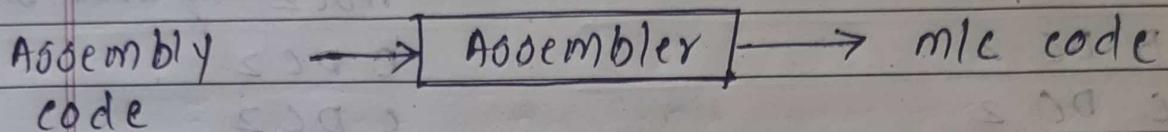
with arguments

loop add areg, a

change it to a call

* Assembler

- Assembler is a pgm for converting instructions written in low level assembly code into relocatable mc code & generating along information for the loader.



- It generates instructions by evaluating the mnemonics (symbols) in operation field & find the value of symbol and literals to produce mc code.
- We can create an assembly lang. code using a compiler or a programmer can write it directly.
- Most programmers use high level langs but, when more specific code is required, assembly lang. is used.
- It uses opcode for the instructions.
- An opcode basically gives info. about the particular instruction.
- The symbolic representation of opcode (mc level inst.) is called mnemonics.

For example ADD A,B

Here, ADD is the mnemonic that tells the processor that it has to perform addition function. A & B are the operands.

Types of Assembler

Assemblers generate instruction, on the basis of a no. of phases used to convert to mic code, assemblers have two types:

1) One-pass Assembler :-

These assemblers perform the whole conversion of assembly code to mic code in one go.

2) Multi-pass / Two-pass Assembler

These assemblers first process the assembly code & store values in the opcode table & symbol table. And then in the second step, they generate the mic code using these tables.

a) Pass 1 -

- symbol table & opcode tables are defined.
- keep the record of the location counter
- Also, processes the pseudo instructions.

b) Pass 2 -

- Finally, converts the opcode into the corresponding numeric opcode.
- Generates mic code according to values of literals & symbols.

opcode table -

They store the value of mnemonics & their corresponding numeric values.

symbol table -

They store the value of programming lang. symbols used by the programmer, & their corresponding numeric values.

Location counter -

It stores the address of the location where the current inst. will be stored.

* Linkers & loaders

Translation Hierarchy

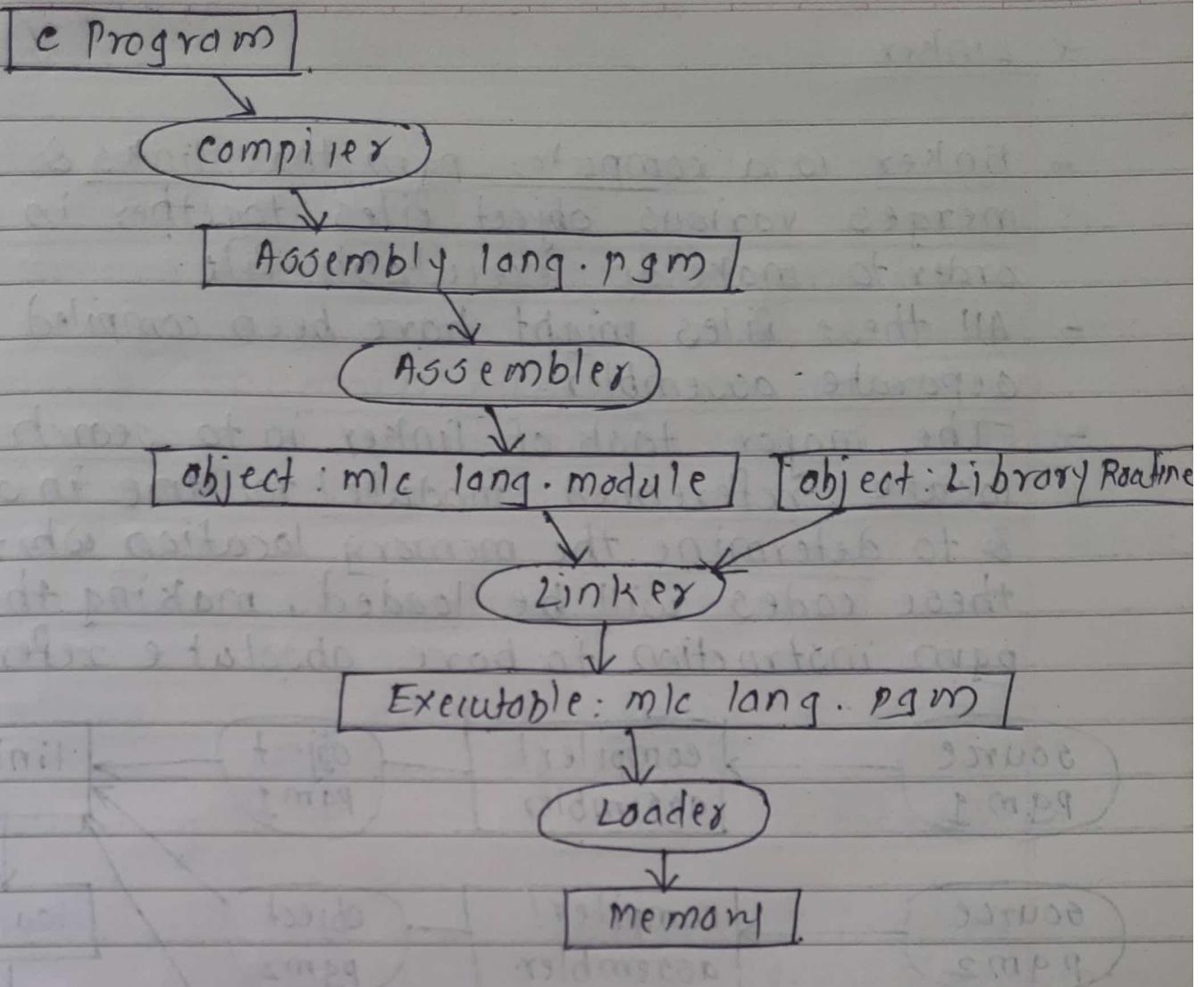
1) Compiler -

Translates high level language pgm into assembly language.

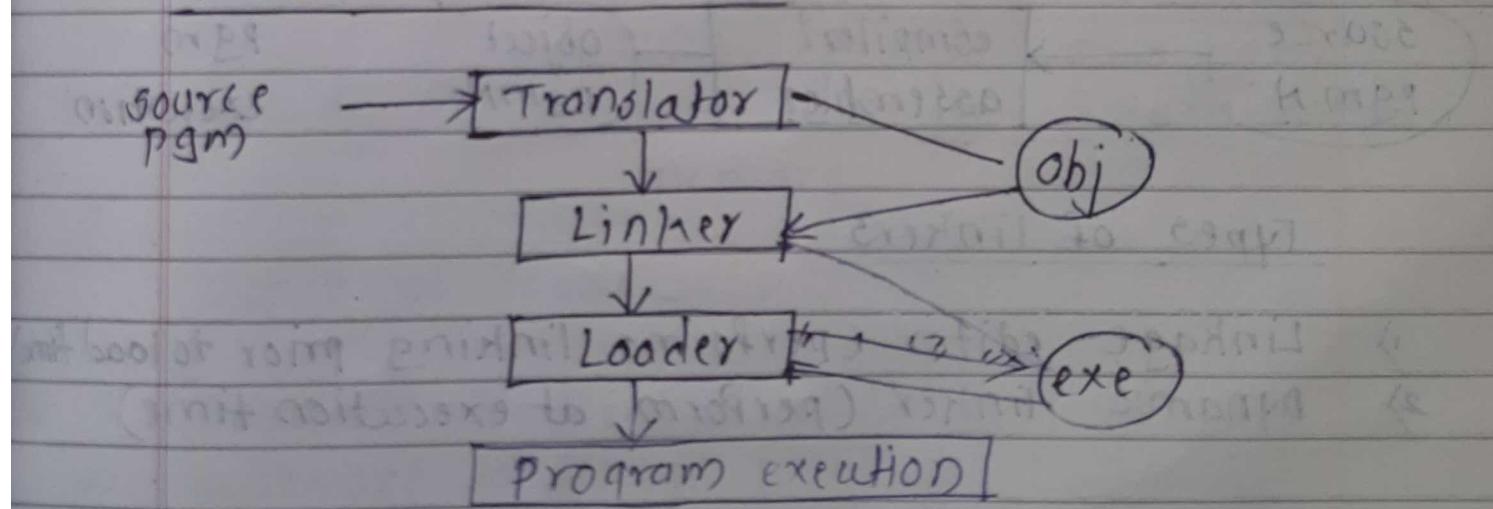
2) Assembler -

converts assembly language pgm into object files.

- object file contains a combination of mc instructions, data & information needed to place instructions properly in memory.



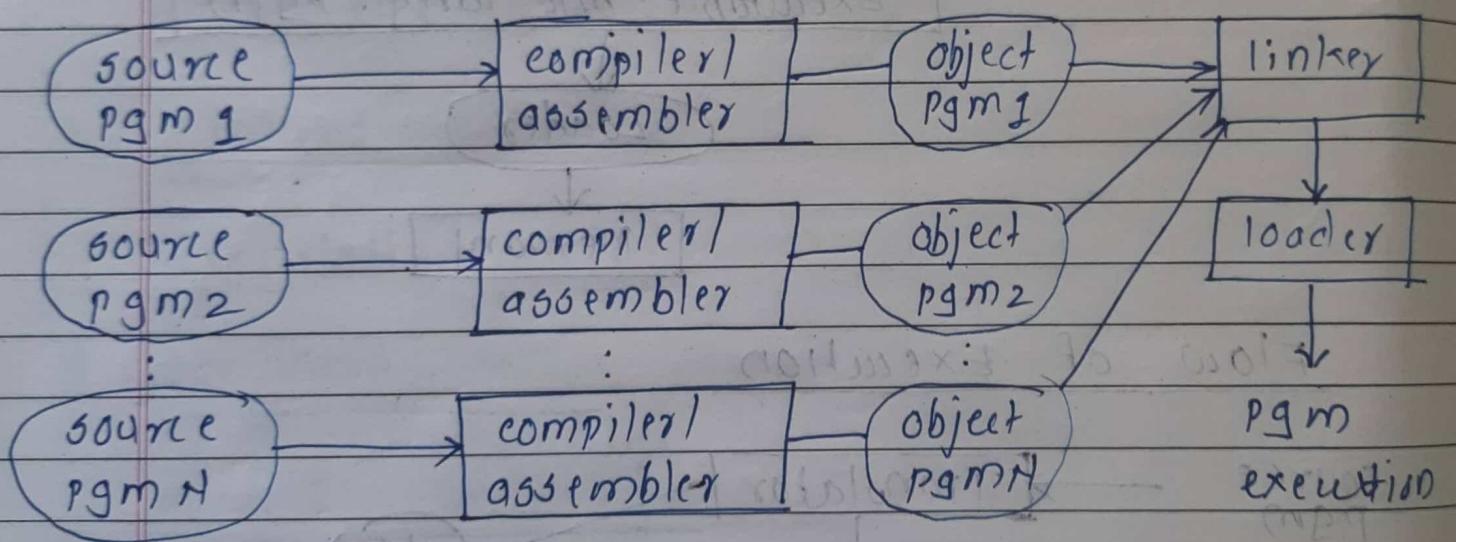
Flow of Execution



printf is a fn. that is not in your code but in an object library. printf calls some target function to add your message. The linker extracts object module from libraries & add them to your executable.

* Linker

- Linker is a computer pgm that links & merges various object files together in order to make an executable file.
- All these files might have been compiled by separate assembler.
- The major task of linker is to search & locate referenced module/routine in a pgm & to determine the memory location where these codes will be loaded, making the pgm instruction to have absolute references.



Types of linkers

- 1) Linkage editor (performs linking prior to loading)
- 2) Dynamic linker (performs at execution time)

linking is a process of collecting & maintaining piece of code & data into a single file.

* Loader

- loader is a pgm which accepts the lbp as linked modules & loads them into main memory for execution by the computer.
- Loaders load or copies pgm from secondary storage to main memory for execution.
- Loaders can also relocate virtual addresses with physical addresses.

Functions of loaders

Allocation - allocated space in memory for the pgms.

Linking - symbol resolution bet. object modules.

Relocation - Adjust address dependent locations of address constants i.e. assign load addresses to different parts of a pgm.

loading - store mc inst. & data into memory.

Types of loaders

- 1) Absolute loader
- 2) Bootstrap loader
- 3) Relocating loader (Relative loader)
- 4) Linking loader

Absolute loader -

It loads the object code to specified locations in the memory. At the end the loader jumps to the specified address to begin execution of the loaded program.

Bootstrap loader -

When a computer is first turned on or restarted a special type of absolute loader, called bootstrap loader is executed. This bootstrap loads the first pgm to be run by the computer - usually an operating system. The bootstrap itself begins at address 0.

Relocating loader -

- loaders that allow for pgm relocation are called relocating loaders or relative loaders.
- The execution of the object pgm is done using any part of the available & sufficient memory.
- The object pgm is loaded into memory wherever there is a room for it.
- The actual starting address of the object pgm is not known until load time.
- Relocation provides the efficient sharing of the mc with larger memory & when several independent programs are to be run together.
- It also supports the use of subroutine libraries efficiently.

Linking loaders -

- Performs all linking & relocation at load time.
The other alternatives:

linkage editors -

- Produces a linked version of the pgm - called a load module or an executable image.
- This load module is written to a file or library for later execution.
- The linked pgm produced is suitable for processing by a relocating loader.
- Using linkage editor, an absolute object program can be created. if starting address is already known.

Dynamic linking -

- This scheme postpones the linking functions until execution.
- A subroutine is loaded & linked to the rest of the program when it is first called.
- It is usually called dynamic linking, dynamic loading or load on call.

Macro Example (Tut)

MACRO

INCR &M, &I, ® = AREG

MOVER ®, &M

ADD ®, &I

MOVE M ®, &M

MEND

MACRO

CALC &X, BY=B, &OP=MULT

MOVER AREG &X

```

&OP AREG, &Y
MOVEM AREG, &X
MEND
START 100
READ N1
CALC A, OP=ADD
ADD AREG, 21
LDA CREG, 100
SUB CREG, A
INCR B, REG=CREG, I=A
A DS 1
B DS 1
END

```

* Single and Two pass assembler

Assembly lang -

Assembly lang is a kind of low level programming lang., which uses symbolic codes or mnemonics as instruction.

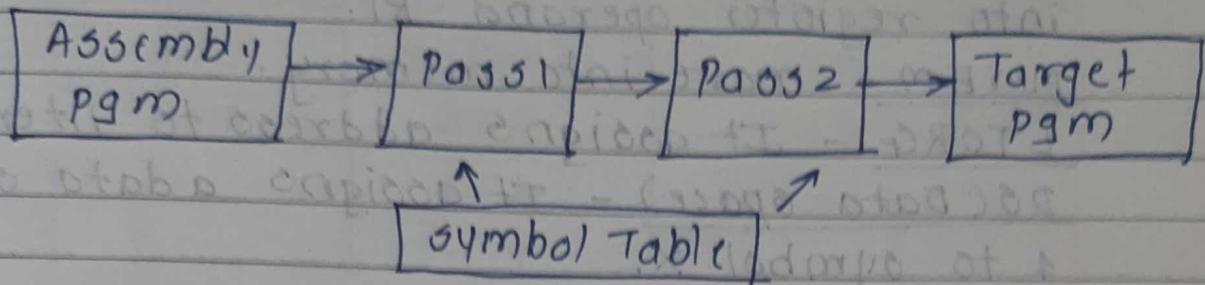
Eg. - ADD, SUB, LDA, STA

Elements of Assembly lang.

- 1) Mnemonic operation code - (micro inst.)
ADD, SUB, MOVE
- 2) Symbolic operands - (variables, symbolic name)
ADD R1, R2, R3
- 3) Data Declaration -
NUM1 03
NUM1 0011H

Assembly processes

convert .asm file into .obj file



Pass 1

complete scan .asm file. Find all labels, instructions & calculating corresponding address.

Pass 2

convert all the instructions into mc lang.

Format

Symbol Table:

stored all the info. of assembly lang. data, variables, instructions, addresses etc.

Example of Assembly lang:

Assembly Program

Label	opcode	(value)	operands	cradle (add?)
JOHN	START	200		
	MOVER	R1, = "3"	200	
	MOVEM	R2, X	201	
L1	MOVER	R2, = "2"	202	
	LTORG		203	
X	DS		204	
	END		205	

START - This inst. starts the execution of pgm.

MOVER - It moves the content of literal (= 'a') into register operand R1.

MOVEM - move into memory.

LTORG - It assigns address to literals.

DSE (Data space) - It assigns a data space of 1 to symbol X.

Design of working of Assembler

1) Analysis Phase (Phase 1 of Assembler) -

- To build symbol table for synthesis phase to process.
- Determine address of each symbols called as memory allocation.
- location counter used to hold address of next inst.
- isolated label, mnemonic opcode, operands, constants etc.
- validate meaning of address of each inst.

2) Synthesis Phase (Phase 2 Assembler) -

- use data structure generated by analysis phase.
- To build mic inst. For every assembly stmt as per mnemonic code & three address allocation.
- Synthesis mic inst. as per source code.

Pass 1 Assembler :

- 1) Symbol Table (ST or SYMTAB) : store value or address assign to the label.

label	address
JOHN	200
L1	202
x	204

- 2) Literal Table (LT or LITAB) : store each literal or constants (= '3') with its location.

Index	literal	Address
0	= '3'	200
1	= '2'	202

- 3) Operation code Table (optab) : store mnemonic operation code with there opcodes & length. on which element action is performed

Inst.	opcode	length (bytes)	
MOVER	3	2	-int
MOVEM	X	1	-char
MOVER	2	2	-int

Pass 2 Assembler :

- Assembled using symbol table.
- generates mic code for inst.
- convert mnemonic operation code, symbol table operands with there equivalent mic code.
- convert data constants to internal mic representations.
- convert complete pgm into obj file.

Example:

START	200	800	1000
MOVER	AREG1 = "5"	at	origin
MOVEM	AREG, X	800	1000
L,	MOV R BREG, = '2'	805	1005
ORIGIN	L1+3	805	11
LTORG		805	15
X	DS 1		

END

comrba	(0x00)	x000
805	'0' =	0
805	'5' =	1

Notes: (dotgo) 800 is 00000000
80000000 is 00000000000000000000000000000000

(certid)	80000000000000000000000000000000	00000000000000000000000000000000
5.	8	MOVE
1	X	MOVE
5	5	MOVE

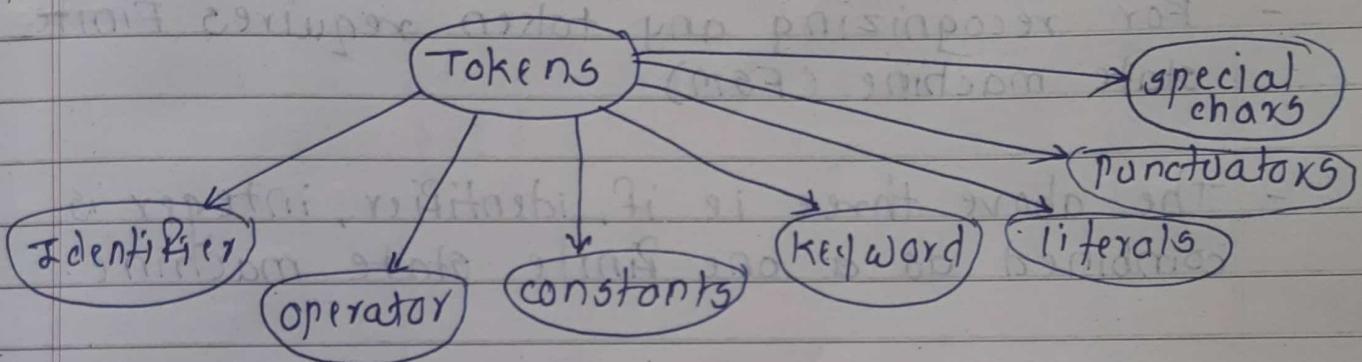
Unit 2: Lexical Analysis & Introduction of Syntax Analysis

Introduction to compiler, Phases & Passes, Bootstrapping, Role of lexical analyzer, specification & Recognition of tokens, LEX/FLEX, Expressing syntax, Top-down Parsing, Predictive Parsers, Implementing scanners, operator precedence parsers.

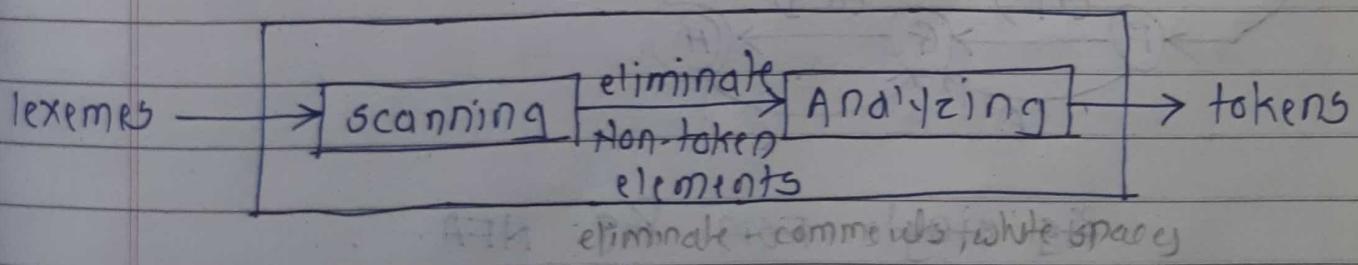
* Lexical Analyzer

Features

- Scans the pure HLL code line by line
- Tokens lexemes as i/p & produces Tokens.



- Tokens are classified into above 7 categories.
- In lexical analyzer two functions take place.
 - Scanning
 - Analyzing



Let see how Analyzing works.

c - Tokens

- if : $\rightarrow A \xrightarrow{i} B \xrightarrow{f} C$
- Identifier : $\rightarrow D \xrightarrow{(a-z|A-Z|-)} E \xrightarrow{(a-z|A-Z|0-g)}$
- Integer : $\rightarrow F \xrightarrow{(+|-)} G \xrightarrow{(0-g)} H$
 - \leftarrow any digit (single digit)
 - \leftarrow any digit (two or more digits)
 - \leftarrow slight transition from F to G
- For recognizing any token requires Finite state machine (FSM)
- The above three ie if, identifier, integer is combined as a one Finite state machine.

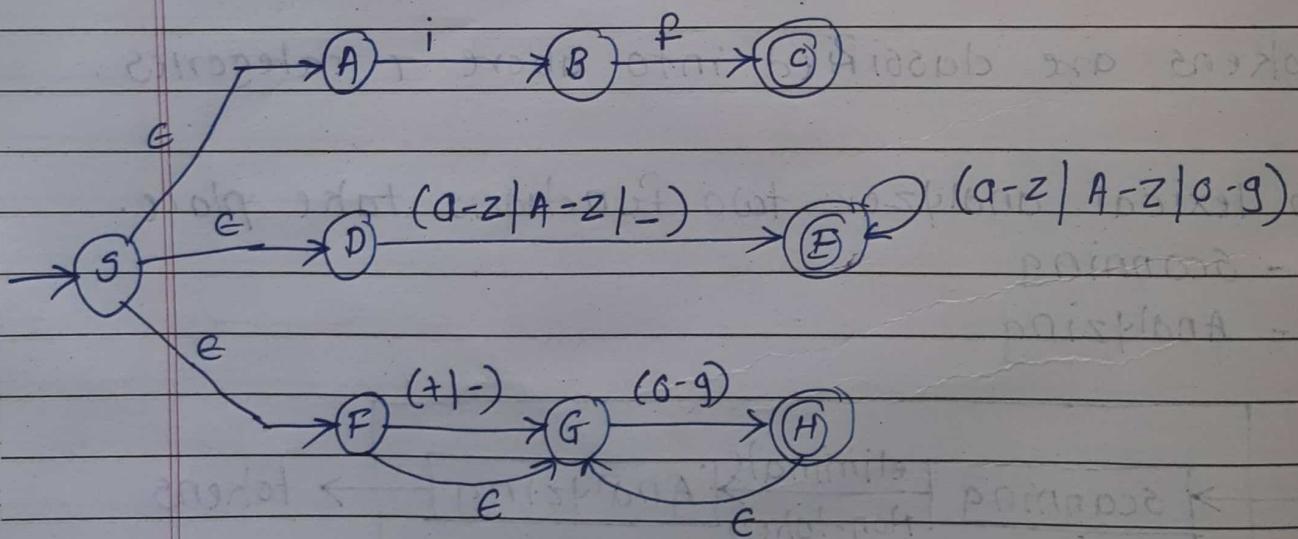


Fig: single FSM Fig: NFA.

- How the above will combine into a single FSM.

First

$$4) E \rightarrow TE' \quad \text{step 2 add id, (,)}$$

$$E' \rightarrow *TE | E \quad *, E$$

$$T \rightarrow FT' \quad \{ p, d, & id, (,) \} \text{ will be}$$

$$T' \rightarrow \epsilon | +FT' \quad \epsilon, +$$

$$F \rightarrow id | (E \quad \rightarrow \text{a instance of id, (,) on follow}$$

FOLLOW()

$\text{Follow}(A)$ contains set of all terminals present immediately in right of A.

Rules:

Follow of start symbol is $\$$

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

$$2) S \rightarrow A C D E$$

$$C \rightarrow a b$$

$$\text{Follow}(A) = \text{First}(C) = \{ a, b \}$$

at right side of A there is C which is non terminal so calculate First of that variable.

$$\text{Follow}(D) = \text{Follow}(S) = \{ \$ \}$$

after D there is nothing so write ϵ , so it will point to start symbol.

