

I N D E X

NAME: Sindhu S Shastri STD.: _____ SEC.: _____ ROLL NO.: _____ SUB.: _____ Compiler Design notes

Cloud Computing Compiler Design notes

Compiler Design

- ✓ Introduction : phases of compilation and Overview. Lexical Analysis (scanner) :
Regular Languages, finite Automata , Regular Expressions, relating regular Expressions and Finite Automata, Scanner generator (lex, flex)
Definitions
- ✓ Syntax Analysis (Parser) : Context-Free Languages and Grammars, Push-down Automata , LLL(1) grammars and top-down parsing, operator Grammars , LR(0), SLR(1), LR(1), LALR(1) grammars and bottom-up parsing, ambiguity and LR Parsing, LALR(1) parser generator (yacc, bison).
- ✓ Semantic Analysis : Attribute grammars, syntax directed definition, evaluation and flow of attribute in a Syntax.
- ✓ Symbol Table : Basic structure, symbol attributes and management
Run-time environment : Procedure activation, Parameters
Passing, value return, memory allocation, scope.
- ✓ Intermediate Code Generation : Translation of different language features, different types of intermediate forms .
- ✓ Code optimization : Control flow, data-flow dependence etc.
local optimization, global optimization, loop optimization, peephole optimization.
- Architecture dependent code optimization : instruction Scheduling
(for pipeline), loop optimization (for cache memory), Register allocation
long code generation.

Introduction

higher level

Compiler: It translates the code written in one language to some other language without changing the meaning of the program. (Machine level instructions)

It is also expected that a compiler should make the target code efficient and optimized in terms of time & space.

Compiler Design:

The principles of CD provides an in-depth view of translation & optimization process. It covers basic translation mechanism and error detection & Recovery.

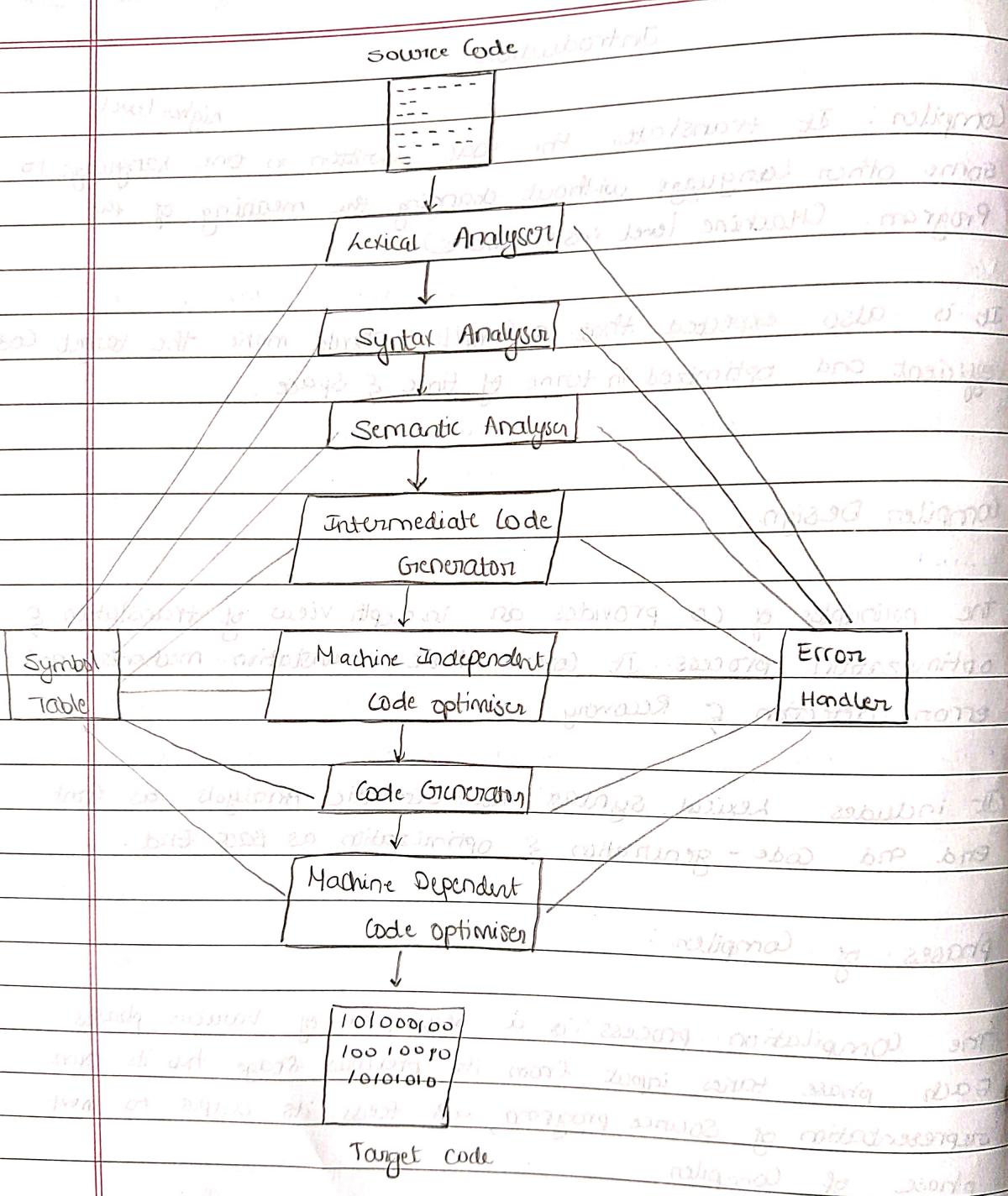
It includes Lexical, Syntax and Semantic Analysis as front end and Code-generation & optimization as Back End.

phases of Compiler:

The compilation process is a sequence of various phases. Each phase takes input from its previous stage, has its own representation of source program, and feeds its output to next phase of compiler.

Analysis Phases - Lexical, Syntax & Semantic

Synthesis Phases - code optimizer, code generation, intermediate code generator



∴ phases of compilation

Lexical Analysis:

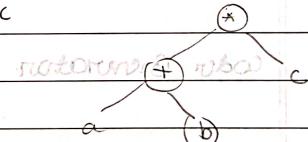
Tokenization and Symbolic

- The first phase of scanner works as a Text Scanner.
- This phase scans the Source Code as a Stream of characters and converts it into meaningful lexemes. Lexical Analysis represents lexemes as tokens in the form of $\langle \text{token-name}, \text{attribute-value} \rangle$.

Ex: $\text{z} = 10$ → Tokens: z - identifier $=$ - assignment operator 10 - Number

Syntax Analysis: noops jrrtq: 10 - Number

- It is also called as Parsing. It takes tokens produced by lexical analysis as input and generates a parse tree (Syntax tree).
- Token arrangements are checked against the source code Grammar, i.e., parser checks if the expression made by tokens is syntactically correct.



Semantic Analysis

- It checks whether the parse tree constructed follows the rules of language. Ex: Assignment of values b/w compatible data types, and adding string to an integer.
- It also keeps track of identifiers, their types & expressions; whether identifiers are declared before use or not etc.
- It produces an annotated Syntax tree as an output.

float x = 30.2;

float y = x * 30;

typecast from integer 30 to

float = 30.0.

Intermediate Code Generation

Algorithm based

- After Semantics Analysis, the compiler generates m intermediate code of m for the target machine.
- It t represents a m program for some abstract machine. It is in b/w the high-level language & machine language. This informs - dict code should be generated in such a way that it makes it easier to be translated into the target machine code.
$$\text{total} = \text{count} + \text{state} * 5$$
$$t1 := \text{int_to_float}(5) \quad t3 := \text{count} + t2$$
$$t2 := \text{state} * t1 \quad \text{total} := t3$$

Code Optimization

- Code optimization of the intermediate code is done.
- It removes unnecessary code lines, and arranges the sequence of statements in order to speed up the program execution without wasting resources (CPU, Memory).

Code Generation

- It takes the optimized representation of intermediate code & maps it to the target machine language. The code generator translates the intermediate code into a sequence of relocatable machine code.
- Sequence of instructions of machine code performs the task as the intermediate code would do.

$a = b + 60.0$

```
MOV B, R1      MOV #16B0, R2
MUL R0, R1, R2
ADD R1, R2
MOV R2, A
```

Symbol table: a ~~pair~~ with ~~name~~ ~~value~~ ~~etc.~~ ~~etc.~~ ~~etc.~~ ~~etc.~~

- It is a data-structure maintained throughout all the phases of a Compiler.
- All the identifier's name along with their types are stored here.
- It makes it easier for Compiler to quickly search the identifier record and retrieve it.
- It is also used for scope Management.

Regular language:

It is a Language that can be expressed with a Regular Expression / a DFA or NFA.

- + It is a set of strings which are made up of characters from a specified alphabet or set of symbols.

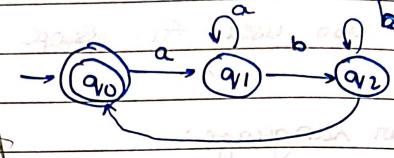
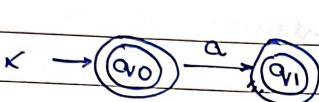
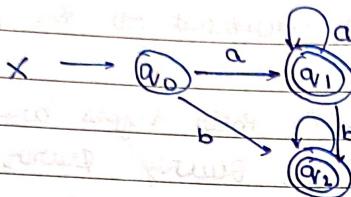
Finite Automata

FA is the simplest machine to recognize patterns.

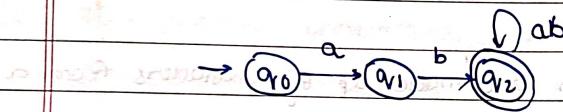
- + It is an abstract machine with 5 tuples.
- $\textcircled{0}$ - Start State
- $\textcircled{0}$ - End State
- Σ : set of Input Symbols
- Q : Finite set of States
- δ : Transition Function

$\delta: Q \times \Sigma \rightarrow Q$

1. Draw a DFA to input the string $a^n b^n$, $n > 0$



2. Draw a DFA to input a string $(ab)^n$



Regular Expression:

It can be recursively defined as follows :

- \emptyset is a regular expression denoting an empty language.
- ϵ is a regular expression indicating the language containing an empty string.

→ R is a regular expression containing only terminals.

→ $R_1 R_2$ is a regular expression denoting the language $L_{R_1} \cup L_{R_2}$.

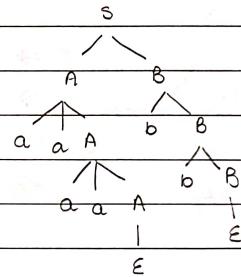
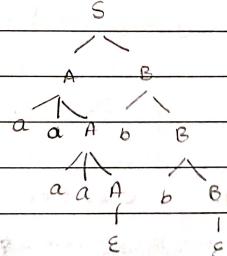
- R^* corresponds to Language L_R^* .
- $R.S$ corresponds to Language $L_R L_S$
 - R^+ corresponds to Language L_R^* .

26/10/21

CFG - A tuple $G = (V, T, P, S)$ $T \rightarrow \text{lower case}$ $V \rightarrow \text{Upper case}$ Derivations \rightarrow leftmost

Rightmost

- 1) $S \rightarrow AB$ Obtain String aaaaabb
 $A \rightarrow aaA | \epsilon$
 $B \rightarrow bB | \epsilon$

LMD : $S \Rightarrow AB$ $\Rightarrow aaA B$ $\Rightarrow aaaaAB$ $\Rightarrow aaaa\epsilon B$ $\Rightarrow aaaabB$ $\Rightarrow aaaabbB$ $\Rightarrow aaaaabb\epsilon$ RMD : $S \Rightarrow AB$ $\Rightarrow A bB$ $\Rightarrow A b bB$ $\Rightarrow A b b\epsilon$ $\Rightarrow aaA bb$ $\Rightarrow aaaaA bb$ $\Rightarrow aaaa\epsilon bb$ $\Rightarrow aaaa bb$ 

Ambiguous Grammar. Not Ambiguous.

