

COMPILER DESIGN NOTES

SEM-6

R-MADHULLIKA

106121071

CSE-A'25

Lexical Analysis :

Scan the given input
line by line
char by char

Scan left to right of find words in the sentence

Valid statement : I is in CSPCG2
acc to lex analyzer a = b + c;

In terms of compiler, var is called an identifier

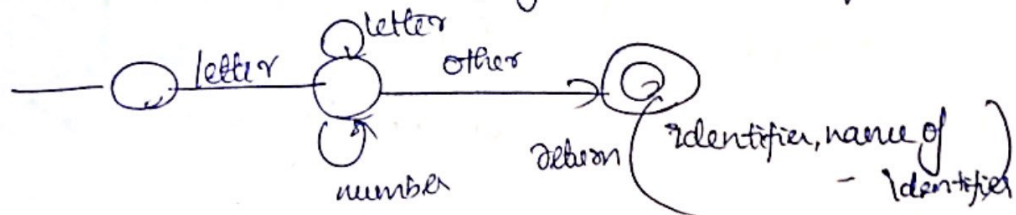


identifier :

letter (letter) number

When a letter/number is not scanned, due to white space or eof then the analyzer will stop and take var as everything before it.

Doesn't decide on anything as correct/incorrect → next phase will just scan, create the words & give it to next phase.



If not an identifier or a white space, it needs to be scanned and given a tokenization

The pointer needs to be moved back

whatever scanned extra is given back to the input & scanned sequentially

Hence,

the unnecessary spaces will be removed & only information is passed to next phase

Syntax Analyzer:

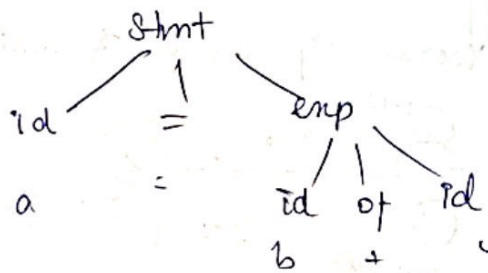
whether syntax is fine w.r.t the written language format of the information

not look at the meaning, just check whether the specifications are correct.

meaning - semantics \times doesn't look
syntax - structure

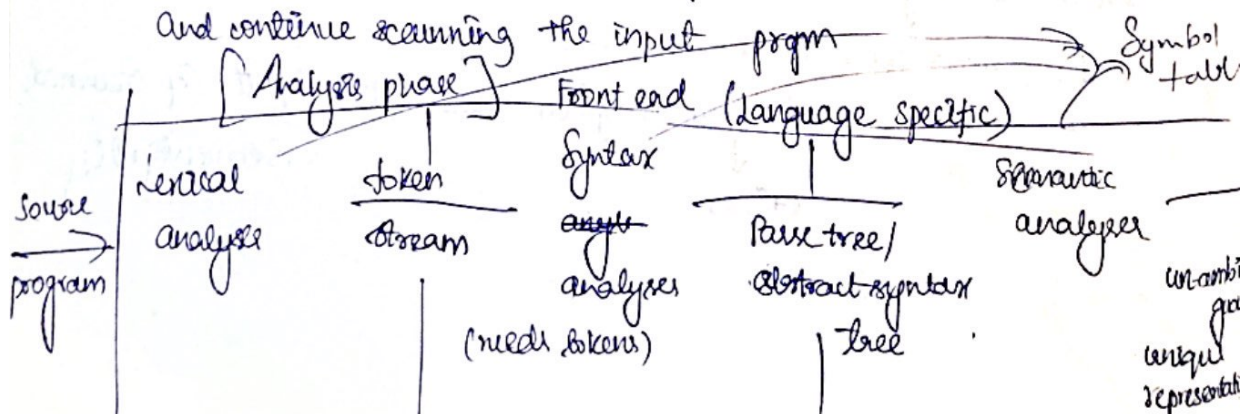
parse tree - 3 line structure

ex: stmt \rightarrow id = exp
exp \rightarrow id of id



When an error is detected, it reports & skips the error

and continue scanning the input program



Semantic analyser :

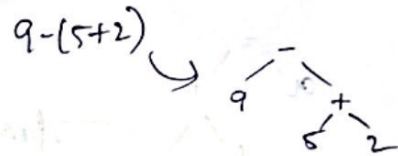
The type checking can be done but not all meanings

produce statement with only one interpretation without ambiguity

1. If C_1 then if C_2 then else C_3 } dangling else problem
 } ambiguous statements

$9 - 5 + 2$ → not ambiguous since $*$ has higher precedence

but $-, +$ has same precedence ∴ it can be $(9 - 5) + 2$



the semantic analyser

should be fed the correct representation

dead code part which will never be executed

Optimiser : (optional phase) → can be skipped

removes extra parts of the program

↓
dead code

Common subexpression elimination

- some expressions which are repeated

$$X = Y + Z$$

$$W = Y + Z + K$$

Compiler sees

$Y + Z$ calculated twice

$$W = X + Z$$

a variable which is of no use in the calculation of loop, the declaration goes out of loop by compiler