3) $S \rightarrow aSbS \mid bSaS \mid \$ \mid \epsilon$

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

Follow never contain ϵ

4) $S \rightarrow AaAb \mid BbBa$

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

$$\text{Follow}(A) = \{a, b\}$$

$$\text{Follow}(B) = \{b, a\}$$

5) $S \rightarrow ABC$

$$A \rightarrow DEF$$

$$B \rightarrow \epsilon$$

$$C \rightarrow \epsilon$$

$$D \rightarrow G$$

$$E \rightarrow \epsilon$$

$$F \rightarrow \epsilon$$

$$\text{Follow}(A) = \text{First}(B) = \text{First}(C) = \{\phi\} =$$

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

6) $E \rightarrow E + T \mid T$

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

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

After eliminating left recursion the grammar is:

$E \rightarrow TE'$

Follow

 $E' \rightarrow +TE'|E$

$\Rightarrow \emptyset$ (Follow of start)

 $T \rightarrow FT'$

\Rightarrow if $A \rightarrow \alpha B\beta$ - Follow of B

 $T' \rightarrow *FT'|E$

$\Rightarrow A \rightarrow \alpha B$, First of B is except first of B

 $F \rightarrow (\underline{E})+id$

\Rightarrow Follow(B) is Follow(A)

First(E) = First(T) = First(F) = { (, id } \Rightarrow First of A

First(E') = { +, E }

First(T) = First(F) = { (, id }

First(T') = { *, E }

First(F) = { (, id } \Rightarrow Follow of B will be same

Follow(E) = { \$,) } \Rightarrow E is start symbol go to \$

Follow(E') = { *, E } \Rightarrow { \$,) } \Rightarrow After E is nothing so Follow(E') is Follow(E)

Follow(T) = { +, \$,) } \Rightarrow First of T is + \Rightarrow First of T \Rightarrow First of T' \Rightarrow First of E \Rightarrow Follow(E)

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

Follow(F) = { +, *, \$,) } \Rightarrow First of T' + Follow(T)

For Follow(E) check on right side where the E appear & taken the right of E
ie (E) =)

Follow(E') = Follow(E) union : () , *

7) $s \rightarrow iETs | iETsEs | a$

$E \rightarrow b$

after eliminating left factoring, \leftarrow

$s \rightarrow iETsS' | a$

$S' \rightarrow eS | e$ follow & form two shift item -

$E \rightarrow b$

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

parsible top's rule.

$$A \rightarrow \alpha$$

$\Rightarrow m[A, \alpha] = A \rightarrow \alpha = \alpha$ is
print of α .

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

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

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

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

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

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

$\Rightarrow m[A, b] = A \rightarrow \alpha - \text{if } e$
is in $\text{First}(A)$ & b is
in $\text{Follow}(A)$

* Parsing Table (For eq. 6f) substitute all production rules in table
Predictive Parsing Table : Find first of right hand side of production

N.T.	+	*	() { }	id	\$
E			$E \rightarrow TE'$		$E \rightarrow TE'$
E'	$E' \rightarrow +TE'$				
T			$T \rightarrow FT'$		$T \rightarrow FT'$
T'	$T' \rightarrow E$	$T' \rightarrow *FT'$			$T' \rightarrow G$
F			$F \rightarrow (E)$		$F \rightarrow id$

$\Rightarrow E \rightarrow TE' \Rightarrow \text{First}(T) \Rightarrow \{+, id\}$ so in (+, id)
continue - we write $E \rightarrow TE'$

$\Rightarrow E' \rightarrow +TE'/e \Rightarrow \text{First}(+) = + \Rightarrow \text{Follow}(E') \Rightarrow \$,$

$\Rightarrow T \rightarrow FT' \Rightarrow \text{First}(F) \Rightarrow \{, id\}$

$\Rightarrow T' \rightarrow *FT'/e \Rightarrow \text{First}(*) \Rightarrow \text{Follow}(T') \Rightarrow +, \$,$

* LL(1) : Grammar parsing tables $\Rightarrow F \rightarrow (E) + id \Rightarrow$
① ② $\text{First}(F) = \{, id\}$

$$S \rightarrow (L) | q$$

$$L \rightarrow SL'$$

$$L' \rightarrow E | ., SL'$$

- First find out first & follow for all non-terminals.

- Introduce new initial state 3 & have slight transitions to all other states.
- This one is NFA, so this is an epsilon NFA
- NFA is conceptual means we can't implement it.
- For implementation we will derive an equivalent Deterministic finite automata (DFA)
- The below is equivalent DFA

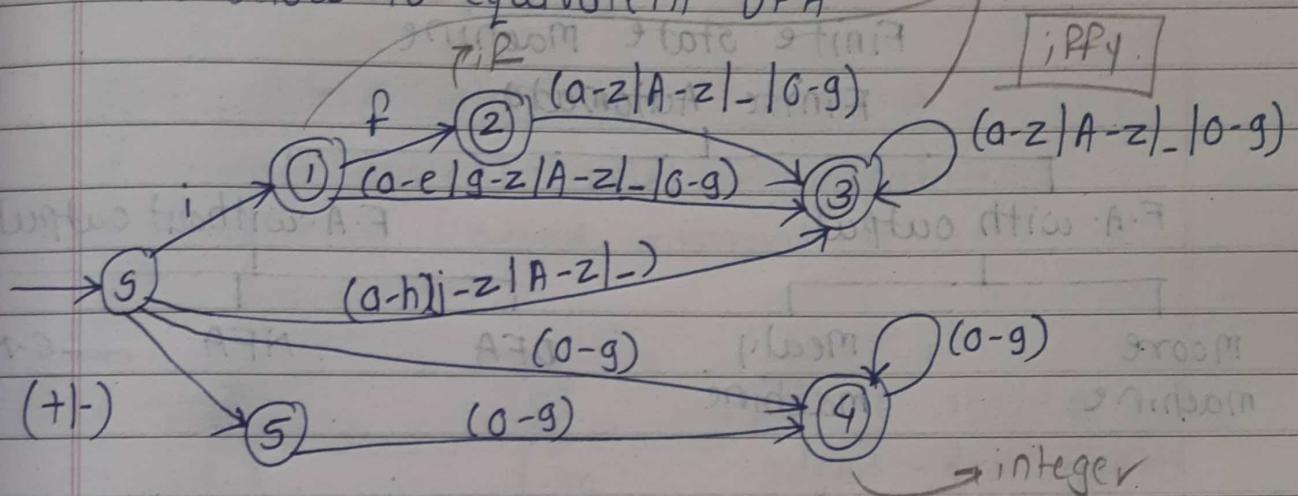
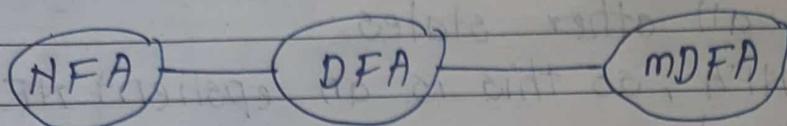


Fig: DFA (AF) automata. string

- $s - 1 - 2 \Rightarrow if$
- $s - 1 - 2 (if) - 3 (iff) - 3 (iffy)$
- $s - 1 - 3$ (except i)
- $s - 3$ (except i)
- $s - 4$ (any single digit)
- $s - 4 - 4$ (self loop) (two or more digits)
- lexical analyzer takes lexemes as if & produces tokens - which uses DFA for pattern matching.



* Finite Automata

Finite state machine
Finite Automata

F.A. with output
Moore machine

F.A. without output
DFA

- Finite automata (FA) is the simplest mc to recognize patterns.
- It has set of states and rules for moving from one state to another but it depends upon the applied input symbol.

7

Features of automata:

1. Input
2. Output
3. States of automata
4. State relation
5. Output relation

It consists of following

Q : finite set of states

Σ : set of input symbols

q_0 : initial state

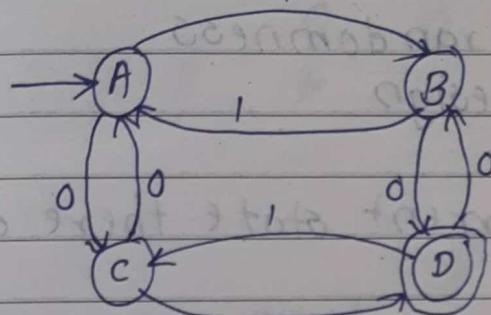
F : set of final states

δ : transition function

$$\{ Q, \Sigma, q_0, F, \delta \}$$

DFA - Deterministic Finite Automata

- It is the simplest model of computation
- It has very limited memory



where,
circle - states
edges - transitions
0,1 - inputs

Every DFA can be defined using 5 tuples

$$(Q, \Sigma, q_0, F, \delta)$$

Q - set of all states

Σ - inputs

q_0 - start state / initial state

F - set of final states

δ - transition function from $Q \times \Sigma \rightarrow Q$

$$Q = \{ A, B, C, D \}$$

$$\Sigma = \{ 0, 1 \}$$

$$q_0 = A$$

$$F = \{ D \}$$

$$\delta =$$

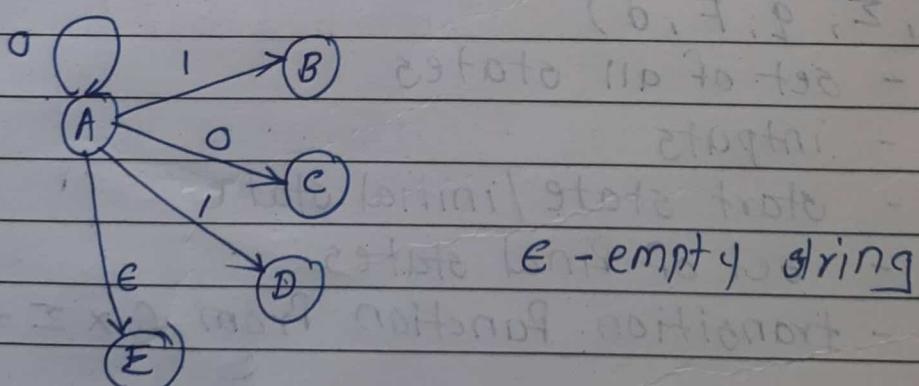
For δ

	0	1
A	C	B
B	D	A
C	A	D
D	B	C

when we are in state A
if my input is 0 then we
will go to state C & if
my input is 1 we will go
to state B, likewise.

NFA - Non-deterministic Finite automata

- In DFA, given the current state we know what the next state will be
- It has only one unique next state
- It has no choices or randomness.
- It is simple & easy to design
- In NFA, given the current state there could be multiple next states.
- The next state may be chosen at random.
- All the next states may be chosen in parallel



- NFA is also define with 5 tuples $(Q, \Sigma, q_0, F, \delta)$

Q - set of all states

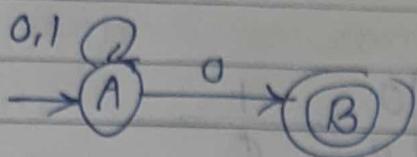
Σ - input

q_0 - start or initial state

F - set of final states

δ - transition function that maps to

$$Q \times \Sigma = 2^Q$$



$$Q = \{A, B\}$$

$$\Sigma = \{0, 1\}$$

$$q_0 = A$$

$$F = B$$

$$\delta = ?$$

	0	1	
A	A	B	A state A on getting 0
B	B	A	A state B on getting 1

$A \times 0 \rightarrow A$ state A on getting 0
 $A \times 1 \rightarrow B$ state A on getting 1
 $B \times 0 \rightarrow \emptyset$ does not go anywhere
 $B \times 1 \rightarrow \emptyset$ represented by \emptyset

$A \rightarrow A, B, AB, \emptyset$ for state A getting 0, either 0 or 1 it

if there are 2 states there are 4 possibilities

For 3 it is 8 possibilities

* Conversion of NFA to DFA

- Every DFA is an NFA, but not vice versa, but there is an equivalent DFA for every NFA.

DFA

$$\delta = Q \times E \rightarrow Q$$

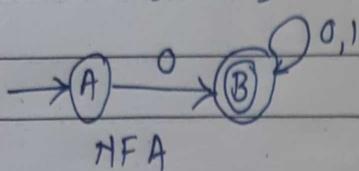
NFA

$$\delta = Q \times \Sigma \rightarrow 2^Q$$

$$NFA \cong DFA$$

$L = \{ \text{set of all strings over } (0,1) \text{ that starts with '0'} \}$

$$\Sigma = \{ 0, 1 \}$$



	0	1
A	B	\emptyset
B	B	(\emptyset, A)

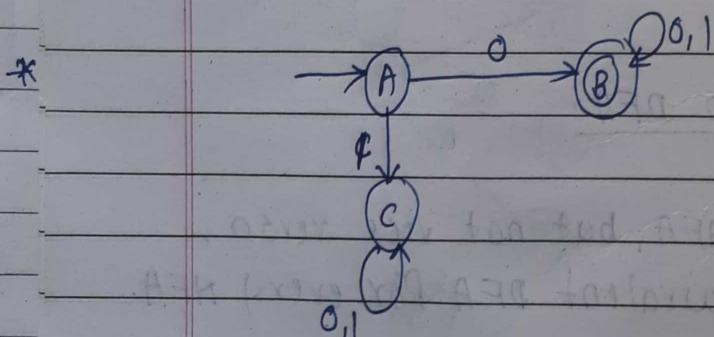
* Now convert this state transition diagram to equivalent DFA

DFA -

	0	1	
A	B	C	C - dead state
B	B	B	trap state
C	C	C	

* we can't indicate \emptyset in DFA means a dead state can't indicate so we introduce a new state i.e. C.

state diagram for DFA



Feature:

- 3) Removes comments & whitespaces from the Pure HLL code.

Types of comments in C

// single line comment

/* multi line comment */

- while scanning lexical analyzer detects comments
- If it's a single line comment till enter encounters it ignores that line & if its multiline till the next */ it will ignore all the text.
- The lexical analyzer classify the comments as non token elements & eliminates those.

For eg. int v1; /* it is a comment */ pune;

int v1; pune;

/*

int v1 pune;

- after scanning by lexical analyzer the comment will be removed, the blank white space is put between the comment.
- During syntactic analysis an error is generated, there pune is reported as undeclared variable.
- The lexical analyzer predicts the 2 different tokens (ie v1 pune).

whitespaces

types

' ' space

'\t' horizontal tab

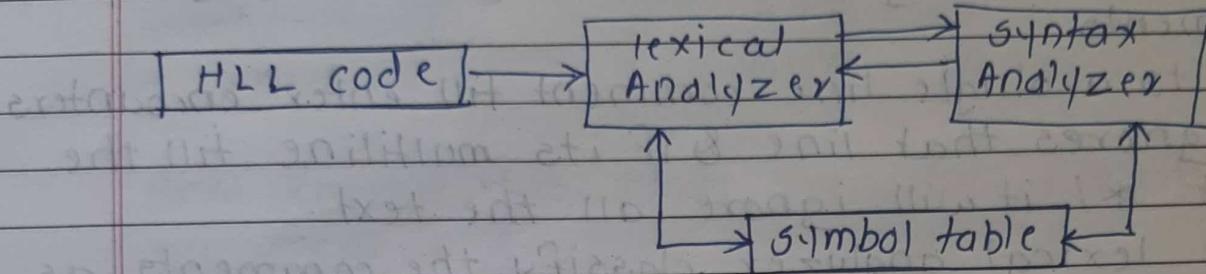
'\n' newline

'\v' vertical tab

'\f' form feed

'\r' carriage return

- all whitespaces are non-token elements.
- 4) Feature
Helps in macro expansion in the Pure HLL code.



* Specification & Recognition of Tokens

- Recognition of tokens by lexical analyzer.
- It is done to separate out different tokens.
- To understand the concept, let us consider that we have certain regular expression such as:
 - if → if
 - then → then
 - else → else
 - relop → < | ≤ | ≥ | = | = | <= | - relation operator
 - id → letter (letter/digit)* identifier (* - closer dual value)
 - num → digit* (digit)* ? (ε (+|-)) ? digit + ?
 - delim → blank / tab / newline

* Specification of Token

There are 3 specifications of tokens:

- 1) strings
- 2) languages
- 3) Regular expression

String & languages

- An alphabet or character class is a finite set of symbols.
- A string over an alphabet is a finite sequence of symbols drawn from that alphabet.
- A language is any countable set of strings over some fixed alphabet.
- In language theory, the terms 'sentence' & 'word' are often used as synonyms for string.
- The length of string s , usually written $|s|$, is the no. of occurrences of symbols in s .
- For example, banana is a string of length six. The empty string, denoted ϵ , is the string of length zero.

Operations on string

1. A prefix of string is any string obtained by removing zero or more symbols from the end of string s . For eg., ban is a prefix of banana.
2. A suffix of string s is any string obtained by removing zero or more symbols from the beginning of s . For eg., nana is a suffix of banana.
3. A substring of s is obtained by deleting any prefix & any suffix from s . For eg., non is a substring of banana.
4. The proper prefixes, suffixes, substrings of a string s are those prefixes, suffixes & substrings, respectively of s that are not ϵ or not equal to s itself.
5. A subsequence of s is any string formed by deleting zero or more not necessarily consecutive

positions of s.

6. For eg., baan is a subsequence of banana.

Operations on languages

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

- 1) Union
- 2) Concatenation
- 3) Kleene closure
- 4) Positive closure

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

1) Union : $L \cup S = \{0,1,a,b,c\}$

2) Concatenation : $L \cdot S = \{0a, 1a, 0b, 1b, 0c, 1c\}$

3) Kleene closure : $L^* = \{\epsilon, 0, 1, 00, \dots\}$

4) Positive closure : $L^+ = \{0, 1, 00, \dots\}$

Regular Expression

- Each regular expression r denotes a language $L(r)$.
- Here are the rules that define the regular expressions over some alphabet Σ & the languages that those expressions denote.
- 1) ϵ is regular expression, & $L(\epsilon)$ is $\{\epsilon\}$, ie, the lang. whose sole member is the empty string.

- 2) If 'a' is a symbol in Σ , then 'a' is regular expression & $L(a) = \{a\}$, ie, the lang. with one string, of length one, with 'a' in its one position.
- 3) Suppose r & s are regular expressions denoting the langs $L(r)$ & $L(s)$. Then $r|s$ is a regular expression denoting the lang $L(r) \cup L(s)$.
- b) $(r)s$ is a regular expression denoting the lang $L(r)L(s)$.
- c) $(r)^*$ is a regular expression denoting $(L(r))^*$.
- d) (r) is a regular expression denoting $L(r)$.
- 4) The unary operator * has highest precedence & is left associative.
- 5) Concatenation has second highest precedence & its left associative.
- 6) | has lowest precedence & is left associative.

Regular Set

- A lang. that can be defined by a regular expression is called a regular set. If two regular expressions r & s denote the same regular set, we say that they are equivalent & write $r = s$.
- There are a no. of algebraic laws for regular expressions that can be used to manipulate into equivalent forms.

For instance, $r|s = s|r$ is commutative; $r(s|t) = (rs)|t$ is associative.

Regular definitions

- Giving names to regular expressions is referred to as a regular definitions. If Σ is an alphabet

of basic symbols, then a regular definition is a sequence of definitions of the form

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

$$d_n \rightarrow r_n$$

1. Each d_i is a distinct name.

* 2. Each r_i is a regular expression over the alphabet $\Sigma \cup \{d_1, d_2, \dots, d_i\}$.

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

id \rightarrow letter (letter | digit)*

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

How to write regular expression?

1) *, +, {} \Rightarrow

These symbols act as repeaters & tell the computer that the preceding character is to be used for more than just one time.

2) \ - escape character

3) : - any character

4) \d - digit

5) \D - not a digit

6) \w - word character

7) \W - not a word character

8) \s - whitespace

9) \S - not whitespace

10) \b - word boundary

11) \B - not a word boundary

12) ^ - beginning of string

13) \$ - end of string

- [] - matches chars in brackets
- [^] - matches chars not in brackets
- | - either or
- () - capturing group
- * - 0 or more
- + - 1 or more
- ? - 0 or 1
- { } - exact no. of characters
- [min, max] - range of no. of chars.

* Specification of Tokens

* Recognition of Tokens

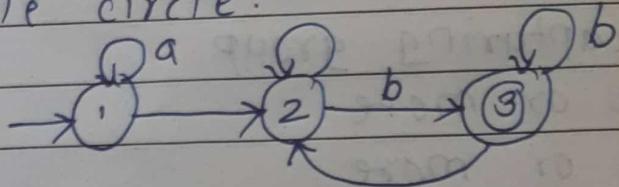
- Tokens can be recognized by Finite automata.
- A finite automata is a simple idealized machine used to recognize patterns within input taken from some character set.
- The job of FA is to accept or reject an i/p depending on whether the pattern defined by the FA occurs in the i/p.
- There are two notations for representing finite automata they are:

- 1) Transition Diagram
- 2) Transition Table

- Transition diagram is a directed labeled graph in which it contains nodes & edges.
- Nodes represents the states & edges represents the transition of a state.

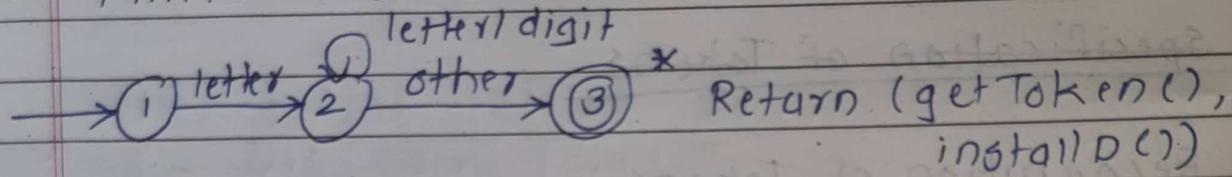
- Every transition diagram is only one initial state represented by an arrow mark (\rightarrow) & zero or more final states are represented by double circle.

eg. -

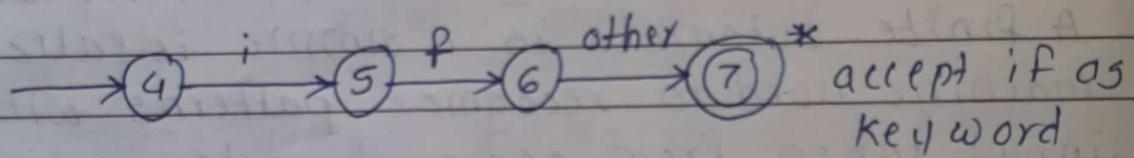


where state '1' is initial state & state 3 is final state.

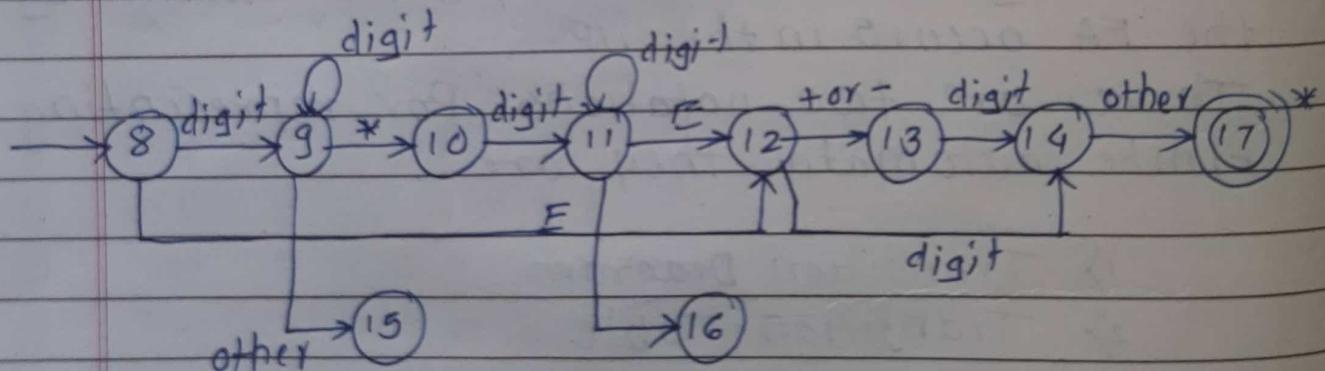
- Finite Automata for recognizing identifiers.



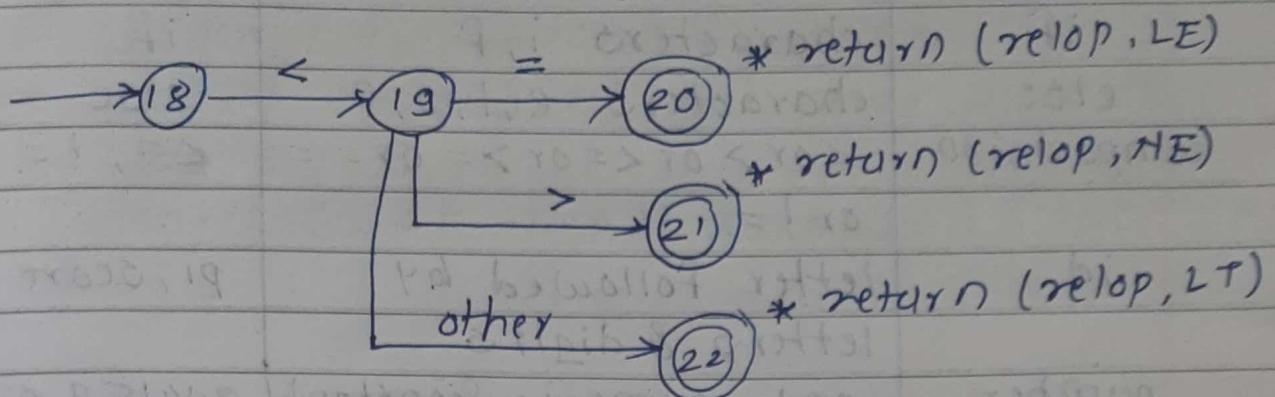
- Finite Automata for recognizing keywords



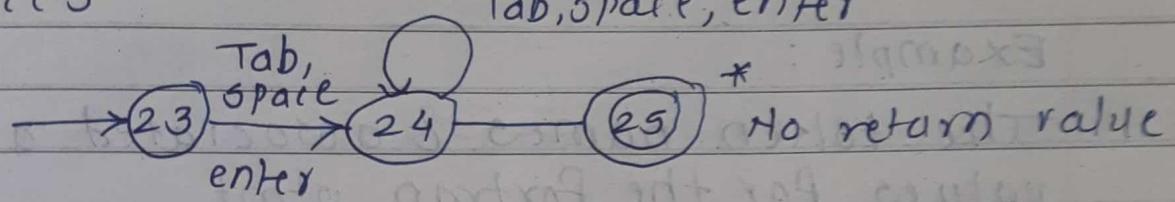
- Finite automata for recognizing numbers



- Finite automata for relational operators



- Finite automata for recognizing white spaces



* Tokens, Pattern, Lexemes

- A token is a pair consisting of a token name & an optional attribute value.
- A pattern is a description of the form that the lexemes of a token may take.
- A lexeme is a sequence of characters in the source program that matches the pattern for a token & is identified by the lexical analyzer as an instance of that token.

For eg :

`printf("Total = %d\n", score);`
both printf & score are ~~term~~ lexemes
matching the pattern for token id &
"Total = %d\n" is a lexeme matching
literal.

Token	informal description	sample lexem
if	characters i, f	if
else	characters e, l, s, e	else
comparison	< or > or <= or >= or == or !=	<=, !=
id	letter followed by letters & digits	pi, score, 02
number	any numeric constant	3.14159, 0, 6.02e
literal	anything but ", surrounded by "s	"core dumped"

Example :

The token names & associated attribute values for the fortran stmt

E = m * c ** 2

<id, pointer to symbol table entry for E>

<assign-op>

<id, pointer to symbol table entry for m>

<mult-op>

<id, pointer to symbol table entry for c>

<exp-op>

<number, integer value 2>

Example :

Divide the following C++ pgm into appropriate lexemes, which lexemes should get associated lexical values? What should those values be?

```

float limitedsquare(x) float x {
/* returns x-squared, but never more than
100 */
return (x <= -10.0 || x >= 10.0) ? 100 : x*x;
}

```

* Examples: Specification of tokens

Describe the languages denoted by the following regular expressions:

- 1) $a(a|b)^*a$
- string of a's & b's that start & end with a.
- 2) $((\epsilon | a)b^*)^*$
- string of a's & b's.
- 3) $(a|b)^*a(a|b)(a|b)$
- string of a's & b's that the character third from the last is a.
- 4) $a^*ba^*ba^*ba^*$
- string of a's and b's that only contains three b.
- 5) $((aa|bb)^*((ab|ba)(aa|bb)^*(ab|ba)(aa|bb)^*))^*$
- string of a's and b's that has an even number of a and b.