- 2) Obtain CFG & tree for $3 * 2 + 5$ using

$S \rightarrow ASB | + - * / \% \mid () \mid E$

$A \rightarrow C \mid 011 \dots 19$

$B \rightarrow 011 \dots 19 \mid)$

OR

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

$E \rightarrow E \mid 011 \mid 2 \dots 19$

RMD:

$3 * 2 + 5$

$S \Rightarrow ASBS$

$\Rightarrow 3SBS$

$\Rightarrow 3 * B \mid 3 * SBS$

$\Rightarrow 3 * A \Rightarrow 3 * BS$

$\Rightarrow 3 * 2S$

$\Rightarrow 3 * 2 ASBS$

$\Rightarrow 3 * 2 E SBS$

$\Rightarrow 3 * 2 + BS$

$\Rightarrow 3 * 2 + 5S$

$\Rightarrow 3 * 2 + 5E$

$E \Rightarrow E * E$

$\Rightarrow 3 * E$

$\Rightarrow 3 * E + E$

$\Rightarrow 3 * 2 + E$

$\Rightarrow 3 * 2 + 5$

RMD:

$E \Rightarrow E * E$

$\Rightarrow E * E + E$

$\Rightarrow E * E + 5$

$\Rightarrow E * 2 + 5$

$\Rightarrow 3 * 2 + 5$

E

$E \times E$

$3 + E + E$

$2 + 5$

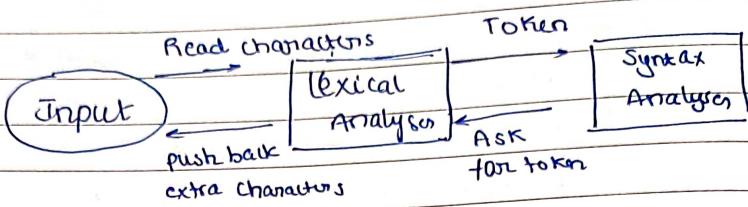
$$1+2*3 - (5*6)/7 + 5$$

minimum 2. modern formular

LNO: $E \Rightarrow E+E$ $E+E+E$ E
 $\Rightarrow 1+E$ $E+E+E$ E
 $\Rightarrow 1+E+E$ $E+E+E-E$ E
 $\Rightarrow 1+2*E$ $E+E+E-E+E+E$ $E \times E$
 $\Rightarrow 1+2*E-E$ $E+E+E-(E)/E+E$ $E-E$
 $\Rightarrow 1+2*3-E$ $E+E+E-(E)/E+E$ $E-E$
 $\Rightarrow 1+2*3-E/E$ $E+E+E-(E)/E+E$ E/E
 $\Rightarrow 1+2*3-(E)/E$ $E+E+E-(E)/E+E$ E/E
 $\Rightarrow 1+2*3-(E \times E)/E$ $E+E+E-(E)/E+E$ $(E) E+E$
 $\Rightarrow 1+2*3-(5*E)/E$ $E+E+E-(E)/E+E$ $E \times E$
 $\Rightarrow 1+2*3-(5*6)/E$ $E+E+E-(E)/E+E$ $E \times E$
 $\Rightarrow 1+2*3-(5*6)/7+E$ $E+E+E-(E)/E+E$ $E \times E$
 $\Rightarrow 1+2*3-(5*6)/7+5$

RNO: $E \Rightarrow E+E$ E times 2
 $\Rightarrow E+E \times E$ $E+E$
 $\Rightarrow E+E \times E-E$ $E \times E$
 $\Rightarrow E+E \times E-E/E$ $E+E$
 $\Rightarrow E+E \times E-E/E/E$ $E-E$
 $\Rightarrow E+E \times E-E/E/E/E$ $E-E$
 $\Rightarrow E+E \times E-E/E/E/E$ $E+E$
 $\Rightarrow E+E \times E-(E)/E+E$ $E+E$
 $\Rightarrow E+E \times E-(E \times E)/E+E$ $E \times E$
 $\Rightarrow E+E \times E-(E \times E)/E+E$ $E \times E$
 $\Rightarrow E+E \times E-(E \times G)/E+E$ $E \times E$
 $\Rightarrow E+E \times E-(S \times G)/E+E$ $S \times G$
 $\Rightarrow E+E \times 3-(S \times G)/E+E$ $S \times G$
 $\Rightarrow E+E \times 3-(S \times G)/E+E$ $S \times G$
 $\Rightarrow 1+2*3-(S \times G)/E+E$ $S \times G$

Lexical Analyser (Scanner)



- It makes modified source code from language Preprocessors.
 - It breaks syntaxes into series of tokens, removing any whitespace or comments in the source code.
 - It can be implemented using DFA.
 - The output is a sequence of tokens that is sent to parser for Syntax Analysis.
- TOKEN: A sequence of characters that can be treated as a unit in grammar of programming languages.

Ex: Keywords, identifiers, operators, separators.

Non-tokens: Comments, preprocessor directive, macros, blanks, tabs, newline etc.

○ Lexeme: The sequence of characters matched by a pattern to form the corresponding token. On a sequence of input characters that comprises a single token is called lexeme.

Ex: "float", "=" , "_" , "273" etc.

How it Works

- It tokenizes, that is Tokenization takes place to divide the program into valid tokens.
 - Remove white space characters
 - Remove Comments
 - It also provides help in generating error by providing row & column numbers.

Suppose we pass a statement,

$a = b + c$ through lexical Analyzer.

It will generate tokens sequence like

$\text{id} = \text{id} + \text{id}$ where id refers to its variable in
Symbol table

Ex: int main() { } valid tokens are: int, main, {, }

1/ 2 variables

int a, b;

$$Q \approx 10;$$

return 0;

— 1 —

'int', 'main' '(', ')'

'i', 'int', 'a', '(', 'b', ')', ',', '.', ','

'a', 'z', '10', ';' , 'return', '0'

“I am not a person who can be easily swayed by others.”

100

Representing valid tokens of a language in regular expression

- x^* - zero / more occurrence of x
 - x^+ - One / more occurrence of x

- o $x?$ means at most one occurrence of x . $\{ \in, x \}$
- o $[a-z]$ - lower case alphabets
- o $[A-Z]$ - upper case alphabets
- o $[0-9]$ - all digits

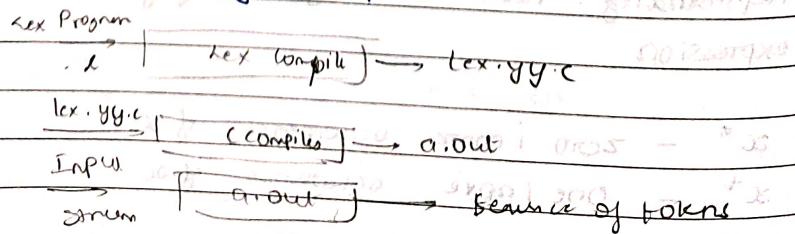
Decimal = (sign)? (digit)⁺ '.' (digit)*
 identifier = (letter) (letter | digit)

→ Lex

- o Lex is a program that generates lexical Analyser. It is used with YACC parser generator.
- o The lexical Analyser is a program that converts input stream into sequence of tokens.

The Function of Lex

- Lexical Analyser creates a program in lex language, then Lex compiler runs the program and produces a C Program lex.yy.c
- Finally C compiler runs lex.yy.c program and produces an object program a.out
- a.out is a lexical Analyser that transforms an input stream into sequence of tokens.



Lex file Format

It is separated by three sections by `/* */` delimiters.

`/*`

definition section

`/* */`

The first section is optional

`/* */`

2 rules 3

`/* */`

Else Subroutines from Lex

`}`

A rule contains a pattern and an action

`p1 {action1}``p2 {action2}``p3 {action3}``...``pn {actionn}`

→ Definitions include declarations of constant, variable and headerfiles.

→ Rules : define the statement of form

`p1 {action1}`
`p2 {action2}`
`p3 {action3}`
`...`
`pn {actionn}`

where `pi` describes regular expression and `actioni` describes the actions what lexical analyzer should take when pattern `pi` matches a lexeme.

User Subroutine : are auxiliary procedures needed by the actions.

Special characters

.	any character except newline
\.	literal '.'
\n	newline
\t	tab
^	beginning of line
\$	end of line

Operations

? zero or one copy of the expression

* zero or more copies

+ one or more copies

a|b a or b

(ab)+ one or more copies of ab

(abc)*

abc*

"abc"

"abc*"

abc*

a(bc)*

a(bc)?

Character class

[abc]

one of : a b c

[a-z]

any LC a to z

[a-zA-Z]

one of : a - z

[-az]

one of :- a z

$[a \wedge b]$

a b

$[A-Z \ A-Z \ 0-9]^+$

one / more alphanumeric

$[\backslash t \ \backslash n]^+$

whitespace

$[^ab]$

except : a b

$[a^b]$

one of : a ^ b / 0 1 records

$[a|b]$

one of : a | b

Predefined Variables in lexer

int yytext (void) { Call to invoke lexer, returns token

char *yytext pointer to Matched String

yylen length of Matched String

yyval value associated with token

int yywrap(void) wrapup, return 1 if done,
 0 if not done

FILE *yyout Output file

FILE *yyin Input file

ECHO write Matched String.

Syntax Analysis / Parsing

→ Limitations of lexical analysis is that it cannot check the Syntax of a given sentence due to limitations of Regular expressions such as balancing parenthesis.

→ This phase uses Context-Free Grammar recognized by Push down Automata.

CFG is superset of Regular grammar.

→ CFG is helpful tool in describing the Syntax of programming languages.

Context-Free Grammar:

A Grammar $G = (V, T, P, S)$ is said to be Context-Free Grammar if all the productions are of the form

$$A \rightarrow \alpha \quad \text{where } \alpha \in (V \cup T)^*$$

A is a non-terminal.

ϵ can appear on right hand side of any production.

It is also a 4 tuple denoted by

$$G = (V, T, P, S)$$

V = Set of Variables / non-terminals

T = Set of terminals

P = Set of Productions

S = start symbol.

Derivation: The process of obtaining a string of terminals from non-terminals from start symbol by applying some set of productions is called derivation.

- At each step, only one variable is replaced using a production
- It starts with start symbol S and ends in string of terminals

Right-most Derivation: The process of obtaining a string of terminals from a sequence of replacements such that only rightmost variables is replaced at each and every step.

left-most Derivation: The process of obtaining a string of terminals by from a sequence of replacements such that only left-most non-terminal is replaced at each and every step.

Parse tree (Derivation Tree)

Let $G = (V, T, P, S)$ be a context-free grammar. The tree is derivation tree if:

- The Root has label S
- Every vertex has label which is in $(V \cup T^*)$
- Every leaf node has label from T and interior vertex has label from V
- If a vertex is labelled A and x_1, x_2, \dots, x_n are all children of A from left, then $A \rightarrow x_1 x_2 \dots x_n$ is the production.

Ambiguous Grammar

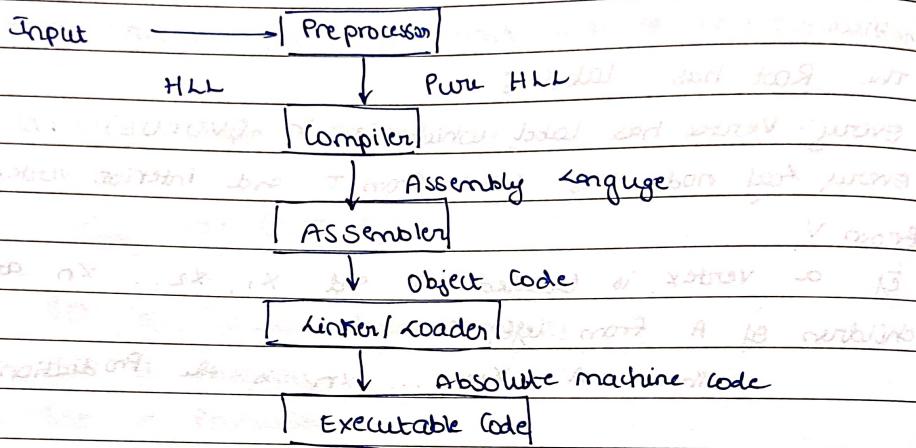
Let $G = (V, T, P, S)$ be a context-free grammar. Grammar G

is said to be ambiguous iff there exists one string $w \in T^*$ for which two or more different parse trees exist by applying either LR(0) or NLR(0).

Language Processing in Compiler Design

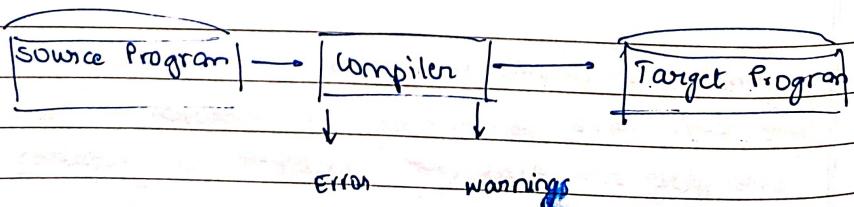
The higher level languages are fed into a series of devices and operating system (OS) components to obtain the desired code that can be used by the machine.