reduce the computation effort in the loop

Code motion

copy propagation

constant folding

strength reduction

$$X = 5 * 2 \rightarrow X = 10$$

whenever constant calculation comes up compiler puts the result into var

$$X = X * 2$$

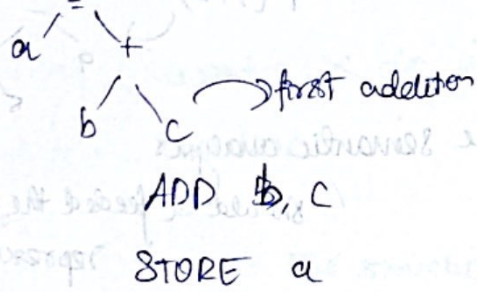
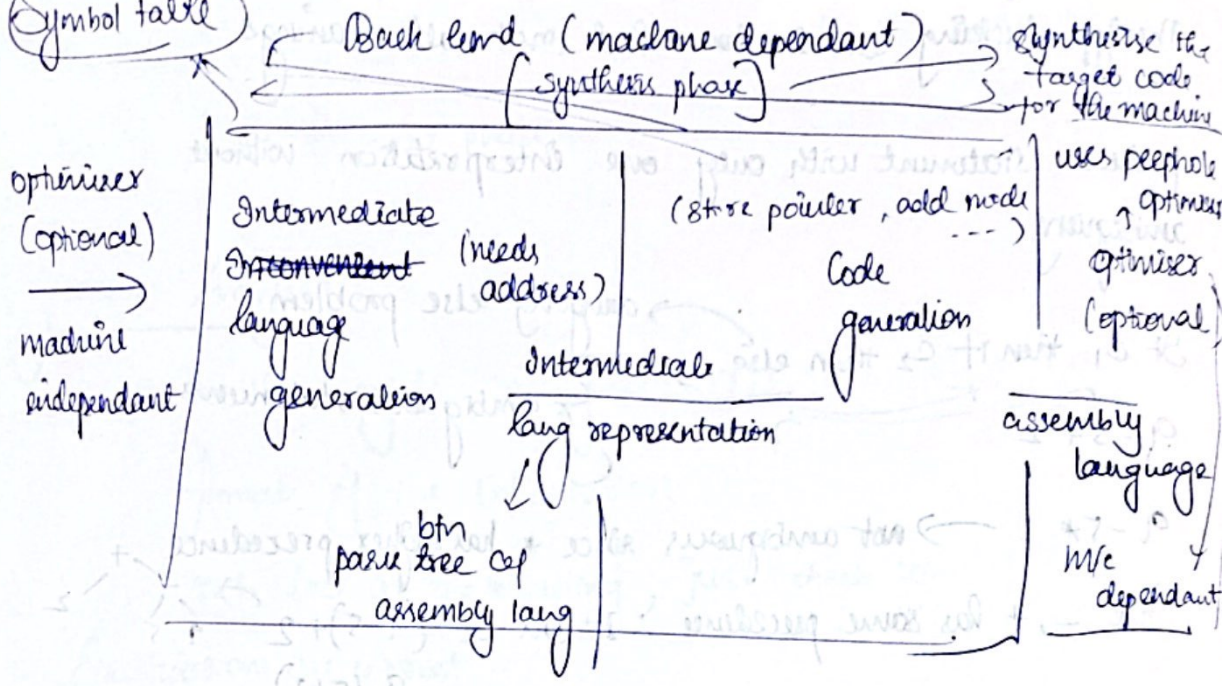
↳ entire operation

can be replaced by

$$X \ll 1$$

$$X + X$$

Symbol table



postfix, prefix - stack notation are intermediate languages

st, a[i] → [Add (a[i]) + 1 location]

redundancy occurs with repeated calculations
↓
can occur in code generation stage

one optimiser atleast is recommended

int count = 0;

<id, count>

The information produced by compiler needs to be stored
∴ the code can be executed only after passing other phases (comp, assembly)

Needs to be stored:

Type of token

name " "

-type information

amount of storage needed,

Address,

pointer

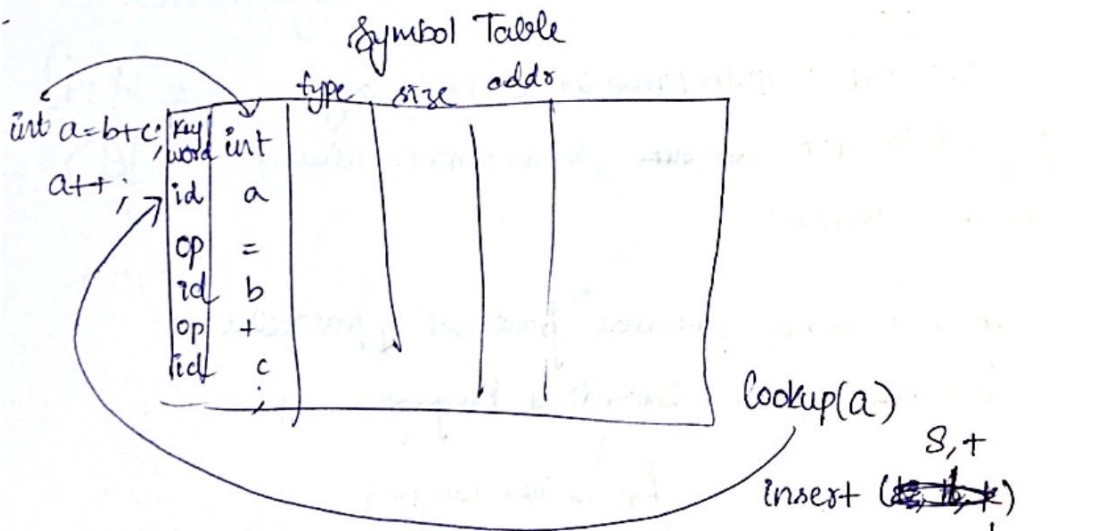
target language
→ code

(H.T, array)