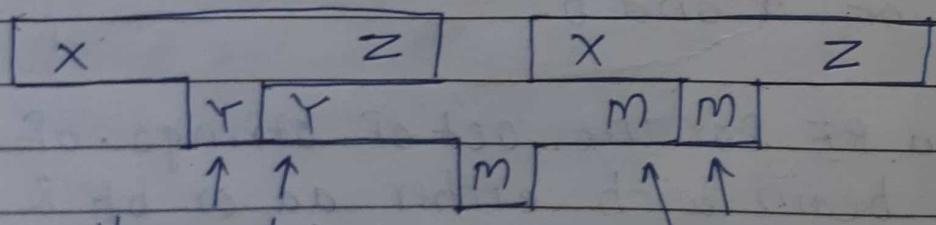
Write a RE for the set of strings of a's & b's that begin with either aa or bb & have an optional c at the end.

$\Rightarrow (aa|bb)(ab)^*c? / (aa|bb)[ab]^*c?$

* Bootstrapping in compiler design

- Bootstrapping is a process in which simple lang. is used to translate more complicated program which in turn may handle for more complicated pgm.
- This complicated pgm can further handle even more complicated pgm & so on.
- Writing a compiler for any high level lang. is complicated process. It takes lot of time to write a compiler from scratch.
- Hence simple lang. is used to generate target code in some stages.
- Suppose we want to write a cross compiler for new language X. The implementation lang. of this compiler is say Y, and the target code being generated is in language Z.
- That is, we create XYZ.
- Now if existing compiler Y runs on machine M & generates code for M then it is denoted as YM.
- Now if we run XYZ using YM then we get a compiler XMZ; that means a compiler for source lang. X that generates a target code in lang. Z & which runs on machine M.

Example:

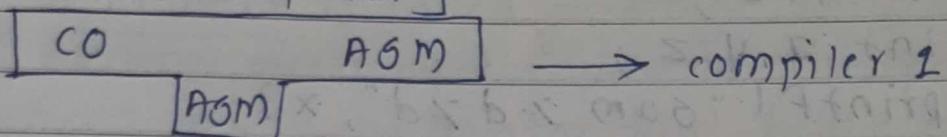


these two
langs must be same

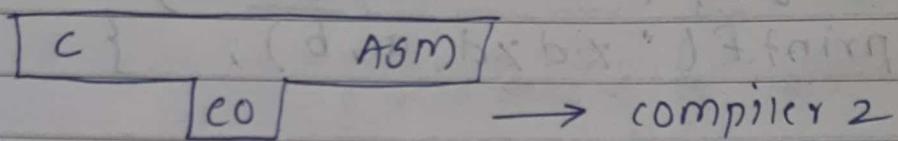
these two
langs must be same

compiler which takes C lang. & generates an assembly lang. as an o/p with the availability of a mic of assembly lang.

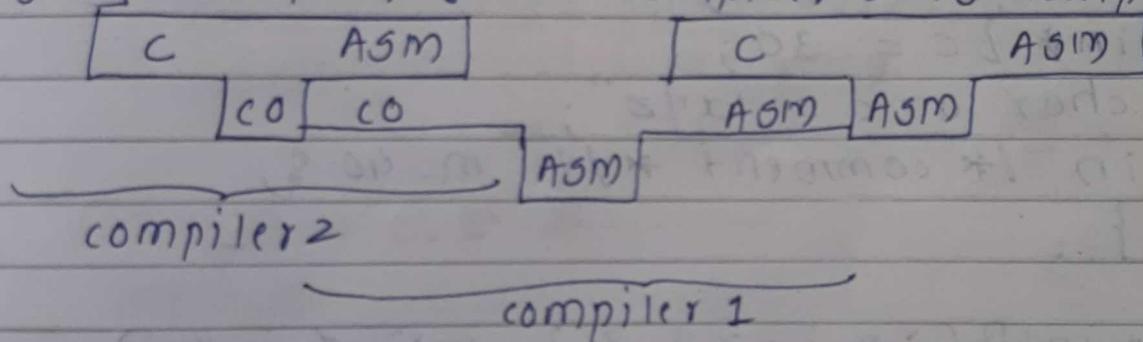
Step 1: First we write a compiler for a small of C in assembly lang.



Step 2: Then using with small subset of C ie CO, for the source lang: C the compiler is written



Step 3: Finally we compile the second compiler, using compiler 1 the compiler 2 is compiled.



Step 4: Thus we get a compiler written in ASM which compiles C & generates code in ASM.

Find the no. of tokens in the given pgm

1) `int main() {`

$x = 4 + z;$; 11

$\text{int } x, y, z;$; 18

`printf("sum %d %d", x);`

{ 26 20 21 22 } 26

* 2) `main() {`

$a = b + + - - + + + = ;$; 13

`printf("%d %d, a, b);` ; 16

⇒ stop here only

3) `main() {`

`int a = 10;` ; 13

`char b = "abc";` ; 14

`int c = 30;` ; 18

`char arr[] = "xyz";` ; 26

`/* comment */` ; 22

m = 40.5; ; 28

{ 33 ; 29 30 31 32 }

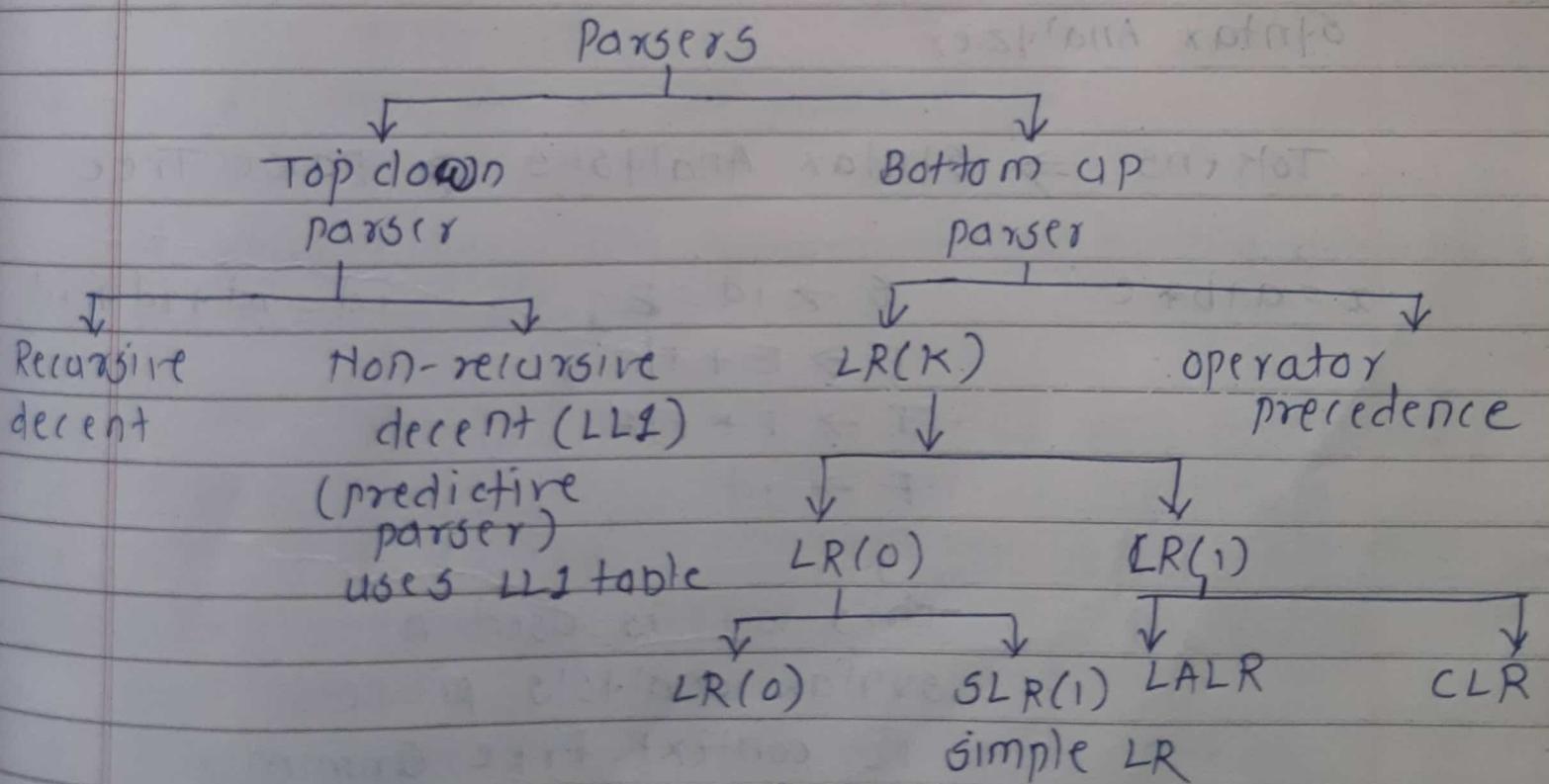
* 4) `printf("%i = %d, &i = %x", &i, &i);`

count ⇒ 10

* Parses

Parsing :

- The process of transforming the data from one format to another is called parsing. This process can be accomplished by the parser. The parser is a component of the translator that helps to organise linear text structure following the set of defined rules which is known as grammer.
- It is a process of deriving string from a given grammer.
- Stream of tokens given as a input to parser (ie to syntax analyzer).
- whether the token belong to proper grammar or not.
- It creates parse tree (syntax tree, derivation tree).
- It uses CFG (context free grammer).



- L2(1) - (2) left to right
 (2) left most derivation
 (1) look at one symbol

LR(0), LR(1) - find the canonical item

LR(0) - (2) - left to right scan
 (R) - right most derivation

LALR - look ahead LR

CLR - canonical LR

operator precedence - operates on unambiguous grammar & ambiguous grammar

* Formal Grammer

Syntax Analyzer

Tokens \Rightarrow Syntax Analysis \Rightarrow parse Tree

$$x = a + b * c \quad S \rightarrow id = E ; \quad id = id + id * id$$

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

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

$$F \rightarrow id$$

\uparrow
 this one is used in

Syntax analysis phase

i.e. context free Grammar

((CFG))

- CFG is a grammar of formal language.
- For the grammar we have set of non-terminals (variables) and terminals (constants).

non-terminals - noun, verb, adjective, verb phrase.

Terminals - constants, words in sentences.

$\langle S \rangle \rightarrow \langle \text{noun} \rangle \langle \text{verb phrase} \rangle$

$\langle \text{verb phrase} \rangle \rightarrow \langle \text{verb} \rangle \langle \text{adjective} \rangle$

$\langle \text{noun} \rangle \rightarrow \langle \text{noun-word} \rangle$

$\langle \text{verb} \rangle \rightarrow \langle \text{verb-word} \rangle$

$\langle \text{adjective} \rangle \rightarrow \langle \text{adjective-word} \rangle$

So, the formal grammar can be described by using 4 tuples $\langle N, T, P, S \rangle$

N - non-terminals

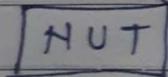
T - terminals

P - specific set of rules

S - sentence

A phrase structure grammar (or simply grammar) is (N, T, P, S) where

- 1) N is a finite, non-empty set of non-terminals.
- 2) T is a finite, non-empty set of terminals.
- 3) $N \cap T = \emptyset$ (Set of non-terminals & terminals don't have anything in common).
- 4) S is a special non-terminal (i.e. SEN), called start symbol.
- 5) P is a finite set whose elements are of the form, $\alpha \rightarrow \beta$



production rules

* must have at least one non-terminal.

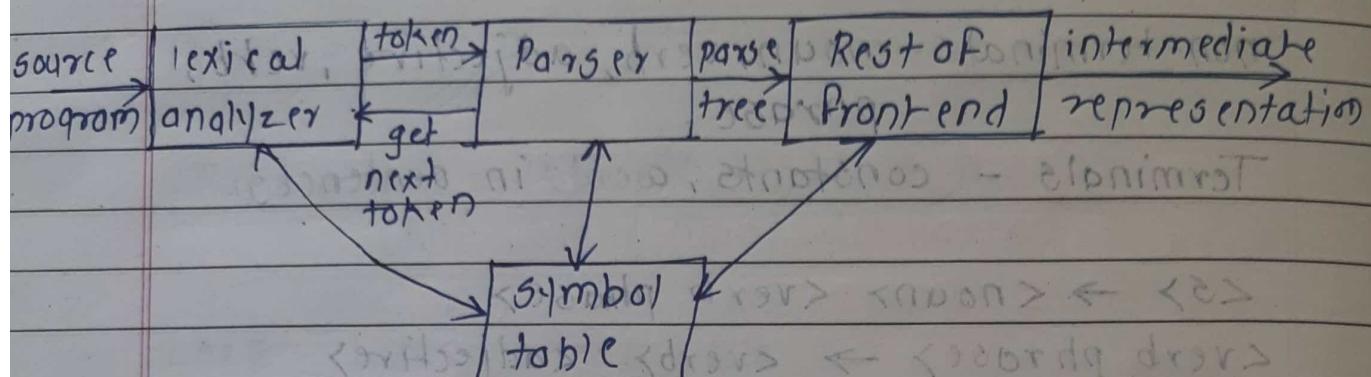


Fig: Position of parser in complex models

Context free Grammar

1. $E \rightarrow E + E$
2. $E \rightarrow E * E$
3. $E \rightarrow id$

Grammar can be defined as 5 tuples

(N, T, P, S)

N = E

T = +, *, id

P = all sentences

S = E (start symbol)

CFG (Type 2) can be written as $\Phi = T^* N$

Type 2 : i) $A \rightarrow \alpha \beta \gamma \dots \eta$ (A E H O S C I E P
 $\alpha \in V^*$)

$A \in N$	$\beta, \gamma, \dots, \eta \in V^*$
$\alpha \in V^*$	

$\alpha \in V^*$ - any combination of terminals & nonterminals including ϵ

Suppose from the above we want to derive $id + id * id$

left most derivation

$$E \rightarrow E + E$$

$$\rightarrow id + E$$

$$\rightarrow id + E * E$$

$$\rightarrow id + id * id$$

right most derivation

$$E \rightarrow E + E$$

$$\rightarrow E + E * E$$

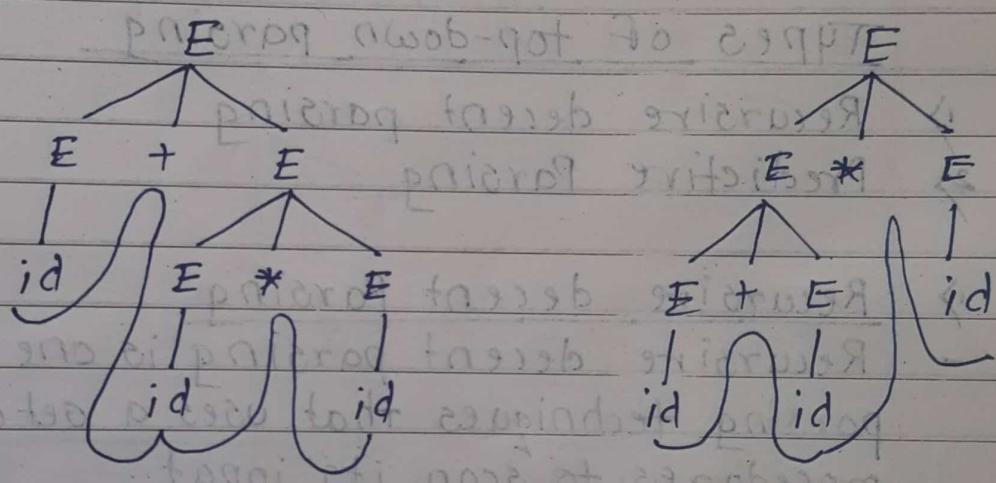
$$\rightarrow id + id * id$$

$$\rightarrow E + E * id$$

$$\rightarrow E + id * id$$

$$\rightarrow id + id * id$$

Parse tree derivation:



In the above example, if we obtain more than one left most derivation, more than one right most derivation, more than one parse tree derivation then the grammar is ambiguous.

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow id$$

$$id + id * id$$

$$: 1970$$

* Top-Down Parser

- A parser can start with the start symbol & try to transform it to input string.
- Example: LL Parsers

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

TYPES OF top-down parsing

- 1) Recursive decent parsing
- 2) Predictive Parsing

b) Recursive decent Parsing

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

Backtracking example:

consider the grammar $G: S \rightarrow CAD$

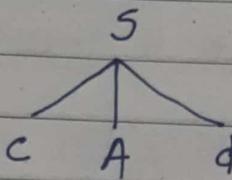
$$A \rightarrow ab|a$$

& the input string $w = cad$.

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

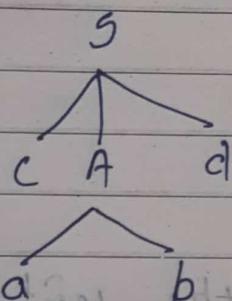
Step 1:

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



Step 2:

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



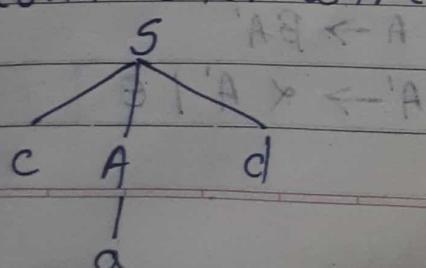
Step 3:

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

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

Step 4:

Now try the second alternative for A.



A left recursive grammar can cause a recursive decent parser to go into an infinite loop. Hence elimination of left recursion must be done before parsing.

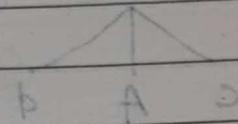
Example of recursive decent parsing:

consider the grammar for arithmetic expressions

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

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

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



After eliminating the left recursion the grammar becomes

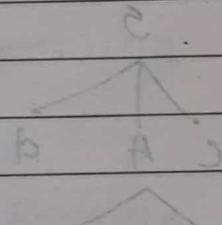
$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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



How to eliminate the left recursion from the Grammer?

A Grammer $G(V, T, P, S)$ is left recursive if it has a production in the form

$$A \rightarrow A\alpha \mid B$$

The above grammer is left recursive because the left of production is occurring at a first position on the right side of production. It can eliminate left recursion by replacing a pair of production with

$$A \rightarrow BA'$$

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

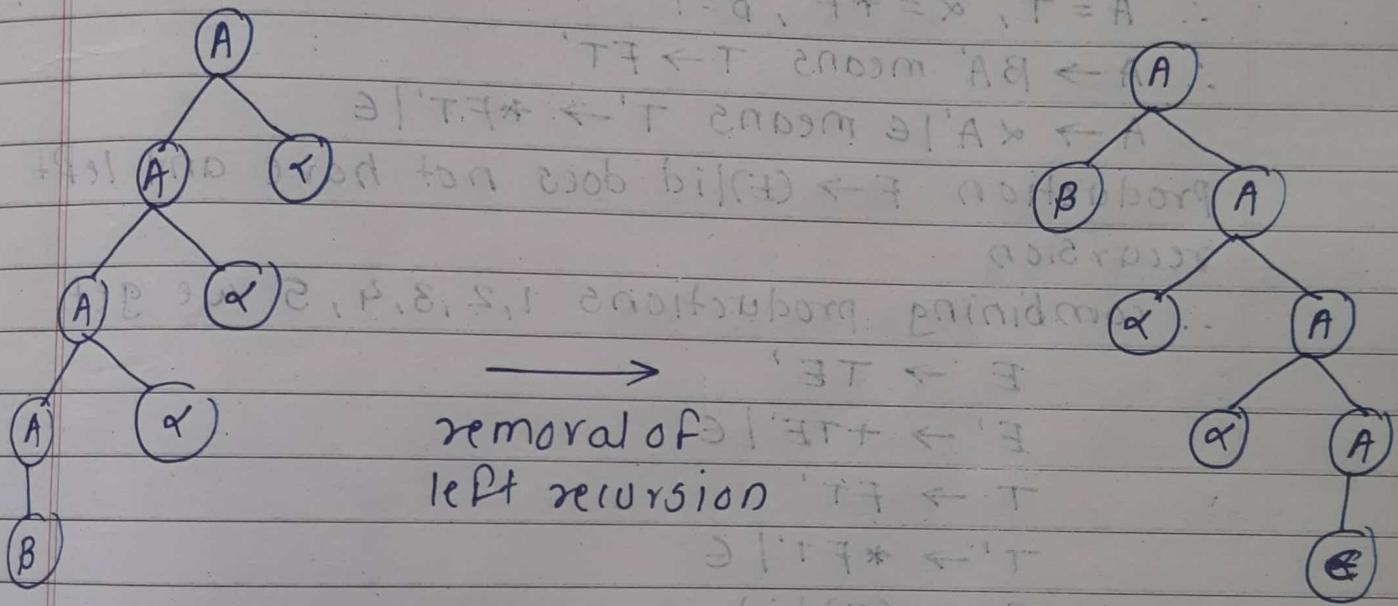
left recursion - A production of grammar is said to be left recursion if the leftmost variable of its RHS is same as variable of its LHS.

Elimination of left recursion

left recursion can be eliminated by introducing new non-terminal A such that

$$A \rightarrow A\alpha | B \xrightarrow{\text{removal of left recursion}} A \rightarrow BA'$$

Example In left recursive grammar, expansion of A will generate Aa, Aaa, Aaaa at each step, causing it to enter into an infinite loop



Example :

consider the left recursion from the grammar

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

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

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

solution :

comparing $E \rightarrow E + T | T$ with $A \rightarrow A\alpha | B$

E	\rightarrow	E	$+ T$	$ $	T
A	\rightarrow	A	α	$ $	B

$\therefore A = E, \alpha = +T, B = T \in F$ to minimize

$A \rightarrow A \alpha | B$ is change to

$A \rightarrow BA'$ &

$A' \rightarrow \alpha A' | \epsilon$

$\therefore A \rightarrow BA'$ means $E \rightarrow TE'$

$A' \rightarrow \alpha A' | \alpha$ means $E' \rightarrow +TE' | \epsilon$

comparing $T \rightarrow T * F | F$ with $A \rightarrow A \alpha | B$

$\therefore A = T, \alpha = *F, B = F$

$\therefore AA \rightarrow BA'$ means $T \rightarrow FT'$

$A \rightarrow \alpha A' | \epsilon$ means $T' \rightarrow *FT' | \epsilon$

Production $F \rightarrow (E) | id$ does not have any left recursion

\therefore combining productions 1, 2, 3, 4, 5 we get

$E \rightarrow TE' \quad \leftarrow$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) | id$

Example: most ambiguous \Rightarrow not obvious

$S \rightarrow a | 1 | (T)$

$T \rightarrow T, S | S$

solution :

$S \rightarrow a | 1 | (T)$

$S T \Rightarrow ST'$

$T' \rightarrow , ST' | \epsilon$

T	+	T+	E	→	E
S	!	X	A	↑	A

Example :

$$E \rightarrow E(T) | T$$

$$T \rightarrow T(F) | F$$

$$F \rightarrow id$$

Solution :

$$E \rightarrow TE'$$

$$E \rightarrow (T)E' | \epsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow (F)T' | G$$

$$F \rightarrow id$$

* left Factoring (left Recursion)

$$A \rightarrow \alpha B_1 | \alpha B_2 | \alpha B_3$$

(grammars with common prefixes)

- This kind of grammar creates a problematic situation for top-down parsers.
- Top down parsers can not decide which production must be chosen to parse the string in hand. To remove this confusion, we use left factoring.
- Left factoring is a process by which the grammar with common prefixes is transformed to make it useful for top down parsers.

In left factoring -

1) we make one production for each common prefixes.

2) The common prefix may be a terminal or a non-terminal or combination of both.

3) Rest of derivation is added by new productions.

Example :

$$P \rightarrow P + Q \mid Q$$

Here P is similar to generate non-terminal
so,

$$\begin{array}{l} P \rightarrow P + Q \mid Q \\ \downarrow \quad \downarrow \quad \downarrow \mid \downarrow \\ A \rightarrow A \mid B \end{array}$$

* so if we will observe the eq. ① its left recursive, so we can eliminate it

IF there is any production

$$\begin{array}{l} A \rightarrow \alpha B_1 \mid \alpha B_2, \text{ it can be rewritten as} \\ A \rightarrow \alpha A' \\ A' \rightarrow B_1 \mid B_2 \end{array}$$

Example :

$$G : S \rightarrow iEsi \mid iEsse \mid a$$

$$E \rightarrow b$$

* left factored, this grammar becomes

$$S \rightarrow iEsi \mid a$$

$$S' \rightarrow iE \mid es \mid e$$

$$E \rightarrow b$$

PAGE NO.	8958
DATE	/ /

Example :

$$A \rightarrow aAB/aBC/aAC$$

Left Factored - (Step 1)

$$A \rightarrow aA'$$

$$A' \rightarrow AB/BC/AC$$

Step 2 :

$$A \rightarrow aA'$$

$$A' \rightarrow AD/BC$$

$$D \rightarrow B/C$$

Example :

$$S \rightarrow bSSaaS/bSSaSb/bSb/a$$

Step 1 :

$$S \rightarrow bSS'/a$$

$$S' \rightarrow SaA/b$$

$$A \rightarrow aS/bS$$

Example :

$$S \rightarrow a/ab/abc/abcd$$

Step 1 :

$$S \rightarrow aS'$$

$$S' \rightarrow b/bc/bcd/\epsilon$$

Step 2 :

$$S \rightarrow aS'$$

$$S' \rightarrow bA/\epsilon$$

$$A \rightarrow c/cd/\epsilon$$

capital letter denotes variables
small letter denotes terminals

PAGE NO.	11
DATE	

Step 3:

$$S \rightarrow aS' \quad S' \rightarrow bA/bE \quad A \rightarrow cB/cE \quad B \rightarrow d/E$$

$$S' \rightarrow bA/bE \quad A \rightarrow cB/cE$$

$$A \rightarrow cB/cE$$

$$B \rightarrow d/E$$

$$S \rightarrow aS' \quad S' \rightarrow bA/bE \quad A \rightarrow cB/cE \quad B \rightarrow d/E$$

* Predictive Parser (LL1)

Steps:

- * 1. Remove left recursion, left factoring, ambiguity
- 2. Use first() & follow() method for all non-terminals
- 3. construct ^{predictive} Parsing table
- 4. Stack implementation
- 5. Parse tree generation the given input string using stack & parsing table.

* FIRST()

First & follow sets are needed so that the parser can properly apply the needed production rule at the correct position.