This is known as a Language Processing System.



- Preprocessor:** It includes all the header files and also evaluates whether a macro is included. And expands the macro by its contents. It takes source code as input and produces modified source code as output.
- It also known as Macro evaluate.

Compiler :



- It takes the modified source code as input and converts it into assembly language.

Assembler

- It takes Assembly language and converts it into machine code instructions as output.

Linker:

A linker or link editor is a program that takes a collection of objects and combines them into an executable program.

Loader:

The Loader keeps the linked program in the main memory.

Executable Code

It is the low level and machine specific code that machine can easily understand.

Compiler

- Compiler takes entire program as input
- Intermediate code object code is generated
- Conditional control statements are executed first

Interpreter

- Interpreter takes single instruction as input
- No intermediate object code is generated
- Conditional control statements are executed slower

- Memory Requirement is More. • Memory Requirement is slow.
 - Program need not be Compiled every time. • Every time, higher level program is converted into lower level program.
 - Errors are displayed after entire program is checked. • Errors are displayed for every instruction.

Context-Free Language:

- It is a language generated by Context-Free Grammar.
- They have many applications in programming languages, in particular, most arithmetic expressions are generated by Context-Free Grammar.

Ex: $\{a^n b^n \mid n \geq 0\}$

Context-Free Grammar:

A Grammar $G = (V, T, P, S)$ is Said to be Context-Free if all the productions are of the form $A \rightarrow \alpha \beta \gamma$ where $\alpha, \beta, \gamma \in (V \cup T)^*$ and A is a non-terminal.

E can appear on RHS of any production.

It can be implemented through Push-Down Automata.

Push Down Automata.

It is a finite automata with extra memory called Stack which helps to recognize Context-Free Grammar.

$$P = \{Q, \Sigma, \Gamma, q_0, Z, F, \delta\} \quad 7\text{-tuple}$$

Q = Set of all states

q_0 = initial state

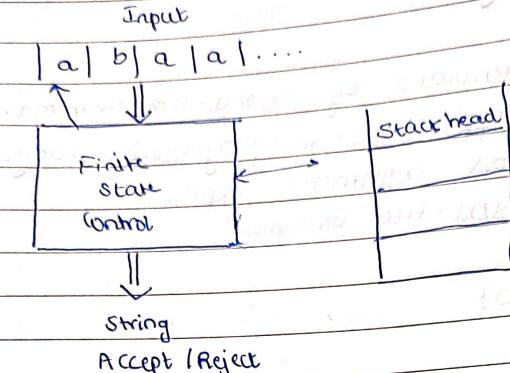
Σ = Set of input Symbol

Z = initial Pushdown Symbol

Γ = Set of pushdown symbols

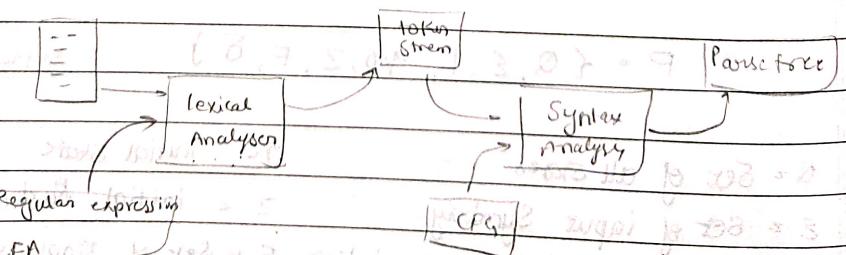
F = Set of Final States

δ = Transition Function



Syntax Analyser

- It is the second phase of the Compiler design process in which the given input string is checked for confirmation of rules and structure of formal grammar.
- It analyses the syntactical structure and checks if the given input has correct syntax or not.
- The Syntax Analyser analyzes the source code (token stream) against the production rules to detect any errors in the code. The output is a parse tree.
- Parser parses through the whole code even if some errors exist in the program and uses error recovering strategies.



First and Follow Sets:

These sets can provide the actual position of any terminal in the derivation. This is done to create the parsing Table.

First Set:

This set is created to know what terminal symbol is derived by in the 1st position by a Non-terminal.

Ex: $A \rightarrow tB$ $t \in \text{First}(A)$ i.e., A derives t (terminal) in the 1st position.

Algorithm for calculating First set.

→ If A is a terminal, then $\text{First}(A) = A$

- $A \rightarrow \epsilon$, then $\text{First}(A) = \{\epsilon\}$
- $A \rightarrow a$, then $\text{First}(A) = \{a\}$
- $A \rightarrow xyz$, then $\text{First}(A) = \text{First}(x) = \{a\}$

$$x \rightarrow a \quad \text{First}(x) = \{a\}$$

$$\{ \} \cup A \rightarrow abc \quad \text{then } \text{First}(A) = \{a\}$$

$$\{ \} \cup A \rightarrow aA \quad \text{then } \text{First}(A) = \{a\}$$

* Problems

$$s \rightarrow abc \mid def \mid ghi$$

$$s \rightarrow abc$$

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

$$s \rightarrow def$$

$$s \rightarrow ghi$$

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

$$A \rightarrow a b l c \quad \text{First}(A) = \{a, b, c\}$$

$$B \rightarrow b \quad \text{First}(B) = \{b\}$$

$$C \rightarrow d \quad \text{First}(C) = \{d\}$$

$$\text{First}(S) = \{\text{First}(A), g, j\}$$

3)

$$S \rightarrow A B C \mid \dots \quad (A \text{ first } \Rightarrow S \leftarrow A \rightarrow \dots)$$

$$A \rightarrow a b l \epsilon \quad \text{First}(A) = \{a, b, \epsilon\}$$

$$B \rightarrow c l d l \epsilon \quad \text{First}(B) = \{c, d, \epsilon\}$$

$$C \rightarrow e l f l \epsilon \quad \text{First}(C) = \{e, f, \epsilon\}$$

$$\text{First}(S) = \text{First}(A)$$

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

$$= \{a, b, c, d, \epsilon\} \text{First}(C)$$

$$= \{a, b, c, d, e, f, \epsilon\}$$

4)

$$E \rightarrow T E l$$

$$E' \rightarrow * E' \mid \epsilon$$

$$T \rightarrow F . T l$$

$$T' \rightarrow + F T' \mid \epsilon$$

$$F \rightarrow id \mid C E$$

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

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

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

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

$$\text{First}(F) = \{id, ()\}$$

- 5) $S \rightarrow aBDg$ $\text{First}(S) = \{a\}$
 $B \rightarrow CC$ $\text{First}(B) = \{C\}$
 $D \rightarrow EF$ $\text{First}(D) = \text{First}(E) = \{g, E, f\}$
 $E \rightarrow g|E$ $\text{First}(E) = \{g, E\}$
 $F \rightarrow f|E$ $\text{First}(F) = \{f, E\}$

Follow sets

We calculate what terminal symbol immediately follows a non-terminal in Production Rules.

→ If S is start Symbol, $\text{Follow}(S) = \{\text{terminal}\}$

→ $\text{Follow}(S) = \{\$\}$ since $\{\$\} = \{\text{terminal}\}$

→ $A \rightarrow ACD$, $\text{Follow}(A) = \text{First}(C)$, $\text{Follow}(C) = \text{First}(D)$

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

→ $S \rightarrow asbs$ $\text{Follow of } S = \{\$, b\}$ and $\text{First}(S) = \{\$,\}$

→ $A \rightarrow ACD$ $\text{Follow}(D) = \text{Follow}(A)$

→ $S \rightarrow ACD$ $\text{First}(C) = \text{First}(A)$ $S \rightarrow asbs$

$C \rightarrow a|b$ $\text{First}(C) = \{a, b\}$ $S \rightarrow bsas$

$S \rightarrow \epsilon$

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

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

$\text{Follow}(C) = \text{First}(D)$

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

3) $S \rightarrow AaAb \mid BbBa$

$A \rightarrow E$

$B \rightarrow E$

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

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

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

$S \rightarrow aBDh$	$\text{First}(S) = \{a\}$	$b \rightarrow a$
$B \rightarrow CC$	$\text{First}(B) = \{c\}$	$c \rightarrow b$
$C \rightarrow bC \mid \epsilon$	$\text{First}(C) = \{b, \epsilon\}$	$d \rightarrow \epsilon$
$D \rightarrow EF$	$\text{First}(D) = \{e, f, \epsilon\}$	$e \rightarrow g$
$E \rightarrow gIE$	$\text{First}(E) = \{g, \epsilon\}$	$f \rightarrow h$
$F \rightarrow f \mid E$	$\text{First}(F) = \{f, \epsilon\}$	

 $\text{Follow}(S) = \{\$\}$ $\text{Follow}(B) = \text{First}(D) = \{e, f, \epsilon\}$ with $[E]$ chuck removed. $\text{Follow}(D) = \{h\}$ $\text{Follow}(C) = \text{Follow}(B) = \{g, f, h\}$ with $[E]$ removed. $\text{Follow}(E) = \text{First}(F) = \{f\}$ with $\epsilon \in \text{Follow}(D) = \{f, h\}$ $\text{Follow}(F) = \text{Follow}(D) = \{h\}$ with $\epsilon \in \text{Follow}(D)$

Types of Parser

TOP - DOWN:

1. LR(0) parser (deterministic) = Bottom up parser
 2. LR(1) parser (non-deterministic) = Predictive parser (LL(1))

	Parsing	Direction	
1. LR(0)	Top-down	Bottom-up	LR-parser (operator)
2. LR(1)	Parsing	Parsing	LR(0) parser (precedence parser)

Top - Down Parser:

- It generates parse tree for the given input string with the help of grammar productions by expanding the non-terminals i.e., it starts from the start symbol and ends on the terminals.
- It uses left-most Derivation.

Top-down Parser is further classified into

- Recursive descent parser.

- Predictive / non-Recursive descent parser.

Recursive Descent Parser (RDP):

→ It is also known as Brute force parser or the backtracking parser.

It basically generates parse tree by using brute force and backtracking.

$$\text{1) } A \rightarrow abC \mid abD \mid aAD$$

$$B \rightarrow bB \mid \epsilon$$

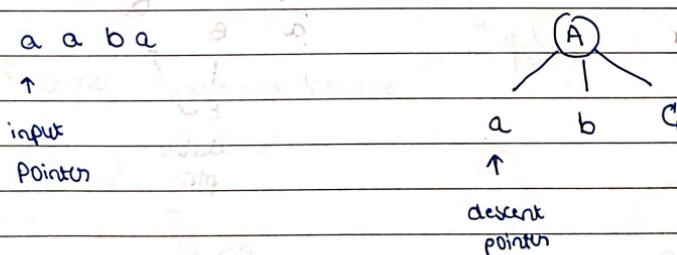
$$C \rightarrow d \mid \epsilon$$

$$D \rightarrow a \mid b \mid \epsilon$$

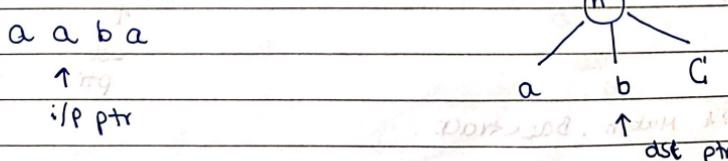
Given input string: aaba

We consider 2 pointers : input pointer and descent pointer.

Construct the Parse tree.



(try match). Increment i/p ptr and dst pointer by 1

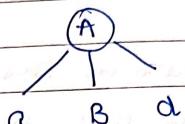


Doesn't Match. Backtrack.