Whenever a data structure is used to store, the problem is

when a value is changed, it needs to be changed at cell instances each time by the compiler

Hence a ^{global} table for all phases, where every phase can store its structure in the table



Compiler gives the line number of the error to debug the code

token along w/ addr
token id " "
string " "

S - seq of char - a

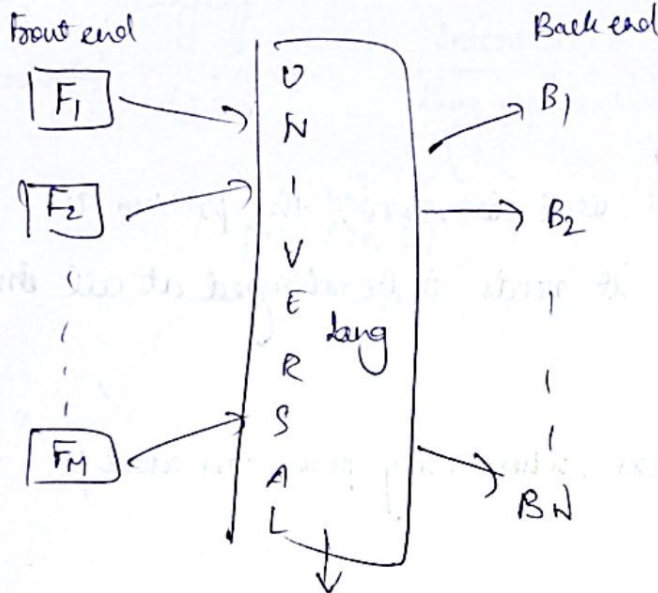
t - token - id

Many logical errors cannot be detected by compiler

$M \times N$

no. of prog lang * no. of architecture - not possible

" + " "



UCOL

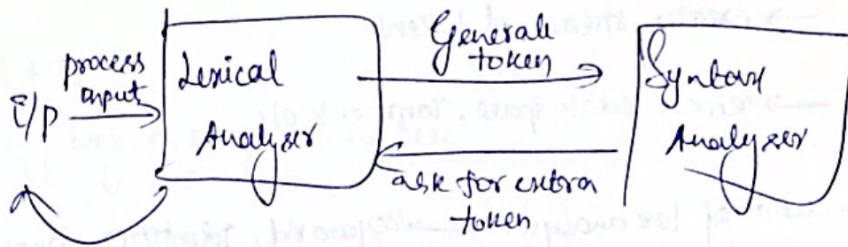
Universal computer/compiler oriented lang

proposed in 1961 ~~but not~~, its an abstract idea and never materialized

- compiler to map particular front end of particular back end of has a intermediate language
↓
for similar language

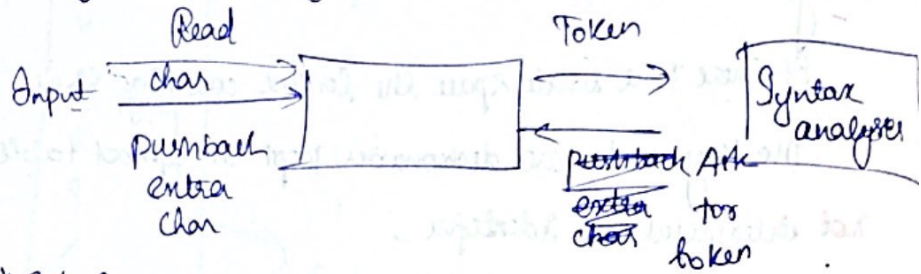
Test suite:
↓
costly operation

no. of specifically written large programs to check the entire program test



extra character
to be put back
to the input

30/11/24 Implementing Lex. Analyser:



Let

$$E = m * C + C$$

input ~~loop~~ output operation which requires disc access are costly
Hence A block of data is stored in a buffer.

The whole $E = m * C + C$ will be in the same block

one var - Start of lexeme

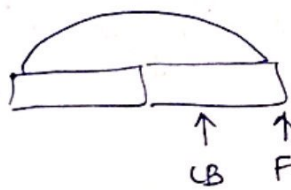
another var - to indicate the end

Suppose the lexeme starts at the end of a buffer

Hence lookahead - req extra info

another buf stores extra char

When buffer ends, we need to reload the buffer



Hence a circular loop is implemented

→ even then the lexeme is not decided to be

If there are more than 2 buffer of lexeme length, then this strategy fails

Block size = 4096 bytes

We need to check:

1. whether the char is end of buffer
2. the char which ends the lexeme

2 io operations

∴ use a special char which marks the end of char

When a variable eof is loaded, it either marks the end of buf or the end of lexeme.

eof - either not need to do anything or load a buffer

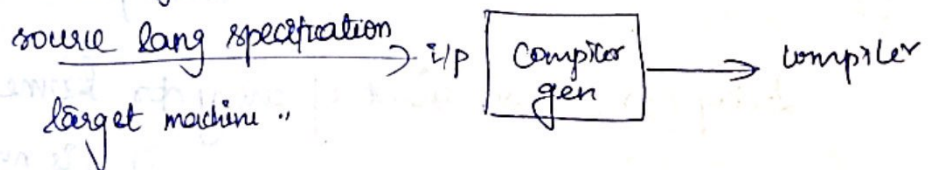
Implementation approaches

→ To write an assembly lang for the io

when to start, read, end - - -

→ take 2 arrays, 2 pointers and implement in high level lang

→ Tools for compiler generator - Lex ⇒ Flex (Fast Lex)



Lex (lex analyser generator)

YACC

(yet another compiler compiler)

Tools give a particular structure to use it.