Example:

$$A \rightarrow abc/deP/ghi$$

Then we have,

$$\text{FIRST}(A) = \{a, d, g\}$$

Rules:

1. For a production rule $X \rightarrow \epsilon$,

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

2. For any terminal symbol 'a',

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

3. For production rule

$$X \rightarrow Y_1 Y_2 Y_3$$

Then,

calculating $\text{FIRST}(X)$

- If $\epsilon \notin \text{First}(Y_1)$, then $\text{First}(X) = \text{First}(Y_1)$
- If $\epsilon \in \text{First}(Y_1)$, then $\text{First}(X) = \{\text{First}(Y_1) - \epsilon\} \cup \text{First}(Y_2 Y_3)$

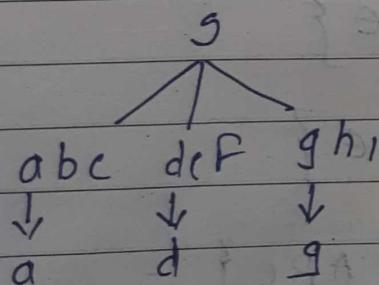
calculating $\text{First}(Y_2 Y_3)$

- If $\epsilon \notin \text{First}(Y_2)$, then $\text{First}(Y_2 Y_3) = \text{First}(Y_2)$
- If $\epsilon \in \text{First}(Y_2)$, then $\text{First}(Y_2 Y_3) = \{\text{First}(Y_2) - \epsilon\} \cup \text{First}(Y_3)$

- $\text{First}(A)$ contains all terminals present in first place of every string derived by A.

Examples :

1) $S \rightarrow abc \mid def \mid ghi$



$$\text{First}(S) = \{a, d, g\}$$

first (terminal) = terminal
First (ϵ) = ϵ

PAGE No.	/ / /
DATE	/ / /

2) $S \rightarrow ABC \mid ghi \mid jkl$

$A \rightarrow alb \mid c$

$B \rightarrow b$

$D \rightarrow d$

First(D) = {d}

First(B) = {b}

First(A) = {a, b, c}

first(S) = first{ABC} \rightarrow First{A} = a, b, c

\therefore First(S) = {a, b, c, g, j}

3) $S \rightarrow ABC$

$A \rightarrow alb \mid \epsilon$

$B \rightarrow clc \mid \epsilon$

$C \rightarrow cfc \mid \epsilon$

First(C) = {e, f, \epsilon}

First(B) = {c, d, \epsilon}

First(A) = {a, b, \epsilon}

First(S) = First{ABC} =

First{a, b, c, d, e, f, \epsilon}

$$\text{First}(S) = \{ (, a \}$$

$$\text{First}(L) = \text{First}(S) = \{ (, a \}$$

$$\text{First}(L') = \{ \epsilon,) \}$$

$$\text{Follow}(S) = \{ \$, (,),) \}$$

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

$$\text{Follow}(L') = \text{Follow}(L) = \{) \}$$

	($\$$	ϵ	a)	$\$$
S	$S \rightarrow (L)$	ϵ	ϵ	$S \rightarrow a$	ϵ	ϵ
L	$L \rightarrow SL'$	ϵ	ϵ	$L \rightarrow SL'$	ϵ	ϵ
L'		$L \rightarrow \epsilon$			$L' \rightarrow SL'$	ϵ

Fill the productions in parse table.

First production

$S \rightarrow (L)$ - calculate the First of right side. (First is ())

$S \rightarrow a$ - First is a

$L \rightarrow SL'$ - First(S) = {(, a)}

$L' \rightarrow \epsilon$ - $L' \rightarrow \epsilon$ - Follow(L') = {)}

$L' \rightarrow , SL'$ $\Rightarrow ,$

In the parse table if a column contains more than one entry then it is not a LL(1) grammar.

The above is LL(1) grammar.

\$ - represents the bottom of stack

E - represents start symbol

PAGE No.	/ / /
DATE	/ /

continue -

$$w = id * id + id$$

we can parse this string by using above predictive parsing table

stack implementation

stack

read E

From stack

id

From input

parsing table

so we get

get

E → TE'

we get

E → T'F

so we get

E → T'id

E → T'*

E → T+F*

E → T'F*

E → T'id

E → T'

E →

\$ E

\$ E'TC

\$ E'T'F

\$ E'T'id

\$ E'T'

\$ E'T'F*

\$ E'T'F

\$ E'T'id

\$ E'T'

\$ E'T+F*

\$ E'T'F*

\$ E'T'id

\$ E'T'

\$ E'

\$

id * id + id \$

id + id \$

id + id \$

E → TE'

T → FT'

F → id

T' → F FT'

F → id

T' → e

T' → e

E' → E

so we get E → TE' & E → T'F

- \$ represents the bottom of stack & E is the start symbol

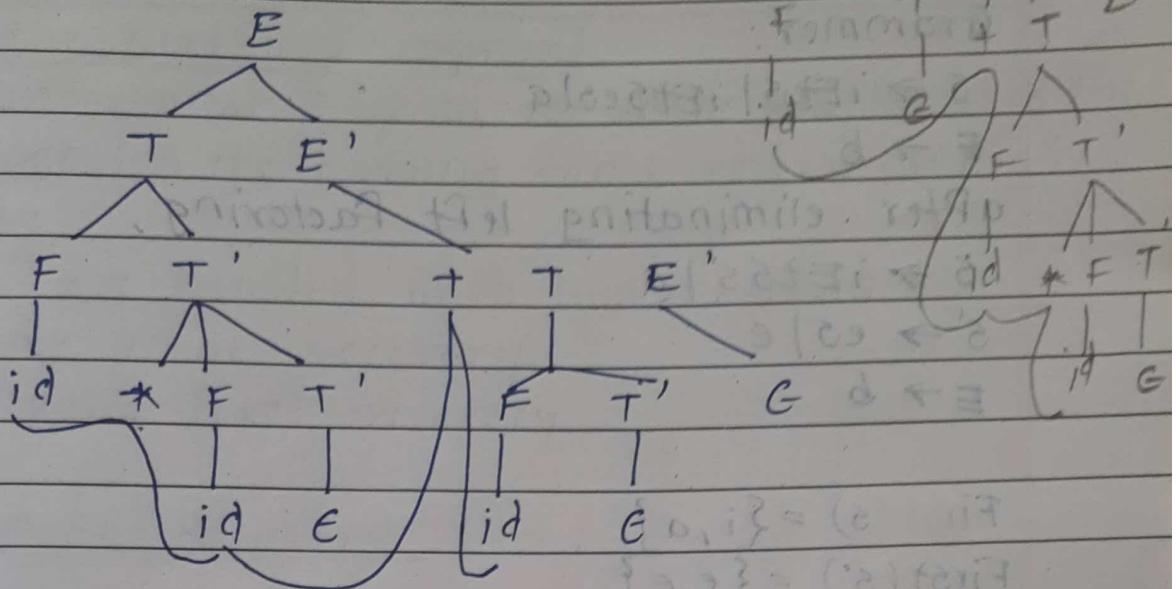
- w (ie input string) followed by \$

- if stack contains nonterminal element & input contains terminal symbol then refer predictive parsing table

- so E → TE' that production is return in O/P & instead we can replace the E by its production ie E → TE

in reverse order ie E' T in stack.

Parse Tree :



From parse Tree : $id * id + id$

- Non-terminal on top of stack & input contains terminal symbol then we have to refer to predictive parsing table.
- E & id gives some production from table that will be written to output
- How stack is replace with output from reverse reverse order.
- if terminal symbol on top of stack & that match with top of input then, we pop the terminal symbol from top of stack & remove the terminal from input buffer.

→ Now $T id \rightarrow$ so in output $T \rightarrow FT'$ & in stack $E' T' F$

- Now $F id \rightarrow$ so $F \rightarrow id \rightarrow$ so in output $F \rightarrow id$ & in stack pop F & put id.
- Now stack & input both containing terminal symbol ie id so remove from both

* LL(1) grammar Example

: 9/9/2019

Grammer:

$$S \rightarrow iES | iESe | a$$

$$E \rightarrow b$$

After eliminating left factoring,

$$S \rightarrow iESe | a$$

$$S' \rightarrow eS | e$$

$$E \rightarrow b$$

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

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

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

$$\text{Follow}(S) = \{b, e\}$$

$$\text{Follow}(S') = \{b, e\}$$

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

Paring tables

Non-terminal	i	a	b	e	S'	E
S	$i \rightarrow S$	$a \rightarrow S$	$b \rightarrow S$	$e \rightarrow S$	$S' \rightarrow eS$	$E \rightarrow b$
S'				$e \rightarrow S$	$S' \rightarrow e$	
E			$b \rightarrow E$			

since there are more than one production, the grammer is not LL(1) grammer.

* Error recovery in Predictive Parsing

(PPT)

- An error is detected during predictive parsing when
- The terminal on the top of the stack does not match the next input symbol.
- When non-terminal A is on top of the stack, a is the next symbol & the parsing table entry $m(A, a)$ is empty.
- Errors are handled with two methods
 - Panic mode recovery
 - Phrase level recovery

Panic Mode Recovery

- Skip symbols on the input until a token in a selected set of synchronizing tokens appear instead of follow we can place synch word in parsing table to indicate synchronizing tokens obtained from the Follow set of non-terminals.

Rules :

1. IF parse looks up entry $m[A, a]$ & finds it blank then the input symbol a is skipped.
2. IF the entry is synch then the nonterminal on top of stack is popped in an attempt to resume parsing.
3. If a token on the top of the stack does not match the input symbol, then we pop the token from the stack.

Example :

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

	First	Follow
F	{(, .id}	{+, *,), \$}
T	{(, .id}	{+,), \$}
E	{(, .id}	{), \$}
E'	{+, E}	{), \$}
T'	{*, E}	{+,), \$}

Non

Input symbol

Terminal	id	+	*	()	?	\$
E	$E \rightarrow TE'$				$E \rightarrow TE'$	synch	synch
E'		$E' \rightarrow +TE'$				$E' \rightarrow E$	$E' \rightarrow E$
T	$T \rightarrow FT'$	synch			$T \rightarrow FT'$	synch	synch
T'		$T' \rightarrow E$	$T' \rightarrow *FT'$			$T' \rightarrow E$	$T' \rightarrow E$
F	$F \rightarrow id$	synch	synch	$F \rightarrow (E)$		synch	synch

- Follow(E) is {(), \$} so put synch in parsing table.

stack	Input	Action
\$ E	id * + id \$	error, skip)
E TE' from \$ E	id * + id \$	id is in FIRST(E)
\$ TE' *	id * + id \$	(OK)
FT'E'	id * + id \$	
^{id & id match} so pop id id T'E'	id * + id \$	(OK)
T'E'	* + id \$	
^{+ & + match} so pop * FT'E'	* + id \$	(OK)
FT'E'	+ id \$	Error, M(F, +) = synch
T'E' \$	business + id \$	F has been popped
^{TE'} more to id \$		
TE'	id \$	
		like wise it goes continue on checking error

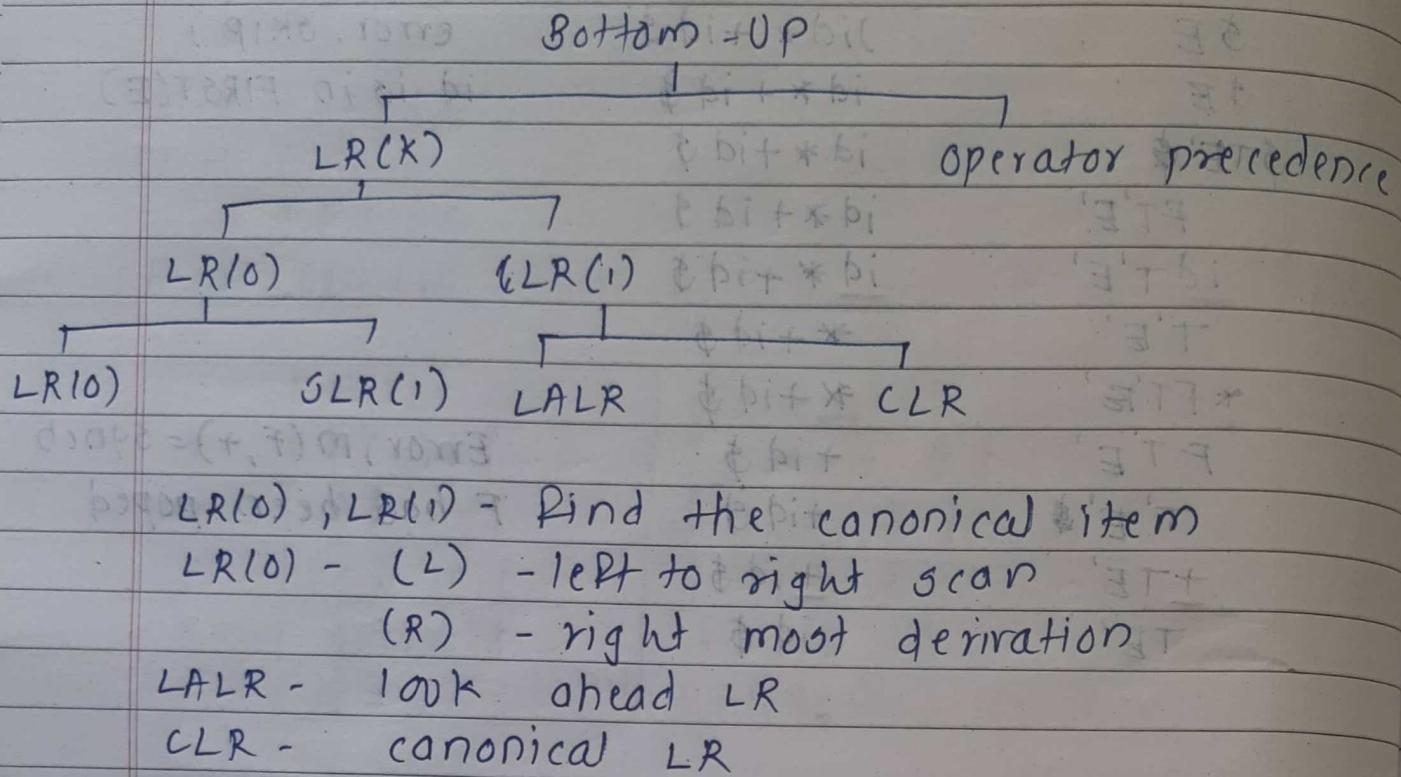
- Panic mode recovery does not address the important issue of error messages.

Phrase Level Recovery

- This involves, defining the blank entries in the table with pointers to some error routines which may
 - change, delete or insert symbols in the iIP
 - May also pop symbols from the stack.

whatever may be the action, care should be taken that it should not lead the parser into infinite loop.

* Bottom - Up Parsing



$LR(0), LR(1)$ - Find the canonical item

$LR(0)$ - (L) - left to right scan

(R) - right most derivation

LALR - look ahead LR

CLR - canonical LR

- Bottom up parse corresponds to the construction of a parse tree for an input string beginning at the leaves & working up towards the root (the top).

$id * id \quad F * id \quad T * id \quad T * F$

$\quad \quad | \quad \quad | \quad \quad | \quad \quad |$

id

$F \quad op \quad F \quad op \quad id$

$| \quad \quad | \quad \quad | \quad \quad |$

id

$id \quad id$

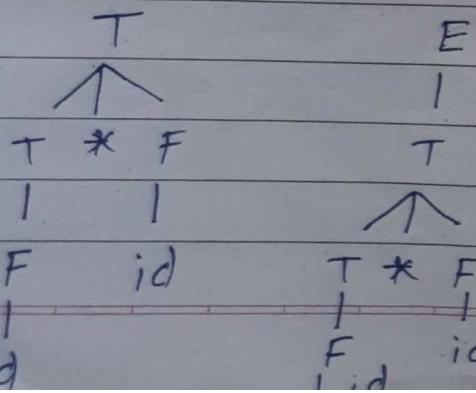


Fig. A bottom-up
parse for $id * id$

- bottom-up parsing as the process of reducing a string w to the start symbol of the grammar. At each reduction step, a specific substring matching the body of a production is replaced by the nonterminal at the head of that production.
- The key decisions during bottom-up parsing are about when to reduce & about what production to apply, as the parse proceeds.

$id \neq id, F \neq id, T \neq id, T \neq F, T, E$

- The goal of bottom up parsing is therefore to construct a derivation in reverse. The following derivation corresponds to the parse in

$$E \Rightarrow T \Rightarrow T * F \Rightarrow T * id \Rightarrow F * id \Rightarrow id \neq id$$

- this derivation is in fact a rightmost derivation.

* Handle Pruning

- Bottom up parsing during a left to right scan of the input constructs a right most derivation in reverse.
- A handle is a substring that matches the body of a production & whose reduction represents one step along the reverse of rightmost derivation.

Right sentential Form	Handle	Reducing production
$id_1 * id_2$	id_1	$F \rightarrow id$
$F * id_2$	F	$T \rightarrow F$
$T * id_2$	id_2	$F \rightarrow id$
$T * F$	$T * F$	$E \rightarrow T * F$

Examples:

- * 1. For the grammar $S \rightarrow 0S1|01$
indicate the handle in each of the
following right sentential forms:
 a) 000111
 b) 00S11
- 2) For the grammar $S \rightarrow SS + 1SS * 1a$
indicate the handle in each of the
following right sentential forms.
 a) SSS+a*+
 b) SS+a*a+a+
 c) aaa*a*a++
- 3) Give bottom up parses for the following
input strings & grammars
 a) The input 000111 for grammar $S \rightarrow 0S1|01$
 b) The input aaa*a+a+ for grammar
 $S \rightarrow SS + 1SS * 1a$

* Shift Reduce Parsing

- It is a type of bottom up parsing

Stack

\$

\$ s

Input

w \$

\$

Actions :

- 1) Shift : Parser shifts 0 or more i/p symbols until handle B appears on top of stack.
- 2) Reduce : B is reduced to left hand side of the production
- 3) Accept : announces successful compilation.
- 4) Error : calls an error recovery routine.

$A \rightarrow B + C D$ $C \rightarrow a b c$ with handle cab on top of stack. When handle on top of stack it must be reduced

Example :

$S \rightarrow CC$

$C \rightarrow cCd$

i/p string w: cdcd

stack

\$

\$ c

\$ cd

\$ CC

\$ C

\$ Cc

\$ Ccd

\$ Ccc

input

cdcd \$

cd \$

cd \$

cd \$

d \$

\$

\$

\$

Action

shift

shift

reduce by $C \rightarrow d$

reduce by $C \rightarrow CC$

shift

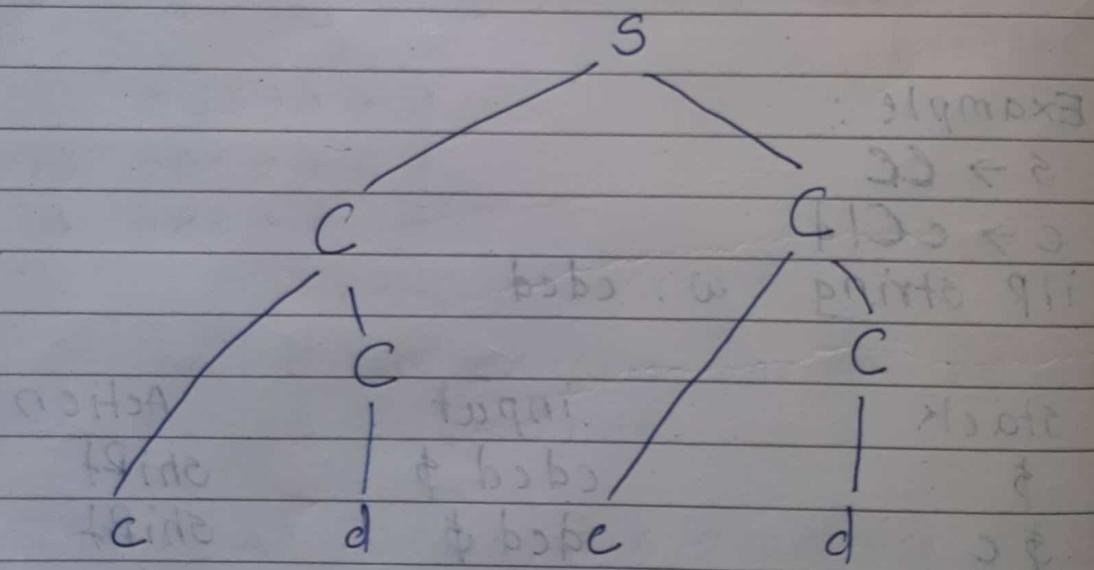
shift

reduce by $C \rightarrow d$

reduce $C \rightarrow CC$

- c shifted on stack
 - check whether c is handle or not (ie to check whether c is the right side of any production)
 - we can reduce if we will found the handle in production rule.

\$ CC reduce by $S \rightarrow C$
\$ Samos accept
How,
LLP is cdccl, How we can construct the
parse tree in bottom-up manner.



PAGE NO.	
DATE	11

Examples:

consider the grammar & perform the shift reduce parsing for the given i/p string

$$1) S \rightarrow S + S$$

$$S \rightarrow S * S$$

$$S \rightarrow id$$

$$i/p \Rightarrow id + id + id$$

$$2) E \rightarrow 2E2$$

$$E \rightarrow 3E3$$

$$E \rightarrow 4$$

$$i/p \Rightarrow 32423$$

$$3) S \rightarrow (L)1q$$

$$L \rightarrow L, S | S$$

$$i/p \Rightarrow (a, (a, a))$$

$$4) S \rightarrow TL$$

$$T \rightarrow int | float$$

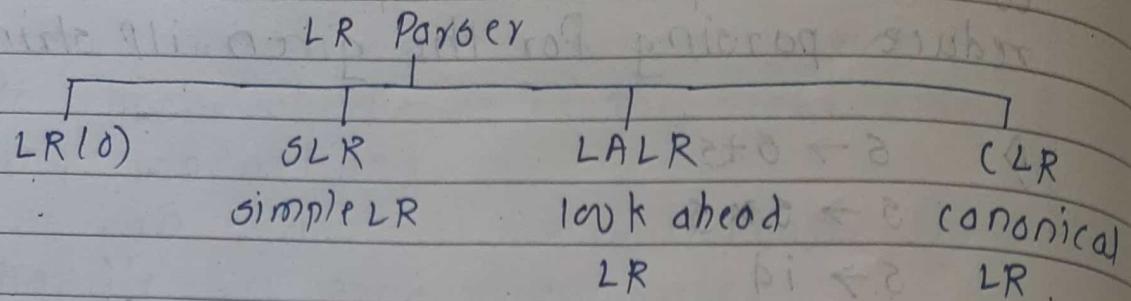
$$L \rightarrow L, id | id$$

$$i/p \Rightarrow int id . id$$

$$5) S \Rightarrow 0S0 | 1S1 | 2$$

$$i/p \Rightarrow 10201$$

* LR parsers



- An efficient bottom up syntax analysis technique that can be used to parse a large class of CFG is called LR(k) parsing.
 - L - left to right scanning
 - R - constructing rightmost derivation in reverse
 - K - no. of input symbols.
when K is omitted it is assumed to be 1.

Advantages :

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

DisAdvantages :

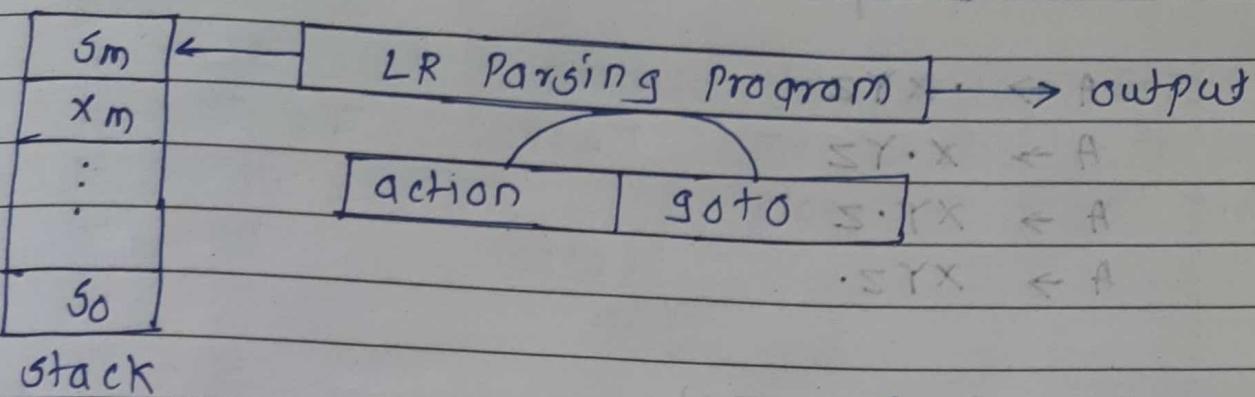
It is too much of work to construct a LR parser by hand for a programming language grammar. A specialized tool, called a LR

parser generator, is needed. Example YACC.

* LR Parsing Algorithm

Input

$[a_1 \mid a_2 \mid \dots \mid a_i \mid \dots \mid a_n \mid \$]$



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

Action : The parsing pgm determines s_m , the state currently on top of stack & a_i , the current iIP symbol. It then constructs action $[s_m, a_i]$ in the action table which can have one of four values :

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

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

LR(0) items:

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

$$A \rightarrow \cdot XYZ$$

$$A \rightarrow X \cdot YZ$$

$$A \rightarrow XY \cdot Z$$

$$A \rightarrow XYZ \cdot$$

* LR(0) Parsing Table

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

$$T \rightarrow id$$

(3)

- To make LR(0) parsing table we use LR(0) canonical items.

- How to make canonical item.
For a given grammar make it augmented grammar.

For eg.