Using Backtracking,

a a b a

i/p ptr
↑



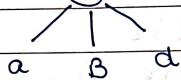
dst ptr
↑

(They Match). Increment the pointers.

a a b a

n
i/p ptr
↑

2nd a is not matched



b
B

dst
ptr
↑

Doesn't Match. Backtrack.

a a b a

↑

i/p ptr = (A)



dst
ptr
↑

a a b a

i/p
ptr
↑

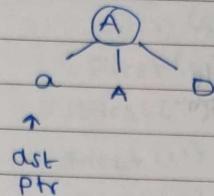


dst
ptr
↑

Doesn't Match. Backtrack.

a a ba

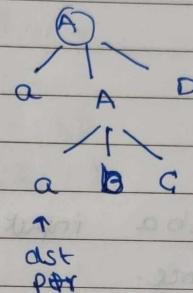
↑
i/p
ptr



Matches, increment ptrs.

a a ba

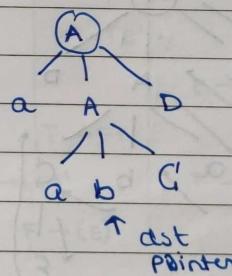
↑
i/p
ptr



Matches, Increment ptrs

a a b a

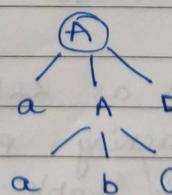
↑
i/p
ptr



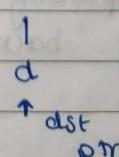
Matches, Increment pointers

a a b a

↑
i/p
ptr



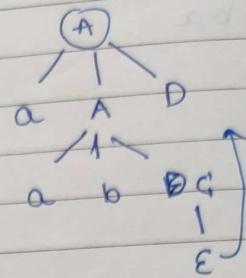
Doesn't Match, Backtrack.



a a b a

↑ i/p
ptr

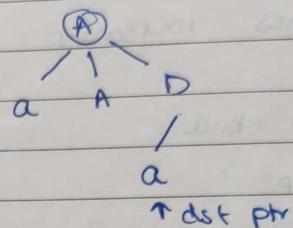
G → E



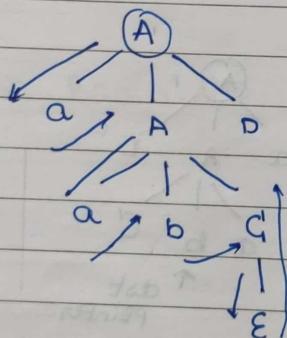
a a b a

↑ i/p ptr

Matches.



Hence, aaba input string is accepted using the below parse tree.



A → aAD
a → abc
c → ε
D → a

Non-Recursive Descent Parser: L-left to Right, L-left Derivation
(1) - using 1 input symbol

It is also known as LL(1) parser, Predictive parser or without Backtracking parser and Dynamic Parser.
It uses parsing table to generate the parse tree instead of backtracking.

$$1) E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

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

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

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

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

$$\text{First}(F) = \{id, ()\}$$

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

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

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

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

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

follow, first of
by fill table

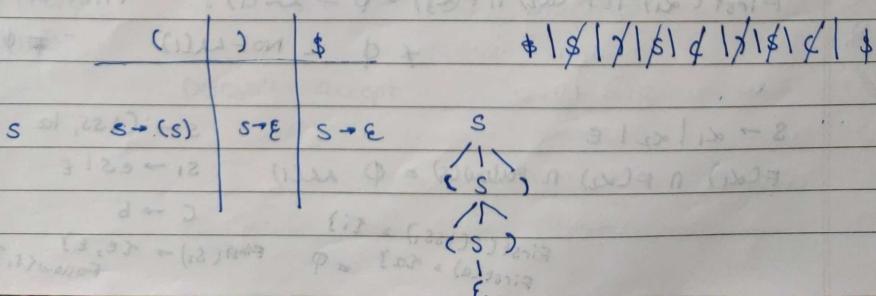
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$	
E'		$E' \rightarrow +TE'$			$E \rightarrow E$
T		$T \rightarrow FT'$		$T \rightarrow FT'$	
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$	$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$			$F \rightarrow (E)$	

$$2) S \rightarrow CS \mid \epsilon$$

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

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

$$() \$$$



Stack Implementation

 $w = id * id + id$

Stack	AI/P	O/P
\$ E	$id * id + id \$$	$E \rightarrow T E'$
\$ E' T	$id * id + id \$$	$T \rightarrow FT'$
\$ E' T' F	$id * id + id \$$	$F \rightarrow id$
\$ E' T' id	$id * id + id \$$	$pop(id)$
\$ E' T'	$* id + id \$$	$T' \rightarrow * FT'$
\$ E' T' F *	$* id + id \\$	$Pop(*)$
\$ E' T' F	$id + id \$$	$F \rightarrow id$
\$ E' T' id	$id + id \$$	$pop(id)$
\$ E' T'	$+ id \$$	$T' \rightarrow E$
\$ E'	$+ id \$$	$E' \rightarrow + TE'$
\$ E' T'	$+ id \\$	$pop(+)$
\$ E' T	$id \$$	$T \rightarrow FT'$
\$ E' T' F	$id \$$	$F \rightarrow id$
\$ E' T' id	$id \$$	$pop(id)$
\$ E' T'	$\$$	$T' \rightarrow E$
\$ E'	$\$$	$E' \rightarrow E$
\$	$\$$	$pop(\$)$

To check LLL(1) or not.

Ex: $S \rightarrow asa | bsa | c$ $S \rightarrow \alpha_1 | \alpha_2 | \alpha_3 - \text{NO } E \text{ production}$ $\text{First}(asa) = \{a\}$ $\text{First}(bsa) = \{b\}$ $\text{First}(\alpha_1) \cap \text{First}(\alpha_2) \cap \text{First}(\alpha_3) = \emptyset - \text{LLL}(1)$ $\text{First}(c) = \{c\}$ $\neq \emptyset - \text{Not LLL}(1)$ $\neq \emptyset : \text{LLL}(1)$ $S \rightarrow \alpha_1 | \alpha_2 | \epsilon$ $F(\alpha_1) \cap F(\alpha_2) \cap \text{Follow}(S) = \emptyset \text{ LLL}(1)$ $s \rightarrow i(tss, 1a)$ $s_1 \rightarrow es | \epsilon$ $\text{First}(i(tss)) = \{i\}$ $\text{First}(a) = \{a\} = \emptyset$ $c \rightarrow b$ $\text{First}(s_1) \rightarrow \{e, \epsilon\}$ $\text{Follow}(s, \gamma) = \text{First}(S) \cup \{\epsilon\}$

Q) $A \rightarrow abG \mid abd \mid aAD$ $\text{First}(A) = \{a\}$
 $B \rightarrow bB \mid \epsilon$ $\text{First}(B) = \{b, \epsilon\}$
 $G \rightarrow aLE$ $\text{First}(G) = \{a, L, E\}$
 $D \rightarrow aLB \mid \epsilon$ $\text{First}(D) = \{a, b, \epsilon\}$

$\text{Follow}(A) = \{\$, \text{First}(D)\} = \{\$, a, b\}$
 $\text{Follow}(G) = \text{Follow}(A) = \{\$, a, b\}$
 $\text{Follow}(B) = \{a, \text{Follow}(B)\} = \{a\}$
 $\text{Follow}(D) = \text{Follow}(G) = \{\$, a, b\}$

	a	b	d	\$	
A	aAD				
	$A \rightarrow abG$	$B \rightarrow bB$	$B \rightarrow \epsilon$		
B		$B \rightarrow bB$	$B \rightarrow \epsilon$		
C	$G \rightarrow \epsilon$	$G \rightarrow \epsilon$	$G \rightarrow a$	$G \rightarrow \epsilon$	
D	$D \rightarrow a$	$D \rightarrow b$		$D \rightarrow \epsilon$	
	$D \rightarrow \epsilon$	$D \rightarrow \epsilon$			

Stack

IIP

OIP

\$A	a b a	$A \rightarrow aAD$
\$DAg	* a b a	pop(a)
\$DA	a b a	$A \rightarrow a b d$
\$DdBga	a b a	pop(a)
\$DdB	b a	$B \rightarrow bB$
\$DBb	b a	pop(b)
\$DB	a	

Doesn't accept.

Bottom-up Parser

→ Bottom up Parser is the Parser which generates the parse tree for given input string with the help of grammar productions by compressing the non-terminals i.e., it starts from non-terminals and ends on the start symbol.

It uses reverse of Right-most derivation.

There are 2 classifications

• LR parser:

It generates parse tree for given string using unambiguous grammar. It follows reverse of Right-most derivation.

LR Parser is of 4 types:

- a) LR(0)
- b) SLR(1)
- c) LALR(1)
- d) CLR(1)

- LR(0) - Left to Right and Reverse of Rightmost Derivation.
(0) - doesn't care about lookahead items.

Rules for Construction.

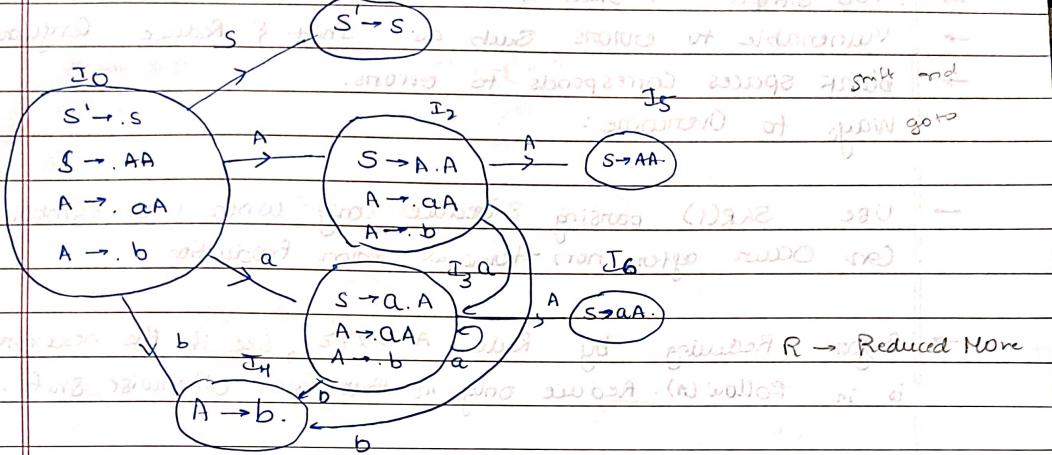
- 1) Construct Context-Free Grammar.
- 2) Check for ambiguity.
- 3) Add a Augment Production.
- 4) Draw the DFA / DFD.
- 5) Create a Collection of Canonical LR(0) items.
- 6) Construct LR(0) Parsing Table.
- 7) Check for Validity with Sample String.

LR(0) items are

my item with ':' in front
of them.

$s' \rightarrow s'$

canonical LR(0) items

 $S \rightarrow AA$ $S \rightarrow AA \text{ or } S \rightarrow A \text{ or } S \rightarrow \epsilon$ $A \rightarrow aA \mid b$ $A \rightarrow aA \text{ or } A \rightarrow b$ $A \rightarrow b$ $A \rightarrow aA \text{ or } A \rightarrow b$ 

Parsing Table

States	Actions			Go To		
$S \rightarrow \epsilon$	a	b	\$	Accept	A	S
I_0	S_3	S_4		I_2	I_1	
I_1		B		Accept		I_5
I_2	S_3	S_4			I_5	S
I_3	S_3	S_4		I_2	I_6	A
I_4	R_3	R_3	R_3	I_2	I_2	
I_5	R_1	R_1	R_1	I_2	I_2	A
I_6	R_2	R_2	R_2			A

see stack implementation

0 | x | x | x | x | A | * | b | x | b | A | x | x | b | A | A

1 | S | 1 | Accept.

A → aA (2)

A → b (3)

B → AA (1)

Disadvantages of LR(0) Parser

- It doesn't lookahead of characters
- Too simple for real parser
- Vulnerable to errors such as Shift & Reduce Conflicts.
- Blank spaces corresponds to errors.

Ways to Overcome:

- Use SLR(1) parsing : reduce only when next symbol can occur after non-terminal from Production.
- Before Reducing by Rule $A \rightarrow x y z$, see if the next token is in $\text{Follow}(A)$. Reduce only in that case. Otherwise, Shift.

SLR(1):

The main difference is the Parsing Table only.

