

Would it be better if the phases of the compiler are combined into a single phase?

Adv

* By Grouping the Diff. phases into some no. of passes reducing the time in generating target code.

* How the Grouping is done?

Combining the phases into a pass in such a way that the operations in each ~~every phase of a compiler are grouped into~~ are incorporated during the pass.

g) tokens are translated directly into intermediate code.

* Minimizing no. of passes improves time efficiency.

Differentiate compilers from interpreters.

COMPILER	INTERPRETER
→ performs the translation is done as whole	→ performs statement by statement translation.
→ Execution is faster	→ It is slower.
→ Request more memory	Efficient memory utilization.
→ Debugging a code is harder	Debugging is quite easier.
→ Both are generating intermediate object code	It stops translation which first error occurred.
Ex: C, C++	Ex: Python, Ruby, Perl, PHP.

An arithmetic expression with unbalanced parenthesis is lexical or syntax error. Comment on it.

2) Syntactical Error.	Missing semicolon ; unbalanced parenthesis ;
In parser Error handler	- parser Get the present error of and report to parser in clearly and accuracy
Four Common Recovery Error Statistics	Implement parser.
* Panic mode	
* Statement level recovery.	
* Error productions	
* Global corrections.	

A compiler that translates a high-level language into another high-level language is called a source- to-source translator. What advantages are there to use C as a target language for a compiler?

➲ 1.1.4

A compiler that translates a high-level language into another high-level language is called a *source-to-source* translator. What advantages are there to using C as a target language for a compiler?

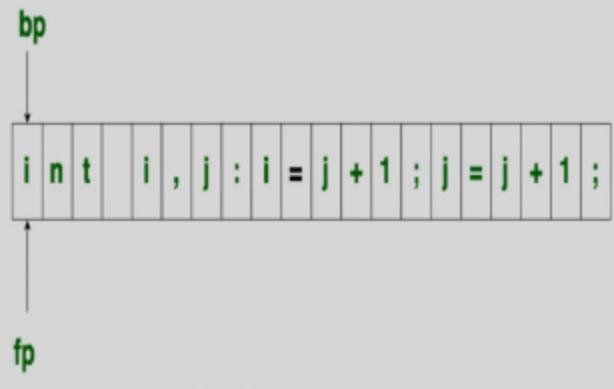
➲ Answer

For the C language there are many compilers available that compile to almost every hardware.

Point out why is buffering used in lexical analysis? What are the commonly used buffering methods?

Input Buffering in Compiler Design

The lexical analyzer scans the input from left to right one character at a time. It uses two pointers begin ptr(bp) and forward to keep track of the pointer of the input scanned.



TYPES

1. One buffer scheme
2. Two buffer scheme

Obtain the regular expressions for the following sets:

- a) The set of all strings over {a, b} beginning and ending with 'a'.
- b) The set of all strings over {b} such that string belong to { b^2, b^5, b^8, \dots }.

③ (a) the set of all strings over $\{a, b\}$ beginning and ending with 'a'

$$a(a|b)^*a$$

(b) the set of all strings over $\{b\}$ such that string belongs to $\{b_2, b_5, b_8, \dots\}$

$$bb(bbb)^*$$

Comment on the efficiency of the compiler if the number of passes in compilation is increased

Adv

* By Grouping the Diff. phases into some no. of passes reducing the time in generating target code.

* How the Grouping is done?

Combining the phases into a pass in such a way that the operations in each ~~every phase of a compiler are grouped into~~ are incorporated during the pass.

g) tokens are translated directly into intermediate code.

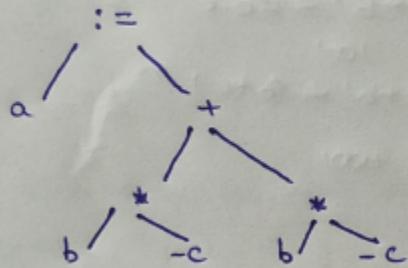
* Minimizing no. of passes improves time efficiency.

What are the various representations of Intermediate languages?

- ↳ 4) Intermediate Code Generation
 - ↳ It produce the intermediate representation of source program
 - ↳ It forms the
 - postfix notation
 - Three Address code (max operands - only 3)
 - Syntax Tree.

Construct the syntax tree for the following assignment statement: $a := b^+ - c + b^+ - c$.

② construct syntax tree $a := b^+ - c + b^+ - c$.



Identify the type of error.

$y = 3 + * 5;$

int a = "hello";

Justify this with the suitable points..

④ $y = 3 + * 5;$

int a = "hello"

Justify this with the suitable points

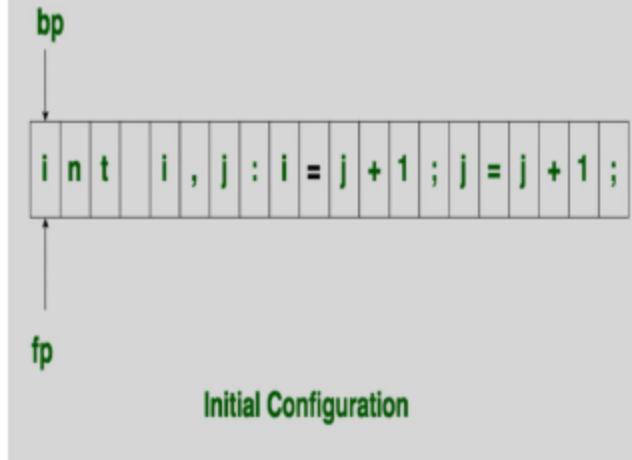
$y = 3 + * 5 \rightarrow$ Syntax error

int a = "hello" \rightarrow Semantic error

Point out why is buffering used in lexical analysis? What are the commonly used buffering methods?

Input Buffering in Compiler Design

The lexical analyzer scans the input from left to right one character at a time. It uses two pointers begin ptr(bp) and forward to keep track of the pointer of the input scanned.



TYPES

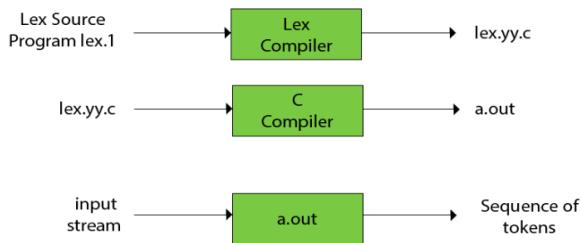
1. One buffer scheme
2. Two buffer scheme

Describe the language denoted by the regular expression $a^*ba^*ba^*ba^*$.

How can a lexical analyzer be constructed using Lex tool?

The function of Lex is as follows:

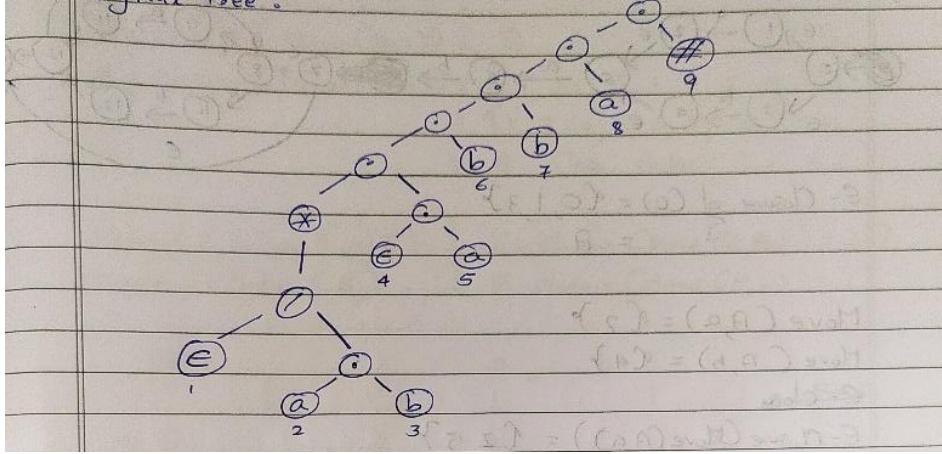
- o Firstly lexical analyzer creates a program lex.1 in the Lex language. Then Lex compiler runs the lex.1 program and produces a C program lex.yy.c.
- o Finally C compiler runs the lex.yy.c program and produces an object program a.out.
- o a.out is lexical analyzer that transforms an input stream into a sequence of tokens.



Draw the syntax tree for the augmented regular expression $(\epsilon \mid ab)^*(\epsilon a)bb\#$

3. Given RE: $(\epsilon \mid a \cdot b)^* \cdot (\epsilon \cdot a) \cdot b \cdot b \cdot a \cdot \#$

Syntactic Tree :



Mention the basic issues in parsing.

Issues in Parsing:

- Variable redeclaration
- Variable initialization before use
- Datatype mismatch for an operation
- Specification of syntax
- Representation of input after parsing.

What is the role of the parser in a compiler model?

The parser obtains a string of tokens from the lexical analyzer and verifies that the string can be generated by the grammar for the source language. It must report any syntax errors if occurs.

Consider the grammar $g = (V, T, P, S)$. Here $V = \{S, N, N_p, ADJ\}$ and $T = \{\text{and, eggs, ham, pencilgreen, cold, tasty}\}$. The set contains the following rules:

$N_p \mid N_p \text{ and } N_p \quad N \mid \text{eggs} \mid \text{ham} \mid \text{pencil}$

$N_p \mid ADJ \mid N_p \quad ADJ \mid \text{green} \mid \text{cold} \mid \text{tasty}$

$N \mid N$

Show that the grammar is ambiguous by constructing two different parse trees for the sentence "green eggs and ham".

2.2 Ambiguity

Consider the following grammar $\mathcal{G} = (V, \Sigma, R, S)$. Here

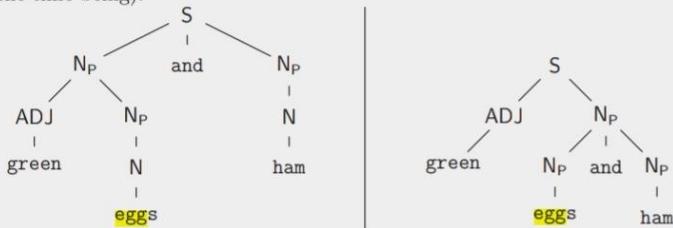
$$V = \{S, N, N_P, ADJ\} \quad \text{and} \quad \Sigma = \{\text{and}, \text{eggs}, \text{ham}, \text{pencilgreen}, \text{cold}, \text{tasty}, \dots\}.$$

The set R contains the following rules:

5

- $N_P \rightarrow N_P \text{ and } N_P$
- $N_P \rightarrow ADJ \ N_P$
- $N_P \rightarrow N$
- $N \rightarrow \text{eggs} \mid \text{ham} \mid \text{pencil} \mid \dots$
- $ADJ \rightarrow \text{green} \mid \text{cold} \mid \text{tasty} \mid \dots$
- \dots

Here are two possible parse trees for the string green **eggs** and ham (ignore the spacing for the time being).



The two parse trees group the words differently, creating a different meaning. In the first case, only the **eggs** are green. In the second, both the **eggs** and the ham are green.

A string w is **ambiguous** with respect to a grammar \mathcal{G} if w has more than one possible parse tree using the rules in \mathcal{G} .

Most grammars for practical applications are ambiguous. This is a source of real practical issues, because the end users of parsers (e.g. the compiler) need to be clear on which meaning is intended.

Eliminate Left recursion for the following grammar.

$S \rightarrow aB \mid aC \mid Sd \mid Se$

$B \rightarrow bBc \mid f$

$C \rightarrow g$

Da

Pa

12)

$g \rightarrow aB \mid ac \mid Sd \mid Se$

$B \rightarrow bBc \mid f$

$C \rightarrow g$

solution

$g \rightarrow aBS' \mid acS'$

$S' \rightarrow dS' \mid eS' \mid e$

$B \rightarrow bBc \mid f$

$C \rightarrow g.$

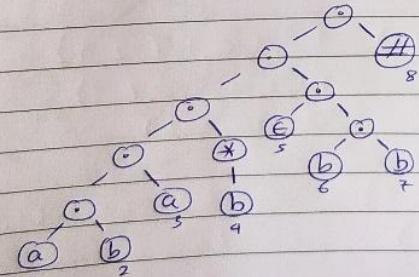
Write short notes on LEX.

Lex is a program that generates lexical analyzer. It is used with YACC parser generator. The lexical analyzer is a program that transforms an input stream into a sequence of tokens. It reads the input stream and produces the source code as output through implementing the lexical analyzer in the C program.

Draw the syntax tree for the augmented expression $(ab)ab^*(\epsilon bb)\#$

Given RE: $(a.b).a.b^*.(e.b.b)^*$

Syntax Tree:



Write a context free grammar that generates the set of all strings of a's and b's which is a palindrome.

Let P be language of palindromes with alphabet {a, b}. One can determine a CFG for P by finding a recursive decomposition.

CFG is:

$$P \rightarrow aPa \mid bPb \mid a \mid b \mid \epsilon$$

How will you eliminate left factor in a grammar?