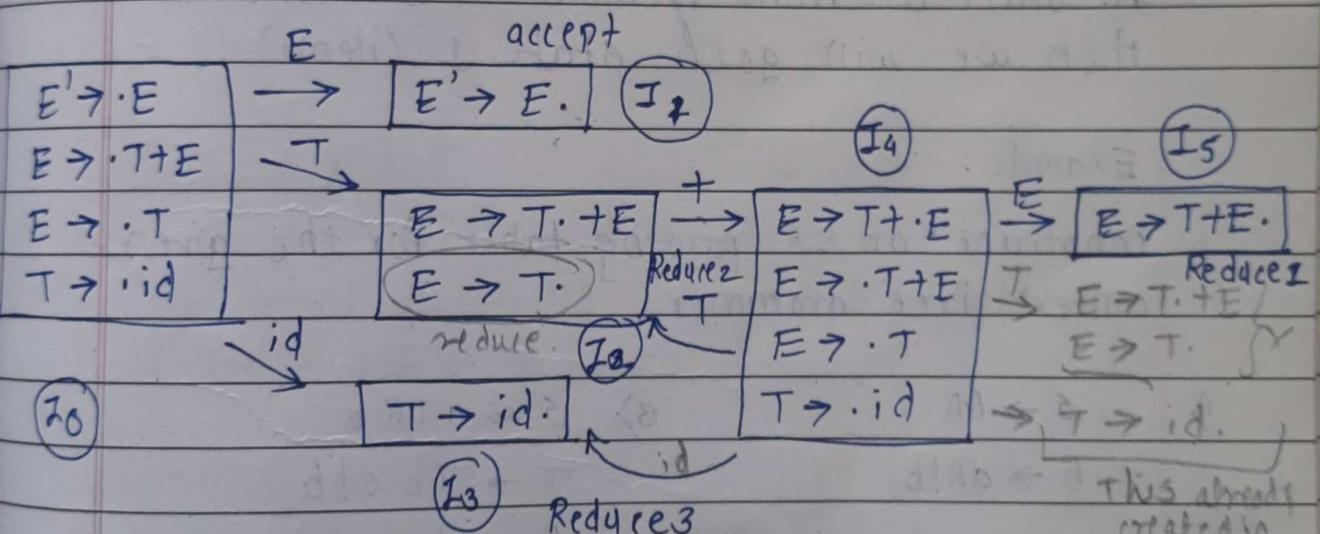
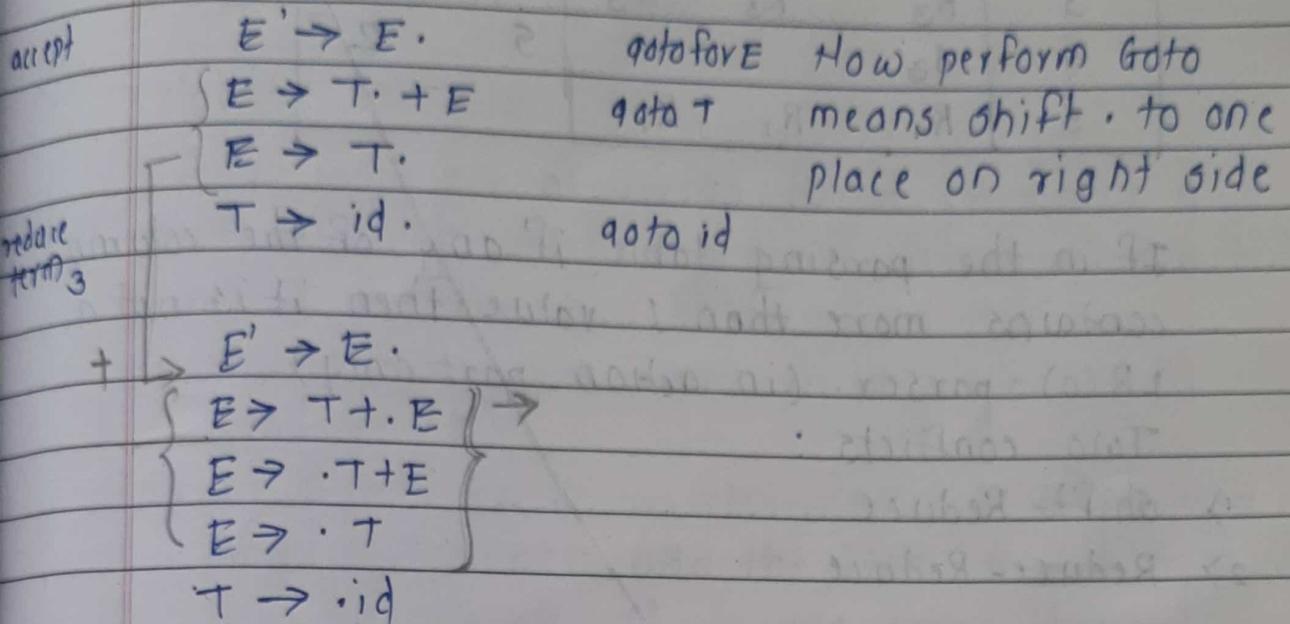
$$E' \rightarrow \cdot E$$

Meaning of $\cdot \Rightarrow$ basically we are using bottom-up technique so \cdot represents what we have seen & what we have not seen.

For eg. $\cdot - id + id \cdot id$ means parser had check with $id +$ & remain to parse id which is on right hand side.

item ₀	$E \rightarrow \cdot E$
item ₁	$E \rightarrow \cdot T + E$
	$E \rightarrow \cdot T$
	$T \rightarrow \cdot id$

→ now we will write the production of E from grammar & before that put a \cdot symbol (closure property)



where, $I_0, I_1, I_2, I_3, I_4, I_5 \Rightarrow$ are item₀, item₁, item₂, item₃, item₄, item₅ to

Fig: canonical LR(0) collection for a given grammar.

state		Action	GOTO
0	S_3	$\text{id} + \$$	E $\rightarrow T$
1	-	- Accept	$b \rightarrow T$
2	$R_2 \text{ or } S_4 / R_2$	R_2	
3	R_3	R_3	
4	S_3	empty	$S \rightarrow 2$
5	$R_1 \text{ or } R$, $R_2 \text{ or } P$		$E \rightarrow T$

* If in the parsing table if any of the columns contains more than 1 value then it is not a LR(0) parser. (in action part only).

Two conflicts :

- 1) Shift- Reduce
- 2) Reduce- Reduce

In state 0 from item 0 if we consider more then we will goto state 1. (item 1)

Example :

Construct an LR parsing table for the given context free grammar

$$\begin{array}{ll} 1) S \rightarrow AA \quad b \rightarrow T \\ A \rightarrow aA \mid b & \end{array} \quad \begin{array}{l} 3) S \rightarrow aTRe \\ T \rightarrow Tb \mid c \mid b \\ R \rightarrow d \end{array}$$

$$2) E \rightarrow E + T \quad E \rightarrow T \quad E \rightarrow F$$

$$4) S \rightarrow L = R \mid R \quad T \rightarrow T * F \quad L \rightarrow * R lid$$

$$T \rightarrow F \quad F \rightarrow (E) \quad R \rightarrow L ason : p is$$

$$F \rightarrow id$$

* SLR(1)

1. Find LR(0) items.
 2. completing the closure.
 3. compute $\text{goto}(I, x)$, where, I is set of items & x is grammar symbol.
- closure operation:

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

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

Goto Operation:

$\text{Goto}(I, x)$ is defined to be the closure of the set of all items $[A \rightarrow \alpha \cdot x \cdot B]$ such that $[A \rightarrow \alpha \cdot x \cdot B]$ is in I .

continue with previous example of LR(0), \Rightarrow

- same till the items are found.

state	id	+	\$	Goto
0	s_3			1 2
1				Accept
2		s_4	R_2	
3		R_3	R_3	5 2
4	s_3			R_1
5				

- Goto for terminals is as it is from the LR(0) parsing table.
- Shift & accept is as it is.
- For reduce =, r_2 written in complete row. In SLR we can't write reduction in incomplete row, we can find the Follow of left hand side.

For r_2 reduction we

will go to production

$$E \rightarrow T + E \quad (2)$$

$$E \rightarrow \cdot T$$

so now calculate Follow(E)

& write r_2 reduction in Follow(E)

so,

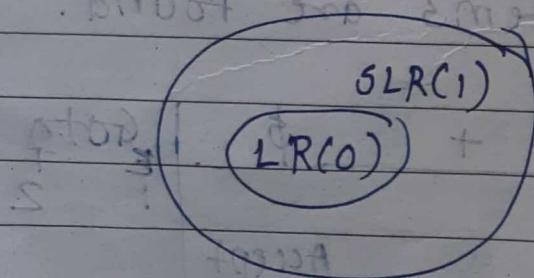
$$\text{Follow}(E) \Rightarrow \$$$

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

so, at place of \$ we can write r_2

$$R_3 \text{ For } \text{Follow}(T) \Rightarrow +, \$$$

so, each column value contains either shift or reduce so, SLR reduces the conflicts of shift Reduce so this grammar is SLR(1)



- The G which is accepted by LR(0) is also accepted by SLR but the G which is not

(DAS) < RAS
 (bottom left 660) > RAS
 accepted by LR(0) is also accepted by SLR, so
 the SLR is more powerful than LR.

- * construct SLR parsing for the following grammar:

$$G : E \rightarrow E + T / T$$

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

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

also shows the action performed for
 $id + id \Rightarrow id$

- * show that the following grammar

$$S \rightarrow AaAb / BbBa$$

$$A \rightarrow G$$

$$B \rightarrow e$$

is LL(1) but not SLR(1)

- * $G : S \rightarrow SAIA$

$$A \rightarrow a$$

is SLR(1) but not LL(1)

- * The following is an ambiguous grammar:

$$S \rightarrow AS / b$$

$$A \rightarrow SA / a$$

construct for this grammar its collection of sets of LR(0) items. If we try to build an LR-parsing table for the grammar, there are certain conflicting actions. What are they? Suppose we tried to use the parsing table by nondeterministically choosing a possible action whenever there is a conflict. Show all the possible sequences of actions on input abab.

LALR > LR(1)
 CLR (add lookahead symbol)

* Canonical LR Parsing tables

- CLR uses LR1 canonical item
- makes full use of the lookahead symbols.

1. First Augment the start symbol, for start symbol the lookahead symbol is added as \$
2. Open the variable i.e. nonterminal symbol.
3. Then add first of that variable to remaining opening production so in below eg. first of S is \$ so add \$ as a lookahead symbol to remaining productions.

$G:$

$$S \rightarrow \overset{①}{aAd} \mid \overset{②}{bBd} \mid \overset{③}{aBe} \mid \overset{④}{bAe}$$

$$A \rightarrow c$$

$$B \rightarrow c$$

$$\begin{aligned} S' &\rightarrow \cdot S, \$ \\ S &\rightarrow \cdot aAd \\ S &\rightarrow \cdot bBd \end{aligned}$$

$$S \rightarrow \cdot aBe$$

$$S \rightarrow \cdot bAe$$

so, Goto (J_1, S)

$S' \rightarrow \cdot S, \$$	$\xrightarrow{S} S \rightarrow S \cdot, \$$	* Accept
$S \rightarrow \cdot aAd, \$$	$\xrightarrow{a} S \rightarrow a \cdot Ad, \$$	
$S \rightarrow \cdot bBd, \$$	$\xrightarrow{} S \rightarrow a \cdot Be, \$$	
$S \rightarrow \cdot aBe, \$$	$A \rightarrow \cdot c, d$	(First of $A \Rightarrow d$)
$S \rightarrow \cdot bAe, \$$	$B \rightarrow \cdot c, e$	(First of $B \Rightarrow c$)

\xrightarrow{b}

$$\begin{aligned} S &\rightarrow b \cdot Bd, \$ \\ S &\rightarrow b \cdot Ae, \$ \\ B &\rightarrow \cdot c, d \\ A &\rightarrow \cdot c, e \end{aligned}$$

Go to S' , \Rightarrow * - indicates reduce

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

$S \rightarrow a \cdot Ad, \$$	$\xrightarrow{A} (S \rightarrow a A \cdot d, \$)$	$\xrightarrow{d} S \rightarrow a Ad \cdot, \$$
$S \rightarrow a \cdot Be, \$$	$\xrightarrow{B} (S \rightarrow a B \cdot e, \$)$	$\xrightarrow{e} S \rightarrow a Be \cdot, \$$
$A \rightarrow \cdot c, d$	$\xrightarrow{c} A \rightarrow c \cdot, d$	
$B \rightarrow \cdot c, e$	$\xrightarrow{B} c \cdot, e$	*

$S \rightarrow b \cdot Bd, \$$	$\xrightarrow{B} (S \rightarrow b B \cdot d, \$)$	$\xrightarrow{d} S \rightarrow b Bd \cdot, \$$
$S \rightarrow b \cdot Ae, \$$	$\xrightarrow{A} (S \rightarrow b A \cdot e, \$)$	$\xrightarrow{e} S \rightarrow b Ae \cdot, \$$
$B \rightarrow \cdot c, d$	$\xrightarrow{c} B \rightarrow c \cdot, d$	
$A \rightarrow \cdot c, e$	$\xrightarrow{A} c \cdot, e$	*

Total \Rightarrow 14 items.

Parsing table of CLR \Rightarrow same as the construction of SLR

- only the difference the reduce will written in at lookahead symbol.
- It is CLR grammar

* LALR (lookahead LR)

- Follows LR(0) canonical items.
- LALR is less powerful than CLR.
- A grammar where the LALR parsing table has no multiply represented entries are said to be LALR(1) or LALR grammar.

working of LALR parser

context Free grammar



construction of set of LR(0) items with look ahead

- (a) construct Augmented Grammar
- (b) Find closure
- (c) Find goto

Merge the items or states with same first component but different look ahead

construct Parsing Table



Parse the input string

String accepted String not accepted

Example :

- 1) $S \rightarrow CC$
- 2) $C \rightarrow CC$
- 3) $C \rightarrow d$

Augmented Grammar

$S' \rightarrow S$

$S \rightarrow .CC$

$C \rightarrow .CC$

$C \rightarrow .d$

LR(1) Items

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

$\text{Goto}(I_2, C)$

$S \rightarrow \cdot CC, \$$

$I_5 : C \rightarrow \cdot d, cld$

$C \rightarrow \cdot CC, cld$

$S \rightarrow CC \cdot, \$$

$C \rightarrow \cdot d, cld$

$\text{Goto}(I_0, S)$

$\text{Goto}(I_2, C)$

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

$I_6 : C \rightarrow C \cdot C, \$$

$C \rightarrow \cdot CC, \$$

$C \rightarrow \cdot d, \$$

$\text{Goto}(I_0, C)$

$I_2 : S \rightarrow CC, \$$

$\text{Goto}(I_4, d)$

$C \rightarrow \cdot CC, \$$

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

$C \rightarrow \cdot d, \$$

$\text{Goto}(I_0, C)$

$\text{Goto}(I_5, C)$

$I_3 : C \rightarrow C \cdot C, cld$

$I_8 : C \rightarrow CC \cdot, cld$

$C \rightarrow \cdot CC, cld$

$\text{Goto}(I_6, C)$

$C \rightarrow \cdot d, cld$

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

$\text{Goto}(I_0, d)$

$I_{14} : C \rightarrow d \cdot, cld$

- In this we check the code items are common & look ahead is different, if code items are common then we can combine the code items & merge the look ahead (or perform union op. with look ahead symbol).

- As we have observe the above grammar the I_3 & I_6 are common so we combine I_3 & I_6

~~I_{3+6}~~ : $C \rightarrow C \cdot C, C/d/\$$
 $C \rightarrow \cdot CC, C/d/\$$
 $C \rightarrow \cdot d, C/d/\$$

~~I_{4+7}~~ , $I_{47} : C \rightarrow d \cdot, C/d/\$$

$I_{8g} : C \rightarrow CC \cdot, C/d/\$$

LALR parsing table :

- Instead of I_3 we can write I_{36} similarly instead of I_4 write I_{47} & I_8 write I_{8g}
- From CLR parsing table we can construct LALR parsing table.

CLR Parsing table :

				goto
	C	d	\$	5 6
0	S_3	S_4		1 2
1			accept	
2	S_6	S_7		5
3	S_3	S_4		8
4	R_3	R_3		
5			R_1	
6	S_6	S_7		9
7			R_3	
8	R_2	R_2		
9			R_2	

- From CLR we can construct LALR parsing table
- From CLR 0 - c \Rightarrow S₃ so in LALR

page no.	write S ₃ /6
DATE	1/1

LALR Parsing Table

	C	d	\$	go to non terminal	
	terminal action			S	C
0	S ₃ /6	S ₄ /7		1	2
1			accept		
2	S ₃ /6	S ₄ /7			5
36	S ₃ /6	S ₄ /7			89
47	r ₃	r ₃	r ₃		
5			r ₁		
89	r ₂	r ₂	r ₂		

All entries are single entries so its LALR grammar.

Examples:

1. show that the following grammar

$$S \rightarrow Aa \mid bAc \mid dc \mid bda$$

$$A \rightarrow d$$

is LALR(1) but not SLR(1).

2. show that the following grammar

$$S \rightarrow Aa \mid bAc \mid Bc \mid bBa$$

$$A \rightarrow d$$

$$B \rightarrow d$$

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

* Operator Precedence Parser

- An efficient way of constructing shift reduce parser is called operator precedence parsing
- Operator precedence parser can be constructed from a grammar called operator grammar.
- There are three disjoint precedence relations namely
 - $<$ - less than
 - $=$ - equal to
 - $>$ - greater than
- The relations give the following meaning:
 - $a < b$ - a yields precedence to b
 - $a = b$ - a has same precedence as b
 - $a > b$ - a takes precedence over b .

1) Stack Input
 \$ w \$

2) Terminal on top of the stack is a symbol pointed to by input pointer is b .

(i) $a < b$ or $a = b$

shift

(ii) if $a > b$

pop terminal symbol from stack until top stack terminal is $<$ to top of stack the terminal most recently popped.

Reduce

- if blank entries in relation table then syntax error occur in grammar

PAGE NO.	/ /
DATE	

$$S \rightarrow a \mid \uparrow \mid (T)$$

$$T \rightarrow T, S \mid S$$

operator precedence relation table

	a	\uparrow	()	,	\$
a				>	>	>
\uparrow				>	>	>
(<	<	<	=	<	
)				>	>	>
,	<	<	<	>	>	
\$	<	<	<			Accept

stack

\$

\$ (

\$ (a

Input

(a, a) \$

a, a) \$

, a) \$

Action

\$ < (, shift

(< a, shift)

a > , pop a

once pop a now
terminal symbol
on top of stack is
< so we check
the relation b/w
top of stack &
most recently pop
symbol ie (< a
so we stop popping
until relation is <
so a is handled now
a is right hand
side of production
so, S \Rightarrow a so its
to reduce

reduce by S \Rightarrow a
reduce T \Rightarrow S

\$ (S

\$ (T

S (T,

S (T, a

S (T, S

, a) \$

a) \$

) \$

) \$

(< , shift

, < a shift

a > , pop a

reduces S \Rightarrow a

) > , pop T, S

If we pop the terminal symbol & if there is non-terminal symbol also then pop both i.e. terminal & nonterminal symbols.

PAGE NO. _____
DATE 15/1

$\$ (T)$	$) \$$	$(=)$ shift
$\$ (T)$	$\$$	$) > \$$ POP T
$\$$	$\$$	now top of stack is $($ & current top is T so relation $(> T$ so reduce $S \rightarrow CT$)
$\$ S$	$\$$	

* Example:

Table : operator-precedence relations

	$+$	$-$	$*$	$/$	\uparrow	id	$($	$)$	$\$$
$+$	$>$	$>$	$<$	$<$	$<$	$<$	$<$	$>$	$>$
$-$	$>$	$>$	$<$	$<$	$<$	$<$	$<$	$>$	$>$
$*$	$>$	$>$	$>$	$>$	$<$	$<$	$<$	$>$	$>$
$/$	$>$	$>$	$>$	$>$	$<$	$<$	$<$	$>$	$>$
\uparrow	$>$	$>$	$>$	$>$	$<$	$<$	$<$	$>$	$>$
id	$>$	$>$	$>$	$>$	$>$			$>$	$>$
$($	$<$	$<$	$<$	$<$	$<$	$<$	$<$	$=$	
$)$	$>$	$>$	$>$	$>$	$>$			$>$	$>$
$\$$	$<$	$<$	$<$	$<$	$<$	$<$	$<$		

G :

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

i/P $\Rightarrow id + id * id$

Rules :

Rule 1 :

- If precedence of b is higher than precedence of a, then we define $a < b$.
- If precedence of b is same as precedence of a, then we define $a = b$.
- If precedence of b is lower than precedence of a then we define $a > b$.

- * If operators θ_1 & θ_2 are of equal precedence then make
 - $\theta_1 > \theta_2$ & $\theta_2 > \theta_1$, if operators are left associative
 - $\theta_1 < \theta_2$ & $\theta_2 < \theta_1$, if right associative

Rule 2 :

- An identifier is always given the higher precedence than any other symbol.
- \$ symbol is always given the lowest precedence.

Rule 3 :

- IF two operators have the same precedence, then we go by checking their associativity.

Example : $E \rightarrow EAElid$

$A \rightarrow + * id$

equivalent operator precedence grammar is

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

iIP $\Rightarrow id + id * id$

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

$$2) S \rightarrow (L) \mid a$$

$$L \rightarrow L, S \mid S$$

$$iIP \Rightarrow (a, (a, a))$$

	a	()	,	\$
a	>	>	>	>	>
(<	>	>	>	>
)	<	>	>	>	>
,	<	<	>	>	>
\$	<	<	<	<	

* Semantic Analysis:

Type checking, symbol table management, symbol table structure and managing lexical tokens, stack of symbols or tokens.

Unit 4:

Syntax directed Translation & Intermediate code Generation:

Syntax directed definition, bottom-up evaluation, Intermediate representation, & intermediate code generation: Quadraples, Error detection & recovery, lexical phase errors, syntactic phase errors, semantic errors, more about translation: Array references in arithmetic expressions:

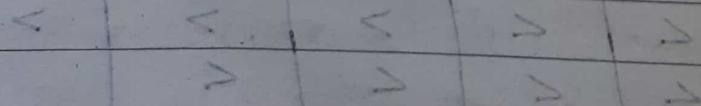
$$b_1 * b_1 + b_1 \leftarrow q_1$$

* Syntax directed Definition:

- SDD is a kind of abstract specification. It is generalization of CFG in which each grammar production $x \rightarrow q$ is associated with it a set of production rules of the form $s = f(b_1, b_2 \dots b_k)$ where s is attribute obtained from Function f .
- The attribute can be a string, number, type or a memory location.
- Semantic rules are fragments of code which are embedded usually at the end of production & enclosed in curly braces. ($\{ \}$).

Example:

$$E \rightarrow E_1 + T \quad \{ \quad E \cdot \text{val} = E_1 \cdot \text{val} + T \cdot \text{val} \quad \}$$



Annotated Parse Tree -

- The Parse Tree containing the values of attributes at each node for given IP string is called annotated or decorated parse tree.
- SDD specifies the values of attributes by associating semantic rules with grammar production.
- For example, an infix to postfix translator might have a production & rule

Production

semantic Rule

$$E \rightarrow E_1 + T \quad E\text{-code} = E_1\text{-code} || T\text{-code} \# '+'$$

- Both E & T have a string valued attribute code.

A syntax directed translation scheme embeds pgm fragments called semantic actions within production bodies, as in

$$E \rightarrow E_1 + T \{ \text{print } '+' \}$$

- By convention, semantic actions are enclosed within curly braces. (if curly braces occur as grammar symbols, we enclose them within single quotes.)
- A SDD is a CFG together with attributes & rules. Attributes are associated with grammar rules symbols & rules are associated with productions.
- If x is symbol & a is one of its attributes, then we write $x.a$ to denote the value of a at a particular parse tree by records or objects.
- Attributes may be of any kind : numbers, types,

table references, string for instance.

* Types of Attribute

- 1) Inherited attribute
- 2) Synthesized attribute

* Synthesized attribute:

- A synthesized attribute for a nonterminal A at a parse tree node N is defined by a semantic rule associated with the production at N . Note that the production must have A as its head. A synthesized attribute at node N is defined is defined only in terms of attribute values at the children of N & at N itself.
- These attributes get values from the attribute values of their child nodes.

e.g. -

$S \rightarrow ABC$

if S is taking values from its child nodes (A, B, C) then it is said to be synthesized attribute, as the values of ABC are synthesized to S .

* Inherited Attribute:

- An inherited attribute for a nonterminal B at a parse tree node N is defined by a semantic rule associated with the production at the parent of N .
- Note that the production must have B as a symbol in its body. An inherited attribute at node N is defined only in terms of attribute values of N 's parent, N itself & N 's siblings.

Inherited attributes can take values from parent & sibling siblings.

e.g. -

$$S \rightarrow ABC$$

A can get values from S, B & C. B can take values from S, A & C. Likewise, C can take values from S, A & B.

Example :

Production

$$1) L \rightarrow E_n$$

$$2) E \rightarrow E_1 + T$$

$$3) E \rightarrow T$$

$$4) T \rightarrow T_1 * F$$

$$5) T \rightarrow F$$

$$6) F \rightarrow (E)$$

$$7) F \rightarrow \text{digit}$$

Semantic Rules

$$L \cdot \text{val} = E \cdot \text{val}$$

$$E \cdot \text{val} = E_1 \cdot \text{val} + T \cdot \text{val}$$

$$E \cdot \text{val} = T \cdot \text{val}$$

$$T \cdot \text{val} = T_1 \cdot \text{val} * F \cdot \text{val}$$

$$T \cdot \text{val} = F \cdot \text{val}$$

$$F \cdot \text{val} = E \cdot \text{val}$$

$$F \cdot \text{val} = \text{digit} \cdot \text{lex val}$$

fig ① Syntax directed definition of a simple desk calculator

- In SDD, each of the nonterminals has a single synthesized attribute, called val.
- Terminal digit has a synthesized attribute lexval, which is an integer value returned by the lexical analyzer.
- An SDD that involves only synthesized attributes is called s-attributed. The above fig. has this property.
- In an s-attributed SDD, each rule computes an attribute for the nonterminal at the head

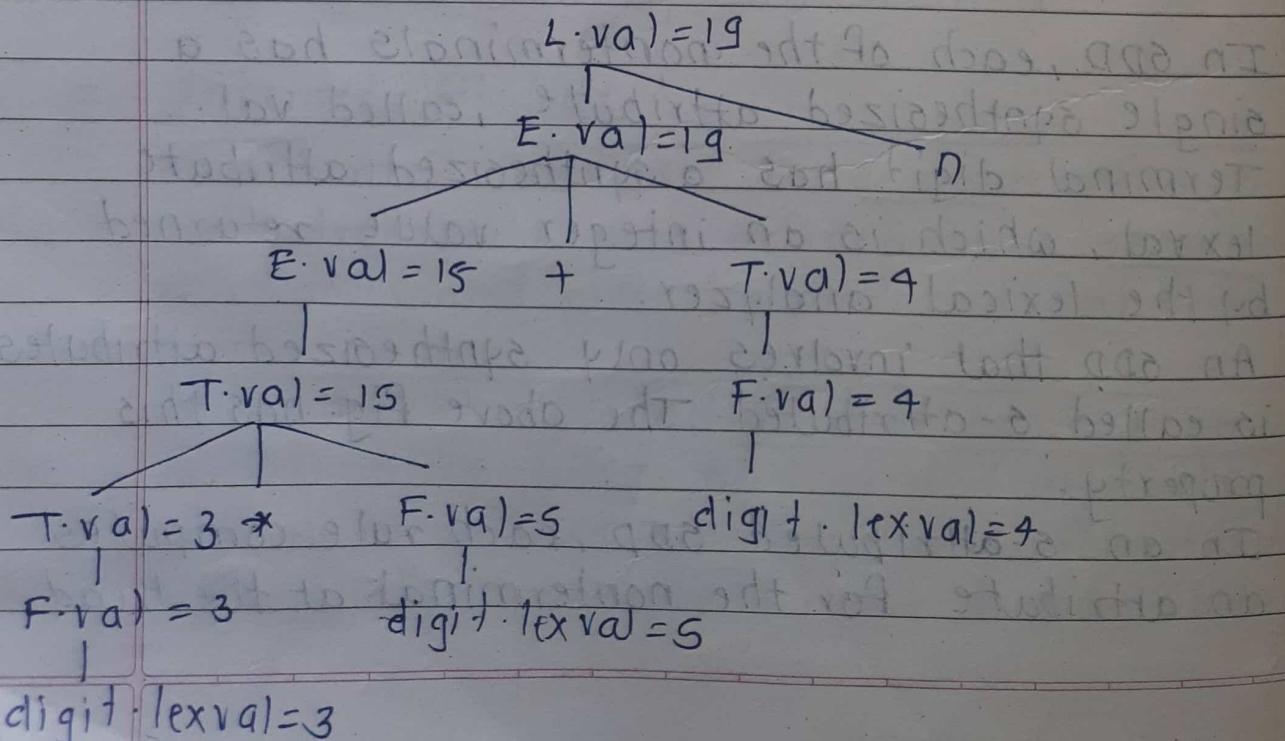
of a production from attributes taken from the body of the production.