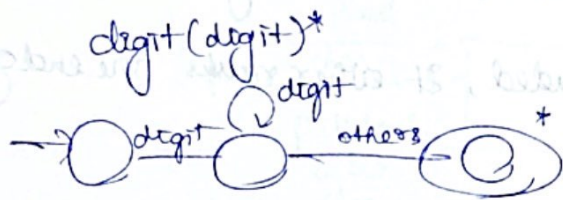
Target of any compiler is to move ~~any~~ as much closer to assembly lang as possible

Construct a lexical analyser

The way of specifying the specs is the challenge
Returns tokens to the syntax analyser
A finite set of tokens is defined

The code given is trying to detect whether the lexeme is a number
(digit)⁺

while(1) → Inf loop cos our lexeme can have any no. of char



return(token,
Token value is returned as lexeme
tokenval - lexeme associated w token
not the num

Symbol table:

Stores information of every token

Insert(s, t): :: same lexeme s and token t and
return pointer

Lookup(s): return index of entry for lexeme s or 0
if s is not found

the structure of a symbol table should be

token	Value	type	Address
id			
number			
keyword			
operator			
separator			

If there are
m tokens
need log₂ m bits

can take 4 bits
∴ 15 tokens can be represented

32 bytes - general max length of a variable

$$(k+32) \times N$$

max allowed
 $K =$ size of the token

$N =$ no of variables defined

size of 2 col

Ex: 4 tokens var size

00 id 32 B

01 num 16 B

10 op 8 B

11 key word 32 B

\therefore the size is wasted
 allocate memory dynamically



define what is a letter & digit

id = letter (letter/digit)* \rightarrow letter = A|B|...|Z|a|b|...|z

digit = 0|1|...|9

defining a regular expression with a name - regular definitions

later we can replace the R.E with the name

Ex:

country \rightarrow (digit)*

area \rightarrow (digit)+

exchange \rightarrow (digit)+

phone \rightarrow (digit)+

fax no. \rightarrow country - area - exchange - phone

Ex:

mail id

letter \rightarrow a|b|...|z|A|B|...|Z

name \rightarrow (letter)+

Address \rightarrow name @ name . net

Floating point

digit \rightarrow 0|1-...|9 e \rightarrow empty string
digits \rightarrow digit⁺
optfraction \rightarrow . digits | e
optexponent \rightarrow (E(+|-|e) digits) | e
number \rightarrow digits optfraction optexponent

123.45
123.45E
12E-3
0.012

? \rightarrow optional

Use of shorthand notation

digit \rightarrow 0-9
digits \rightarrow digit⁺
number \rightarrow digits (. digits)? (E[+-]? digits)?

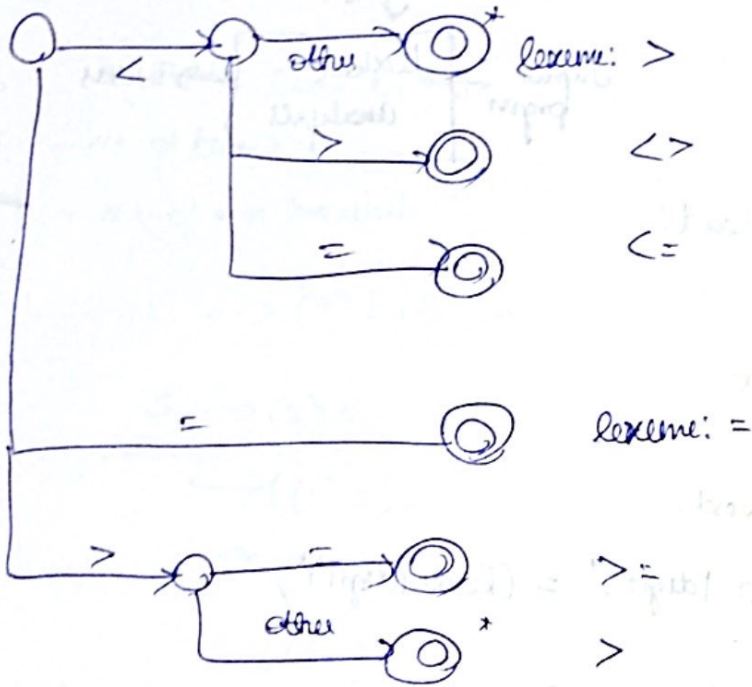
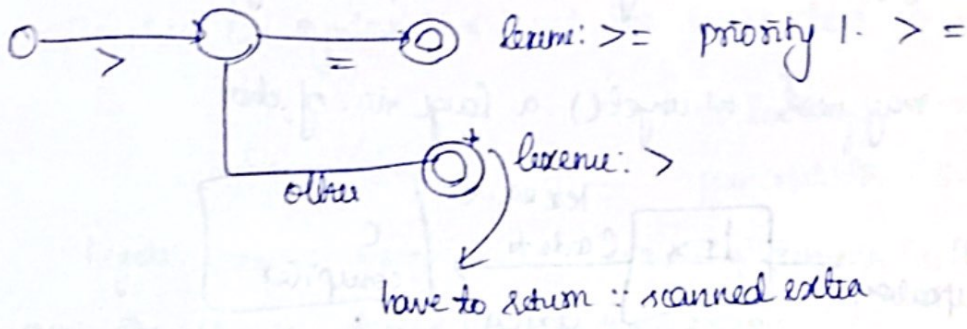
Implementation of specification