States	Actions			Goto
	a	b	\$	
I ₀	S ₃	S ₄		2
I ₁			Accept	
I ₂	S ₃	S ₄		5
I ₃	S ₃	S ₄		6
I ₄	R ₃	R ₃	R ₃	
I ₅			R ₄	
I ₆	R ₂	R ₂	R ₂	

$I_H - R_3$ FOLLOW(LHS)

$A \rightarrow b.$ FOLLOW(A) = FIRST(A) & FOLLOW(S) = $\{\$, a, b\}$

$I_S - R_1$

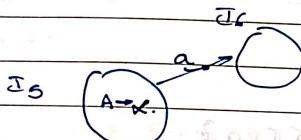
$S \rightarrow AA.$

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

$I_G - R_2$

$A \rightarrow aA$

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



5	s_0, s_1, s_2, s_3
---	----------------------

a, b, c

R1, R2, R3, R4

$S \rightarrow dA | ab$

FIRST(S) = d, a, b \Rightarrow FOLLOW of all

$A \rightarrow bA | c$

FIRST(A) = b, c \Rightarrow LR(0) \Rightarrow

$B \rightarrow bB | c$

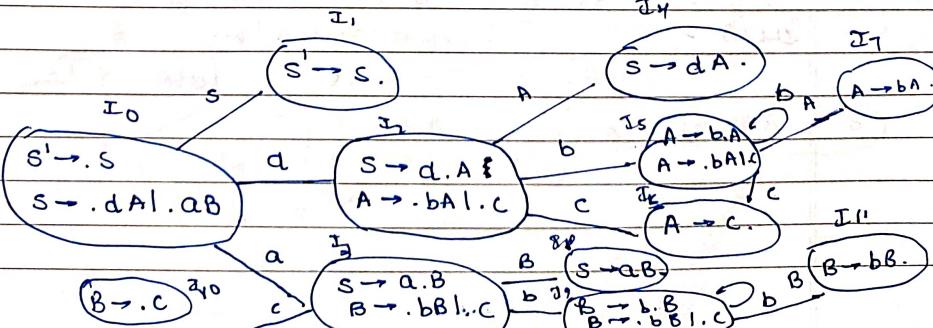
FIRST(B) = b, c \Rightarrow SLR(1) \Rightarrow

$S' \rightarrow S$

$S' \rightarrow dA | ab$

$A \rightarrow bA | c$

$B \rightarrow bB | c$

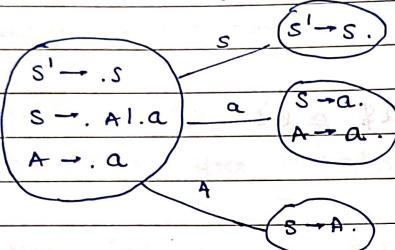


2) $S \rightarrow A \mid a$ $\text{First}(S) = \text{First}(A) \cap \text{LL}(1)$ \times
 $A \rightarrow a$ $\text{First}(A) = \{a\}$ $\text{Follow}(A) = \{\$ \}$ $\text{LR}(0)$ \times
 $\text{SLR}(1)$ \times

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

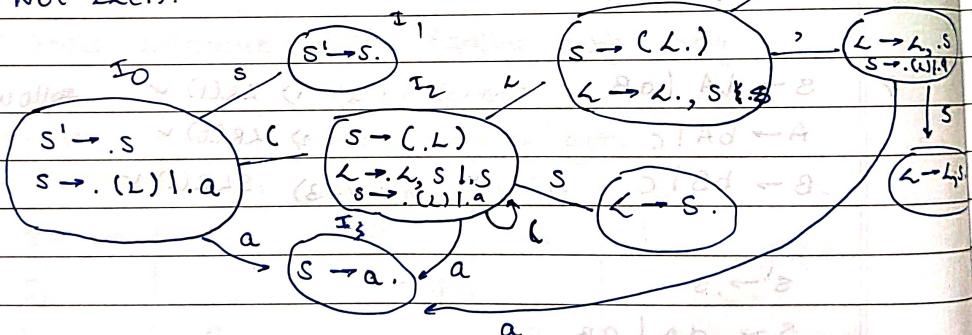
$$\text{First}(A) = \{a\} \quad \text{Both same. Not LL(1)}$$

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



3) $S \rightarrow (L) \mid a$ $\text{First}(S) = \{(, a\}$
 $L \rightarrow L, S \mid S$ $\text{First}(L) = \text{First}(L) \cap \text{First}(S) = \{(, a\}$

i) NOT LL(1).



$\text{LR}(0) \leftarrow S \rightarrow \text{First}(S) = \{ \mid a \}$

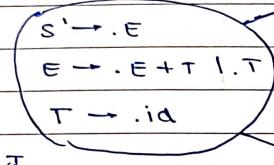
$\text{SLR}(1) \leftarrow \text{First}(L) = \text{First}(L) \cap \text{First}(S)$

LOOK ahead of the symbol

left

$$4) E \rightarrow E + T \mid T$$

$$T \rightarrow id$$



$$\text{First}(E) = \text{First}(e) \cap \text{First}(T)$$

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

$$+ \quad (E \rightarrow E + T \mid T \rightarrow .id) \rightarrow (E \rightarrow E + T)$$

$$5) S \rightarrow T + E \mid T$$

$$LL(1)$$

$$\Rightarrow S \rightarrow i \text{ (CSS, 1st)}$$

$$T \rightarrow id$$

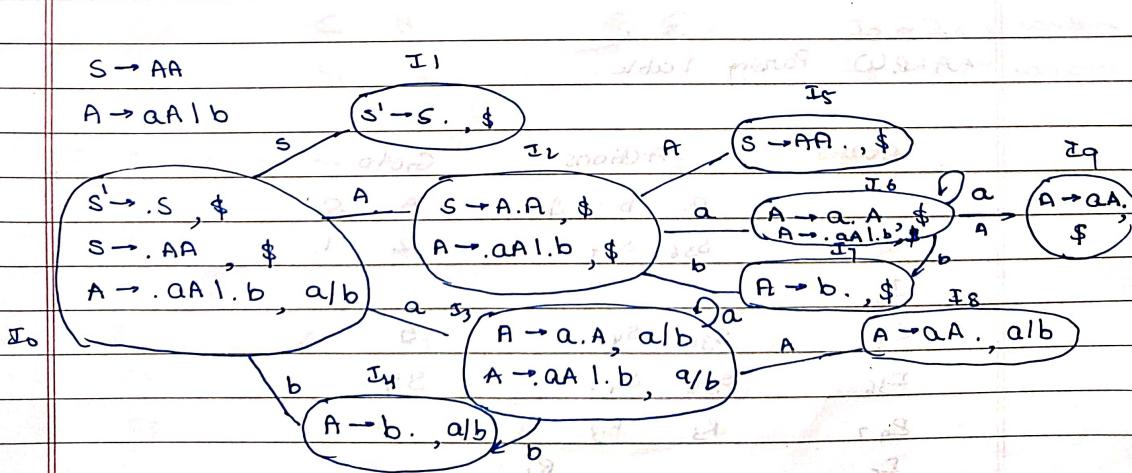
$$\text{First}(E) = \text{First}(T) = \{id\}$$

$$\text{First}(T) = \{id\}$$

$$id \cap id \neq \emptyset \text{ NOT LL(1)}$$

CLR Parser: Canonical Lookahead

$$\text{LA LR}(1) - \text{CL LR}(1) = \text{uses LR}(1) item} = \text{LR}(0) \text{ item} + \text{lookahead.}$$



CLR(G) Parsing Table

states	Actions			Goto
	a	b	\$	
I ₀	S ₃	S ₄		A → S
I ₁				
I ₂	S ₆	S ₇		5
I ₃	S ₃	S ₄		8
I ₄	R ₃	R ₃		
I ₅				(R ₁ , I ₂)
I ₆	S ₅	S ₇		9
I ₇				R ₃
I ₈	R ₂	R ₂		(S ₃ T → (T) T → T → (T) T → T)
I ₉				R ₂

I₃ and I₆ - CombineI₁ and I₇ - CombineI₈ and I₉ - Combine

NLR(G) Parsing Table

states	Actions			Goto
	a	b	\$	
I ₀	S ₃	S ₄	d. IAD → A	2 1
I ₁				
I ₂	S ₃ G	S ₄ G		S
I ₃ G	S ₃ G	S ₄ G	d. A → A	8 9
R ₄ G	R ₃	R ₃	R ₃	(d. d. d. A → A)
I ₅				R ₁

I_{gg} $R_L \quad R_2 \quad R_2$

After grammar derivation

4

27

 $E \rightarrow BB$ (1) $\rightarrow J_1$ $B \rightarrow CB, Id$

(2) (3)

E

B

 $E' \rightarrow E, \$$ $E \rightarrow .BB, \$$ $B \rightarrow .CB, cl/d$ $B \rightarrow .d, cl/d$

c

E'

B

E'

B

E'

B

E'

B

I₀

d

I₁

d

I₂

c

I₃

d

I₄

d

I₅

d

I₆

d

I₇

d

I₈

d

I₉

d

I₁₀

d

I₁₁

d

I₁₂

d

I₁₃

d

I₁₄

d

I₁₅

d

I₁₆

d

I₁₇

d

I₁₈

d

I₁₉

d

I₂₀

d

I₂₁

d

I₂₂

d

I₂₃

d

I₂₄

d

I₂₅

d

I₂₆

d

I₂₇

d

I₂₈

d

I₂₉

d

I₃₀

d

I₃₁

d

I₃₂

d

I₃₃

d

I₃₄

d

I₃₅

d

I₃₆

d

I₃₇

d

I₃₈

d

I₃₉

d

I₄₀

d

I₄₁

d

I₄₂

d

I₄₃

d

I₄₄

d

I₄₅

d

I₄₆

d

I₄₇

d

I₄₈

d

I₄₉

d

I₅₀

d

I₅₁

d

I₅₂

d

I₅₃

d

I₅₄

d

I₅₅

d

I₅₆

d

I₅₇

d

I₅₈

d

I₅₉

d

I₆₀

d

I₆₁

d

I₆₂

d

I₆₃

d

I₆₄

d

I₆₅

d

I₆₆

d

I₆₇

d

I₆₈

d

I₆₉

d

I₇₀

d

I₇₁

d

I₇₂

d

I₇₃

d

I₇₄

d

I₇₅

d

I₇₆

d

I₇₇

d

I₇₈

d

I₇₉

d

I₈₀

d

I₈₁

d

I₈₂

d

I₈₃

d

I₈₄

d

I₈₅

d

I₈₆

d

I₈₇

d

I₈₈

d

I₈₉

d

I₉₀

d

I₉₁

d

I₉₂

d

I₉₃

d

I₉₄

d

I₉₅

d

I₉₆

d

I₉₇

d

I₉₈

d

I₉₉

d

I₁₀₀

d

I₁₀₁

d

I₁₀₂

d

I₁₀₃

d

I₁₀₄

d

I₁₀₅

d

I₁₀₆

d

I₁₀₇

d

I₁₀₈

d

I₁₀₉

d

I₁₁₀

d

I₁₁₁

d

I₁₁₂

d

I₁₁₃

d

I₁₁₄

d

I₁₁₅

d

I₁₁₆

d

I₁₁₇

d

I₁₁₈

d

I₁₁₉

d

I₁₂₀

d

I₁₂₁

d

I₁₂₂

d

I₁₂₃

d

I₁₂₄

d

I₁₂₅

d

I₁₂₆

d

I₁₂₇

d

I₁₂₈

d

I₁₂₉

d

I₁₃₀

d

I₁₃₁

d

I₁₃₂

d

I₁₃₃

d

I₁₃₄

d

I₁₃₅

d

I₁₃₆

d

I₁₃₇

d

I₁₃₈

d

I₁₃₉

d

I₁₄₀

d

I₁₄₁

d

I₁₄₂

d

I₁₄₃

d

I₁₄₄

d

I₁₄₅

d

I₁₄₆

d

I₁₄₇

d

I₁₄₈

d

I₁₄₉

d

I₁₅₀

d

I₁₅₁

d

I₁₅₂

d

I₁₅₃

d

I₁₅₄

d

I₁₅₅

d

I₁₅₆

d

I₁₅₇

d

I₁₅₈

d

I₁₅₉

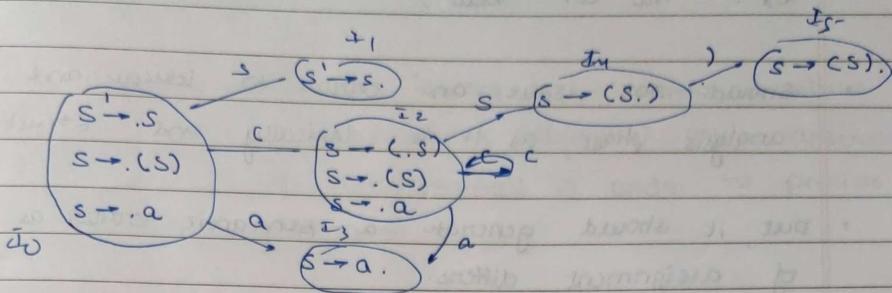
d

I₁₆₀

d

I₁₆₁

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



LALR : Same as CLR with one difference, i.e., if 2 states differ only in look ahead then combine those states.

After minimisation if Parsing Table has no conflict, the grammar is LALR.