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

$$A \rightarrow A\alpha \mid \beta$$

The above Grammar 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

$$A \rightarrow \beta A'$$

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

Draw transition diagrams for predictive parsers for the grammar

$$E \quad E + T \mid T$$

$$T \quad T^* F \mid F$$

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

9. Given:

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

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

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

LR \rightarrow left Recursion
LF \rightarrow Left Factoring

Removing LR using LF,

i) $E \rightarrow TE'$

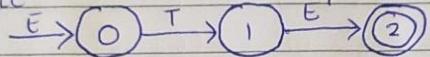
ii) $E' \rightarrow +TE' \mid \epsilon$

iii) $T \rightarrow FT'$

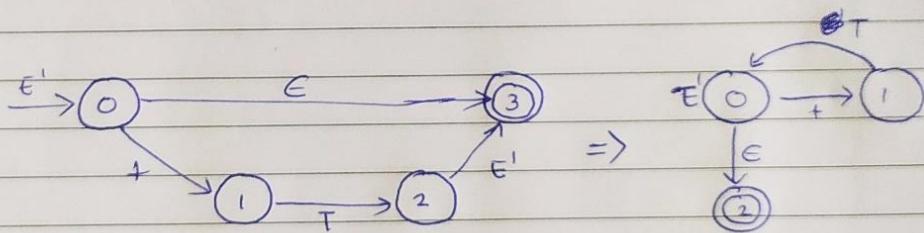
iv) $T' \rightarrow *FT' \mid \epsilon$

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

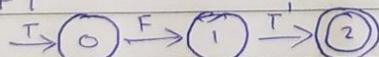
vi) $E \rightarrow TE'$



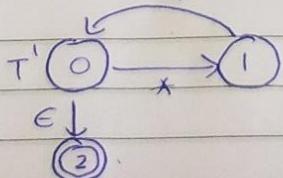
vii) $E' \rightarrow +TE' \mid \epsilon$



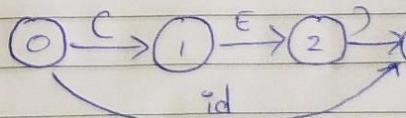
viii) $T \rightarrow FT'$



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



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



Consider the grammar $g = (V, T, P, S)$. Here $V = \{S, N, N_p, ADJ\}$ and $T = \{\text{and}, \text{eggs}, \text{ham}, \text{pencil green}, \text{cold}, \text{tasty}\}$. The set contains the following rules:

Np \quad Np and Np N \quad eggs | bread | pencil |

Np ADJ Np ADJ boiled | cold | tasty|

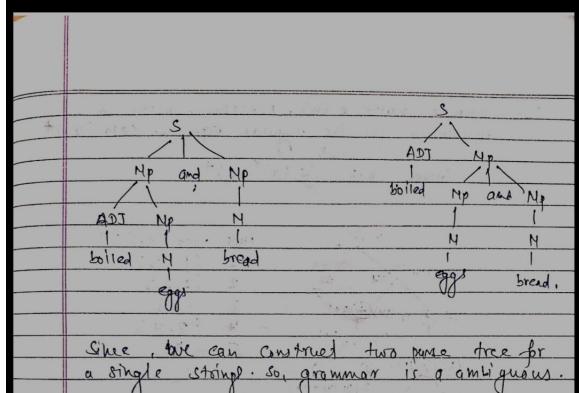
Np N

Is this grammar ambiguous for the sentence "boiled eggs and bread".

11. Consider the grammar given below:
V = {S, N, Np, ADJ} and T = {and, eggs, ham, pencil, green, cold, tasty}. The set contains the following rules.
 $Np \rightarrow Np \text{ and } Np$ $N \rightarrow \text{eggs} \mid \text{bread} \mid \text{pencil}$
 $Np \rightarrow \text{ADJ } Np$ $\text{ADJ} \rightarrow \text{boiled} \mid \text{cold} \mid \text{tasty}$
 $Np \rightarrow N$

Is this grammar ambiguous for the sentence "boiled eggs and bread".
So?

Scanned with CamScanner



What is a front end and back end of a compiler? What are the advantages of breaking up the compiler functionality into these two distinct stages?

Obtain the regular expressions for the following having $\Sigma = \{0,1\}$

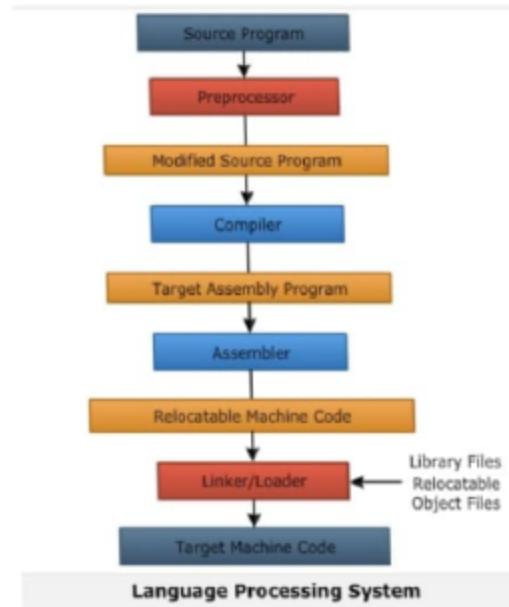
- strings of 0's and 1's beginning with 0 and ending with 1.
- $\{x \mid x \text{ contains an even number of 0's or an odd number of 1's}\}$

(a) strings of 0's and 1's beginning with 0 and ending with 1.
 $0(011)^*1$

(b) $\{x \mid x \text{ contains an even number of 0's or an odd number of 1's}\}$
 $(00)^+ / 1(11)^*$

Depict diagrammatically and explain how a language is processed?

Language Processing System



Language Processing System

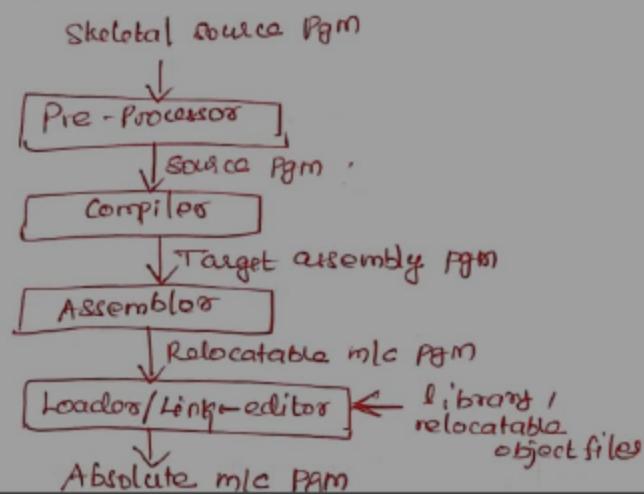
PreProcessor

- > A Source pgm may be divided into modules stored in separate file. The task of collecting ~~macro~~ the source pgm is entrusted to a separate pgm called Preprocessor.
- > It also expands macros into source language statement.

Scanned by CamScanner

Compilers

Compiler is a pgm that takes source pgm as input & produces assembly language pgm as output.



Assembler

It is a pgm that converts assembly language pgm into m/c lang. pgm. It produces re-locatable m/c code as its output

Loader and Link editor

- > The relocatable m/c code has to be ~~linked~~ with other re-locatable object files & library files into the code that actually runs on the machine
- > the linker resolves external memory addresses, where the code in one file may refer to a location in another file
- > The loader puts together the entire executable object files into memory for execution

Preprocessor:

It is a program that reads the source code and prepares it for the translator
The result of preprocessing is called the compiler

Roles of Preprocessor:

Copies all library functions of a file into a source code
Expand macros into a source code

Example: #include<stdio.h> here # is a processor directive. It copies all library functions of stdio.h (standard input output) file into a program.

Compiler:

Compiler is a program that reads a program in high level language known as source code as its input and converts it into an equivalent program in low level language known as object code as its output.

Assembler

The compiler may produce an assembly language as easier to produce as output and is easier to debug .The assembly language is then processed by a program called assembler that produce relocatable machine code as its output.

- **Relocatable files are not human readable format**

Linker/ Loader

The process of linking user functions and system functions to a final executable code known as linker

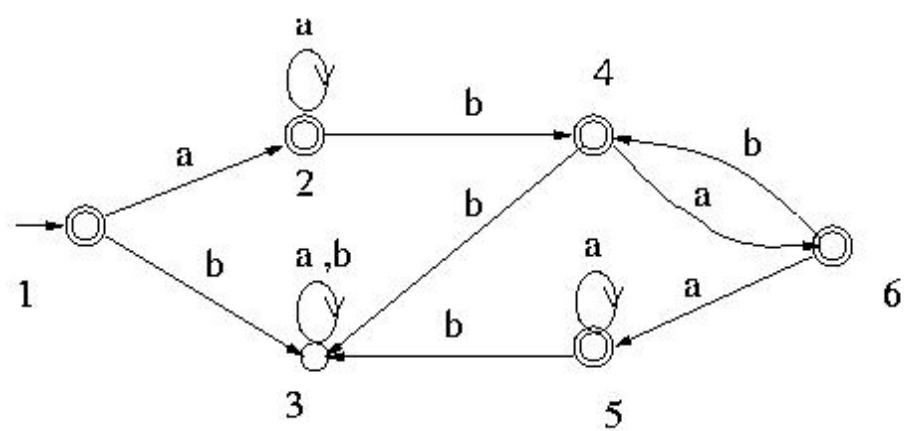
Loader:

Once the program is linked, it is ready for execution. The process of loading the program into primary memory for execution is known as loader

- **The loader finds sufficient place for the program to be loaded for execution.**



minimize the number of states of the following DFA.



Transition Table

states	a	b
* 1	2	3
* 2	2	4
* 3	3	3
* 4	6	3
* 5	5	3
* 6	5	4

$$\pi_0 = \{ q_1^0, q_2^0 \}$$

where, q_1^0 is set of all final states

$q_2^0 = Q - q_1^0$, Q is a set of all the states in DFA

$$q_1^0 = \{ 1, 2, 4, 5, 6 \}$$

$$q_2^0 = \{ 3 \}$$

$$\pi_0 = \{ \underbrace{\{ 1, 2, 4, 5, 6 \}}_1, \underbrace{\{ 3 \}}_2 \}$$

	q.3	2	4	5	6
9	1	1	1	1	1
5	2	1	2	2	1

$$\therefore \pi_1 = \{ \underbrace{\{ 3 \}}_2, \underbrace{\{ 1, 4, 5 \}}_3, \underbrace{\{ 2, 6 \}}_4 \}$$

Again;

	1	2	4	5	6
9	4	4	4	3	3
5	2	3	2	2	3

	1	2	4	5	6	7
9	4	4	4	3	3	
b	2	3	2	2	3	

Scanned with CamScanner

$$\therefore \pi_2 = \left\{ \frac{\{3\}}{2}, \frac{\{1,4\}}{5}, \frac{\{2\}}{6}, \frac{\{5\}}{7}, \frac{\{6\}}{8} \right\}$$

	1	2	4	5	6
9	6	6	8	7	7
b	2	5	2	2	5

$$\therefore \pi_3 = \left\{ \frac{\{3\}}{2}, \frac{\{1\}}{9}, \frac{\{2\}}{6}, \frac{\{4\}}{10}, \frac{\{5\}}{7}, \frac{\{6\}}{8} \right\}$$

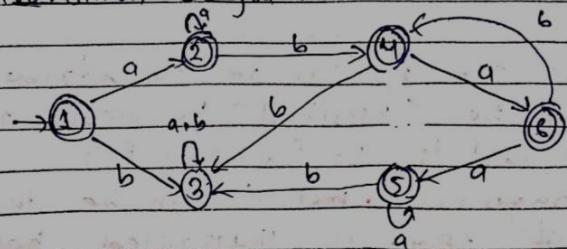
	1	2	4	5	6
9	6	6	8	7	7
b	2	10	2	2	10

$$\therefore \pi_4 = \left\{ \{3\}, \{1\}, \{2\}, \{4\}, \{5\}, \{6\} \right\}$$

Transition Table

states	input a		input b	
	1	2	3	4
1	2		3	
2		2		4
3		3		3
4		6		3
5		5		3
6		5		4

Transition Diagram:



This is the required minimized DFA.

To verify the syntax of programming language constructs, compilers use CFG. Why not Regular expression?. Justify with examples.

Regular Expressions are capable of describing the syntax of Tokens. Any syntactic construct that can be described by Regular Expression can also be described by the Context free grammar.

grammar that's why we are using Regular Expression.

There are several reasons and they are:

Regular Expressions	Context-free grammar
Lexical rules are quite simple in case of Regular Expressions.	Lexical rules are difficult in case of Context free grammar.
Notations in regular expressions are easy to understand.	Notations in Context free grammar are quite complex.
A set of string is defined in case of Regular Expressions.	In Context free grammar the language is defined by the collection of productions.
It is easy to construct efficient recognizer from Regular Expressions.	By using the context free grammar, it is very difficult to construct the recognizer.
There is proper procedure for lexical and syntactical analysis in case of Regular Expressions.	There is no specific guideline for lexical and syntactic analysis in case of Context free grammar.
Regular Expressions are most useful for describing the structure of lexical construct such as identifiers, constant etc.	Context free grammars are most useful in describing the nested chain structure or syntactic structure such as balanced parenthesis, if else etc. and these can't be define by Regular Expression.