1. Write RE for lexemes of each token
number \rightarrow digit⁺
identifier \rightarrow letter (letter (digit))⁺
2. Construct ϵ matching all
lexemes of all tokens
 $R = R_1 + R_2 + R_3 + \dots$
3. Let input be $x_1 \dots x_n$
check $x_1 \dots x_n \in L(\epsilon)$
4. $x_1 \dots x_n$

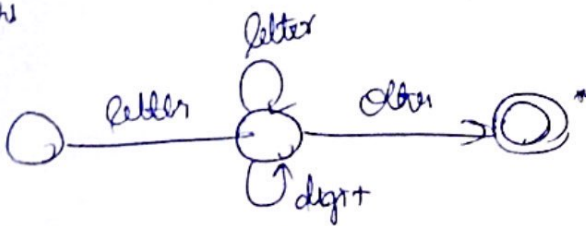
(checking if any part of it
matches with RE)

Relational operators

relOp \rightarrow $\langle | \leq | = | \neq | \geq | = \rangle$



Identifiers



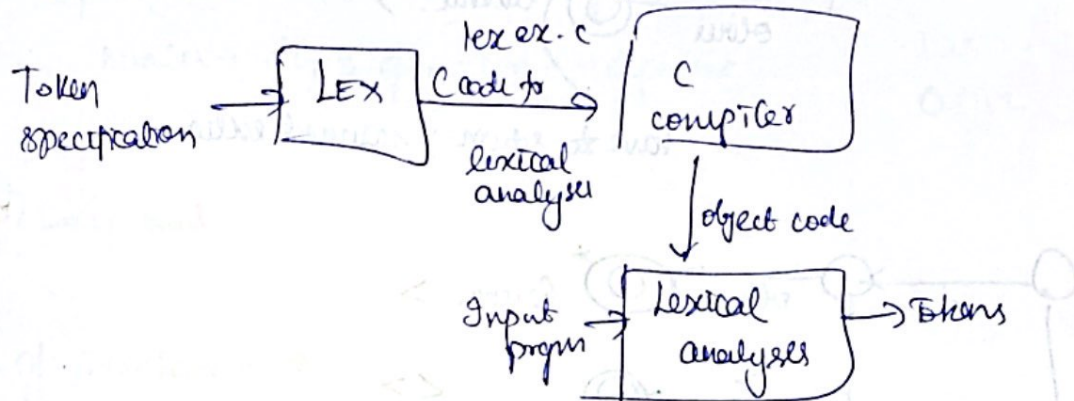
White spaces



For each state, add switch cases

Complexity of transition diagram inc the prone of code.

Tradeoff: may need to `ungetc()` a large no. of char



Implementing lookahead:

DO 10 | = 1, 25

DO 10 | = 1, 25

Spec for DO as keyword:

DO / (letter | digit)* = (letter | digit)*,

Structure of LEX prog:

declaration

translation rule

auxiliary func

Translation rule

Pattern {Action}

`installID()` returns a pointer to symbol table placed in `yglobal`

`ytext`: pointer to beginning of lexeme

`yylen`: length of lexeme

`yglobal` - global variable

SYNTAX ANALYSIS

- Task of this phase is to check the ^{syntax of the} grammar

- To check the syntax, we need the parse tree of labels drive PDA
 ↓
 to look up, more powerful than the RE

Regular grammar cannot count, counting lemma is needed to count the opening, closing braces, indentation

S - ~~symbol~~ start symbol

T - set of terminals

NT - set of non terminals

$$S \rightarrow (S) S | \epsilon$$

$$S \rightarrow (S) S$$

$$\rightarrow ((S) S) S$$

$$\rightarrow (((S) S) S) S$$

$$\rightarrow ((((S) S) S) S) S$$

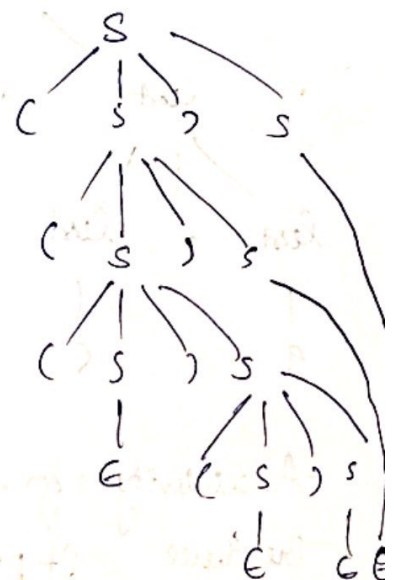
$$\rightarrow (((((S) S) S) S) S) S) S$$

$$\rightarrow ((((((S) S) S) S) S) S) S) S$$

((()))

((() ()))