- An s -attributed SDD can be implemented in conjunction with an LR parser.

Annotated Parse Tree

- A parse tree showing the values of its attributes is called an annotated parse tree.
- construction of annotated parse tree:
 - Before we can evaluate an attribute at a node of a parse tree, we must evaluate all the attributes upon which its value depends.
 - For example: if all attributes are synthesized, then we must evaluate the val attributes at all of the children of a node before we can evaluate the val attribute at the node itself.
 - With synthesized attributes, we can evaluate attributes in any bottom-up order, such as that of postorder traversal of parse tree.

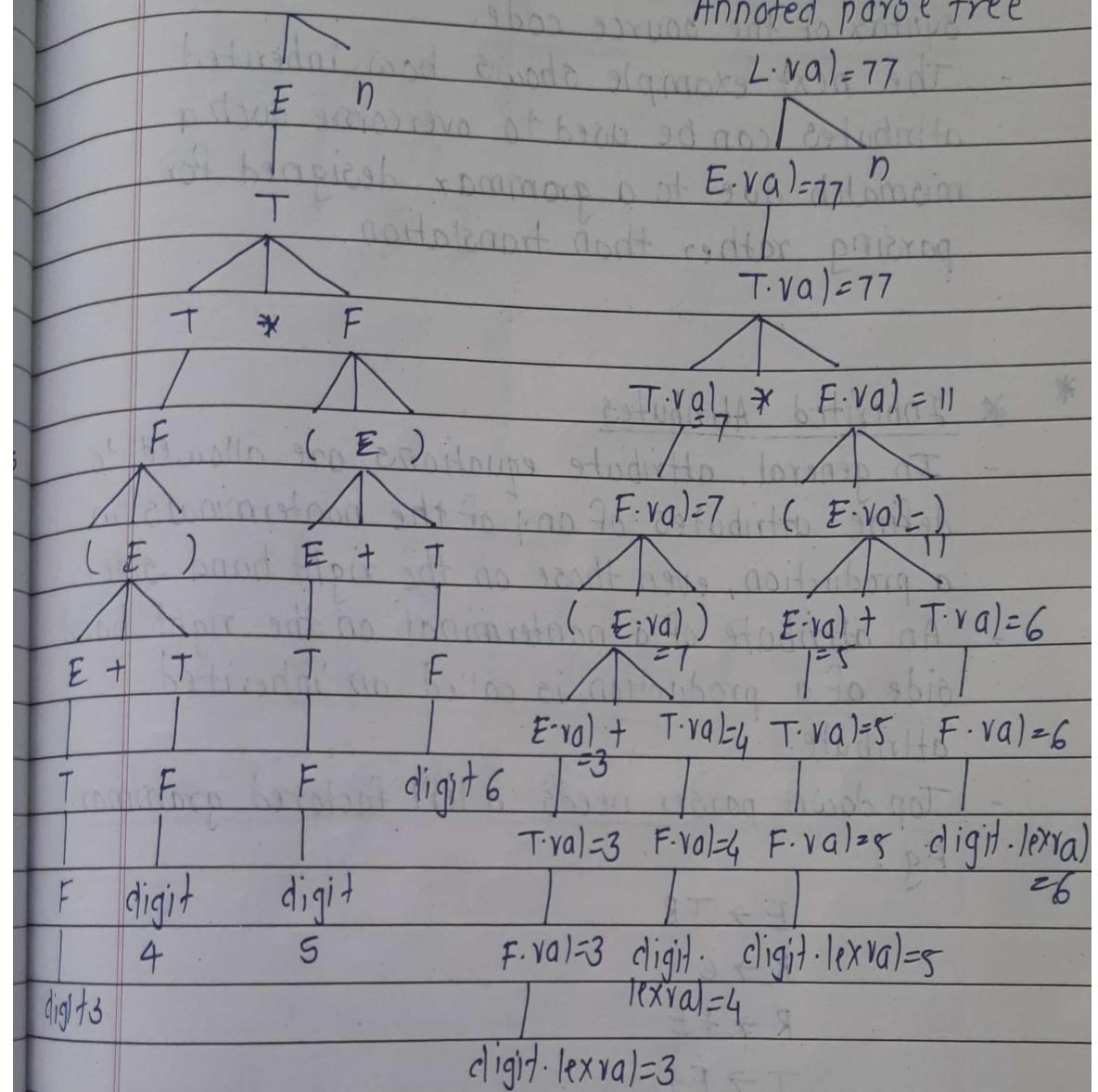
Example: $3 * 5 + 4 \pi$



$$*(3+4)* (5+6) n$$

L

Annotated parse tree



$$* 1 * 2 * 3 * (4 + 5) \triangleright$$

$$*(9+8)*(7+6)+5)*4n$$

- Inherited attributes are useful when the structure of a parse tree does not match the abstract syntax of the source code.
- The next example shows how inherited attributes can be used to overcome such a mismatch due to a grammar designed for parsing rather than translation.

* * Inherited Attributes

- In general, attribute equations are allowed to define attributes of any of the nonterminals in a production, even those on the right hand side. An attribute of a nonterminal on the right hand side of a production is called an inherited attribute.
- Top down parser needs a left factored grammar.

Eg :

$$E \rightarrow TR$$

$$R \rightarrow E$$

$$R \rightarrow +E$$

$$T \rightarrow FS$$

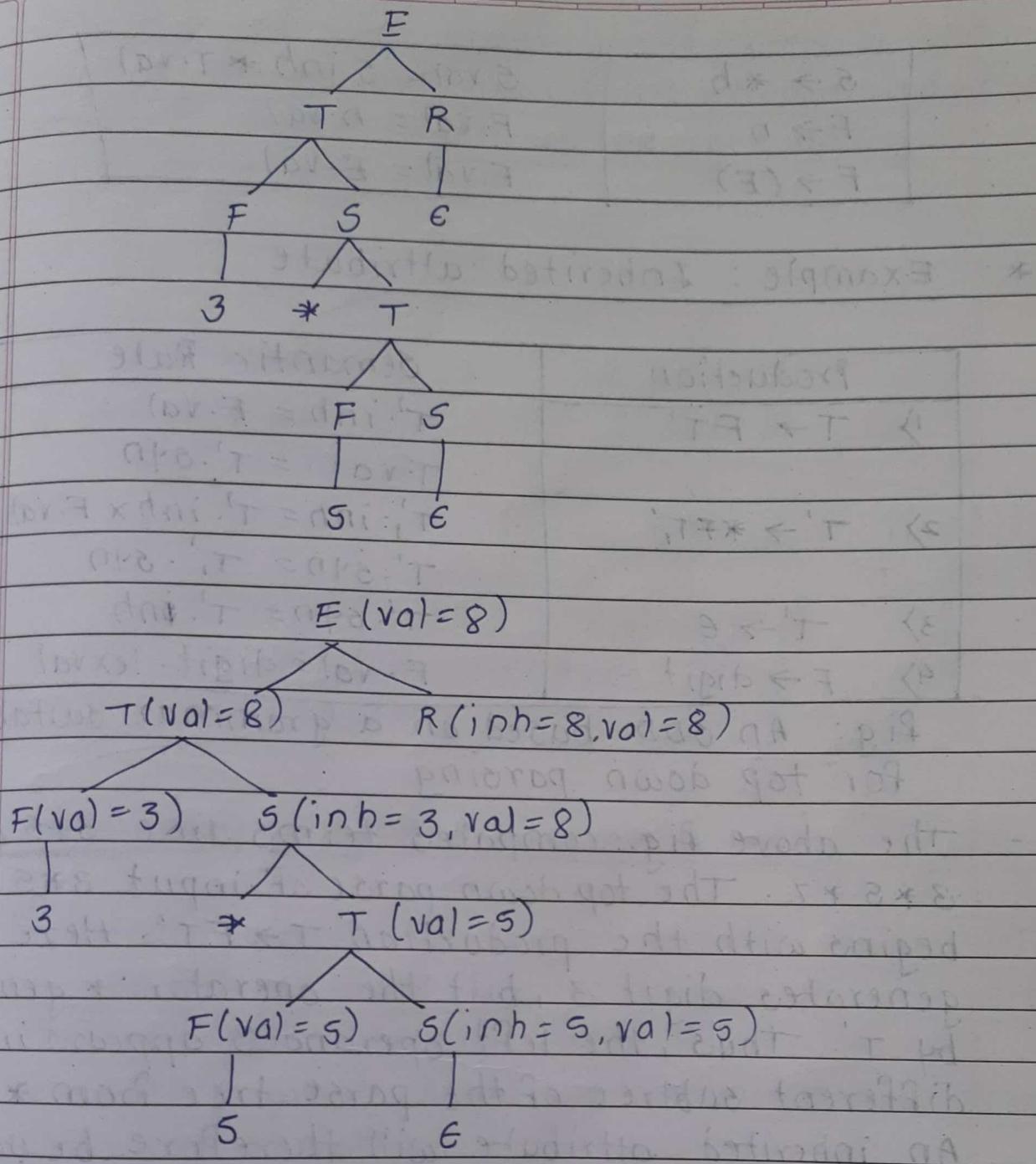
$$S \rightarrow E$$

$$S \rightarrow *T$$

$$F \rightarrow n$$

$$F \rightarrow (E)$$

Suppose that our goal is to associate an attribute of each E , T & F node that is the value of the expression.



Production

$$E \rightarrow TR$$

$$R \rightarrow E$$

$$R \rightarrow +E$$

$$T \rightarrow FS$$

$$S \rightarrow E$$

Semantic Rule

$$E \cdot \text{val} = R \cdot \text{val}$$

$$R \cdot \text{inh} = T \cdot \text{val}$$

$$R \cdot \text{val} = R \cdot \text{inh}$$

$$R \cdot \text{val} = R \cdot \text{inh} + E \cdot \text{val}$$

$$T \cdot \text{val} = S \cdot \text{val}$$

$$S \cdot \text{inh} = F \cdot \text{val}$$

$$S \cdot \text{val} = S \cdot \text{inh}$$

$S \rightarrow * h$	$S \cdot \text{val} = S \cdot \text{inh} * T \cdot \text{val}$
$F \rightarrow n$	$F \cdot \text{val} = n \cdot \text{val}$
$F \rightarrow (E)$	$F \cdot \text{val} = E \cdot \text{val}$

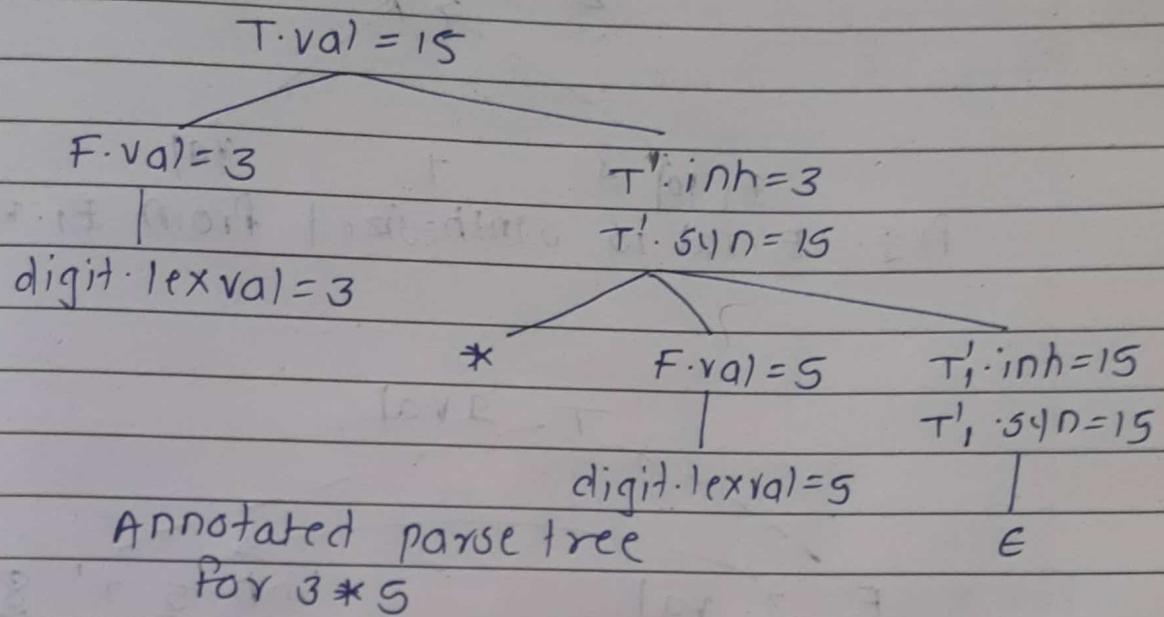
* Example : Inherited attribute

Production	Semantic Rule
1) $T \rightarrow FT'$	$T' \cdot \text{inh} = F \cdot \text{val}$
2) $T' \rightarrow *FT_1$	$T \cdot \text{val} = T' \cdot \text{syn}$
3) $T' \rightarrow E$	$T'_1 \cdot \text{inh} = T'_1 \cdot \text{syn} \times F \cdot \text{val}$
4) $F \rightarrow \text{digit}$	$T'_1 \cdot \text{syn} = T'_1 \cdot \text{inh}$
	$F \cdot \text{val} = \text{digit} \cdot \text{lexval}$

Fig: An SDD based on a grammar suitable for top down parsing

- The above fig. computes terms like $3*5$ & $3*5*7$. The top down parse of input $8*5$ begins with the production $T \rightarrow FT'$. Here F generates digit 3, but the operator $*$ generated by T' . Thus, the left operand 3 appears in a different subtree of the parse tree from $*$. An inherited attribute will therefore be used to pass the operand to the operator.
- The semantic rule are based on the idea that the left operand of the operator $*$ is inherited. More precisely, the head T' of the production $T' \rightarrow *FT_1$ inherits the left operand of $*$ in the production body.
- Given a term $x * y * z$, the root of the subtree for $y * z$ inherits x . Then root of subtree

for $*z$ inherits the value of $x * y$ & so on. once all the factors have been accumulated, the result is passed back up the tree using synthesized attributes.



* Dependency Graphs (Evaluation order of SDD)

- Dependency graph are useful tool for determining an evaluation order for the attribute instances in given parse tree.
- Dependency graph helps us determine how those values can be computed.
- Dependency graph depicts the flow of information among the attribute instances in a particular parse tree.

Example :

Production
 $E \rightarrow E_1 + T$

Semantic Rule
 $\text{Eval} = E_1 \cdot \text{val} + T \cdot \text{val}$

As a convention, we shall show the parse tree edges as dotted lines, while the edges of the dependency graph are solid.

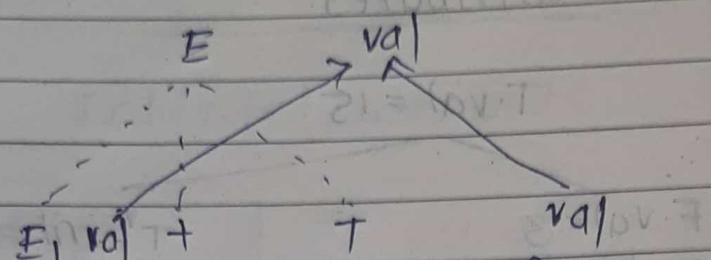


Fig. E.val is synthesized from E1.val & T.val.

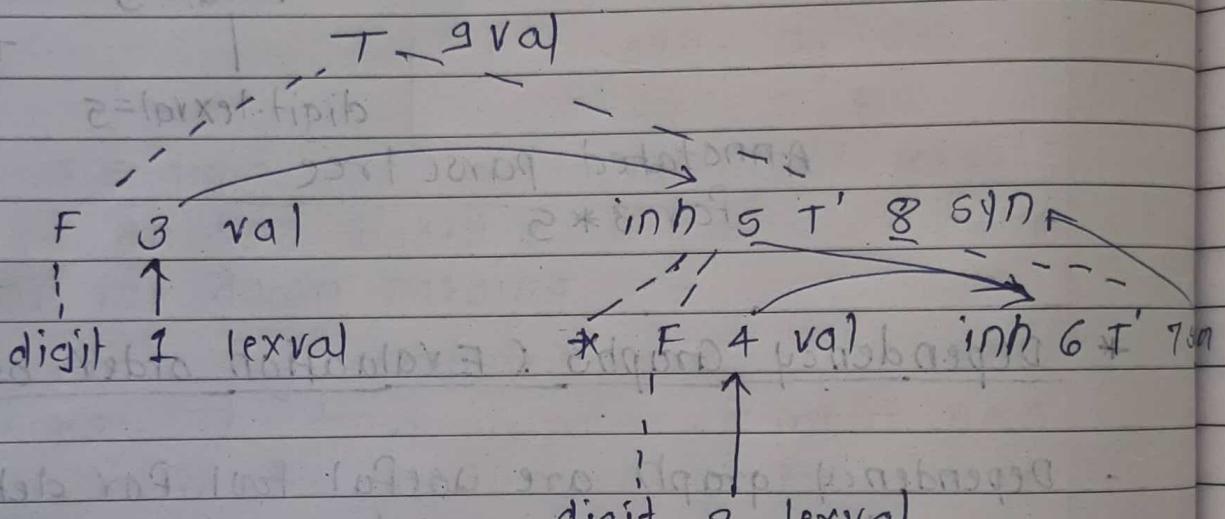


Fig: Dependency graph for annotated parse tree
For 3*5

* s-Attributed SDD

- If every attribute is synthesized, then an SDD is called s-attributed SDD.

- If the value of parent nodes depends upon the value of child nodes, the s-attributed SDD is evaluated in bottom-up parsing.

The right most place of RHS holds the semantic action.

- Fig ① is s-attributed SDD.

- * L-attributed SDD
- If an attribute of an SDD is synthesized or inherited with some restriction on inherited attributes, it can inherit values from left siblings only. It is known as L-attributed SDD.
- Attributes of this SDD are evaluated by depth first & left to right parsing methods.

Eg. -

$$x \Rightarrow AB C$$

$$\{ B.P = x.P, B.P = A.P, B.P = C.P \}$$

This is not L-attributed SDD.

Eg:

Production

$$T \Rightarrow FT'$$

$$T' \Rightarrow *FT';$$

semantic rule

$$T'.inh = F.val$$

$$T'.inh = T.inh \times F.val$$

Example :

Any SDD containing the following production & rules cannot be L-attributed.

Production

$$A \Rightarrow BC$$

semantic rule

$$A.s = B.b;$$

$$B.i = F(C.c, A.s)$$

* SDD with simple type declarations :

Production

$$1) D \rightarrow TL$$

$$2) T \rightarrow int$$

$$3) T \rightarrow float$$

$$4) L \rightarrow L, id$$

semantic Rule

$$L.inh = T.type$$

$$T.type = integer$$

$$T.type = float$$

$$L.inh = L.inh$$

$$addType(id.entry, L.inh)$$

$$addType(id.entry, L.inh)$$

$$5) L \rightarrow id$$

Fig. SDD for simple type declarations

- Fig. takes a simple declarations D consisting of a basic type T followed by a list L of identifiers.
- The purpose of L.inh is to pass the declared type down the list of identifiers, so that it can be added to the appropriate symbol table entries.
- Production 4 & 5 also have a rule in which a function addType is called with two arguments
 - 1) id.entry , a lexical value that points to a symbol table object.
 - 2) L.inh , the type being assigned to every identifier on the list .

Example :

input string : float id1, id2, id3

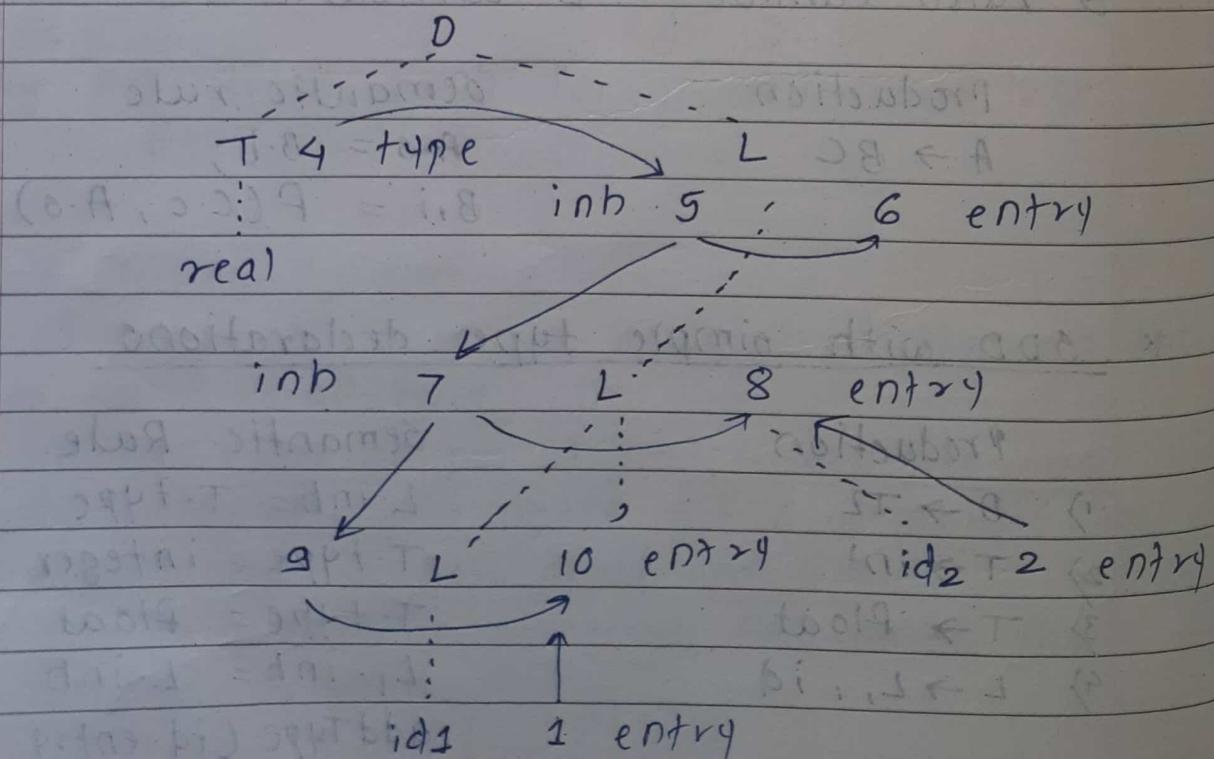


Fig. dependency graph for declaration float id1, id2, id3

* For the given grammar give annotated parse tree for following expression

a) int a, b, c

b) float w, x, y, z

Production

$$1) D \rightarrow TL$$

$$2) T \rightarrow \text{int}$$

$$3) T \rightarrow \text{float}$$

$$4) L \rightarrow L, id$$

$$5) L \rightarrow id$$

Semantic rule

$$L.\text{inh} = T.\text{type}$$

$$T.\text{type} = \text{integer}$$

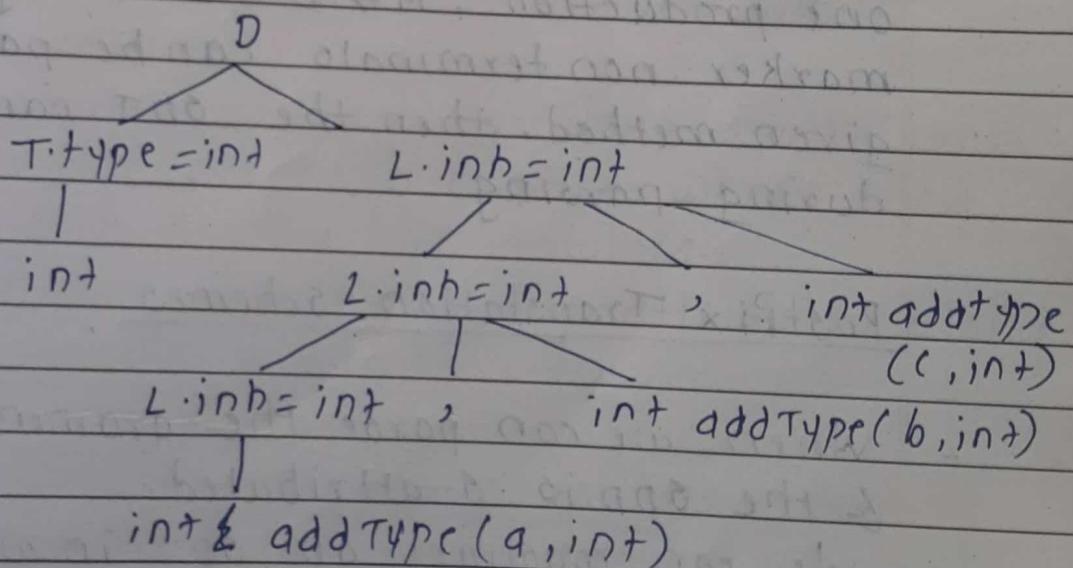
$$T.\text{type} = \text{float}$$

$$L_i.\text{inh} = L.\text{inh}$$

$$\text{addType(id.entry, } L.\text{inh})$$

$$\text{addType(id.entry1, } L.\text{inh})$$

a)



* Syntax Directed Translation Schemes:

- The SDT's are implemented during parsing, without building a parse tree.
- classes of SDD.
- 1) The underlying grammar is LR parseable & the SDD is S-attributed.
- 2) The underlying grammar is LL parseable, & the SDD is L-attributed.
- * SDT's that can be implemented during parsing can be characterized by introducing distinct marker nonterminals in place of each embedded action, each marker m has only one production, $m \rightarrow e$. If the grammar with marker nonterminals can be parsed by a given method, then the SDT can be implemented during parsing.

Postfix Translation schemes

- When we can parse the grammar bottom-up & the SDD is S-attributed.
- We can construct an SDT in which each action is placed at the end of the production & is executed along with the reduction of the body to the head of that production.
- SDT with all actions at the right ends of the production bodies are called postfix SDTs.