Example:

Take the problem of palindrome language, which cannot be described by means of regular expression. That is: $L = \{w | w = w^R\}$ is not a regular language, but it can be described by means of CFG, as illustrated below:-

Scanned with CamScanner

$$G_1 = (V, \Sigma, P, S)$$

where

$$V = \{\emptyset, z, N\}$$

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

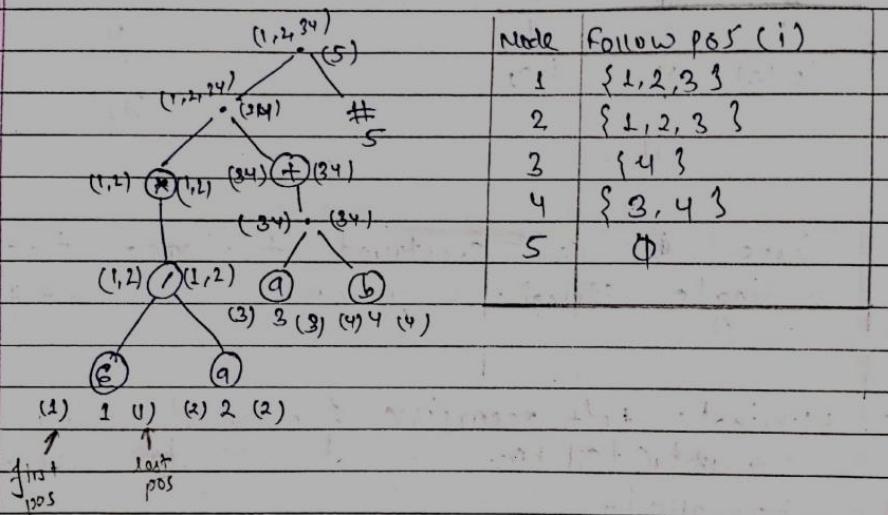
$$P = \{\emptyset \rightarrow z \mid \emptyset \rightarrow N \mid \emptyset \rightarrow \epsilon \mid z \rightarrow 0z0 \mid N \rightarrow 1z1\}$$

This grammar describes palindrome languages, such as: 1001, 110011, 1111 etc.

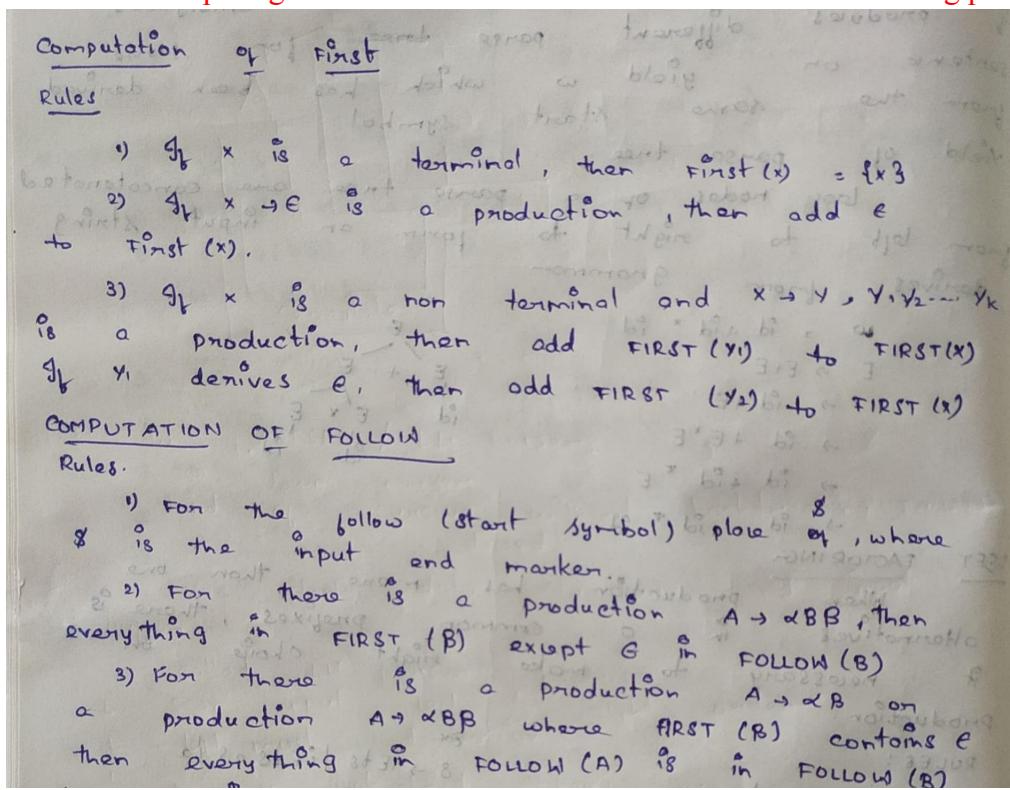
Compute nullable(n), FIRSTPOS, LASTPOS and FOLLOWPOS for the regular expression $(\epsilon \mid a)^*(ab)^+$

13. Compute nullable (η), FIRSTPOS, LASTPOS and FOLLOWPOS for the regular expression $(\epsilon | a)^*(ab)^*$.

Given, RE: $(\epsilon | a)^*(ab)^*$



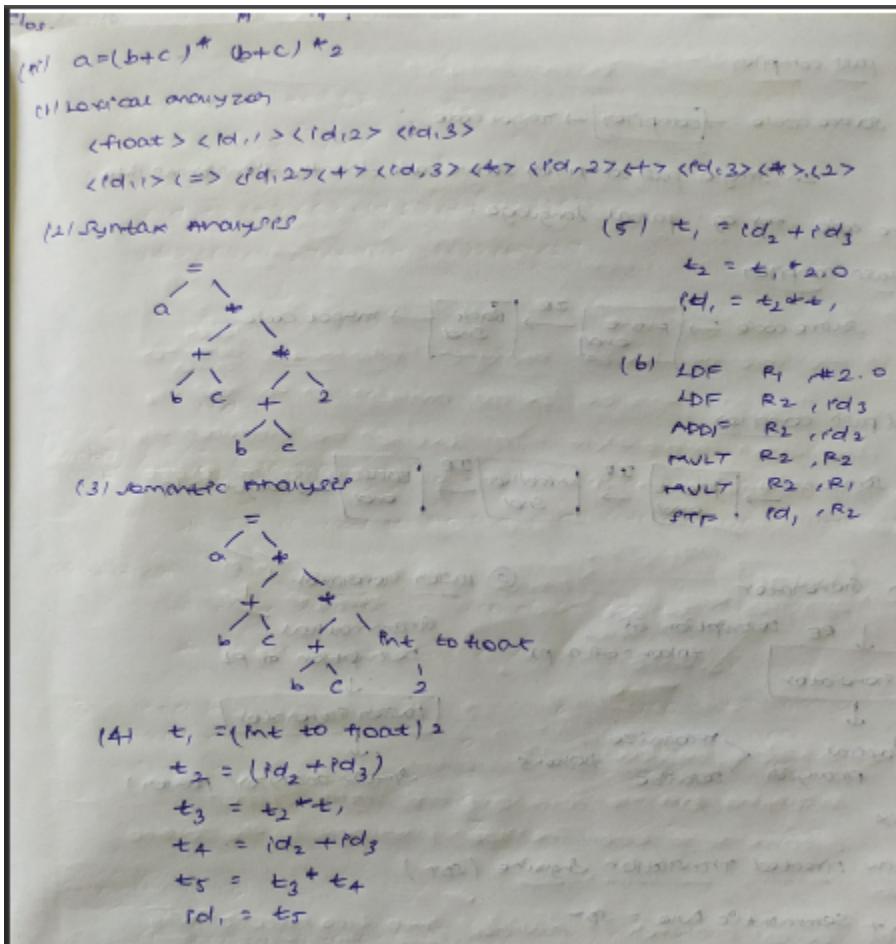
Write the rules for computing the FIRST and FOLLOW functions for constructing predictive parsing table.



(i)

(ii)

Generate the target code for the following code fragment. Assume $a=(b+c) * (b+c) *$ as input to the lexical analyzer and try to show the intermediate outputs of all the phases from lexical analyzer to code generator.



How do you speed up the reading the source program? What are the specialized algorithms used by the compiler? Did all of them read the source program with same time or they differ? Justify with suitable claims and figures.

Can you identify the lexemes that make up the tokens in the following program segment? Write the corresponding tokens, lexemes and pattern.

```
void swap(int i, int j)
{   int t;   t=i;   i=j;   j=t;
```

Do compiled languages beat interpreted languages in performance? What is the difference between those languages? Justify this with appropriate examples.

Construct minimized Deterministic Finite Automata to accept the regular expression : $(0+1)^* (00+11) (0+1)^*$.

Illustrate the output of each phase of compilation while translating the following assignment statement $i = i*70+j+2;$.

How to solve the source program to target machine code by using language processing system?

Can you identify and list the systems/tools to help with the compiler-writing process are often been referred to as compiler-compilers, compiler-generators or translator-writing systems? Why are they called

compiler-compilers? Write the representation and implementation techniques? What flexibility do they provide?

(f) Compiler construction tools:

- The compilers writer (still developing)
many modern software development environments
tools like
 - * language editors;
 - * debuggers
 - * version managers
 - * profilers
 - * test harnesses etc.
- More specialized SW tools have been
created to implement phases of a compiler
- These tools use specialized languages
use more sophisticated Algo's.
- The most successful tools are those
that hide the details of the
generation of the Algorithm &
can easily integrated into the
remainder of the compiler.

(b) Some commonly used compiler constructs

① Parser Generators

Automatically produce Syntax Analyzer
from a grammatical description of a (PL)
Programming Language.

④ scanner Generators

Generates lexical Analyser from a regular expression description of the tokens of a language

⑤ Syntax-directed translation engines

[Parse tree \rightarrow (to) \rightarrow Intermediate code],
produce collections of routines for walking
a parse tree & generating intermediate code

⑥ code Generators

produce a code generator from a collection
of rules for translating each operation of
intermediate lang. into the machine lang.
for a target machine.

⑦ Data flow Analysis Engines

\rightarrow It facilitates the gathering of
information about how values are
transmitted from one part of a program
to each other part.

\rightarrow It is a key part of code optimization

⑧ compiler construction tool kits

\rightarrow Provide an integrated set of
routines for constructing various
phases of a compiler.

Construction tools

- They are also known as a **compiler-compilers**, **compiler-generators** or **translator**
- These **tools** use **specific language or algorithm** for specifying and **implementing the component** of the **compiler**

1. Parser generators.
2. Scanner generators.
3. Syntax-directed translation engines.
4. Data-flow analysis engines.
5. Automatic code generators.

Parser Generators

Input: Grammatical description of a programming language

Parser generator takes the grammatical

Description(BNF) of a programming language produces a syntax analyzer

Production Rule	Meaning
$S \rightarrow NP \cdot VP$	A sentence is a noun phrase followed by a verb phrase
$NP \rightarrow DT \cdot NN$	A noun phrase is a determiner followed by a noun
$NP \rightarrow NP \cdot PP$	A noun phrase can also be another noun phrase followed by a prepositional phrase
$VP \rightarrow VV$	A verb phrase is simply an intransitive verb
$VP \rightarrow VV \cdot NP$	A verb phrase is a transitive verb followed by a noun phrase
$VP \rightarrow VP \cdot PP$	A verb phrase can also be another verb phrase followed by a prepositional phrase
$PP \rightarrow IN \cdot NP$	A prepositional phrase is a preposition followed by a noun phrase
$DT \rightarrow \text{the}$	The is a determiner
$NN \rightarrow \text{dog} \mid \text{cat}$	Dog, cat, and tree are nouns
$IN \rightarrow \text{on}$	On is a preposition
$VV \rightarrow \text{chased} \mid \text{searched}$	Both chased and searched are transitive verbs
$VV \rightarrow \text{saw}$	Saw is an intransitive verb

The generated code is a **parser**, which takes a **sequence of characters** and tries to **match the sequence** against the **grammar**

- Antlr and yacc are popular parser generators

Parser generators for C/C++: [Bison \(in HTML\)](#), [Bison](#), [Yacc](#)

Scanner Generators

Input: Regular expression description of the tokens of a language

Output: Lexical analyzers.

Scanner generator generates lexical analyzers from a regular expression

description of the tokens of a language

- Describing tokens(RE)
- Recognizing (or accepted) -via mathematical model called finite-state acceptor

Scanner generators for C/C++: [Flex](#), [Lex](#)

Syntax-directed Translation Engines

Grammar + semantic rule = SDT (syntax directed translation)

- Every non-terminal can get one or more than one attribute based on the type of the attribute
- The value of these attributes is evaluated by the semantic rules associated with the production rule
- $E \rightarrow E+T \{ E.val = E.val + T.val \} PR\#1$
- $E \rightarrow T \{ E.val = T.val \} PR\#2$
- $T \rightarrow T^*F \{ T.val = T.val * F.val \} PR\#3$
- $T \rightarrow F \{ T.val = F.val \} PR\#4$
- $F \rightarrow INTLIT \{ F.val = INTLIT.lexval \} PR\#5$

Data-flow Analysis Engines

Data-flow analysis is a key part of code optimization.

- These are techniques that derive information about the **flow of data** along program execution paths
- There is an infinite number of paths through a program and there is no bound on the length of a path
- Data-flow analysis engine gathers the information on **values transmitted from one part of a program to each of the other** parts
- Highly used for debugging

Automatic Code Generators

Input: Intermediate language.

Output: Machine language.