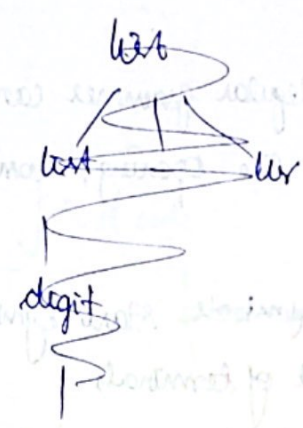
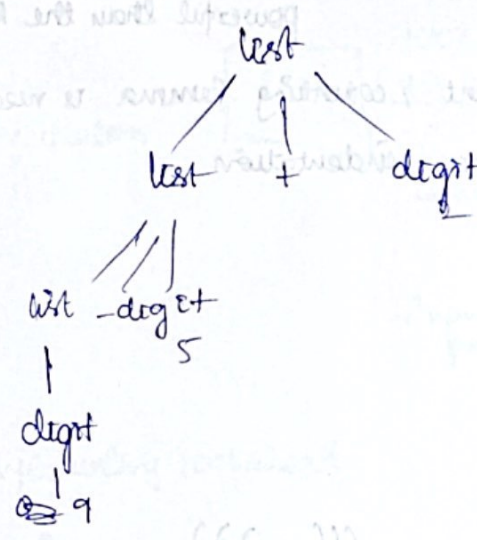
sentential form



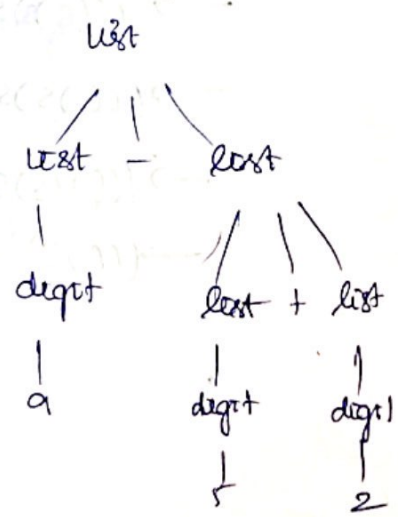
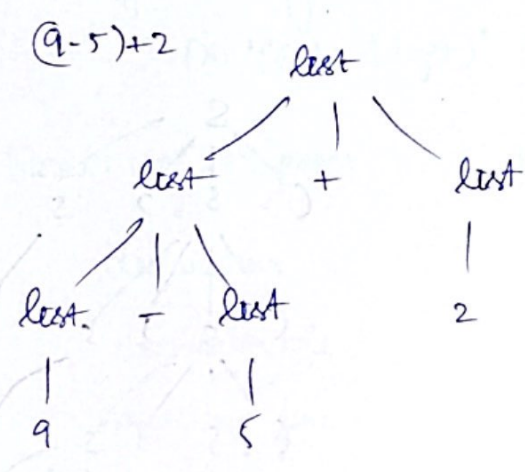
Q - 5 + 2 → parse tree for this using grammar 1, 2

1. $list \rightarrow list + digit \mid list - digit \mid digit$

$digit \rightarrow 0 \mid \dots \mid 9$



2. $list \rightarrow list + list \mid list - list \mid list \mid 0 \mid \dots \mid 9$



- Associativity : for every operator which values to be taken

- Precedence If precedence are same then check for associativity
 right associativity left associativity also cannot be compared.

$$9 \wedge 5 \wedge 2 \rightarrow 9 \wedge (5 \wedge 2) \checkmark$$

$$\downarrow$$

$$(9 \wedge 5) \wedge 2$$

$S \rightarrow SaS | b | d | e$

$\rightarrow SaA$

If $S \rightarrow b | d | e$ the A which is replacing S should produce $b | d | e$

$A \rightarrow b | d | e$

$list \rightarrow list + list \mid list - list \mid$

$list \ 10 - 9$

$digit \rightarrow [0 - 9]$

Another ex:

$S \rightarrow SaS | b$

Removing left recursive grammar

$S \rightarrow SaA | b$

$A \rightarrow b$

Syntax Analyzer :

CFG: $wAx \rightarrow w\beta x$

\uparrow
string of terminals

WLR

DPDAX

$\alpha, \beta \rightarrow$ string of terminals or non-terminal

If $A \rightarrow \beta$

In lex analyzer generated string of tokens \rightarrow terminals

$x = a_1 a_2 \dots a_n$

$w \rightarrow$ empty string / string of terminals

$E \rightarrow E + E \mid E - E \mid rd$

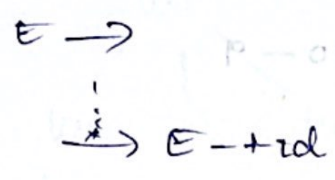
$E \rightarrow E + E$

$\rightarrow rd + E$

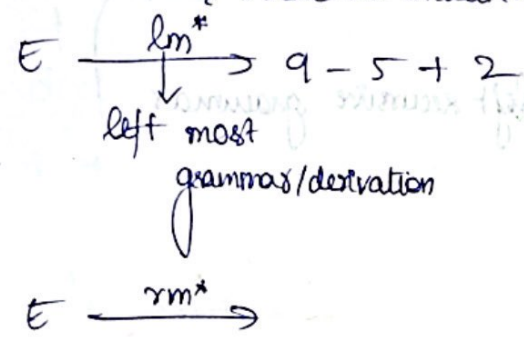
$\rightarrow rd + E - E$

In the parser, we check everytime whether the lexeme is in sentential form or not
 if not valid then report error

Panic mode: to complete whatever in the subsequent part of code
 skip one ~~invalid~~ part and take a valid sentence