LALR may contain RR conflict.

$$LR(0) = n_1$$

$$n_1 = n_2 = n_3 \leq n_4$$

$$SLRF = n_2$$

$$LALR = n_3$$

$$CLR = n_4$$

operator Precedence Parser:

It generates the parse tree from given grammar & string but the only condition is 2 consecutive non-terminals & (E) never appear on RHS of production

$(S -> (S), 1)$

Semantic Analysis

Semantics of a language provide meaning to its constructs, like tokens and syntax structure. Semantics help interpret symbols, their types and their relations with each other.

CFG + Semantic Rules = Syntax Directed Definitions.

Ex : int a = "value";

- should not issue an error in lexical and syntax analysis phase, as it is lexically and structurally correct.
- but, it should generate a semantic error as the type of assignment differs.

These rules are set by the grammar of the language and evaluated in semantic Analysis.

Tasks of Semantic Analysis :

- Scope Resolution : Checks the scope of the variables and also checks the scope Resolution operator.
- Type checking : It checks the compatible types in variables.
- Array-bound checking
- Flow control check : Keeps a check that control structures are used in a proper manner (Ex: no break statement are outside a loop)

Semantic Errors

Some of semantic errors that semantic analyzer recognises are:

- Type Mismatch
- Undeclared variable
- Reserved identifier misuse
- Multiple declaration of Variable in a scope
- Accessing an out of scope variable
- Actual & Formal parameter Mismatch.

Attribute Grammar :

It is a special form of context-free grammar where some additional information (Attributes) are appended to one/more of its non-terminals in order to provide context-sensitive information.

Each attribute has well-defined Domain of values, such as integer, float, character, String and expressions.

Attribute grammar is a medium to provide semantics to the context-free grammar and it can help specify the Syntax and Semantics of programming language.

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

The RHS of CFG contains the semantic rules that specify how the grammar should be interpreted.

Here, the values of Non-terminals $E \& T$ are added together and result is copied to the Non-terminal E .

Semantic attributes may be assigned to their values from their domain at the time of Parsing and evaluated at the time of Assignment.

Based on the way Attributes get their values, there are 2 categories:

- Synthesized attributes
- Inherited attributes

Synthesized Attributes:

These attributes get values from the attribute values of their child nodes.

Ex: $S \rightarrow ABC$

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 .

- Synthesized attributes never take values from their parent nodes or any other sibling nodes.

Inherited Attributes:

Inherited attributes can take value from parent or siblings.

Ex: $S \rightarrow ABC$

'A' can get value from S, B and C .

'B' can take values from S, A , and C .

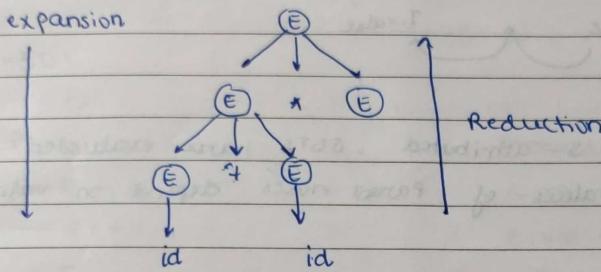
'C' can get values from S, A and B .

- Expansion: when a non-terminal is expanded to terminals as per a Grammatical Rule

- Reduction: when a terminal is reduced to its corresponding non-terminal according to Grammar Rules.

NOTE: Syntax trees are parsed top-down and left to right.

Whenever reduction occurs, we apply its corresponding semantic rules.



Semantic Analysis uses Syntax Directed Translations that perform expansion & reduction actions which facilitate in Semantic Analysis.

Semantic Analyzer receives AST (Abstract Syntax Tree) from Syntax Analysis and attached attribute information with AST, which are called Attributed AST.

Attributes are 2 tuple values $\langle \text{attribute name}, \text{Attribute value} \rangle$

Ex: int value = 5;

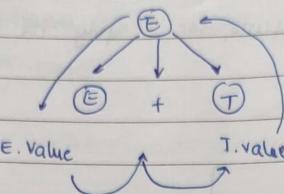
$\langle \text{type}, \text{"integer"} \rangle$

$\langle \text{present value}, \text{"5"} \rangle$

S-Attributed SDT

If an SDT uses only synthesized attributes, it is called as S-attributed SDT.

$$E.\text{Value} = E.\text{Value} \cup T.\text{Value}$$



Attributes in S-attributed SDTs are evaluated in Bottom-up Parsing, as values of Parent nodes depend on values of Child node.

1 - Attributed SDT:

This form of SDT uses both synthesized and inherited attributes with restriction of not taking values from Right siblings.

In 1-attributed SDTs, a Non-terminal can get values from its parent, child and left sibling nodes only.

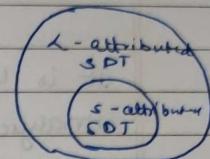
Ex: $S \rightarrow ABC$

- S can take values from A, B and C (synthesized)
- A can take values from only S.
- B can take values from S and A
- C can get values from S, A, B.

Attributes in L-attributed SDTs are evaluated by Depth-First and left to Right Parsing.

Top-down

S-SDT



Production

$L \rightarrow E \sqcap *$

$E \rightarrow EI + T$

$E \rightarrow T$

$T \rightarrow TI * F$

$F \rightarrow F$

$F \rightarrow (E)$

$F \rightarrow \text{digit}$

After

$E \rightarrow$

$T \rightarrow$

$TI \rightarrow$

$F \rightarrow$

$(E) \rightarrow$

$\text{digit} \rightarrow$

Semantic Rules

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

$E.\text{val} = EI.\text{value} + T.\text{value}$

$E.\text{val} = T.\text{value}$

$T.\text{value} = TI.\text{value} + F.\text{value}$

$T.\text{value} = F.\text{value}$

$F.\text{value} = E.\text{value}$

$F.\text{value} = \text{digit. devval}$

L-SDT

Production

Semantic Rules

$D \rightarrow TL$

$L.in = T.type$

$T \rightarrow \text{int}$

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

$T \rightarrow \text{real}$

$\text{real}.type = \text{real}$

$L \rightarrow L id$

$L.in = L.in$

$L \rightarrow id$

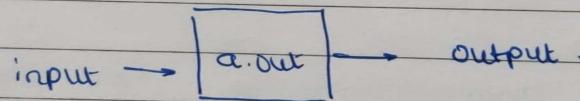
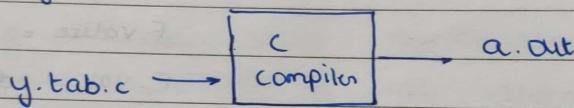
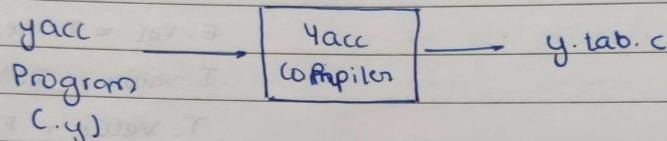
$\text{add type (identity, L.in) }$

Symbol Tabl.

Yacc - Yet another Compiler - Compiler

It is used to produce the source code of the Syntax Analyser of the language produced by LALR(1) grammar.

It is a Parser Generation. The input of Yacc is the rules / grammar along with lexemes / token. And the output is a C program.



Structure of Yacc Program:

It consists of 3 sections:

- Declaration section
- Rules section
- Auxiliary Functions

Declaration

1. I

1.1.

declaration / definitions

Rules

1. }

1. 1.

1. tokens

1. 1.

Rules

1. 1.

Auxiliary Functions

→ declaration Section: This part of YACC has 2 sections; both are optional. The first section has ordinary declarations delimited by '1. I' and '1. }'. Any temporary variable used by the 2nd & 3rd section will be kept in this part.

Ex:

1. I

#include < stdio.h>

1. }

1. token digit

Declaration of grammar tokens also takes part in this section. This part defines the tokens that can be used in later parts.

→ Rules Section: After the first '1. 1.', we place the translation Rules. Every rule has a grammar production and associated semantic action.

<head> ; {Body1>

{semantic action1>}

1 <body2>

{<semantic action2>}

1 <body n>

{<semantic action n>}

EX : S : A \$ B { prints ("Valid"); }

In semantic action, the symbol $\$$ is considered to be an attribute value associated with head Non-terminal.

→ Auxiliary Functions :

It is the last section of the YACC program and it includes the necessary C functions needed for the program to be executed along with few yacc functions such as `yyparse()`, `yyerror` etc.

Meta characters of Yacc

- $\$\$$ - is the attribute value of the left hand production
- $\$i$ - is the attribute value of the RHS of Production.
- `y.tab.h` - is the header file that has to be included so that the yacc program gets to know the tokens used and accessed from the lex program.

Symbol Table

It is an important data structure created and maintained by the compiler in order to keep track of semantics of variables i.e., it stores information about the scope and binding information about names, information about instances of various entities such as variable and function names, classes, objects etc.

- It is built-in lexical and syntax analysis phase.
- The information is collected by the analysis phases of the compiler and is used by the synthesis phases of the compiler to generate code.

Analysis phases : Lexical Analysis, Syntax Analysis, Semantic Analyzer.

Synthesis phases : intermediate code generator, code optimiser, target code generator.

- It is used by the compiler to achieve compile-time efficiency.
- It is used by various phases of compiler as:
 - Lexical Analysis : Creates new table entries in the table, for ex: like entries for about tokens.
 - Syntax Analysis : Adds information regarding attribute type, scope, dimension, line of inference, use in the table.
 - Semantic Analysis : Uses available information in the table to check for semantics i.e., to verify that expressions and assignments are semantically correct and update it.

- Intermediate Code generation: Refers Symbol Table for knowing how much and what type of run-time is allocated and Table helps in adding temporary variable information.
- Code optimization: Uses information present in the Symbol Table for Machine-dependent optimization.
- Target Code Generator: Generates code by using address information of identifiers present in the Table.

Symbol Table Entries: Each Entry in the Symbol Table is associated with attributes that support the compiler in different phases.

Items stored in Symbol Table:

- Variables names and constants
- Procedure and function names
- Literal Constants and Strings
- Compiler generated temporaries
- Labels in source language

Information used by the compiler from Symbol Table:

- Data type and name
- Declaring Procedures
- offset in storage
- If structure or record then, a pointer to structure Table.
- For Parameters, whether parameters Passing by value or by Reference
- Number and Type of arguments passed to Function
- Base address.

operations of Symbol Table - The basic operations defined on a Symbol Table includes:

Operation	Function
allocate	to allocate a new empty symbol table.
Free	to remove all entries and free storage of symbol table
Look up	to search for a name & return pointer to its entry.
Insert	to insert a name in a symbol table and return a pointer to its entry
Set-attribute	to associate an attribute with a given entry
Get-attribute	to get an attribute associated with a given entry.

Implementation of Symbol Table

Commonly used data structures for implementing Symbol Table are:

→ Ordered / Unordered List:

- In this Method, an array is used to store names and associated information.
- A pointer "available" is maintained at end of all stored records and new names are added in the order as they arrive.
- Insertion is fast $O(1)$, but lookup is slow for large tables - $O(n)$ on Average.

- The Advantage is that it takes minimum amount of Space.

→ Linked List:

- A link field is added to each record. A pointer "First" is maintained to point to the first record of the symbol table.
- Insertion is fast $O(1)$, but look up is slow - $O(n)$.

→ Hash Table:

- In Hashing Scheme, 2 tables are maintained - a hash table and a symbol table.
- A Hash Table is an array with an Index Range: $0 \text{ to } \text{table_size} - 1$. These entries are pointers pointing to the names of the symbol table.
- Insertion and Look up can be very fast - $O(1)$.

Search

→ Binary Tree:

- Another approach to implement is binary search tree.
- We add 2 link fields i.e., left and right child.
- Insertion and lookup are $O(\log n)$ on average.

A Symbol Table stores its entry in :

<Symbol name, type, attribute>

Ex: static int stats;

<Stats, int, static>

operations

- insert(): This operation is more frequently used by analysis phase where tokens are identified and names are stored in the table.

This operation is used to add information in the Symbol table about unique names occurring along with their attribute, value, state and scope.

Ex: int a; is processed as

insert(a, int);

- lookUp(): It is used to search a name in Symbol Table to see:

- if the symbol exists in the Table
- if it is declared before its use
- if the name is used in scope
- if symbol is initialised
- if symbol is declared multiple times.

lookUp(Symbol)

If symbol exists in Table, it returns its attributes stored in the Table.