Code-generator takes a collection of rules that define the **translation of each operation of the intermediate language into the machine language** for a target machine.

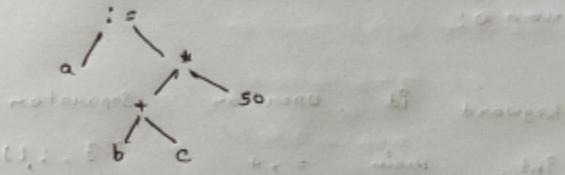
Discuss the working of various phases of Compiler. Interpret the output for each phases for the expression $a := b + c * 50$.

$a := b + c * 50$

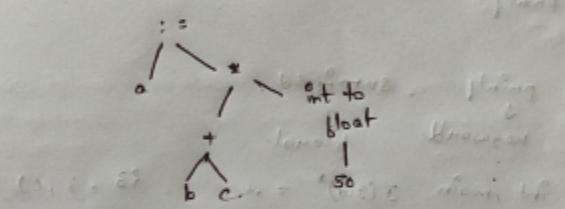
i) lexical.

$\langle \text{id}, 1 \rangle \langle := \rangle \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle * \rangle \langle \text{so} \rangle$

ii) syntax



iii) semantic.



iv) intermediate code generation.

$$\begin{aligned} t_1 &= \text{int to float (50)} \\ t_2 &= \text{id}_3 + \text{id}_2 \\ t_3 &= t_1 * t_2 \\ \text{id}_1 &:= t_3 \end{aligned}$$

v) Code optimization

$$\begin{aligned} t_1 &= \text{id}_3 + \text{id}_2 \\ \text{id}_1 &:= t_1 * 50 \end{aligned}$$

vi) Code Generation.

$$\begin{aligned} \text{LDF } R_1 & \text{ id}_3 \\ \text{ADDF } R_1 & , \text{id}_2 \\ \text{MULF } R_1 & , \# 50 \\ \text{ASSIGN } & \text{id}_1, R_1 \end{aligned}$$

A compiler can choose one of the two options.

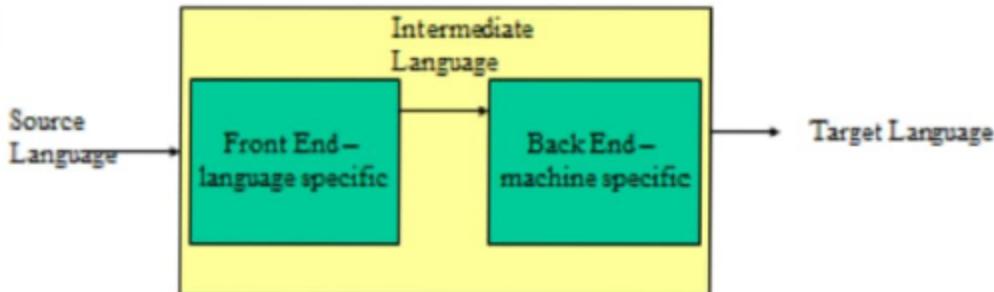
- a) Translate the input source into intermediate code and then convert it to final machine code.
- b) Directly generate the final machine code from the input source.

What is the preferred option and why? Justify this with the suitable claims.

Preferred option is a

Compiler Architecture

In more detail:

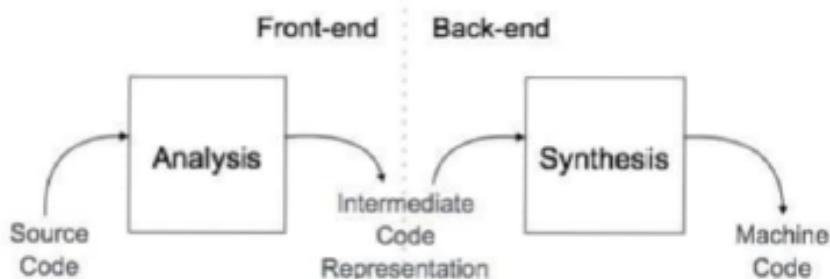


- Separation of Concerns
- Retargeting

A compiler can broadly be divided into two phases based on the way they compile.

Analysis Phase

Known as the front-end of the compiler, the **analysis** phase of the compiler reads the source program, divides it into core parts, and then checks for lexical, grammar, and syntax errors. The analysis phase generates an intermediate representation of the source program and symbol table, which should be fed to the Synthesis phase as input.



Synthesis Phase

Known as the back-end of the compiler, the **synthesis** phase generates the target program with the help of intermediate source code representation and symbol table.

A compiler can have many phases and passes.

- **Pass** : A pass refers to the traversal of a compiler through the entire program.
- **Phase** : A phase of a compiler is a distinguishable stage, which takes input from the previous stage, processes and yields output that can be used as input for the next stage. A pass can have more than one phase.

int fact;

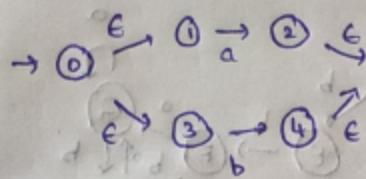
```
int factorial (int n)
{
    int val;
if(n>1){
    val=n*factorial (n-1);
return (val);
}
else
{
return(1);
}
}
int main()
{ printf ("factorial program");
Fact 5=factorial(5);
Printf("fact=5=%d \n", fact5);
}
```

Consider the above program, What are the data structures used to store the information, how the symbol table management is handled?

Describe the need of separating the analysis phase into lexical phase and parsing?

Solve the given regular expression $(a+b) abb (a/b)^*$ into NFA using Thompson construction and then to minimized DFA. . Describe the sequence of moves made by each in processing the input string ababba. Explain in detail the structure of the compilers. Compare language dependent and independent phases

1) $(a+b)^* abb (a/b)^*$



ϵ -closure (0) = {0, 1, 3}

ϵ -closure (1) = {1}

ϵ -closure (2) = {2, 5}

ϵ -closure (3) = {3}

ϵ -closure (4) = {4, 5}

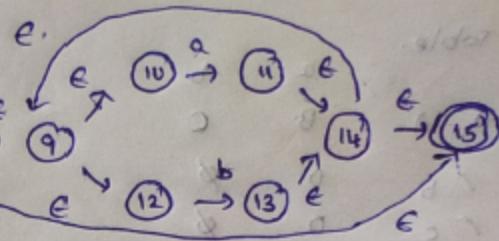
ϵ -closure (5) = {5}

ϵ -closure (6) = {6}

ϵ -closure (7) = {7}

ϵ -closure (8) = {8, 9, 10, 12, 15}

ϵ -closure (9) = {0, 1, 3} = A



ϵ -closure (9) = {9, 10, 12, 15}

ϵ -closure (10) = {10}

ϵ -closure (11) = {11, 14, 15, 9, 10, 12}

ϵ -closure (12) = {12}

ϵ -closure (13) = {13, 14, 9, 10, 12}

ϵ -closure (14) = {14, 9, 10, 12, 15}

ϵ -closure (15) = {15}

$$D(\text{Trans})[A, a] = \{2\} = \{2, 5\} \Rightarrow B$$

$$D(\text{Trans})[A, b] = \{4\} = \{4, 5\} \Rightarrow C$$

$$D(\text{Trans})[B, a] = \{6\} = \{6\} \Rightarrow D$$

$$D(\text{Trans})[B, b] = \{3\}$$

$$D(\text{Trans})[C, a] = \{6\} = \{6\} \Rightarrow D$$

$$D(\text{Trans})[C, b] = \{3\}$$

$$D(\text{Trans})[D, a] = \{3\}$$

$$D(\text{Trans})[D, b] = \{7\} = \{7\} \Rightarrow E$$

$$D(\text{Trans})[E, a] = \{3\}$$

$$D(\text{Trans})[E, b] = \{8\} = \{8, 9, 10, 12, 15\} \Rightarrow F$$

$$D(\text{Trans})[F, a] = \{11\} \Rightarrow G \quad \{11, 12, 15, 9, 10, 12\} \Rightarrow G_1$$

$$D(\text{Trans})[F, b] = \{13\} = \{13, 14, 9, 10, 12\} \Rightarrow H$$

$$D(\text{Trans})[G, a] = \{11\} = \{11, 14, 15, 9, 10, 12\} \Rightarrow G_1$$

$$D(\text{Trans})[G, b] = \{13\} = \{13, 14, 9, 10, 12\} \Rightarrow H$$

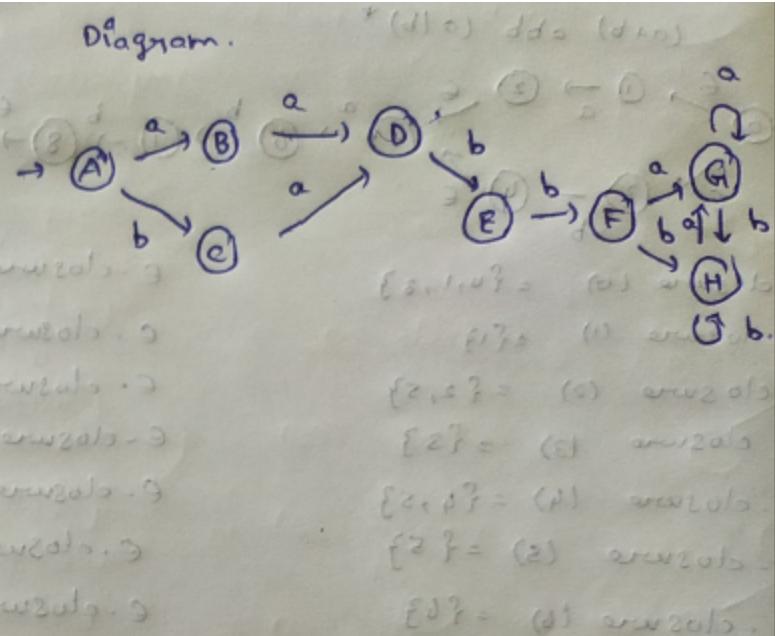
$$D(\text{Trans})[H, a] = \{11\} = \{11, 14, 15, 9, 10, 12\} \Rightarrow G_1$$

$$D(\text{Trans})[H, b] = \{13\} = \{13, 14, 9, 10, 12\} \Rightarrow H.$$

Table.

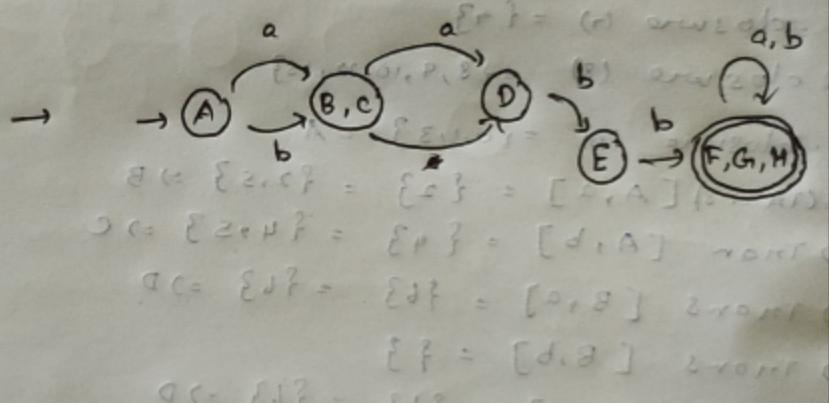
	a	b	
A	B	C	
B	D	\emptyset	
C	D	\emptyset	
D	\emptyset	E	$\{d\} = \{\text{B}\}$ minimal
E	\emptyset	F	$\{e\} = \{\text{C}\}$ minimal
F	G	H	$\{f\} = \{\text{D}\}$ minimal
G	H	I	$\{g\} = \{\text{E}\}$ minimal
H	G, H	I	$\{h\} = \{\text{F}\}$ minimal
			$\{i\} = \{\text{G}\}$ minimal

Diagram.



Minimization.

	a	b
A	B	C
B,C	D	\emptyset
D	\emptyset	E
E	\emptyset	F
F,G,H	G	H



Compare tokens, patterns and lexemes

Lexeme

A lexeme is a sequence of characters that are included in the source program according to the matching pattern of a token. It is nothing but an instance of a token.

Token

The token is a sequence of characters which represents a unit of information in the source program.

Pattern

A pattern is a description which is used by the token. In the case of a keyword which uses as a token, the pattern is a sequence of characters.

Count the number of tokens in the given snippet.

```
int main()
{
    a=10,b=20;
    printf(" sum is : %d", a+b);
    return 0;
}
```

Explain various Errors encountered in different phases of compiler.

LEXICAL ERRORS

- incorrect or misspelled name of some identifiers

Scanned by CamScanner

SYNTACTICAL ERRORS

- Missing Semicolon, unbalanced parenthesis
- When an error is detected, it must be handled by parser to enable the parsing of the rest of the input.
- Most of the errors occurred as syntactic errors
- Role of the Error handler in parser also
 - * Report the presence of errors clearly & accurately
 - * Recover from each error quickly enough to detect subsequent errors
 - * Add minimal overhead to the processing of correcting programs.
- 4 - common Error recovery strategies implemented in the parser to deal with the error in the code
 - * Panic mode
 - * Statement level
 - * Error productions
 - * Global correction

SEMANTIC ERRORS

- Incompatible value assignment
- semantic Analyzer corrects the foll. errors

- * Type mismatch
- * undeclared variable
- * Reserved identifier misuse
- * Multiple declaration of a variable in a scope
- * Accessing an out-of-scope variable
- * Actual & formal Parameter mismatch

Scanned by CamScanner

Lexical phase errors

These errors are detected during the lexical analysis phase. It occurs when compiler does not recognize valid token string while scanning the code.