The given grammar is ambiguous, any derivation has 2 parse trees
 → there are multi steps and not specified



Note: reverse
 NFA → DFA → min DFA

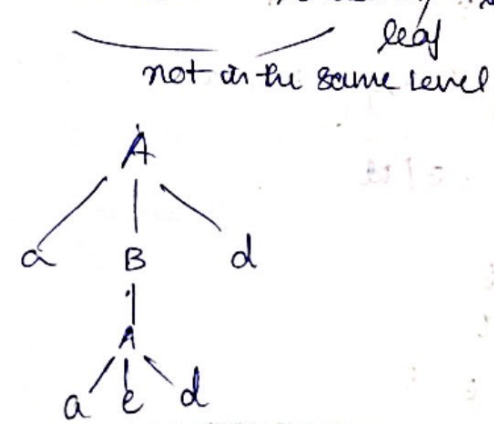
The compiler needs to remove the ambiguity

Top down parsing or Bottom up parsing - generates diff parse trees

In a parse tree,
 start symbol - root

intermediate nodes - non terminals (children) - terminals

$B \rightarrow A$
 $A \rightarrow aBde$
 aaedd



Associativity & precedence

On the lex prog, we define the ops by precedence

*/% * , / left

*/% + , - left

lets take the grammar, $E \rightarrow ETE | E-E | id$
 $E \rightarrow E+E | ETE$

$(9-5)+2$ - left associative

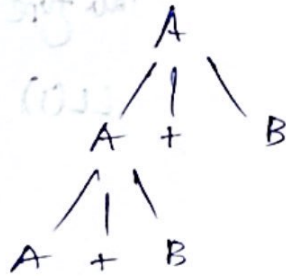
$9+(5+2)$ - right "

$A \rightarrow A+B$

$\rightarrow A+B+B$

$\rightarrow A+B+B+B$

⋮



left heavy / associative

any operator

$A \rightarrow B+A$

Right heavy

needs to be left associative

$E \rightarrow E+T | E-T | id$

$9-5+2$

$F \rightarrow id$

$F \rightarrow ETF | id \rightarrow$ right heavy

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

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

our order of exec is

go to next level

introduce $* , /$
 $F \quad + , -$

go to next level

introduce $+$
 $N-T T$

writing 3 set of prod rules to indicate 3 levels of precedence/associativity

The above is one way of removing ambiguity

Let $9 * 5 - 2$

Top down parsing
 $E \rightarrow E - T$
 $\rightarrow T - T$
 $\rightarrow T * F - T$
 $\rightarrow F * F - T$
 ~~$\rightarrow id * id - id$~~
 $\rightarrow id * F - T$
 $\rightarrow id * id - T$
 $\rightarrow id * id - F$
 $\rightarrow id * id - id$

Bottom up

Lookahead:

look ahead one symbol to decide on the prod rule.

This type of grammar is called

LL(1)

Top-down parsing:

$E \rightarrow ETE | E * E | (E) | E \uparrow E \text{ (id)}$

We have 4 diff levels of precedence here

$E \rightarrow E + F | F$

$F \rightarrow F * T | T$

$T \rightarrow Q \uparrow T | Q$

$Q \rightarrow (E) \text{ (id)}$

This non-ambiguous grammar is
 Not suitable for top down parser

↓
 have specific set of rules

the tree goes deeper as we derive each level deeper we start from left to right that Hence deeper is given more precedence

Left recursion

$$A \rightarrow Ax | b$$

Remove the left recursion

$$A \rightarrow bA$$

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

$$A \rightarrow Ax$$

$$\rightarrow Ax^2$$

$$\rightarrow bx \dots x$$

to end the no. of x's

After removing left recursion,

$$E \rightarrow FE'$$

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

$$F \rightarrow TF'$$

$$F' \rightarrow *TF' | \epsilon$$

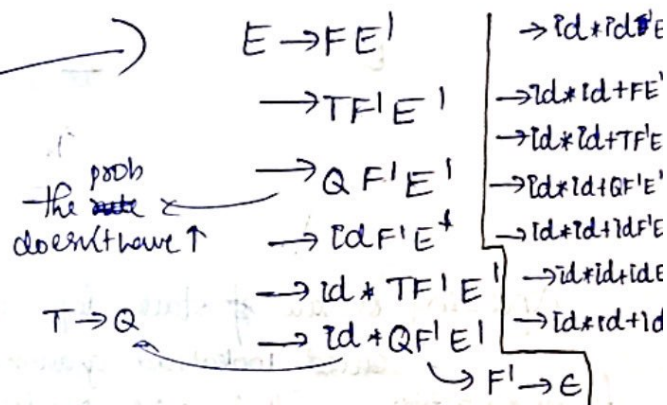
$$T \rightarrow Q$$

$$Q \rightarrow id$$

Let the

$$9 * 5 + 2$$

$$id * id + id$$



Stack

ip

Action

E

id * id + id

E → FE'

↑

pop E

push FE'

but we operate on F

∴ push E' first then F

FE'

id * id + id

F → TF'

↑

TF'E'

id * id + id

T → Q

↑

QF'E'

id * id + id

Q → id

↑

F'E'