Example :

$L \rightarrow E_n$	{ print ($E \cdot \text{val}$); }
$E \rightarrow E_1 + T$	{ $E \cdot \text{val} = E_1 \cdot \text{val} + T \cdot \text{val}$; }
$E \rightarrow T$	{ $E \cdot \text{val} = T \cdot \text{val}$; }
$T \rightarrow T_1 * F$	{ $T \cdot \text{val} = T_1 \cdot \text{val} * F \cdot \text{val}$; }
$T \rightarrow F$	{ $T \cdot \text{val} = F \cdot \text{val}$; }
$F \rightarrow (E)$	{ $F \cdot \text{val} = E \cdot \text{val}$; }
$F \rightarrow \text{digit}$	{ $F \cdot \text{val} = \text{digit} \cdot \text{lexval}$; }

Fig : Postfix SGT implementing the desk calculator

* Parser - stack Implementation of Postfix SGT.

- Postfix SGTs can be implemented during LR parsing by executing the actions when reductions occur.
- The attributes of each grammar symbol can be put on the stack in a place where they can be found during the reduction.
- The best plan is to place the attributes along with the grammar symbols in records on the stack itself.
- Fig shows the parser stack contains records with a field of a grammar symbol (or parser state) & below it a field for an attribute.

x	y	z	state / grammar symbol
x.x	y.y	z.z	synthesized attributes

Fig : Parser stack with a field for synthesized attributes

- Stack manipulation is usually done automatically by the parser.

Product & DOT = { } Actions

$L \rightarrow E_n$	{ print(stack[top-1].val); top = top - 1 }
$E \rightarrow E, + T$	{ stack[top-2].val = stack[top-2].val + stack[top].val; top = top - 2; }
$E \rightarrow T$	
$T \rightarrow T, * F$	{ stack[top-2].val = stack[top-2].val * stack[top].val; top = top - 2; }

$T \rightarrow F$	
$F \rightarrow (E)$	{ stack[top-2].val = stack[top-1].val; top = top - 2; }

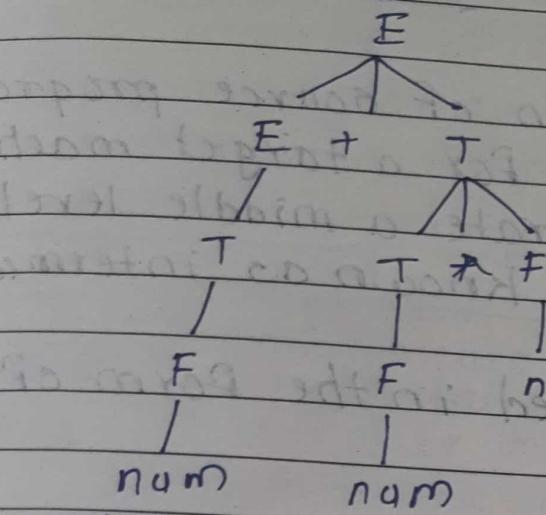
$F \rightarrow \text{digit};$

Fig. Implementing the desk calculator
on a bottom-up parsing stack.

* SDD for infix to postfix translation

Grammar	Syntactic Action	Action
$E \rightarrow E + T$	{ print ("+"); }	
$E \rightarrow T$	{ };	
$T \rightarrow T * F$	{ print ("*"); }	
$T \rightarrow F$	{ };	
$F \rightarrow \text{num}$	{ print (num.val); }	

convert $2 + 3 * 4$ to postfix notation



$2 + 3 * 4$

* Remaining from PPT unit: 3

Intermediate code Generation

- During the translation of source program into the object code for a target machine, a compiler may generate a middle level lang. code which is known as intermediate code.
- It can be represented in the form of :
 - postfix notation
 - syntax tree
 - directed acyclic graph (DAG)
 - three address code
 - quadruples
 - triples

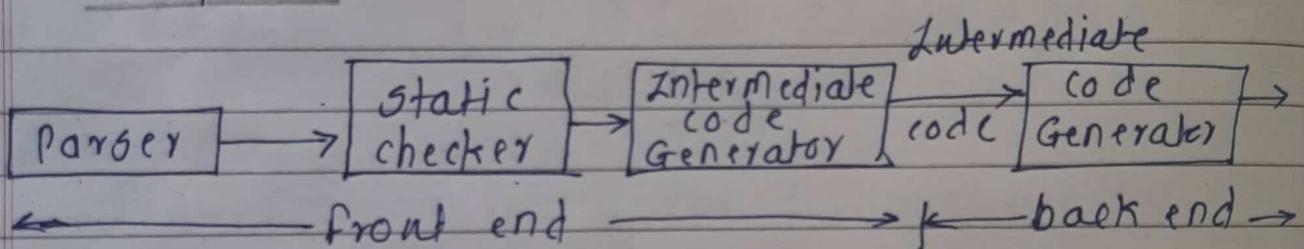


Fig. logical structure of compiler Front end.

- static checking includes type checking, which ensures that operators are applied to compatible operands.
- It also includes any syntactic checks that remain after parsing.
- For example, static checking assures that a break stmt in c is enclosed within a while, for, switch stmt ; an error is reported if such an enclosing statement does not exist.
- Intermediate representations includes

Syntax trees & three address code.

In the process of translating a program in a given source lang into code for a given target mic, a compiler may construct a sequence of intermediate representations.

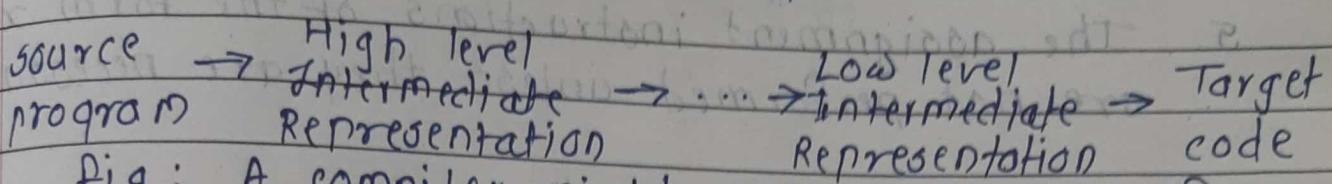


Fig: A compiler might use a sequence of intermediate representations

* Directed Acyclic Graphs

- like a syntax tree for expression, a DAG has leaves corresponding to atomic operands & interior nodes corresponding to operators.
- The difference is that a node N in a DAG has more than one parent if N represents a common subexpression
- Thus a DAG not only represents expressions more succinctly, it gives the compiler important clues regarding the generation of efficient code to evaluate the expression.
- A DAG is constructed for optimizing the basic block.
- A DAG is usually constructed using Three address code.

Construction of DAG (Rules)

1. Interior nodes always represent the operators.
2. Exterior nodes (leaf nodes) always represent

- the names, identifiers or constants.
3. A check is made to find if there exists any node with the same value.
 4. A new node is created only when there does not exist any node with the same value.
 5. The assignment instructions of the form $x := y$ are not performed unless they are necessary.

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

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

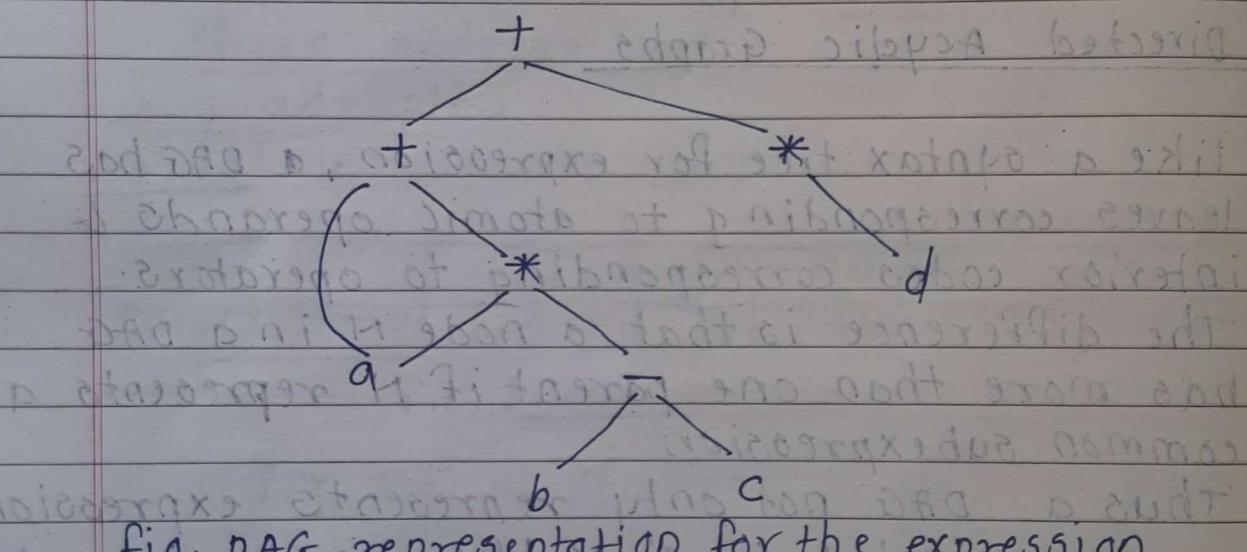


fig. DAG representation for the expression

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

- The leaf for a has two parents, because a appears twice in the exp. The two occurrences of the common subexpression $b - c$ are represented by one node the node labeled $-$.
- That node has two parents, representing its two uses in the subexpression $a * (b - c)$ & $(b - c) * d$.
- Even though b & c appear twice in the complete expression, their nodes each have one parent.

since both uses are in the common subexpression
b-c.

- P₁ = Leaf (id, entry-a)
- P₂ = Leaf (id, entry-a) = P₁
- P₃ = Leaf (id, entry-b)
- P₄ = Leaf (id, entry-c)
- P₅ = Node ('-', P₃, P₄)
- P₆ = Node ('*', P₁, P₅)
- P₇ = Node ('+', P₁, P₆)
- P₈ = Leaf (id, entry-b) = P₃
- P₉ = Leaf (id, entry-c) = P₄
- P₁₀ = Node ('-', P₃, P₄) = P₅
- P₁₁ = Leaf (id, entry-d)
- P₁₂ = Node ('*', P₅, P₁₁)
- P₁₃ = Node ('+', P₇, P₁₂)

fig. steps for constructing the DAG for above fig.

- Here node & leaf return an existing node. We assume that entry-a points to the symbol table entry for a & similarly for the other identifiers.
- When a call to leaf (id, entry-a) is repeated at step 2, the node created by the previous call to is returned, so P₂ = P₁. Similarly, the nodes returned at steps 8 & 9 are the same as those returned at step 8 & 4 (ie P₈ = P₃ & P₉ = P₄). Hence the node returned at step 10 must be the same as that returned at step 5 ie P₁₀ = P₅.

* Value Number method for constructing DAG

- The nodes of syntax tree or DAG are stored in an array of records.
- Each row of the array represents one record & therefore one node. In Fig. (b) leaves have one additional field, which holds lexical value (either symbol table pointer or a constant) & interior nodes have two additional fields indicating the left & right children.

$=$	$89 = (\text{d_prfn}, 6)$	$\text{id} = 89$	$\rightarrow \text{to entry}$
$+$	$10 = (\text{d_prfn}, 6)$	$\text{num} = 10$	$\text{For } i$
$+$	$(89, 89, -3)$	$+1 = 10$	2
$+$	$(\text{d_prfn}, 6)$	$i = 1 = 11$	3
i	$(10, 29, *5)$	$\text{sbch} = 29$	

(a) DAG

(b) Array

Fig. nodes of DAG For $i = i + 10$ allocated in array.

Example: construct DAG for following expression

$$1) b((x+y)-((x+y)*(x-y))) + ((x+y)*(x-y))$$

construct DAG & identify the value numbers for the subexpressions of the following exp., assuming + associates from left.

$$1) a+b+(a+b)$$

$$2) a+b+a+b$$

$$3) a+a+((a+a+a)+(a+a+a+a))$$

* Three Address codes

- In TAC, there is at most one operator on the right side of an instruction.

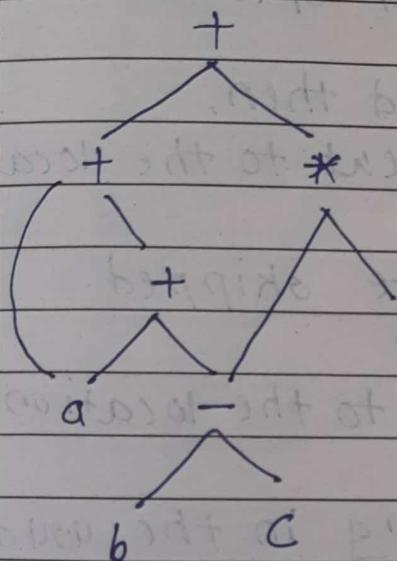
$$\text{Eg: } x + y \rightarrow z$$

So, sequence of three address instructions

$$t_1 = y \rightarrow z$$

$$t_2 = x + t_1$$

where, t_1 & t_2 are compiler generated temporary names.



(a) DAG

$$t_1 = b - c$$

$$t_2 = a * t_1$$

$$t_3 = a + t_2$$

$$t_4 = t_1 * d$$

$$t_5 = t_3 + t_4$$

(b) Three address code

Fig. a DAG & its corresponding three address code.

Instructions to follow for construction of TAC

1. Assignment statement:

$$x = y \text{ op } z \text{ and } x = \text{op } y$$

x, y, z are operands

op is operator

2. copy statement

3. conditional jump

IF $x \text{ relop } y \text{ goto } x$

Here, x & y are operands

~~but x is the target label of the statement.~~

* relop is a relational operator.

if $x \text{ relop } y$ gets satisfied then,

- The control is directly sent to the location specified by label x .
- All stmts in between are skipped.

if $x \text{ relop } y$ fails then,

- The control is not sent to the location specified by label x .
- The next stmt appearing in the usual sequence is executed.

* Example: ~~public void main() { int a = 5; }~~

$$1) a = b + c + d$$

$$t_1 = b + c$$

$$t_2 = t_1 + d$$

$$a = t_2$$

$$2) (a * b) + (c + d) - (a + b + c + d)$$

$$P \leftarrow a \quad P \leftarrow P * b \quad P \leftarrow P + c \quad P \leftarrow P + d$$

$$P \leftarrow P - a \quad P \leftarrow P - b \quad P \leftarrow P - c \quad P \leftarrow P - d$$

$$P \leftarrow P + a \quad P \leftarrow P + b \quad P \leftarrow P + c \quad P \leftarrow P + d$$

3) IF $A < B$ then 1 else 0

\Rightarrow 1) IF ($A < B$) goto (4)

2) $t_1 = 0$

3) goto (5)

4) $t_1 = 1$

5) exit

4) IF $A < B$ and $C < D$ then $t=1$ else $t=0$

\Rightarrow 1) if ($A < B$) goto (3)

2) goto (4)

3) if ($C < D$) goto (6)

4) $t = 0$

5) goto (7)

6) $t = 1$

7) exit

* Quadruples

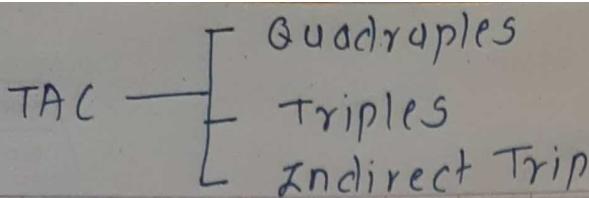
- Three address instructions specifies the components of each type of instruction, but it does not specify the representation of these inst. in DS.

- In compiler these instructions can be implemented as objects or as records with fields for the operator & operands.

- These such representations are called quadruples, triples & indirect triples.

- A quadruple (quad) has four fields i.e. op, arg₁, arg₂, result. op field contains an internal code for the operator.

Rules :



- Instructions with unary operators like $x = \text{minus } y$ or $x = y$ do not use arg2. Note that for a copy stmt like $x = y$, op is =, while for most other operations, the assignment operator is implied.
- Operators like param use neither arg2 nor result.
- conditional and unconditional jumps put the target label in result.

Eg :

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

(-c)	$t_1 = \text{minus } c$	op	arg1	arg2	result
unary operator	$t_2 = b * t_1$	0	minus	c	
	$t_3 = \text{minus } c$	1	*	b	t_1
	$t_4 = b * t_3$	2	minus	c	t_2
	$t_5 = t_2 + t_4$	3	*	b	t_3
	$a = t_5$	4	+	t_2	t_4
(a) TAC		5	=	t_5	a

(b) Quadruples

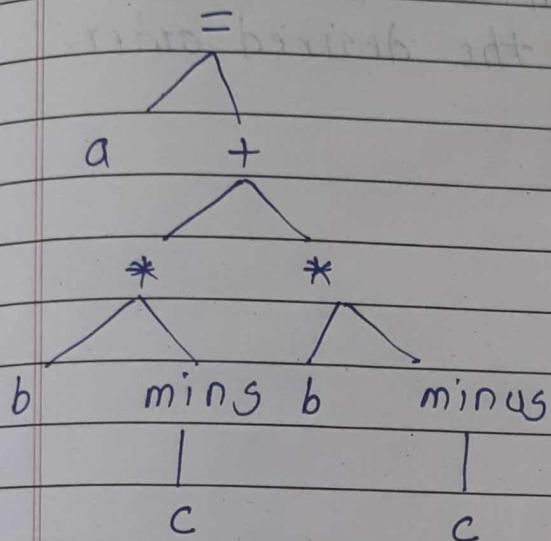
* Triples are objects of formalism changing & nothing is lost.

- A triple has only three fields, which we call op, arg1, & arg2. Note that the result field in above Fig. is used primarily for temporary names.
- Using triples, we refer to the result of an operation $x \text{ op } y$ by its position, rather than

by an explicit temporary name.

Example:

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



(a) Syntax tree

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

(b) Triples

Fig. representation of $a * a * (b - c) + (b - c) * d$

- Thus instead of the temporary t , in above fig. (Quadruples), a triple representation refers to position (0).
- Parenthesized numbers represent pointers into the triple structure itself.
- DAG & triple representations of expressions are equivalent.

In triple representation,

- References to the instructions are made.
- Temporary variables are not used.

* Indirect triples

- It consists of listing of pointers to triples, rather than a listing of triples themselves. For example, let us use an array inst. to list pointers to triples in the desired order.

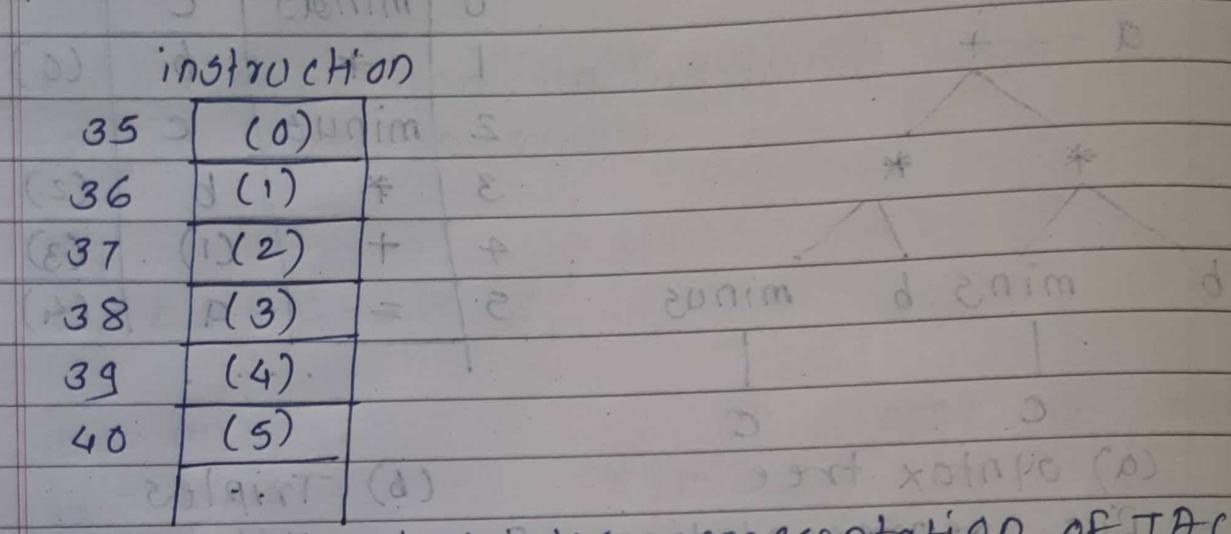


Fig. Indirect triples representation of TAC

Examples

$$1) \quad a+b \times c / e \uparrow f + b \times c$$

TAC :

$$t_1 = e \uparrow f$$

$$t_2 = b \times c$$

$$t_3 = t_2/t_1$$

$$t_4 = b \times \sigma$$

$$t_5 = a + t_3$$

$$t_6 = T_5 + t_4$$

Quadruple

	OP	arg1	arg2	result
(0)	\uparrow	e	f	t1
(1)	x	b	c	t2
(2)	/	t2	t1	t3
(3)	x	b	a	t4
(4)	+	a	t3	t5
(5)	+	t5	t4	t6

Triple

	OP	arg1	arg2
(0)	\uparrow	e	f
(1)	x	b	c
(2)	/	(1)	(0)
(3)	x	b	a
(4)	+	a	(2)
(5)	+	(4)	(3)

Indirect Triple

	statement
35	(0)
36	(1)
37	(2)
38	(3)
39	(4)
40	(5)

	OP	arg1	arg2
(0)	\uparrow	e0	f
(1)	\times	b	e
(2)	\div	c1	(0)
(3)	\times	b	a
(4)	$+$	a	(2)
(5)	$+$	(4)	(3)

Example 3:

Translate the arithmetic exp. into

- a) A syntax tree b) Quadruples
- c) Triples d) Indirect triples

$$1) a + - (b + c)$$

$$2) a = b * - c + b * - c$$

* Types & Declarations

Type checking : uses logical rules to reason about the behavior of a pgm at run time.

Translation Applications: From the type of a name a compiler can determine the storage that will be needed for that name at run time.

* Type Expression

Example :

The array type `int [2][3]` can be read as "array of 2 arrays of 3 integers each" & written as a type expression `array (2, array (3, integer))`.

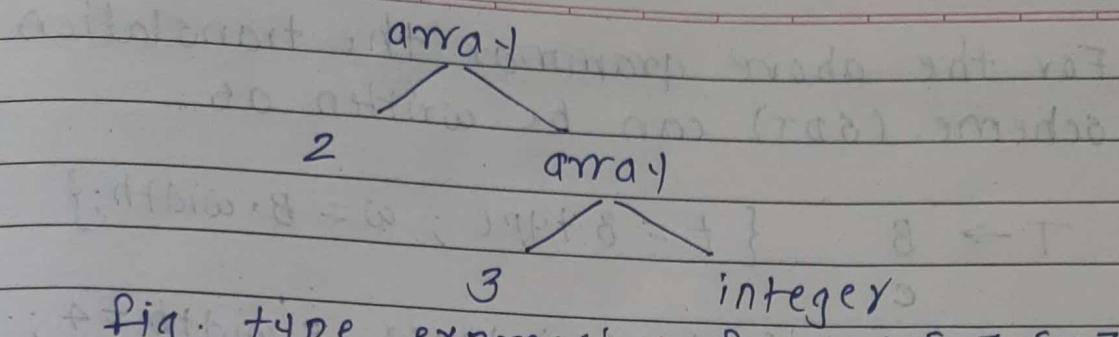


Fig. type expression for `int[2][3]`

* Declarations

- The type & declarations uses a simplified grammar that declares just one name at a time.

The grammar is:

$$D \rightarrow T \text{id} ; D \mid \epsilon$$

$$T \rightarrow B \ C \mid \text{record} \{ \ D \ \}$$

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

$$C \rightarrow \epsilon \mid [\text{num}] C$$

Here, we consider storage layout as well as types. Here,

D - generates the sequence of declarations.

T - generates basic, array or record types.

B - generates basic type int & float.

C - generate string of zero or more integers, each integer surrounded by brackets.

A record-type is a sequence of declarations for the fields of the record; all surrounded by curly braces.

For the above grammar the translation scheme (SDT) can be written as,

$$T \rightarrow B \quad \{ t = B \cdot \text{type} ; w = B \cdot \text{width} ; \}$$

PUSH

$$B \rightarrow \text{int} \quad \{ B \cdot \text{type} = \text{integer} ; B \cdot \text{width} = 4 ; \}$$

$$B \rightarrow \text{float} \quad \{ B \cdot \text{type} = \text{float} ; B \cdot \text{width} = 8 ; \}$$

$$C \rightarrow e \quad \{ c \cdot \text{type} = t ; c \cdot \text{width} = w ; \}$$

$$C \rightarrow [\text{num}] C_1 \quad \{ \text{array}(\text{num} \cdot \text{value}, c_1 \cdot \text{type}) ;$$

$$c \cdot \text{width} = \text{num} \cdot \text{value} \times c_1 \cdot \text{width} ; \}$$

- The SDT computes types & their widths for basic & array types.
- Here type & width are synthesized attributes for each non terminal & two variables t & w to pass type & width info. From B node in a parse tree to the node for production $C \rightarrow e$.
- The body of T production consists of nonterminal B, an action & nonterminal C, which appears on next line.
- The action between B & C gets t to B.type & w to B.width.
- The production for C determines whether T generates a basic type or array type.
- If $C \rightarrow e$ then t becomes c.type & w becomes c.width. otherwise, e specifies an array component.
- The action for $C \rightarrow [\text{num}] C_1$ forms c.type by applying the type constructor array to the operands num.value & c1.type.

* Translation of Expression

Expression : $a = b + -c$

Three address code :

$$t_1 = \text{minus } c$$

$$t_2 = b + t_1$$

$$a = t_2$$

Semantic Rule for above exp. with TAC

Production

$$S \rightarrow id = E ; \quad S\text{-code} = E\text{-code} ||$$

$$E \rightarrow E_1 + E_2$$

$$E\text{-addr} = \text{newTemp}()$$

$$E\text{-code} = E_1\text{-code} || E_2\text{-code} ||$$

$$\text{gen}(E\text{-addr} = 'E_1\text{-addr}' + 'E_2\text{-addr}')$$

$$1 - EP$$

$$E\text{-addr} = \text{newTemp}()$$

$$E\text{-code} = E_1\text{-code} ||$$

$$\text{gen}(E\text{-addr} = '\text{minus}' E_1\text{-addr})$$

$$E\text{-addr} = E_1\text{-addr}$$

$$E\text{-code} = E_1\text{-code}$$

$$\#id \mid id$$

$$E\text{-addr} = \text{top}\cdot\text{get}(id\text{-lexeme})$$

$$E\text{-code} = "1"$$

- Attributes $S\text{-code}$ & $E\text{-code}$ denote the TAC for S & E respectively.
- Attribute $E\text{-addr}$ denotes the address that will hold the value of E .
- consider the last production ie $E \rightarrow id$, when an exp. is a single identifier say x , then x itself holds the value of expression. The semantic rule for this production define $E\text{-addr}$ to point to symbol table entry. Let top denote the current symbol table. Function

top.get retrieves the entry when it is applied to string representation id. lexeme of this instance of id. E.code is set to empty string.