- Spelling error
- Exceeding length of identifier or numeric constants.
- Appearance of illegal characters
- Unmatched string
- Transposition of two characters

Example: 1- **This is a comment */**

Example 2:

```
void main()
{
int x=10, y=20;
char * a;
a= &x;
x= 1xab;
}
```

Error recovery:

Panic Mode Recovery

successive characters from the input are removed one at a time until a designated set of synchronizing tokens is found. Synchronizing tokens are delimiters such as; or

4/5/2022 26

Slide 27 of 37

Syntactic phase errors

These errors are detected during syntax analysis phase. Typical syn-

- Errors in structure
- Missing operator
- Misspelled keywords
- Unbalanced parenthesis

```
Example : swicth(ch)
{
    .....
    .....
}
```

Lexical error: *I took a bite of an appet*

Correction: *I took a bite of an apple.*

Hence, “**Unidentified keyword/identifier**” error occurs.

Semantic errors

-occurs when a statement is syntactically valid, but does not do what the **programmer** intended **without producing error messages**

Example :

- An expression may not be evaluated in the order you expect, yielding an incorrect result
- Undeclared variables(**symbol table entry has to be made for error recovery**)

```
Example : int a[10], b;  
.....  
.....  
a = b;
```

Example: Jack left her homework at home

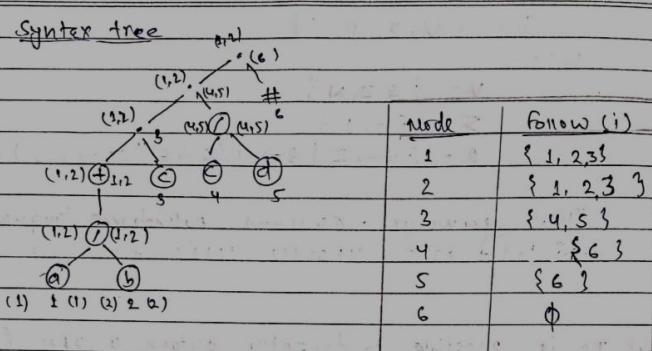
A “type mismatch” between her and Jack; we know they are different people

Obtain the required DFA from the regular expression $(a/b)^* a (a/b)$ using Minimization of DFA using \sqcap new constructions. Write down the suitable algorithm.

Explain in detail the structure of the compilers. Compare language dependent and independent phases.

Construct transition diagram for keyword, identifiers and relational operators

Is it possible to directly obtain a DFA from a regular expression? How is it different from obtaining DFA from an intermediate NFA? Obtain DFA directly for the regular expression $(a/b)^+c(c/d)$.



Now we can construct DFA for the given regular expression. The start state of DFA is firstpos of root.
firstpos of root = {1, 2, 3}.

let the start be A.

$$A = \{1, 2\}$$

consider the symbol 'a' on set A.
position 1 has a transition.

followpos(1)

$$\Rightarrow \{1, 2, 3\}$$

$\Rightarrow B$

consider the input symbol 'b' on set A.
position 2 has a transition.

followpos(2)

$$\Rightarrow \{1, 2, 3\}$$

$\Rightarrow B$

Now set A has no position left for c & d. so,
on input c & d, it will go to dead state (z_d).

Scanned with CamScanner

D.Trans [A, a] = B
D.Trans [A, b] = B
D.Trans [A, c] = z _d
D.Trans [A, d] = z _d

Again,

Consider the symbol 'a' on set B.
position 1 has a transition.

followpos(1)

$\Rightarrow B$

Consider the symbol 'b' on set B.

followpos(2)

$\Rightarrow B$

Consider the symbol 'c' on set B.

followpos(3)

$$\Rightarrow \{4, 5\}$$

$\Rightarrow C$

D.Trans [B, a] = B
D.Trans [B, b] = B
D.Trans [B, c] = C
D.Trans [B, d] = z _d

Again,

Consider the symbol 'c' on set C.

$\text{followpos}(c) \cup$

$\Rightarrow \{c\}$

$\Rightarrow \emptyset$

Consider the symbol d on set C

$\text{followpos}(c) \cup$

$\Rightarrow \{c\} = \emptyset$

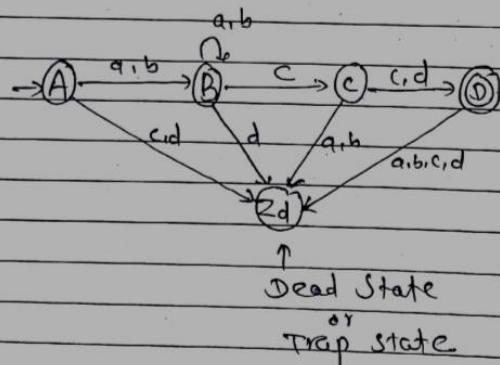
Scanned with CamScanner

D·trans [C, a] = Zd
D·trans [C, b] = Zd
D·trans [C, c] = \emptyset
D·trans [C, d] = \emptyset

Now,

i. On input any symbol on set D, it will go to dead state (Z_d).

DFA Diagram:



29/43

Hence, this is the required DFA for given regular expression.

How do we eliminate ambiguity from the dangling else grammar?

Eliminating Ambiguity.

Date : / /

An ambiguous Gr can be rewritten to eliminate the ambiguity.

Ambiguous Gr cannot distinctively determine which parse tree to select for a sentence or string

Hence disambiguiting rules are used with ambiguous Gr to do away with undesirable parse trees thereby leaving only one tree for a sentence

ex) removing ambiguity from the "dangling else" Gr

$\text{stmt} \rightarrow \text{if expr then stmt}$

| $\text{if expr then stmt else stmt}$

| other

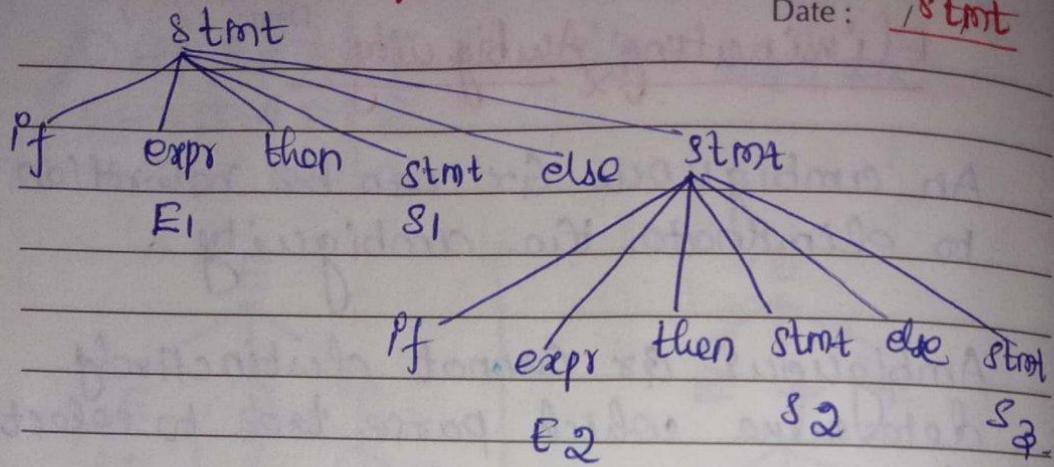
Other \rightarrow any other statement.

Consider,
The compound conditional stmt,

$\text{if } E_1 \text{ then } S_1 \text{ else if } E_2 \text{ then } S_2 \text{ else } S_3$

Parse Tree for the Above conditional

Date : 8/10/2023

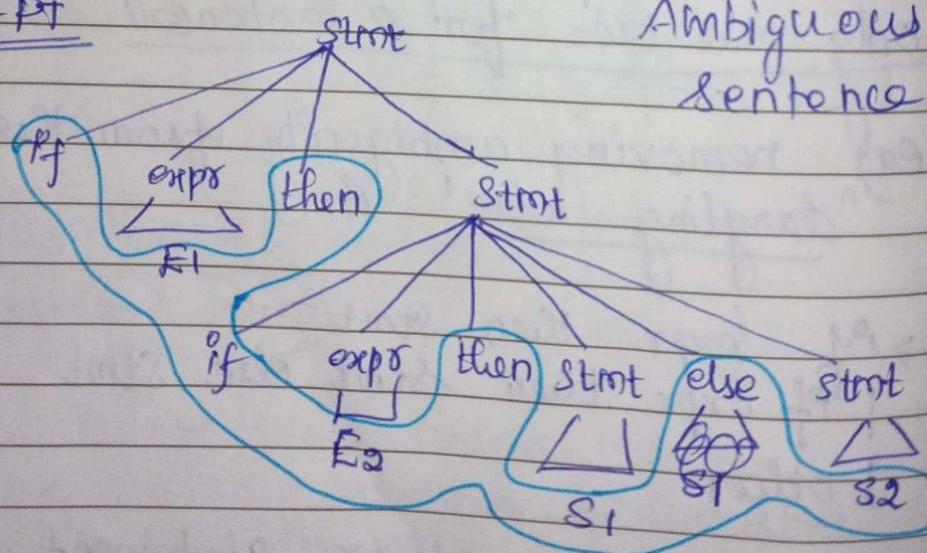


ex: 2

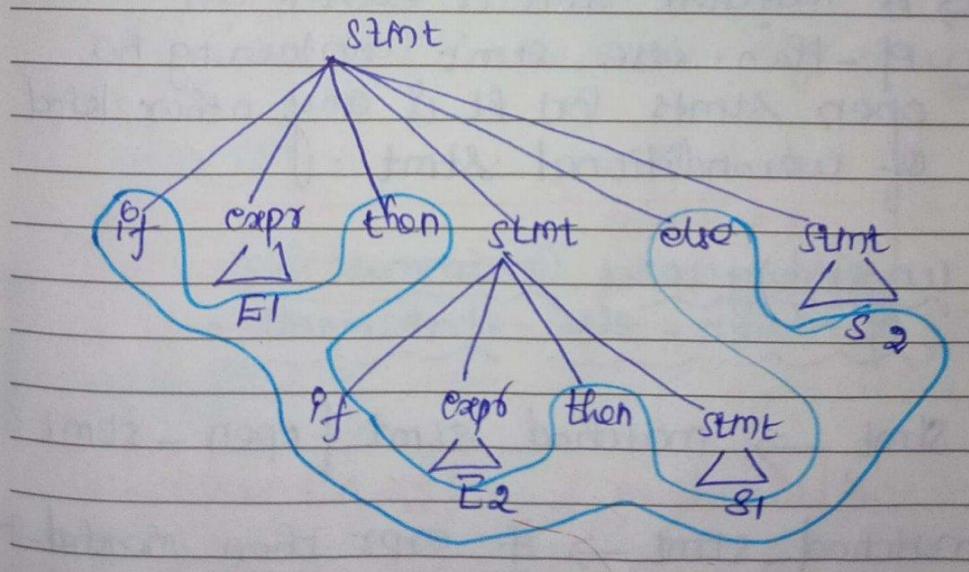
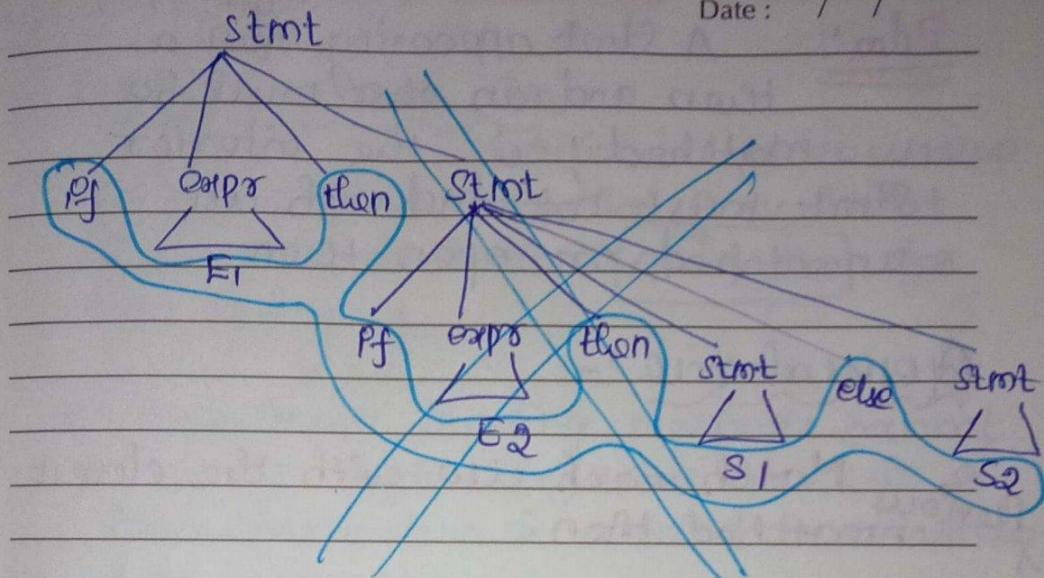
$P_f \ E_1 \ then \ P_f \ E_2 \ then \ S_1 \ else \ S_2$

It has two PT's
Diff PT

\downarrow
Ambiguous sentence



Date : / /



Date : / /

Pdoa: A Stmt appearing b/w a 'then' and an 'else' must be matched; i.e. the interior Stmt must not end with an unmatched or open 'then'.