Scope Management

A compiler maintains 2 types of Symbol Tables:

- a global Symbol Table which can be accessed by all procedures
- a scope Symbol Table that are created for each scope in the program.

Ex: int value = 10;

void pro_one()

{

int one_1;

int one_2;

{

int one_3;

int one_4;

}

int one_5;

void pro_two()

{

int two_1;

int two_2;

{

int two_3;

int two_4;

}

int two_5;

}

Symbol	Type	Scope
Pro_one	Proc	global
Pro_two	Proc	global
value	int	global

Symbol	Type	Scope
One_1	int	Proc para
One_2	int	proc para
One_5	int	proc para

Symbol	Type	Scope
two_1	int	Proc para
two_2	int	proc para
two_5	int	proc para

Symbol	Type	Scope	Scope	Type	Scope
one_3	int	inner	two_3	int	inner
one_4	int	inner	two_4	int	inner

Runtime Environment and Storage Management

By Runtime, we mean a program in execution. Runtime Environment is a state of the Target machine, which may include software libraries, environment variables to provide services to the processes running in the system.

Runtime support system is a package, mostly generated with the executable program itself and facilitates the process of communication b/w the process and the runtime environment. It takes care of the memory allocation and de-allocation while the program is being executed.

Procedure activation

A Program is a sequence of instructions combined into a number of procedures.

Instructions in a procedure are executed sequentially.

A procedure has a start and an end delimiter and everything inside is called body of the procedure.

The execution of procedure is called its activation.

An activation record contains all the necessary information required to call a procedure.

They contain

Temporaries : Stores temporary and intermediate values of ~~an~~ ~~exp~~

Local Data : Stores local data of the called procedure.

Machine Status : Stores machine status such as Registers, Program Counter etc before the procedure is called.

Control Link : Stores the address of activation record of the caller procedure.

Access Link : Stores the information of data which is outside the local scope.

Actual Parameters : Stores actual Parameters i.e., parameters which are used to send input to called Procedure.

Return Value : Stores return value.

whenever a procedure is executed, its activation record is stored on the stack, also known as control stack.

A series of activations are easier to represent in form of a tree, known as activation tree.

Ex : `printf(" Enter Name\n");`

`scanf("%s", Name);`

`Show-data(Name);`

`printf(" Press any key");`

Return value

Actual para
meter

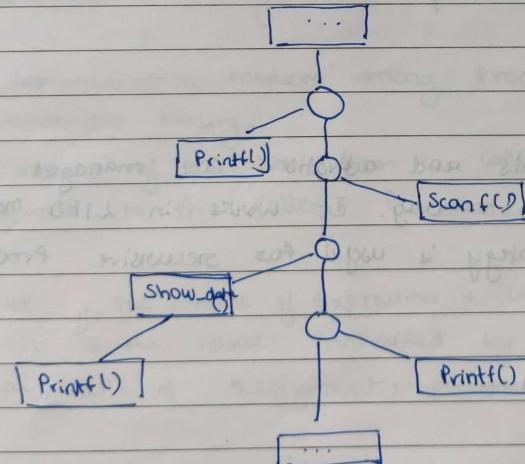
Control
link

Access
type

Machine
status

Local Data

Temporary



Storage Allocation

Runtime environment manages Runtime memory requirements for:

- Code
- Procedures
- Variables

Static Allocation :

In this Allocation Scheme, the compilation data is bound to a fixed location in the memory and it does not change when the program executes.

As the memory requirement and storage locations are known in advance, run-time support package for memory allocation & de-allocation is not required.

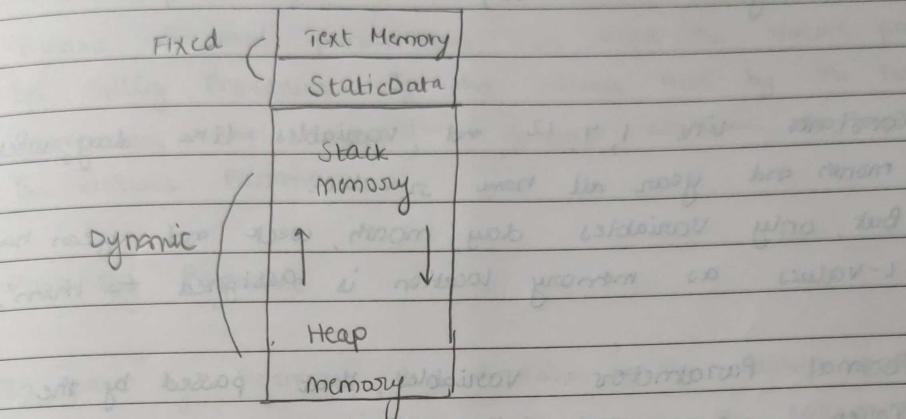
Stack Allocation :

Procedure calls and activations are managed by means of stack memory. It works in LIFO method and this strategy is useful for recursive procedure calls.

Heap Allocation

Variables local to a procedure are allocated & de-allocated only at runtime. Heap memory is used to dynamically allocate memory to the variables and claim it back when not required.

Except static allocation memory area, both stack and heap memory can grow / shrink dynamically.



The text part of code is also allocated a fixed amt of memory.

Parameter Passing

- The communication medium among procedures is known as parameter passing.
- The values of variables from a calling procedure are transferred to the called procedure.

r-value : The value of expression is called its r-value, and it is the value contained by a single variable on the RHS of assignment operation.

l-value : The location of memory (address) where expression is stored is known as l-value of that expression. It appears on the LHS of assignment operation.

Ex:

```

day = 1;
week = day * 7;
month = 1;
year = month * 12;

```

Constants like 1, 7, 12 and variables like day, week, month and year all have no-value.

But, only variables day, month, week and year have l-values as memory location is assigned to them.

Formal Parameters: Variables that passed by the Caller procedure are called Formal Parameters. These variables are declared in the definition of the called fn.

Actual Parameters: Variables whose values are being passed to the called procedure are called actual parameters and are specified in the function call as arguments.

Ex:

```

main()
{
    int x = 10, y = 20;
    add(x, y); // → actual calling
}

void add(int a, int b) // → Formal, called no
{
}

```

```

    printf("sum = %d", (a+b));
}

```

Pass by Value

The calling procedure passes the I-value of the actual parameters and compiler puts that into the called procedure's Activation Record. Formal parameters thus hold the value passed by the calling procedure. If the values held by the Formal Parameters are changed; it should not have no impact on the actual parameters.

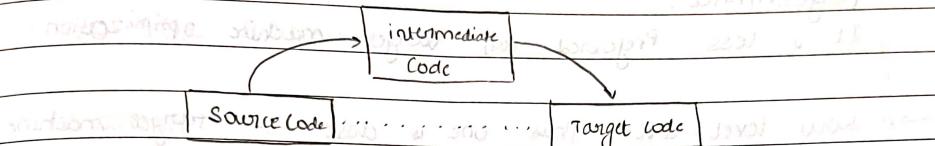
Pass By Reference

In pass by Reference, the I-value of the actual parameter is copied to the activation record of the called procedure. The called procedure now has the address of the actual parameter and formal Parameters refer to the same memory location. If the value pointed by Formal Parameters is changed, then it impacts on Actual Parameter as well.

Pass By Copy - Restore

Intermediate Code Generation

Reason why intermediate code generation is required.



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

Intermediate Representation

They can be represented in a variety of ways and they have their own benefits.

→ High level IR : This is very close to source language itself. They can be easily generated from source code and we can easily apply code modifications to enhance performance.

It is less preferred for target machine optimization.

→ Low level IR : This one is close to target machine, which makes it suitable for register and memory allocation, instruction set selection etc.

It is good for machine-dependent optimizations.

Intermediate code can be either language specific (Byte code for Java) or language independent (3-address code).

The following are commonly used intermediate code representation:

→ Postfix Notation :

The expressions are written in infix notations are

The Postfix Notation for the some expression places the operator after the operands

$a b +$

No parenthesis are needed in Postfix Notation becoz the position and no. of arguments of operators permit only one way to decode postfix expression.

ab - cd + * ab +

$$\text{Ex: } (a-b) * (c+d) + (a-b)$$

(a-b) (c+d)

$$T_1 \quad T_2 \quad T_3 = T_1 T_2 *$$

$$T_4 \quad T_2 = cd +$$

$$ab - cd + * ab - b : \text{Post fix Notation. } T_4 = ab - b$$

$$T_3 + T_4$$

$$T_3 T_4 +$$

$$ab - cd + * ab - +$$

Three - Code Address

A statement that has 3 references i.e., 2 for operands and one for result is known as 3 - address statement.

A sequence of 3 - address statement is known as 3 - address code.

They are of the form $x = y \text{ op } z$ hence x, y and z will have address.

Ex: 3 - address code for $a + b * c / d : (a) - abcd$

$T_1 = b * c$

$T_2 = a + b * c$

$T_3 = T_2 / d$

T_1, T_2 and T_3 are temporary variables

Syntax tree:

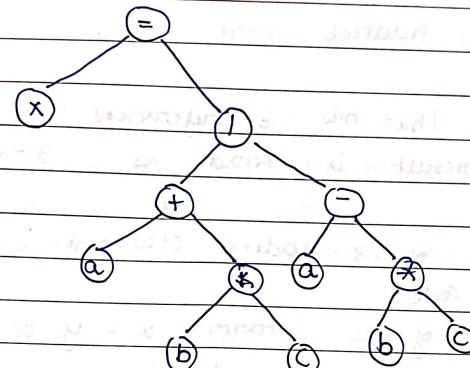
It is nothing more than condensed form of a Parse tree.

The operator and keyword nodes of the parse tree are moved to their parents and a chain of single productions is replaced by single link in syntax tree.

The internal nodes are operators, child nodes are operands.

$$\text{Ex: } x = (a + b * c) / (a - b * c)$$

$$x = (a + (b * c)) / (a - (b * c))$$



Code optimization

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

Objectives of code optimization:

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

when to optimize?

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

Why optimize?

Optimizing an algorithm is beyond the scope of the code optimization phase.

So the program is optimized, it may involve reducing the size of code.

Optimization helps to:

- Reduce the space consumed and increases the speed of compilation.
- Manually performing the optimization is tedious, hence we use a code optimizer.
- An optimised code often promotes re-usability of code.

Types of code optimization

• Machine Independent Optimization:

It attempts to improve the intermediate code to get a better target code as the output. The part of intermediate code which is transformed does not involve any CPU registers or absolute memory locations.

→ Machine Dependent Optimization:

It is done after the target code has been generated and here when the code is transformed, it involves CPU registers and absolute memory references. It put efforts to take max advantage of memory hierarchy.

Different ways of code optimization

1. Compile Time Evaluation:

$$\text{if } A = 2 * (22.0 / 1.0) * 5$$

Performs $2 * (22.0 / 1.0) * 5$ during compile time

$$\text{ii) if } x = 12.4 \text{ then } y = 12.3$$

$$y = x / 2.3$$

Evaluate $x / 2.3$ as $12.4 / 2.3$ at compile time

2. Variable Propagation

$$c = a * b$$

$$x = a$$

After

$$c = a * b$$

till expression contains x

$$d = x * b + 4$$

$$d = a * b + 4$$

Hence, $a * b$ and $x * b$ will be identified as common sub-expression.

3. Dead Code elimination:

variable propagation often leads to making assignment statement into dead code.

$$c = a * b$$

$$x = a$$

till

$$d = a * b + 4$$

Aff

elimination

$$c = a * b$$

till

$$d = a * b + 4$$

4. Code Motion:

- o Create such a modified set of unique statements.
- o Reduce the evaluation frequency of expression.
- o Bring loop invariant statements out of the loop.

$$a = 200 ;$$

while ($a > 0$)

T

$$b = x + y ;$$

if ($a \% b == 0$)

$$a = 200 ;$$

$b = x + y ;$

while ($a > 0$)

$$\leftarrow \text{loop body}$$

if ($a \% b == 0$)

loop body

5. Induction variable and strength Reduction:

- o An induction variable is used for loop form following
- o $i = i + \text{constant}$ \rightarrow high strength with i
- o Strength reduction means replacing the high strength operators by the low strength.

$$i = 1 ;$$

while ($i \leq 10$)

$$\left\{ \begin{array}{l} y = i * 4 ; \\ \end{array} \right.$$

```

i = 1
t = 4
while (t < 40)
    y = t;
    t = t + 4;
}

```

6. Constant folding:

→ It eliminates the expressions that calculate a value before code execution.