- When $E \rightarrow (E_1)$ the translation of E is same as that of subexpression E_1 . Hence, $E\text{-addr}$ equals $E_1\text{-addr}$ & $E\text{-code}$ equals $E_1\text{-code}$.
- The rule $E \rightarrow E_1 + E_2$, generate code to compute the value of E from the values of E_1 & E_2 . values are computed into newly generated temporary names. ie $t = E_1\text{-addr} + E_2\text{-addr}$, where t is new temporary name. $E\text{-addr}$ is set to t .
- Here we used notation gen ($x = 'y' + 'z'$) to represent three address inst. $x = y + z$. expressions appearing in place of variables like x, y & z are evaluated when passed to gen, & quoted string like $'$ are taken literally.

* In SDD, gen builds an inst. & returns it. In translation scheme, gen builds an inst. & incrementally emits it by putting it into the stream.

* Incremental Translation

slide 22

- Instead of building E-code, we arrange to generate only the new three address instructions.
- In the incremental approach, gen not only constructs a three address inst., it appends the inst. to the sequence of inst.'s generated so far.
- The sequence may either be retained in memory for further processing or it may be output incrementally.
- With the incremental approach, the code attribute is not used, since there is a single sequence of instructions that is created by successive calls to gen.

$\$ \rightarrow id = E; \{ gen(top.get(id.lexeme) '=' E.addr); \}$
 $E \rightarrow E1 + E2 \{ E.addr = new Temp();$
 $gen(E.addr = 'E1.addr +' E2.addr); \}$

$1 - E1 \{ E.addr = new Temp();$
 $gen(E.addr = 'minus' E1.addr); \}$

$| (E1) \{ E.addr = E1.addr; \}$

$| id \{ E.addr = top.get(id.lexeme); \}$

Fig : Generating three address code for expressions incrementally

Fig 2

* Translation of Array References

- Let nonterminal L generate an array name followed by a sequence of index expressions:

$$L \rightarrow L(E) \mid id [E]$$

$S \rightarrow id = E ; \{ gen(top.get(id.lexeme)) = 'E.addr' ; \}$
 $| L = E ; \{ gen(L.addr.base ['L.addr']) = 'E.addr' ; \}$

$L \rightarrow id[E] \quad \{ \quad L.array = top.get(id.lexeme);$
 $L.type = L.array.type.elem;$
 $L.addr = newTemp();$
 $gen(L.addr) = E.addr \star$
 $L.type.width); \}$

```
L1(E) { L.array = L1.array;
         L.type = L1.type.elem;
         t = new Temp();
         L.addr = new Temp();
         gen(t '=' E.addr '*' L.type.width);
         gen(L.addr '=' L1.addr + t); }
```

Fig: semantic actions for array references

Non-terminal L has three synthesized attributes:

1. $L\text{-addr}$ - denotes a temporary that is used while computing the offset for the array reference by summing the terms $ij * w_j$.
2. $L\text{-array}$ - is a pointer to symbol table entry for the array name. The base address of the array, $L\text{-array.base}$ is used to determine the actual l-value of an array reference after all index expressions are analyzed.
3. $L\text{-type}$ is the type of the subarray generated by L .

- the semantic action for $S \rightarrow L = E;$ generates an indexed copy inst. to assign the value denoted by expression E to the location denoted by the array reference L . The attribute $L\text{-array}$ gives the symbol table entry for the array.
- The arrays base address - the address of its 0th element - is given by $L\text{-array.base}$.
- Attribute $L\text{-addr}$ denotes the temporary that holds the offset for the array ref. generated by L .
- The location for the array reference is therefore $L\text{-array.base}[L\text{-addr}]$.
- The generated inst. copies the r-value from address $E\text{-addr}$ into the location for L .

Example:

1. Add the to the translation of Fig ① rules for the following productions:
 - a) $E \rightarrow E_1 * E_2$
 - b) $E \rightarrow + E_1$ (unary plus)
2. Repeat above eg. for the incremental translation of Fig. ②

* Flow of control statements :

34,

Grammar:

 $s \rightarrow \text{if}(B)s,$ $s \rightarrow \text{if}(B)s, \text{else } s_2$ $s \rightarrow \text{while}(B)s,$

Nonterminal B represents a boolean expression & terminal s represents a statement.

Production

semantic Rule

 $P \rightarrow s$ $s.\text{next} = \text{newlabel}()$ $P.\text{code} = s.\text{code} \text{ || label}(s.\text{next})$ $s \rightarrow \text{assign}$ $s.\text{code} = \text{assign}.\text{code}$ $s \rightarrow \text{if}(B)s$ $B.\text{true} = \text{newlabel}()$ $B.\text{false} = s_1.\text{next} = s.\text{next}$ $s.\text{code} = B.\text{code} \text{ || label}(B.\text{true})$
 $\text{|| } s_1.\text{code}$

cont'd

Fig: code for if, if-else & while statements

The translation of $\text{if}(B) S_1$ consists of B-code followed by S_1 -code. Within B-code are jumps based on the value of B. If B is true, control flows to the first inst. of S_1 -code, & if B is false, control flows to the inst. immediately following S_1 -code.

- The SDD for above grammar produces three address code for boolean expressions in the context of if, if-else & while statements.

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

- B-true = newlabel()
- B-false = newlabel()
- $S_1\text{-next} = S_2\text{-next} = s\text{-next}$
- $s\text{-code} = B\text{-code}$
- || label(B-true) || S_1 -code
- || gen('goto' s-next)
- || label(B-false) || S_2 -code

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

- begin = newlabel()
- B-true = newlabel()
- B-false = s-next
- $S_1\text{-next} = \text{begin}$
- $S_1\text{-code} = \text{label(begin)} || B\text{-code}$
- || label(B-true) || S_1 -code
- || gen('goto' begin)

$s \rightarrow S_1 S_2$

- $S_1\text{-next} = \text{newlabel()}$
- $S_2\text{-next} = S\text{-next} + 1$
- $S\text{-code} = S_1\text{-code} || \text{label}(S_1\text{-next}) || S_2\text{-code}$

- We assume that newlabel() creates a new label each time it is called, & that label(L) attaches label L to the next three address inst. to be generated.

P \rightarrow S \Rightarrow

The semantic rule associated with this production initialize s.next to a new label. P-code consist of s.code followed by the new label s.next.

assign \Rightarrow

Taken assign in the production S \rightarrow assign is a placeholder for assignment statements. s.code is simply assign.code.

S \rightarrow if(B) S₁ \Rightarrow

(create a new label B.true & attach it to the first three address inst. generated for stmt S₁, Thus jumps to B.true within the code for B will go to the code for S₁. By setting B.false to s.next we ensure that control will skip the code for S₁ if B evaluates to false.)

S \rightarrow if(B) S₁ else S₂ \Rightarrow

the code for the boolean expression B has jumps out of it to the first inst. of the code for S₁ if B is true, & to the first inst of the code for S₂ if B is false.

Further, control flows from both S₁ & S₂ to the three address inst. immediately following the code for S - its label is given by the inherited attribute s.next. An explicit goto s.next appears after the code for S₁ to skip over the code for S₂. No goto is

needed after s_2 , since $s_2 \cdot \text{next}$ is the same as $s_3 \cdot \text{next}$. \Rightarrow - 10

$s \rightarrow \text{while } CB \cdot s_1 \Rightarrow$ - 10

is formed from B-code & s_1 -code. we use a local variable begin to hold a new label attach to the first inst. for this while Stmt, which is also the first inst. for B. The inherited label $s \cdot \text{next}$ marks the inst. that control must flow to if B is False, hence, B-False is set to be $s \cdot \text{next}$. B-true is attached to the first inst. For s_1 , the code for B generates a jump to this label if B is true. After the code for s_1 we place the inst. goto begin, which causes a jump back to the beginning of the code for the boolean expression. $s \cdot \text{next}$ is set to this label begin, so jumps from within s_1 -code can go directly to begin

$s \rightarrow s_1 s_2 \Rightarrow$

consist of the code for s_1 , followed by the code for s_2 . The semantic rule maps the labels. the first inst. after the code for s_1 is the beginning of the code for s_2 , & the inst. after the code for s_2 is also the inst. after the code for s_1 .

Example:

1) Repeat s while B

$S \rightarrow TL$

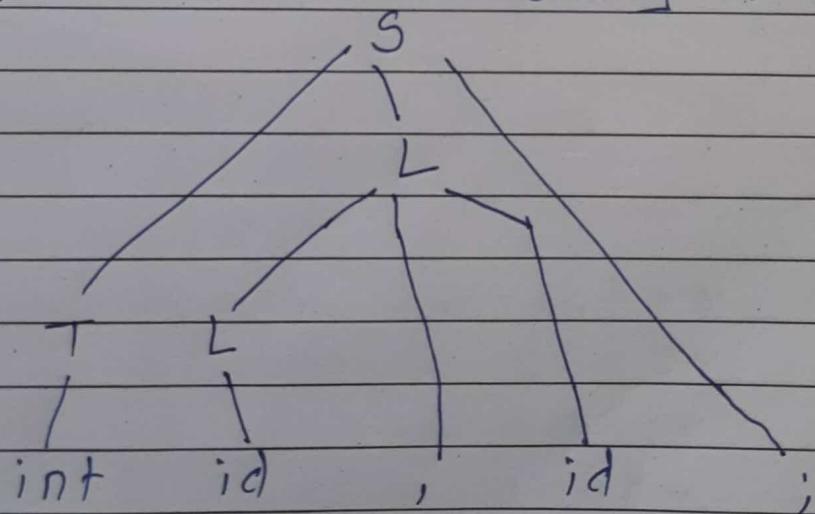
$T \rightarrow int \mid float$

$L \rightarrow L, id \mid id$

int id, id ;

stack	input	Action
\$	int id, id ; \$	shift
\$ int	id, id ; \$	Reduce int $\to T$
\$ T	a x mid, id ; \$	shift
\$ T id	, id ; \$	Reduce id $\to L$
\$ TL	, id ; \$	shift
\$ TL ,	id ; \$	shift
\$ TL , id	; \$	Reduce L, id $\to L$
\$ TL	; \$	Reduce TL $\to S$
\$ S	; \$	shift
\$ S ;	\$	Accept

As ; is the end of string pt is accepted.



Ambiguity in CFG:

* Recursive decent Parsing

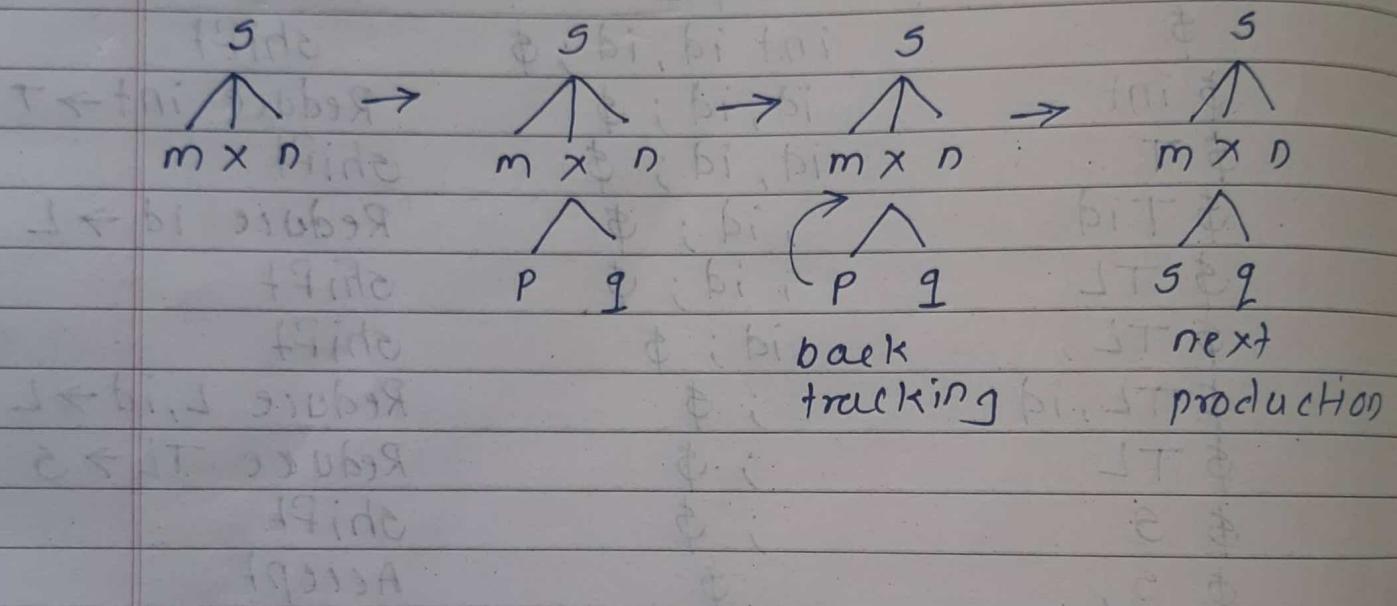
Example :

$S \rightarrow mxn \mid mzn$

$x \rightarrow pq \mid sq$

$z \rightarrow qr$

input $\Rightarrow msqn$



Examples

1. Consider the context free grammar

$$S \rightarrow SS + 1SS * 1a \quad \&$$

the string $aat+a*$

- a) Give the left most derivation for the string
- b) Give the rightmost derivation for the string
- c) Give the parse tree for the string
- d) Is the grammar ambiguous or unambiguous
Justify your answer.
- e) Describe the lang. generated by this grammar.

2. Repeat above for following grammar

a) $S \rightarrow 0S1101$ with string 000111

b) $S \rightarrow +SS1 * SS1a$ $+ * aaa$

c) $S \rightarrow S(S)S1E$ $((())(().$

d) $S \rightarrow S+S|SS|(S)|S*1a$ $(a+a)*a$

e) $S \rightarrow (L)1a$ & $L \rightarrow L, S|S$ $((a,a),a,(a))$

f) $\rightarrow aSbS|bSaS|E$ aabbab.

Solution

a) $aa+a*$

$$S \rightarrow SS *$$

$$\rightarrow S+a* \rightarrow SS+S*$$

$$\rightarrow SS+a* \rightarrow aS+S*$$

$$\rightarrow \quad \quad \quad \rightarrow aa+S*$$

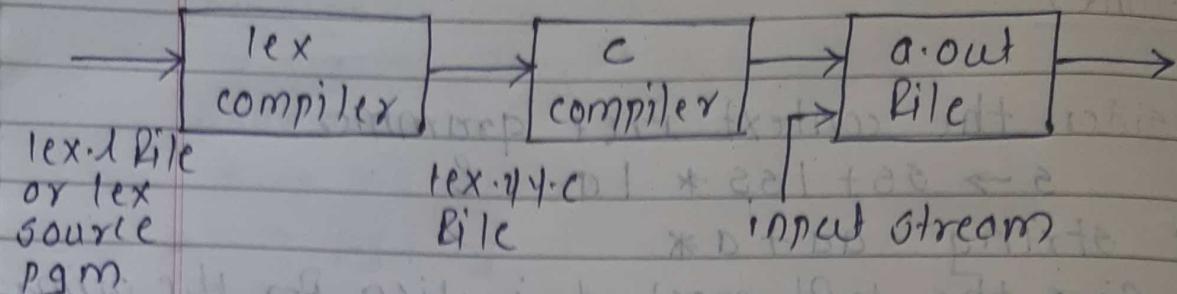
$$\rightarrow aa+a*$$

b) $S \rightarrow SS *$

$$\rightarrow Sa* \rightarrow SS+a* \rightarrow S+a+a* \rightarrow aa+a*$$

c) unambiguous

d) The set of all postfix exp? consist of addition & multiplication



Program Structure

Definition Section -

`y {`

↳ **1>** **2>** **3>** **4>** **5>** **6>** **7>** **8>** **9>** **10>** **11>** **12>** **13>** **14>** **15>** **16>** **17>** **18>** **19>** **20>** **21>** **22>** **23>** **24>** **25>** **26>** **27>** **28>** **29>** **30>** **31>** **32>** **33>** **34>** **35>** **36>** **37>** **38>** **39>** **40>** **41>** **42>** **43>** **44>** **45>** **46>** **47>** **48>** **49>** **50>** **51>** **52>** **53>** **54>** **55>** **56>** **57>** **58>** **59>** **60>** **61>** **62>** **63>** **64>** **65>** **66>** **67>** **68>** **69>** **70>** **71>** **72>** **73>** **74>** **75>** **76>** **77>** **78>** **79>** **80>** **81>** **82>** **83>** **84>** **85>** **86>** **87>** **88>** **89>** **90>** **91>** **92>** **93>** **94>** **95>** **96>** **97>** **98>** **99>** **100>** **101>** **102>** **103>** **104>** **105>** **106>** **107>** **108>** **109>** **110>** **111>** **112>** **113>** **114>** **115>** **116>** **117>** **118>** **119>** **120>** **121>** **122>** **123>** **124>** **125>** **126>** **127>** **128>** **129>** **130>** **131>** **132>** **133>** **134>** **135>** **136>** **137>** **138>** **139>** **140>** **141>** **142>** **143>** **144>** **145>** **146>** **147>** **148>** **149>** **150>** **151>** **152>** **153>** **154>** **155>** **156>** **157>** **158>** **159>** **160>** **161>** **162>** **163>** **164>** **165>** **166>** **167>** **168>** **169>** **170>** **171>** **172>** **173>** **174>** **175>** **176>** **177>** **178>** **179>** **180>** **181>** **182>** **183>** **184>** **185>** **186>** **187>** **188>** **189>** **190>** **191>** **192>** **193>** **194>** **195>** **196>** **197>** **198>** **199>** **200>** **201>** **202>** **203>** **204>** **205>** **206>** **207>** **208>** **209>** **210>** **211>** **212>** **213>** **214>** **215>** **216>** **217>** **218>** **219>** **220>** **221>** **222>** **223>** **224>** **225>** **226>** **227>** **228>** **229>** **230>** **231>** **232>** **233>** **234>** **235>** **236>** **237>** **238>** **239>** **240>** **241>** **242>** **243>** **244>** **245>** **246>** **247>** **248>** **249>** **250>** **251>** **252>** **253>** **254>** **255>** **256>** **257>** **258>** **259>** **260>** **261>** **262>** **263>** **264>** **265>** **266>** **267>** **268>** **269>** **270>** **271>** **272>** **273>** **274>** **275>** **276>** **277>** **278>** **279>** **280>** **281>** **282>** **283>** **284>** **285>** **286>** **287>** **288>** **289>** **290>** **291>** **292>** **293>** **294>** **295>** **296>** **297>** **298>** **299>** **300>** **301>** **302>** **303>** **304>** **305>** **306>** **307>** **308>** **309>** **310>** **311>** **312>** **313>** **314>** **315>** **316>** **317>** **318>** **319>** **320>** **321>** **322>** **323>** **324>** **325>** **326>** **327>** **328>** **329>** **330>** **331>** **332>** **333>** **334>** **335>** **336>** **337>** **338>** **339>** **340>** **341>** **342>** **343>** **344>** **345>** **346>** **347>** **348>** **349>** **350>** **351>** **352>** **353>** **354>** **355>** **356>** **357>** **358>** **359>** **360>** **361>** **362>** **363>** **364>** **365>** **366>** **367>** **368>** **369>** **370>** **371>** **372>** **373>** **374>** **375>** **376>** **377>** **378>** **379>** **380>** **381>** **382>** **383>** **384>** **385>** **386>** **387>** **388>** **389>** **390>** **391>** **392>** **393>** **394>** **395>** **396>** **397>** **398>** **399>** **400>** **401>** **402>** **403>** **404>** **405>** **406>** **407>** **408>** **409>** **410>** **411>** **412>** **413>** **414>** **415>** **416>** **417>** **418>** **419>** **420>** **421>** **422>** **423>** **424>** **425>** **426>** **427>** **428>** **429>** **430>** **431>** **432>** **433>** **434>** **435>** **436>** **437>** **438>** **439>** **440>** **441>** **442>** **443>** **444>** **445>** **446>** **447>** **448>** **449>** **450>** **451>** **452>** **453>** **454>** **455>** **456>** **457>** **458>** **459>** **460>** **461>** **462>** **463>** **464>** **465>** **466>** **467>** **468>** **469>** **470>** **471>** **472>** **473>** **474>** **475>** **476>** **477>** **478>** **479>** **480>** **481>** **482>** **483>** **484>** **485>** **486>** **487>** **488>** **489>** **490>** **491>** **492>** **493>** **494>** **495>** **496>** **497>** **498>** **499>** **500>** **501>** **502>** **503>** **504>** **505>** **506>** **507>** **508>** **509>** **510>** **511>** **512>** **513>** **514>** **515>** **516>** **517>** **518>** **519>** **520>** **521>** **522>** **523>** **524>** **525>** **526>** **527>** **528>** **529>** **530>** **531>** **532>** **533>** **534>** **535>** **536>** **537>** **538>** **539>** **540>** **541>** **542>** **543>** **544>** **545>** **546>** **547>** **548>** **549>** **550>** **551>** **552>** **553>** **554>** **555>** **556>** **557>** **558>** **559>** **560>** **561>** **562>** **563>** **564>** **565>** **566>** **567>** **568>** **569>** **570>** **571>** **572>** **573>** **574>** **575>** **576>** **577>** **578>** **579>** **580>** **581>** **582>** **583>** **584>** **585>** **586>** **587>** **588>** **589>** **590>** **591>** **592>** **593>** **594>** **595>** **596>** **597>** **598>** **599>** **600>** **601>** **602>** **603>** **604>** **605>** **606>** **607>** **608>** **609>** **610>** **611>** **612>** **613>** **614>** **615>** **616>** **617>** **618>** **619>** **620>** **621>** **622>** **623>** **624>** **625>** **626>** **627>** **628>** **629>** **630>** **631>** **632>** **633>** **634>** **635>** **636>** **637>** **638>** **639>** **640>** **641>** **642>** **643>** **644>** **645>** **646>** **647>** **648>** **649>** **650>** **651>** **652>** **653>** **654>** **655>** **656>** **657>** **658>** **659>** **660>** **661>** **662>** **663>** **664>** **665>** **666>** **667>** **668>** **669>** **670>** **671>** **672>** **673>** **674>** **675>** **676>** **677>** **678>** **679>** **680>** **681>** **682>** **683>** **684>** **685>** **686>** **687>** **688>** **689>** **690>** **691>** **692>** **693>** **694>** **695>** **696>** **697>** **698>** **699>** **700>** **701>** **702>** **703>** **704>** **705>** **706>** **707>** **708>** **709>** **710>** **711>** **712>** **713>** **714>** **715>** **716>** **717>** **718>** **719>** **720>** **721>** **722>** **723>** **724>** **725>** **726>** **727>** **728>** **729>** **730>** **731>** **732>** **733>** **734>** **735>** **736>** **737>** **738>** **739>** **740>** **741>** **742>** **743>** **744>** **745>** **746>** **747>** **748>** **749>** **750>** **751>** **752>** **753>** **754>** **755>** **756>** **757>** **758>** **759>** **760>** **761>** **762>** **763>** **764>** **765>** **766>** **767>** **768>** **769>** **770>** **771>** **772>** **773>** **774>** **775>** **776>** **777>** **778>** **779>** **780>** **781>** **782>** **783>** **784>** **785>** **786>** **787>** **788>** **789>** **790>** **791>** **792>** **793>** **794>** **795>** **796>** **797>** **798>** **799>** **800>** **801>** **802>** **803>** **804>** **805>** **806>** **807>** **808>** **809>** **810>** **811>** **812>** **813>** **814>** **815>** **816>** **817>** **818>** **819>** **820>** **821>** **822>** **823>** **824>** **825>** **826>** **827>** **828>** **829>** **830>** **831>** **832>** **833>** **834>** **835>** **836>** **837>** **838>** **839>** **840>** **841>** **842>** **843>** **844>** **845>** **846>** **847>** **848>** **849>** **850>** **851>** **852>** **853>** **854>** **855>** **856>** **857>** **858>** **859>** **860>** **861>** **862>** **863>** **864>** **865>** **866>** **867>** **868>** **869>** **870>** **871>** **872>** **873>** **874>** **875>** **876>** **877>** **878>** **879>** **880>** **881>** **882>** **883>** **884>** **885>** **886>** **887>** **888>** **889>** **890>** **891>** **892>** **893>** **894>** **895>** **896>** **897>** **898>** **899>** **900>** **901>** **902>** **903>** **904>** **905>** **906>** **907>** **908>** **909>** **910>** **911>** **912>** **913>** **914>** **915>** **916>** **917>** **918>** **919>** **920>** **921>** **922>** **923>** **924>** **925>** **926>** **927>** **928>** **929>** **930>** **931>** **932>** **933>** **934>** **935>** **936>** **937>** **938>** **939>** **940>** **941>** **942>** **943>** **944>** **945>** **946>** **947>** **948>** **949>** **950>** **951>** **952>** **953>** **954>** **955>** **956>** **957>** **958>** **959>** **960>** **961>** **962>** **963>** **964>** **965>** **966>** **967>** **968>** **969>** **970>** **971>** **972>** **973>** **974>** **975>** **976>** **977>** **978>** **979>** **980>** **981>** **982>** **983>** **984>** **985>** **986>** **987>** **988>** **989>** **990>** **991>** **992>** **993>** **994>** **995>** **996>** **997>** **998>** **999>** **1000>** **1001>** **1002>** **1003>** **1004>** **1005>** **1006>** **1007>** **1008>** **1009>** **1010>** **1011>** **1012>** **1013>** **1014>** **1015>** **1016>** **1017>** **1018>** **1019>** **1020>** **1021>** **1022>** **1023>** **1024>** **1025>** **1026>** **1027>** **1028>** **1029>** **1030>** **1031>** **1032>** **1033>** **1034>** **1035>** **1036>** **1037>** **1038>** **1039>** **1040>** **1041>** **1042>** **104**

Flex (Fast Lexical Analyzer Generator)

- lex/flex to recognize patterns in text, the pattern must be described by a regular expression
- The input to lex/flex is a machine readable set of regular expressions.
- The input is in the form of pairs of regular expressions & c code, called rules.
- It generates as output a c source file, lex.yyy.c, which defines a routine `yyflex()`.
- This file is compiled & linked with the -lfl library to produce an executable.
- Whenever the executable is run, it analyzes its input for occurrences of the regular expression. Whenever it finds one, it executes the corresponding c code.
- It is a tool/computer pgm for generating lexical analyzers (scanners).
- It is used together with yacc or Bison parser generator.
- Bison produces parser from the input file provided by the user.
- The function `yyflex()` is automatically generated by the flex when it is provided with a .l file & this `yyflex()` function is expected by parser to call to retrieve tokens from current token stream.