General rule:-

1. ^{Previous} Match each 'else' with the closest unmatched 'then'.

↳ A matched Stmt is either an 'if - then - else' Stmt containing no openStmts (or) It is any other kind of unconditional Stmt.

unambiguous Grammar for
if - then - else statements :-

Stmt → matched_stmt | open_stmt

matched_stmt → if expr then matched_stmt
else matched_stmt | other

open_stmt → if expr then Stmt |
if expr then matched_stmt else
open_stmt

Explain the structure of Lex program. Write Lex specifications for a Desk calculator application.

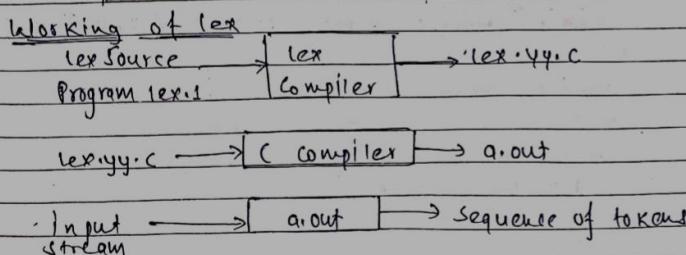
⇒ Lex is a program that generates lexical analyzer. It is used with YACC parser generator. The lexical analyzer is a program that transforms an input stream into a sequence of tokens. It reads the input stream and produces the source code as output through implementing the lexical analyzer into the C program.

Structure of lex program:

A lex program is separated into three sections by % delimiters.

```
{ definitions }  
%%  
{ rules }  
%%  
{ user subroutines }
```

Working of lex



Scanned with CamScanner

lex program to implement a calculator

```
%%  
int op = 0, i;  
float a, b;  
%  
dig [0-9]+([0-9]*\.[0-9]+)  
add "+"  
sub "-"  
mul "*"  
div "/"  
pow "^"  
ln "ln"  
%  
/* digic is a user defined function */
```

```
{dig} {digic();}  
{add} {op=1;}  
{sub} {op=2;}  
{mul} {op=3;}  
{div} {op=4;}  
{pow} {op=5;}
```

```

/* digic is a user defined function */
{dig} {digic();}

{add} {op=1; }

{sub} {op=2; }

{mul} {op=3; }

{div} {op=4; }

{pow} {op=5; }

{ln} {printf ("In the Answer : %f\n", a);}

%{"}

digic()
{
    if (op == 0)
        a = atoi (yytext);
    switch (op)
    {

```

Scanned with CamScanner

```

case 1: a = a+b;
break;
case 2: a = a-b;
break;
case 3: a = a*b;
break;
case 4: a = a/b;
break;
case 5: for (i=a; b>1; b--)
    a = a*i;
break;
}
op = 0;
}

main (int argc, char *argv[])
{
    yylex();
    yywrap();
    return 1;
}

```

Describe the different strategies that a parser can employ to recover from a syntactic error.

Error Recovery Strategies:-

Planning the error handling right from the start can both simplify,

- * Structure of a Compiler
- * Improve its handling of errors

Date : / /

Error handler is a parser 'has
the following goals :-

- * Report the presence of errors clearly & accurately
- * Recover from each error quickly enough to detect subsequent errors.
- * Add the minimal overhead to the processing of correct programs.

Four Error Handling Strategies:-

- * Panic mode
- * phrase Mode Level Recovery
- * Error production
- * Global correction.

① Panic Mode :-

- * easiest method , prevents the parser from infinite loops

Date: / /

* when parser finds an error in the start, it ignores the rest of the start by not processing the input.

* when an err is found, all i/p are discarded until it finds a synchronizing token.

↳ delimiters, semicolon

Adv

* simplicity
* no loops formed

DisAdv

Additional errors
are not checked
since i/p's are skipped

② phrase level Recovery

On discovering an error, Parser performs local correction on the remaining i/p.

e.g)

- * Replacing a Prefix by some str
- * Replacing comma by semicolon
- * Deleting extraneous semicolon
- * Inserting Missing semicolon.

Date :

* by performing local corrections,
allow the parser to continue.

* one wrong correction will lead to an
infinite loop.

Adv

It can correct any
I/P string

DisAdv

correction should
not lead to infinite
loop.

By the error has occurred
at the point of detection

③ Error Production

→ Productions which generate erroneous
constructs are augmented to the Gt by
considering common errors that occur

→ These productions detect the anticipated
errors during parsing.

→ Error diagnostics about the erroneous
constructs are generated by the parser

Consider the following CFG $G = (N=\{S, A, B, C, D\}, T=\{a,b,c,d\}, P, S)$ where the set of productions P is given below:

$S \rightarrow A$
 $A \rightarrow BC \mid DBC$
 $B \rightarrow Bb \mid \epsilon$
 $C \rightarrow c \mid \epsilon$
 $D \rightarrow a \mid d$

Is this grammar suitable to be parsed using the recursive descendent parsing method?

Justify and modify the grammar if needed

Recursive descent parsing is actually a technique.
It cannot handle left-recursion because it
is a top-down parsing technique and also
it should be left-factored grammar.

Eliminating the left recursion:

$$\begin{aligned} S &\rightarrow A \\ A &\rightarrow BC \mid DBC \\ B &\rightarrow B' \\ B' &\rightarrow bB' \mid \epsilon \\ C &\rightarrow c \mid \epsilon \\ D &\rightarrow ald \end{aligned}$$

NOW,

Eliminating left-factoring

Scanned with CamScanner

$$\begin{aligned} S &\rightarrow A \\ A &\rightarrow EBC \\ E &\rightarrow D \mid \epsilon \\ B &\rightarrow B' \\ B' &\rightarrow bB' \mid \epsilon \\ C &\rightarrow c \mid \epsilon \\ D &\rightarrow ald \end{aligned}$$

Hence, this is the required grammar suitable
to be parsed using the recursive descent
parsing method.

Consider the grammar

bexpr \rightarrow bexpr or bterm | bterm
bterm \rightarrow bterm and bfactor | bfactor
bfactor \rightarrow not bfactor | (bexpr) | true | false

- a) Construct a parse tree for the input string **not (true or false)** Show that this grammar generates all boolean expressions.

6:22 PM | 0.0KB/s

LTE 4G+ 12%

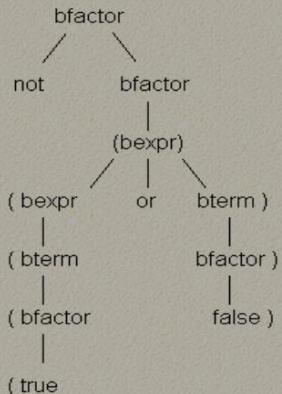
X ⚠ Homework 2
p.web.umkc.edu



bexpr ::= bexpr **or** bterm | bterm
 bterm ::= bterm **and** bfactor | bfactor
 bfactor ::= **not** bfactor | (bexpr) | **true** | **false**

Construct a parse tree for the sentence **not (true or false)**.
 Show that this grammar generates all Boolean expressions

Parse Tree for: **not**
(true or false)



Proof that this language supports all Boolean expressions:

bfactor=> bfactor => TRUE	bfactor=> bfactor => FALSE
bfactor r=>bfactor => not bfactor => not TRUE	bfactor=>bfactor => not bfactor => not FALSE
nor: bfactor => not bfactor => not (bexpr) => not (bexpr or bterm) => not (bterm or bterm) => not (bfactor or bterm) => not (true or bterm) => not (true or bfactor) => not (true or false)	nand: bfactor => not bfactor => not (bexpr) => not (bterm) => not (bterm and bfactor) => not (bfactor and bfactor) => not (true and bfactor) => not (true and false)
xor: bexpr => bexpr or bterm => bterm or bterm => bfactor or bterm => (bexpr) or bterm => (bterm) or bterm => (bterm and bfactor) or bterm => (bfactor and bfactor) or bterm => (true and bfactor) or bterm => (true and false) or bterm => (true and false) or bfactor => (true and false) or (bexpr) => (true and false) or (bexpr or bterm) => (true and false) or (bterm or bterm) => (true and false) or (bfactor or bterm) => (true and false) or (bfactor or bfactor) => (true and false) or (false and bfactor) => (true and false) or (false and true)	xnor: bexpr => bexpr or bterm => bterm or bterm => bfactor or bterm => (bexpr) or bterm => (bterm) or bterm => (bterm and bfactor) or bterm => (bfactor and bfactor) or bterm => (true and bfactor) or bterm => (true and false) or bfactor => (true and false) or (bexpr) => (true and false) or (bexpr or bterm) => (true and false) or (bterm or bterm) => (true and false) or (bfactor or bterm) => (true and false) or (bfactor or bfactor) => (true and false) or (false and bfactor) => (true and false) or (false and true)

NOT is implemented by bfactor

OR is implemented by bterm

AND is implemented by bexpr

•

Illustrate the heuristic techniques for error-recovery in predictive parsing

Consider the grammar

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

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

- i) Construct FIRST and FOLLOW sets
- ii) Construct the predictive parsing table
- iii) Show the moves of the predictive parser for input (a,(a,a))

After remove the left recursion.

$$S \rightarrow (L) a$$

$$L \rightarrow SL'$$

$$L' \rightarrow SL' | \epsilon$$

(i)	Productions	First	Follow
	$S \rightarrow (L) a$	{(, ,)}	{\$, , ,)}
	$L \rightarrow SL'$	{(, ,)}	{,)}
	$L' \rightarrow SL' \epsilon$	{, , } {, }	{,)}

Scanned with CamScanner

(ii) Parse table

Non Terminals	Input symbols			
	()	,	\$
S	$S \rightarrow (L)$		$S \rightarrow a$	
L		$L \rightarrow SL'$	$L \rightarrow SL'$	
L'		$L' \rightarrow E$	$L' \rightarrow SL'$	

(iii) $(a, (a, a))$

stack	input	Production Applied
$\$ S$	$(a, (a, a)) \$$	$S \rightarrow (L)$
$\$ L ($	$(a, (a, a)) \$$	matched (, pop
$\$) L$	$a, (a, a) \$$	$L \rightarrow SL'$
$\$) L' S$	$a, (a, a) \$$	$S \rightarrow a$
$\$) L' a$	$a, (a, a) \$$	matched a, pop
$\$) L'$	$, (a, a) \$$	$L' \rightarrow SL'$
$\$) L' S,$	$, (a, a) \$$	matched S, pop
$\$) L' S$	$(a, a) \$$	$S \rightarrow (L)$
$\$) L') L ($	$(a, a) \$$	matched (, pop
$\$) L') L$	$a, a) \$$	$L \rightarrow SL'$
$\$) L') L' S$	$a, a) \$$	$S \rightarrow a$
$\$) L') L' a$	$a, a) \$$	matched a, pop
$\$) L') L'$	$) \$$	$L' \rightarrow SL'$
$\$) L'$	$) \$$	matched) , pop
$\$) L'$	$\$$	$L' \rightarrow E$
$\$)$	$\$$	matched) , pop
$\$$	$\$$	Accept & successful completion

Scanned with CamScanner

Left factoring is a technique to remove non-determinism. How left factoring does avoid backtracking in grammars? Apply left factoring in the following grammar to make it deterministic.

$$A \rightarrow aBcC \mid aBb \mid aB \mid a$$

$$B \rightarrow \epsilon$$

$$C \rightarrow \epsilon$$

SOL

⇒ Left factoring is a technique to remove non-determinism. sometime we find common prefix in many productions like $A \rightarrow \alpha\beta_1$, $A \rightarrow \alpha\beta_2$ or $A \rightarrow \alpha\beta_3$, where α is common prefix, while processing α we cannot decide whether to expand A by $\alpha\beta_1$ or by $\alpha\beta_2$, so this needs backtracking. To avoid such problem, grammar can be left factoring.

left factoring is a process of factoring the common prefixes of the alternatives of grammar rule. If the production of the form $A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \alpha\beta_3$ has α as common prefix, by left factoring we get the equivalent grammar as

$$A \rightarrow \alpha A'$$

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

So, we can immediately expand A to $\alpha A'$.

⇒ After left factoring grammar is,

$$A \rightarrow \alpha A'$$

$$A' \rightarrow BA'' | \epsilon$$

$$A'' \rightarrow aC | bC | \epsilon$$

$$B \rightarrow C$$

$$C \rightarrow \epsilon$$

Hence, this is the required grammar.

How do SLR(1), LR(1) and LALR(1) methods compare against each other in the process of constructing the parsing table from the grammar?

SLR Parser

SLR represents "Simple LR Parser". It is very easy and cost-effective to execute. The SLR parsing action and goto function from the deterministic finite automata that recognizes viable prefixes. It will not make specifically defined parsing action tables for all grammars but does succeed on several grammars for programming languages. Given a grammar G . It augment G to make G' , and from G' it can construct C , the canonical collection of a set of items for G' . It can construct ACTION the parsing action function, and GOTO, the goto function, from C using the following simple LR Parsing table construction technique. It needed us to understand FOLLOW (A) for each non-terminal A of a grammar.