Ex: $i < 320 * 200 * 32$

$$= 320 * 0 \quad \downarrow \quad 320 * 0 = 0$$

$$i < 0 \quad (0 < 0) \text{ is true}$$

Basic Blocks

Source code generally have a no. of instructions, which are always executed in sequence and are considered as the basic blocks of the code.

Basic Block Identification

They are important for code optimization and code generation.

Source Code

Basic Blocks

w = 0;

x = x + y;

y = 0;

if (x > z)

{

y = z;

x++;

}

else

{

w = x + z;

w = 0;

x = x + y;

y = 0;

if (x > z)

B1

y = x;

x++;

B2

y = z;

z++;

B3

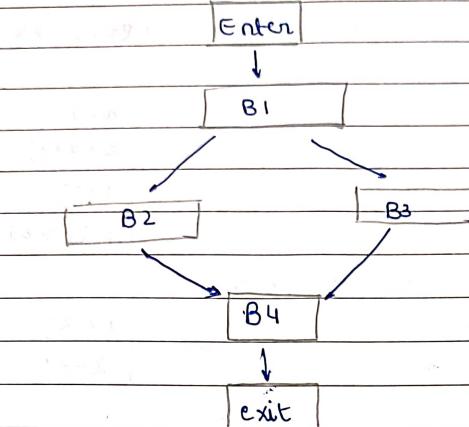
w = x + z;

w =

They play an imp role in identifying variables, which are being used more than once in a single block.

Control flow graph:

Basic blocks in a Program can be represented by means of control flow graphs. It depicts how the program control is being passed among the blocks. It helps in optimization by locating any unwanted loops.



Phases of optimization (Scope of optimization)

Local optimization: The optimization which occurs within the basic block.

Techniques used in local optimization are:

- ## 1) Local common Subexpression elimination

$$\begin{aligned} \Delta t &= y * z + K \quad \text{on} \quad \rightarrow \quad T_1 = y * z \quad \text{is the result} \\ A &= y * z - m \quad \quad \quad x = t_1 + K \end{aligned}$$

- 2) Dead code elimination within the basic block.

These games in a broader sense are called **recreational**.

if (0) all the rights Control would never be given to it

3 printf ("dead");

- ### 3) Constant Propagation

$$\begin{array}{l} x = 10 \\ A = B + 5 \\ z = 10 + y \end{array}$$

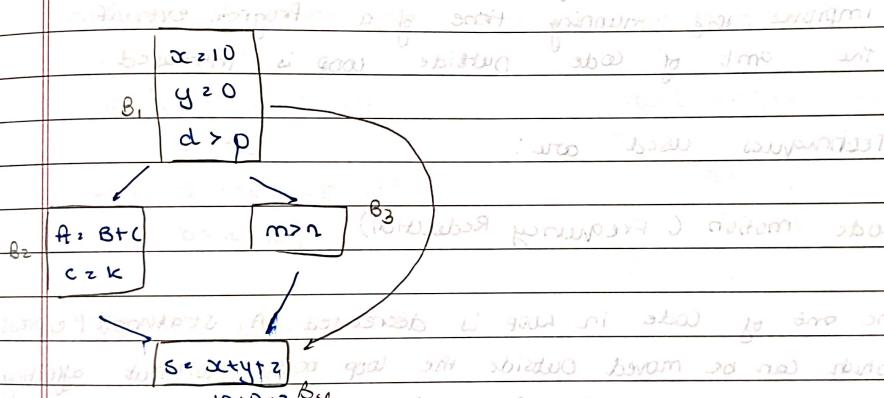
- Algebraic Laws (associative, commutative)
→ Redefining statements.

Global Optimization

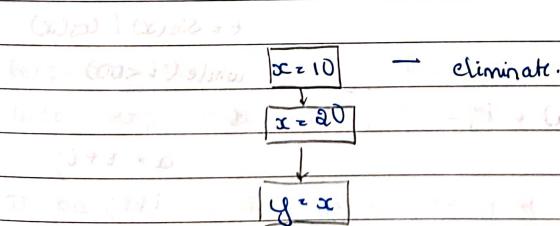
Transformations are applied to large Program segment that includes functions, Procedures and no loops, between the basic blocks.

The optimizations are extended to entire Control flow graph:

i) Constant Propagation b/w basic blocks



ii) Dead Code elimination



Loop Optimization: Comes under Machine Independent Optimization

It is the process of increasing execution speed and reducing the overheads associated with loops.

It plays an imp role in improving Cache Performance and making effective use of Parallel Processing capabilities.

Most execution time of a scientific program is spent on loops.

Decreasing the no. of instructions in an inner loop improves the running time of a program even if the amt of code outside loop is increased.

Techniques used are:

- Code motion (Frequency Reduction)

The amt of code in loop is decreased. A statement / expression which can be moved outside the loop body without affecting the semantics of the program, is moved outside the loop.

`while(i<100)`

$t = \sin(x) / \cos(x)$

`i`

$\xrightarrow{\text{Separate}} \quad \text{while } (i < 100)$

$a = \sin(x) / \cos(x) + i;$

$\quad \quad \quad i$

`i++;`

$a = t + i;$

`}`

`i++;`

→ Loop Unrolling:

It helps to optimize the execution time of a program. we basically remove or reduce iterations.

It increases the program's speed by eliminating loop control instruction and loop test instructions.

```
for (int i=0; i<3; i++)
    printf("Pankaj\n");
printf("Pankaj\n");
printf("Pankaj\n");
printf("Pankaj\n");
```

→ Loop Jamming (Fusion):

It is combining the 2 or more loops in a single loop. It reduces the time taken to compile many loops.

```
for (int i=0; i<5; i++)
    a = i+5;
for (int i=0; i<5; i++)
    b = i+10;
```

Peep-hole Optimization:

The small set of instructions on small part of code is called peephole or window.

Peep-hole optimization is performed on a small part of the code, very small set of instructions in a segment of code.

It basically works on the theory of replacement. Part of code is replaced by shorter and faster code without change in output.

Peephole is a Machine Dependent optimization:

The objectives of Peephole Optimization is to effect local optimization on the code.

- To Improve Performance
- To Reduce memory footprint
- To Reduce code size

Peephole optimization Techniques:

a) Redundant load and store elimination; in this, redundancy is eliminated

Move b, R0

Add c, R0

Move R0, a

Move a, R0

Add e, R0

Move R0, d

Move b, R0

Add f, R0

Move R0, a

Add e, R0

Move a, R0, d

b) Strength Reduction

The operators that consume higher execution time are replaced by the operators consuming less execution time.

$$x^2 \Rightarrow x * x$$

$x * 2 \rightarrow$ left shift \Leftarrow value $\ll 2$ bits are shifted to the left

$x / 2 \Rightarrow$ Right shift \Leftarrow $x \gg 2$, where 2 is shift to the right

3) Simplify Algebraic expressions

$$a = a + 0$$

$$a = a * 1$$

$$a = a / 1$$

~~a = a - 0~~ ~~cancel out 0 with 0~~ ~~is simplifying redundant~~

4) Unreachable code ~~is code that you can never run~~

int add (int x)

{
 return (x+10);
}

printf("Hello");

5) Flow of Control optimization

Mov R1, R2 ; R1 = R2 and now Mov R1, R2

Goto L1 ;

Goto L2 ; R = R

L1 : goto L2

R = R ; R = R

L2 : INCR R

R = R + 1

Instruction scheduling (for Pipeline)

Instruction scheduling is one of the compiler optimization passes that occurs in the back end of the code generation.

Front End → Optimization → Instruction Selection
→ Instruction Scheduling → Pre-Pass

Register Allocation → Post-Pass

Instruction Scheduling - Post-Pass

It is typically for a compiler to have 2 instruction Scheduling passes, one before & one after Register Allocation.

Instruction scheduling

The main purpose of this is to Reorder instructions, in such a way that there is less load when executing in Pipeline.

i.e., if a load is an example of an instruction that has long latency.

Ex: If a write instruction is followed by Read instruction then the CPU has to fetch the value before it writes into it.

A load instruction has typically a latency of few cycles of 3-4 cycles.

A : Load I 5 → m_1

B : Load $\alpha \rightarrow m_2$

C : add $m_1, m_2 \rightarrow m_3$

D : Load E 7 → m_4

E : Load Y → m_5

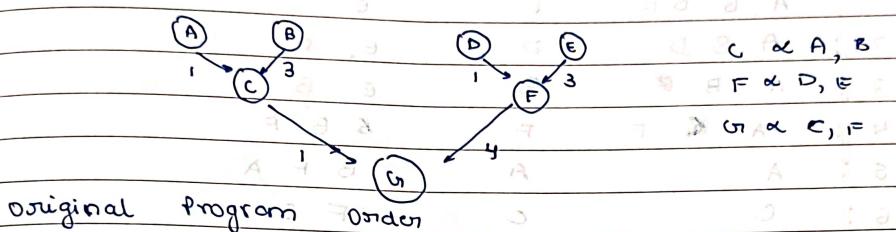
F : mult $m_4, m_5 \rightarrow m_6$

G : add $m_3, m_6 \rightarrow m_7$

If we consider the load latency has 3 cycles and mult latency is 2 cycles, and rest of latency of 1 cycle.

Assuming it is a single issue processor, the ordering of the instructions that done based on the target machine Model.

We construct a Data Dependency graph before scheduling the instructions.



1st cycle : A 2nd cycle : B (takes 3 cycles) 12 : stall

3rd cycle : stall 4th cycle : stall 13 : stall

5th cycle : C In total, it takes 14 cycles.

6th cycle : D

7th cycle : E (takes 3 cycles)

8th cycle : stall

9th cycle : stall

10th cycle : F (takes 4 cycles)

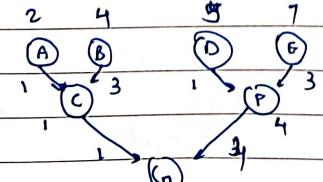
11th : stall

The code optimizer will minimize the no of cycles by scheduling other instructions.

It can be done through a technique called List Scheduling which keeps track of Ready List.

distance

It goes with the critical path, first of all instruction X is the length of the longest path from X to a key instruction 2 4 5 7



	Ready List	INSTR	Active List
1 :	A B D E	E	E
2 :	A B D	D	E, D
3 :	A B F	B	E B
4 :	A F	F	E B F
5 :	A	A	E B F A
6 :	C	C	F C (sort of) LIFO
7 :	Stall		F
8 :	G	G	A (shifted)

Hence the stall is just one after re-ordering and is executed by 8 cycles with the help of pipeline.

Register Allocation

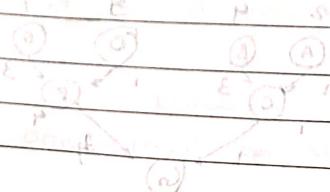
Registers are the fastest locations in the memory hierarchy, but their resource is very limited.

The Register Allocation of a Compiler maps the virtual registers to Physical registers for storage and reuse.

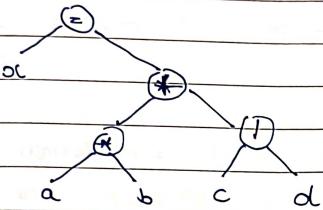
The Register allocation optimization takes into account of the register pressure as they are limited in no. of registers physically present on the machine.

Let us consider an example of an expression

$x = a * b + c / d$



Drawing an Abstract Syntax tree, we get.



The linear Assembly of instructions will be

Load a $\rightarrow r_1$

Load b $\rightarrow r_2$

~~Half Load~~ $r_1, r_2 \rightarrow r_3$

Load c $\rightarrow r_4$

Load d $\rightarrow r_5$

Dir $r_4, r_5 \rightarrow r_6$

Add $r_3, r_6 \rightarrow r_7$

Store $r_7 \rightarrow x$.

The virtual registers required are 7 registers. But

if we are using a compiler which has only 3 physical registers that can be used.

Then we apply register allocation technique to use it in the optimized way.

Assuming that the Target machine has only 3 physical registers P_1, P_2, P_3 . we can see how the optimization technique works

For that we need to keep track of live-range of the time of definition and its usage instruction.

The Main aim of Register allocations is to use min no of physical registers and to reuse registers for optimization.

$$r_1 \rightarrow P_1$$

$$r_2 \rightarrow P_2$$

$$r_3 \rightarrow P_1$$

$$r_4 \rightarrow P_2$$

$$r_5 \rightarrow P_3$$

$$r_6 \rightarrow P_2$$

$$r_7 \rightarrow P_1$$

$$P_1 \rightarrow x$$