CLR Parser

CLR refers to canonical lookahead. CLR parsing uses the canonical collection of LR (1) items to construct the CLR (1) parsing table. CLR (1) parsing table make more number of states as compared to the SLR (1) parsing. In the CLR (1), it can locate the reduce node only in the lookahead symbols.

LALR Parser

LALR Parser is Look Ahead LR Parser. It is intermediate in power between SLR and CLR parser. It is the compaction of CLR Parser, and hence tables obtained in this will be smaller than CLR Parsing Table. For constructing the LALR (1) parsing table, the canonical collection of LR (1) items is used. In the LALR (1) parsing, the LR (1) items with the equal productions but have several look ahead are grouped to form an

individual set of items. It is frequently the similar as CLR (1) parsing except for the one difference that is the parsing table.

The overall structure of all these LR Parsers is the same. There are some common factors such as size, class of context-free grammar, which they support, and cost in terms of time and space in which they differ.

Let us see the comparison between SLR, CLR, and LALR Parser.

SLR Parser	LALR Parser	CLR Parser
It is very easy and cheap to implement.	It is also easy and cheap to implement.	It is expensive and difficult to implement.
SLR Parser is the smallest in size.	LALR and SLR have the same size. As they have less number of states.	CLR Parser is the largest. As the number of states is very large.
Error detection is not immediate in SLR.	Error detection is not immediate in LALR.	Error detection can be done immediately in CLR Parser.

SLR Parser	LALR Parser	CLR Parser
SLR fails to produce a parsing table for a certain class of grammars.	It is intermediate in power between SLR and CLR i.e., $SLR \leq LALR \leq CLR$.	It is very powerful and works on a large class of grammar.
It requires less time and space complexity.	It requires more time and space complexity.	It also requires more time and space complexity.

Consider the grammar

$S \rightarrow S(S)S \mid \epsilon$

Construct the predictive parsing table for the input $() ()$

Consider the following (subset of a) CFG grammar

stmt \mid NIL \mid stmt \mid ';' stmt \mid ifstmt \mid whilststmt \mid stmt

ifstmt \mid IF bexpr THEN stmt ELSE stmt

IF bexpr THEN stmt

whilststmt \mid WHILE bexpr DO stmt

where NIL, \mid , IF, THEN, WHILE and DO are terminals, and "stmt", "ifstmt", "whilststmt" and "bexpr" are non-terminals.

For this grammar answer the following questions:

a) Is it ambiguous? Why? Is that a problem?

b) Rewrite the grammar so as to eliminate ambiguity

Draw the precedence function graph for the following table.

	A	()	,	\$
a			>	>	>
(<	<	=	<	
)			>	>	>
.	<	<	>	>	
\$	<	<			

For the following three address code, give the quadruple representation

t1 := -r

t2 := q*t1

t3 := -r

t4 := s * t3

t5 := t2 + t4

p := t5

Construct LALR(1) parsers for the following grammar and parse the sentence id=id

S → L = R

S → R

L → * R

L → id

R → L

The LR(0) finite-state machine for G has the following states.

State 0:

S' → · S

S → · L = R

S → · R

L → · * R

L → · id

R → · L

State 1:

S' → S ·

State 2:

S → L · = R

R → L ·

State 3:

$S \rightarrow R \cdot$

State 4:

$L \rightarrow * \cdot R$

$R \rightarrow \cdot L$

$L \rightarrow \cdot * R$

$L \rightarrow \cdot id$

State 5:

$L \rightarrow id \cdot$

State 6:

$S \rightarrow L = \cdot R$

$R \rightarrow \cdot L$

$L \rightarrow \cdot * R$

$L \rightarrow \cdot id$

State 7:

$L \rightarrow * R \cdot$

State 8:

$R \rightarrow L \cdot$

State 9:

$S \rightarrow L = R \cdot$

The FIRST and FOLLOW sets are as follows.

$FIRST(S) = \{*, id\}$

$FIRST(L) = \{*, id\}$

$FIRST(R) = \{*, id\}$

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

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

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

Look at state 2. LR(0) item $S \rightarrow L \cdot = R$ calls for a shift on lookahead $=$. But LR(0) item $R \rightarrow L \cdot$ calls for a reduce on lookahead $=$, since $FOLLOW(L)$ contains $=$. That is a parsing conflict.

But, in reality, $=$ can only follow L when L occurs in the context of production $S \rightarrow L = R$, never when L is derived using production $R \rightarrow L$.

So the shift reduce conflict in state 2 can be resolved by choosing the shift action without affecting any parses.

Only one occurrence of each object is allowable at a given moment during program execution. Justify your answer with respect to static allocation

Static Storage and Dynamic Allocation

Dynamic memory is controlled by the **new** and **delete** operators, **not** by **scope** and **linkage** rules.

So, dynamic memory can be allocated from one function and freed from another function.

Although the storage schemes don't apply to dynamic memory, then **do** apply to **automatic** and **static** pointer variables used to keep track of dynamic memory.

Let's look at the following line of code:

```
int *ptr = new int[10];
```

The 40 bytes of memory allocated by **new** remains in memory until the **delete** frees it. But the pointer **ptr** passes from existence when the function containing this declaration terminates.

If we want to have the 40 bytes of allocated memory available from another function, we need to pass or return its address to that function.

On the other hand, if we declare **ptr** with external linkage, the **ptr** pointer will be available by using:

```
extern int *ptr;
```

However, a statement that uses **new** to set **ptr** has to be in a function because static storage variables can only be initialized with constant expressions as shown in the following example.

```

int *ptr;
// initialization with non-const not allowed here
// int *ptr = new int[10];

int main()
{
    ptr = new int[10];
    ...
}

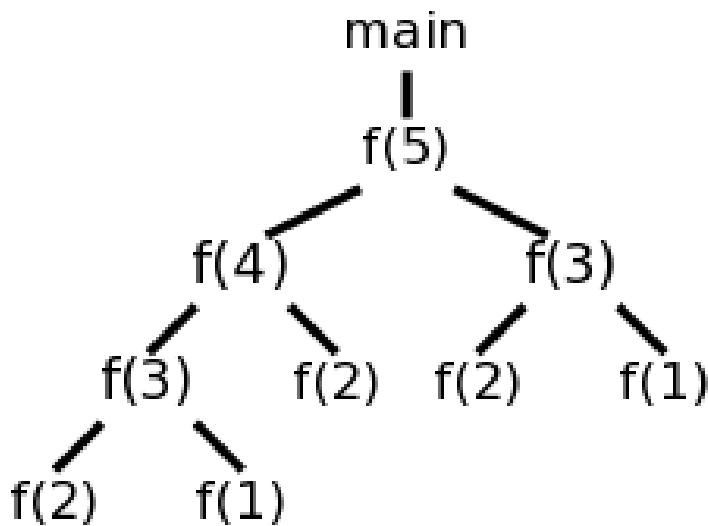
```

Design an activation tree for the Fibonacci sequence 1,1,2,3,5,8, ... defined by $f(1)=f(2)=1$ and, for $n>2$, $f(n)=f(n-1)+f(n-2)$. Consider the function calls that result from a main program calling $f(5)$.

```

system starts main
enter f(5)
    enter f(4)
        enter f(3)
            enter f(2)
            exit f(2)
            enter f(1)
            exit f(1)
        exit f(3)
        enter f(2)
        exit f(2)
    exit f(4)
    enter f(3)
        enter f(2)
        exit f(2)    s
        enter f(1)
        exit f(1)
    exit f(3)
exit f(5)
main ends

```



Consider the following grammar

$S \rightarrow E$
 $E \rightarrow E_1 + T$
 $E \rightarrow T$
 $T \rightarrow T_1 * F$
 $T \rightarrow F$
 $F \rightarrow \text{digit}$

Draw the annotated parse tree to compute s-attribute for the input string $4 * 5 + 6$.

The SDD for the above grammar can be written as follow

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

Let us assume an input string $4 * 5 + 6$ for computing synthesized attributes. The annotated parse tree for the input string is



For computation of attributes we start from leftmost bottom node. The rule $F \rightarrow \text{digit}$ is used to reduce digit to F and the value of digit is obtained from lexical analyzer which becomes value of F i.e. from semantic action $F.\text{val} = \text{digit}.lexval$. Hence, $F.\text{val} = 4$ and since T is parent node of F so, we get $T.\text{val} = 4$ from semantic action $T.\text{val} = F.\text{val}$. Then, for $T \rightarrow T_1 * F$ production, the corresponding semantic action is $T.\text{val} = T_1.\text{val} * F.\text{val}$. Hence, $T.\text{val} = 4 * 5 = 20$

Similarly, combination of $E_1.\text{val} + T.\text{val}$ becomes $E.\text{val}$ i.e. $E.\text{val} = E_1.\text{val} + T.\text{val} = 26$. Then, the production $S \rightarrow E$ is applied to reduce $E.\text{val} = 26$ and semantic action associated with it prints the result $E.\text{val}$. Hence, the output will be 26.

2. Inherited Attributes – These are the attributes which derive their values from their parent or sibling nodes i.e. value of inherited attributes are computed by value of parent or sibling nodes.

Example:

```
A --> BCD { C.in = A.in, C.type = B.type }
```

Computation of Inherited Attributes –

- Construct the SDD using semantic actions.
- The annotated parse tree is generated and attribute values are computed in top down manner.

Example: Consider the following grammar

```

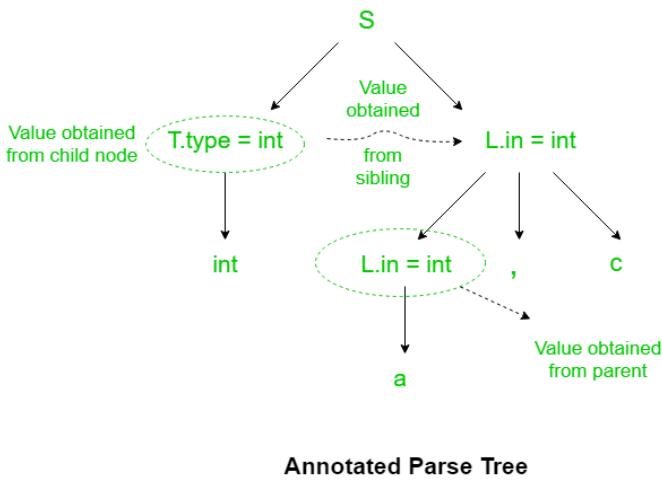
S --> T L
T --> int
T --> float
T --> double
L --> L1, id
L --> id

```

The SDD for the above grammar can be written as follow

Production	Semantic Actions
$S \rightarrow T L$	$L.in = T.type$
$T \rightarrow int$	$T.type = int$
$T \rightarrow float$	$T.type = float$
$T \rightarrow double$	$T.type = double$
$L \rightarrow L_1, id$	$L_1.in = L.in$ $Enter_type(id.entry, L.in)$
$L \rightarrow id$	$Entry_type(id.entry, L.in)$

Let us assume an input string **int a, c** for computing inherited attributes. The annotated parse tree for the input string is



The value of L nodes is obtained from T.type (sibling) which is basically lexical value obtained as int, float or double. Then L node gives type of identifiers a and c. The computation of type is done in top down manner or preorder traversal. Using function Enter_type the type of identifiers a and c is inserted in symbol table at corresponding id.entry.

Write the program for dot product of two vectors. Optimize this code by Eliminating the common sub expressions.

```

def dotProduct(vect_A, vect_B):

    product = 0

    # Loop for calculate dot product
    for i in range(0, n):
        product = product + vect_A[i] * vect_B[i]

    return product
if __name__ == '__main__':
    vect_A = [3, -5, 4]
    vect_B = [2, 6, 5]
    cross_P = []

# dotProduct function call
print("Dot product:", end = " ")
print(dotProduct(vect_A, vect_B))

```

For the following three address code, give the triple and indirect triple implementation

t1 := -r
t2 := q*t1
t3 := -r
t4 := s * t3
t5 := t2 + t4
p := t5

Consider the grammar.

- (1) S → A a
- (2) S → b A c
- (3) S → d c
- (4) S → b d a
- (5) A → d

Parse input string "bdc" constructing LALR Parsing Table. Is the grammar LALR(1)?

Solution

The first number the production as below –

Step1– Construct Augmented Grammar

$$(0) S' \rightarrow S$$

$$(1) S \rightarrow A \ a$$

$$(2) S \rightarrow b \ A \ c$$

$$(3) S \rightarrow d \ c$$

$$(4) S \rightarrow b \ d \ a$$

$$(5) A \rightarrow d$$

Step2– Find Closure & goto. Find the canonical set of LR (1) items for the Grammar.

$I_0 :$	$S \rightarrow . S, \$$ $S \rightarrow . A a, \$$ $S \rightarrow . b A c, \$$ $S \rightarrow . dc, \$$ $S \rightarrow . bda, \$$ $A \rightarrow . d, a$	$I_4 = \text{goto } (I_0, d)$ $S \rightarrow d . c, \$$ $A \rightarrow d . , a$	$I_8 = \text{goto } (I_4, c)$ $S \rightarrow d C . , \$$
I_1	$I_1 = \text{goto } (I_0, S)$ $S' \rightarrow S . , \$$	I_5	$I_5 = \text{goto } (I_2, a)$ $S \rightarrow A a . , \$$
I_2	$I_2 = \text{goto } (I_0, A)$ $S \rightarrow A . a, \$$	I_6	$I_6 = \text{goto } (I_3, A)$ $S \rightarrow b A . c, \$$
I_3	$I_3 = \text{goto } (I_0, b)$ $S \rightarrow b . Ac, \$$ $S \rightarrow b . da, \$$ $A \rightarrow . d, c$	I_7	$I_7 = \text{goto } (I_3, d)$ $S \rightarrow bd . a, \$$ $A \rightarrow d . , c$

In the states, I_0 to I_{10} , no states have a similar first element or core. So, we cannot merge the states. Some states will be taken for building the LALR parsing table.

LALR Parsing Table

States	Action						goto	
	a	b	c	d	\$	S	B	
0		s3		s4		1	2	
1					accept			
2	s5							
3				s7			6	
4	r5		s8					
5								
6			s9					
7	s10		r5					
8					r3			
9					r2			
10					r4			

Parsing of String "bdc"

Stack	Input String	Action
\$ 0	bdc \$	Shift 3
\$ 0 b 3	dc \$	Shift 7
\$ 0 b 3 d 7	c \$	Reduce by A → d
\$ 0 b 3 A 6	c \$	Shift 9
\$ 0 b 3 A 6 c 9	\$	Reduce by S → b Ac

\$ 0 s 1	\$	accept
----------	----	--------

Consider the following grammar:

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

$$T \rightarrow TF \mid F$$

$$F \rightarrow F^* \mid a \mid b$$

Construct the SLR parsing table and also parse the input “a*b+a”

Solution

Step1 – Construct the augmented grammar and number the productions.

$$(0) E' \rightarrow E$$

$$(1) E \rightarrow E + T$$

$$(2) E \rightarrow T$$

$$(3) T \rightarrow TF$$

$$(4) T \rightarrow F$$

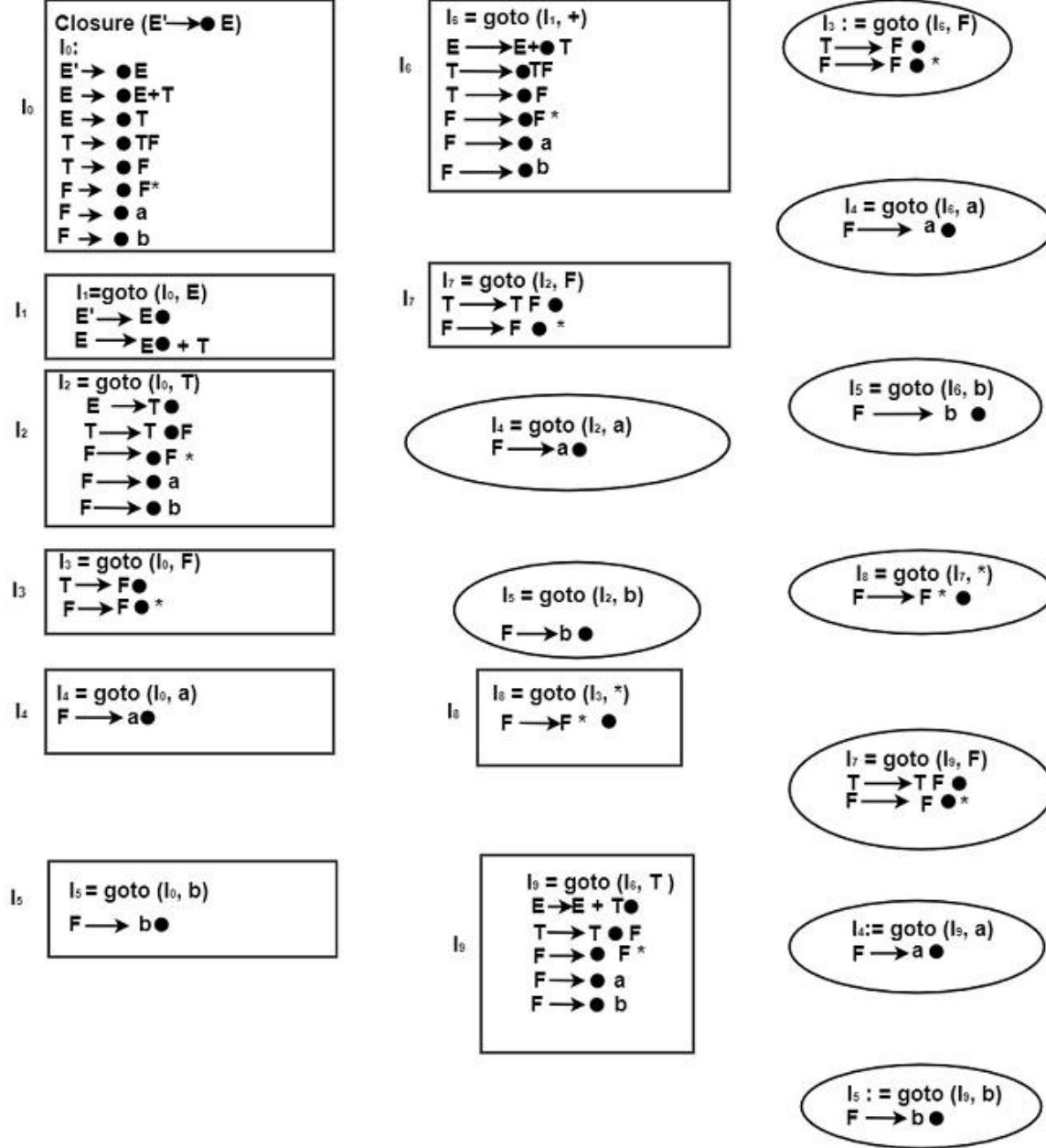
$$(5) F \rightarrow F^*$$

$$(6) F \rightarrow a$$

$$(7) F \rightarrow b.$$

Step2 – Find closure & goto Functions to construct LR (0) items.

Box represents the New states, and the circle represents the Repeating State.



Computation of FOLLOW

We can find out

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

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

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

State	Action					goto		
	+	*	a	b	\$	E	T	F
0			s4	s5		1	2	3
1	s6				accept			
2	r2		s4	s5	r2			7
3	r4	s8	r4	r4	r4			
4	r6	r6	r6	r6	r6			
5	r6	r6	r6	r6	r6			
6			s4	s5		9	3	
7	r3	s8	r3	r3	r3			
8	r5	r5	r5	r5	r5			
9	r1		s4	s5	r1			7

Parsing for Input String a * b + a -

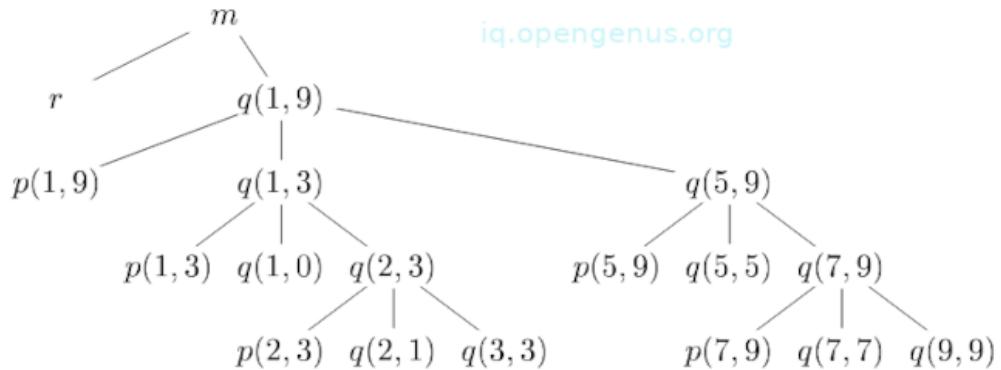
Stack	Input String	Action
0	a * b + a \$	Shift
0 a 4	* b + a \$	Reduce by F → a.
0 F 3	* b + a \$	Shift
0 F 3 * 8	b + a \$	Reduce by F → F *
0 F 3	b + a \$	Reduce by T → F
0 T 2	b + a \$	Shift
0 T 2 b 5	+a \$	Reduce by F → b
0 T 2 F 7	+a \$	Reduce by T → TF
0 T 2	+a \$	Reduce by E → T
0 E 1	+a \$	Shift
0 E 1 + 6	a \$	Shift
0 E 1 + 6 a 4	\$	Reduce by F → a
0 E 1 + 6 F 3	\$	Reduce by T → F
0 E 1 + 6 T 9	\$	Reduce by E → E + T
0 E 1	\$	Accept

Design an activation tree for the following function calls:
enter main()

```

enter readarray()
leave readarray()
enter quicksort(1,9)
  enter partition(1,9)
  leave partition(1,9)
  enter quicksort(1,3)
  .....
  leave quicksort(1,3)
  enter quicksort(5,9)
  .....
  leave quicksort(5,9)
  leave quicksort(1,9)
leave main()

```



Evaluate the execution state of the following C program

```

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

```

into

- A. Syntax tree.
- B. Postfix notation.
- C. 3 address code.

Consider the following grammar

S \rightarrow E
E \rightarrow E1 + T
E \rightarrow T
T \rightarrow T1 * F
T \rightarrow F
F \rightarrow digit

Draw the annotated parse tree to compute s-attribute for the input string 3 * 5 + 5.

Efficient code generation requires the Remember of internal architecture of the target machine. Justify your answer with an Example

Code generation can be considered as the final phase of compilation. Through post code generation, optimization process can be applied on the code, but that can be seen as a part of code generation phase itself. The code generated by the compiler is an object code of some lower-level programming language, for example, assembly language. We have seen that the source code written in a higher-level language is transformed into a lower-level language that results in a lower-level object code, which should have the following minimum properties:

- It should carry the exact meaning of the source code.
- It should be efficient in terms of CPU usage and memory management.

We will now see how the intermediate code is transformed into target object code (assembly code, in this case).

Directed Acyclic Graph

Directed Acyclic Graph (DAG) is a tool that depicts the structure of basic blocks, helps to see the flow of values flowing among the basic blocks, and offers optimization too. DAG provides easy transformation on basic blocks. DAG can be understood here:

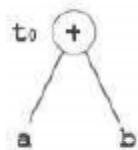
- Leaf nodes represent identifiers, names or constants.
- Interior nodes represent operators.
- Interior nodes also represent the results of expressions or the identifiers/name where the values are to be stored or assigned.

Example:

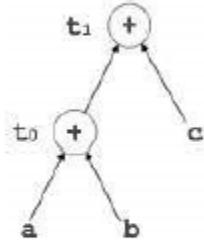
$$t_0 = a + b$$

$$t_1 = t_0 + c$$

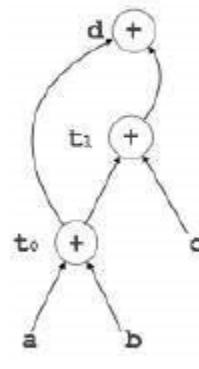
$$d = t_0 + t_1$$



$$[t_0 = a + b]$$



$$[t_1 = t_0 + c]$$



$$[d = t_0 + t_1]$$

Flow of control optimization

There are instances in a code where the program control jumps back and forth without performing any significant task. These jumps can be removed. Consider the following chunk of code: ... MOV R1, R2 GOTO L1 ... L1 : GOTO L2 L2 : INC R1

In this code, label L1 can be removed as it passes the control to L2. So instead of jumping to L1 and then to L2, the control can directly reach L2, as shown below: ... MOV R1, R2 GOTO L2 ... L2 : INC R1

Code Generator

A code generator is expected to have an understanding of the target machine's runtime environment and its instruction set. The code generator should take the following things into consideration to generate the code:

- Target language : The code generator has to be aware of the nature of the target language for which the code is to be transformed. That language may facilitate some machine-specific instructions to help the compiler generate the code in a more convenient way. The target machine can have either CISC or RISC processor architecture.
- IR Type : Intermediate representation has various forms. It can be in Abstract Syntax Tree (AST) structure, Reverse Polish Notation, or 3-address code.
- Selection of instruction : The code generator takes Intermediate Representation as input and converts (maps) it into target machine's instruction set. One representation can have many ways (instructions) to convert it, so it becomes the responsibility of the code generator to choose the appropriate instructions wisely.

Descriptors

The code generator has to track both the registers (for availability) and addresses (location of values) while generating the code. For both of them, the following two descriptors are used:

- Register descriptor : Register descriptor is used to inform the code generator about the availability of registers. Register descriptor keeps track of values stored in each register. Whenever a new register is required during code generation, this descriptor is consulted for register availability.
- Address descriptor : Values of the names (identifiers) used in the program might be stored at different locations while in execution. Address descriptors are used to keep track of memory locations where the values of identifiers are stored. These locations may include CPU registers, heaps, stacks, memory or a combination of the mentioned locations.

Code generator keeps both the descriptor updated in real-time. For a load statement, LD R1, x, the code generator:

- updates the Register Descriptor R1 that has value of x and
- updates the Address Descriptor (x) to show that one instance of x is in R1.

Code Generation

Basic blocks comprise of a sequence of three-address instructions. Code generator takes these sequence of instructions as input.

getReg : Code generator uses getReg function to determine the status of available registers and the location of name values. getReg works as follows:

- If variable Y is already in register R, it uses that register.
- Else if some register R is available, it uses that register.
- Else if both the above options are not possible, it chooses a register that requires minimal number of load and store instructions.