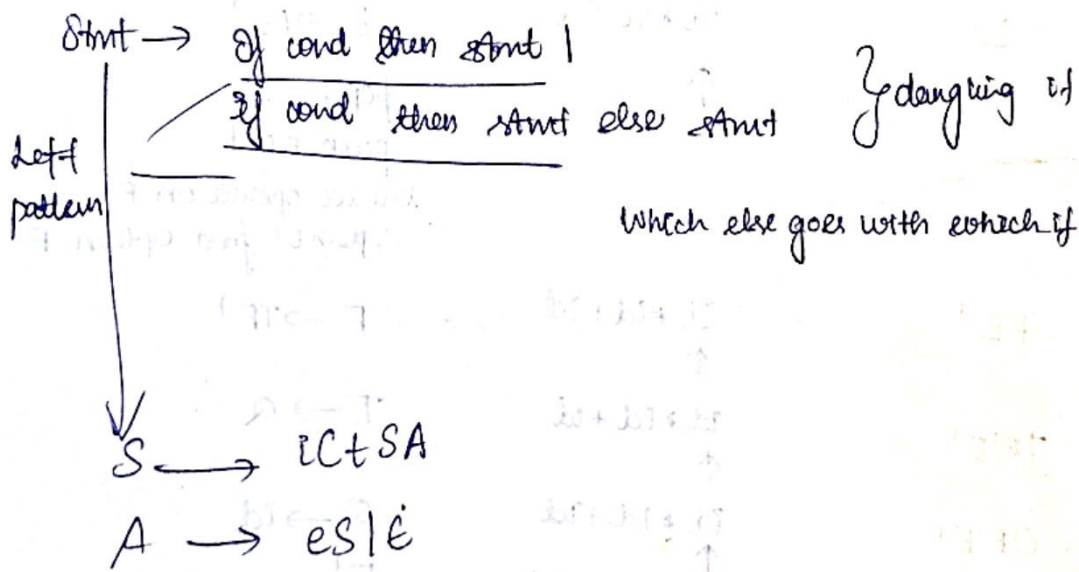
* id + id

F' → *TF'

↑

	Stack	I/P	Action
* is processed	TF'E'	id + id	$T \rightarrow Q$
	QF'E'	↑ id + id	$Q \rightarrow id$
	F'E'	↑ + id	$F' \rightarrow \cancel{+FE}$ E
	E'	↑ + id	$E' \rightarrow +FE'$
+ is processed	FE'	↑ id	$F \rightarrow TF'$
	TF'E'	↑ id	$T \rightarrow Q$
	QF'E'	↑ id	$Q \rightarrow id$
	F'E'	—	$F' \rightarrow E$
	E'	—	$E' \rightarrow E$
	—	—	HALT

Depending on the stack top and the input, the compiler called lookahead grammar (LL(1)) decides on action, left to right production with 1 position lookahead



Every then there is an associated else \rightarrow matched stmt

stmt \rightarrow matched stmt | unmatched stmt

We removed ambiguity, now

apply this condition every else goes with immediate \neq

Matched stmt \rightarrow If cond then matched stmt
else matched stmt | Others

unmatched stmt \rightarrow If cond then matched stmt
else unmatched stmt

extra else

or

some of not

matched with any else

If cond then stmt

associate the else with the nearest If

From the top down parsing,

how to know the input has exhausted (i/o op are costly)

Use some marker to indicate the input has exhausted

also " " " stack " "

X
stack top

a is current element

If $X = a$ (a terminal)

i/p symbol

(string of terminal)

just pop X and move read pointer to the next symbol

stack

input

$a = \$$

$X = \$$

TF

*

E operation performed

We directly skip the * but

it takes one more step, we push *

into stack and now stack top

symbol are equal the we move

the read pointer

a is a terminal \rightarrow pop stack top of
program i/p pointer $\rightarrow X = a$

$F \rightarrow *TF*$

(pop)

X is a non-terminal (x.a)

a is always a terminal \Rightarrow

but x can be non terminal

then apply the rule and reach to the id (terminal)

Recursive descent parsing

A

ANTLR \downarrow
tools to generate recursive descent parsing

\rightarrow backtracking

(not a good approach)

\rightarrow Predictive parsing

predict depending on the input to be processed
 \hookrightarrow lookahead

LL(K) - Generally $K=1$

\downarrow
no. of lookahead symbols needed to decide

no ambiguity

Generate the parse table for every possible state and select the correct.

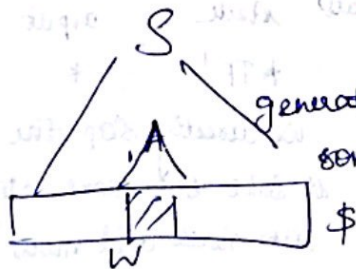
First set
Follow set

First set:

guide us on which prod rule to use

\hookrightarrow end marker following the sentence should be reached

(for start symbol)



generate a string w

some NT 'A' which produced a part of

1. Follow(S) contains \$

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

Follow(A) contains first(B)

~~First~~

If we apply A, what are the possible values \rightarrow first(A)
rules on A,

After applying all prod rules what follows # A gives follow set

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

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

E is our start symbol \therefore any derivation of any N.T derives an input with (or id

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

Note:

If E is the first symbol, then the next symbol acts as first

$X \rightarrow X_1 X_2 \dots X_n$

$\Rightarrow First(X) \neq \epsilon = X_2$

$X_1 \rightarrow \epsilon$

If all $X_1 X_2 \dots X_n \rightarrow \epsilon$ then only $first(x) = \epsilon$

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

3. $A \rightarrow \alpha B \beta$ & $B \rightarrow \epsilon$

β is empty \therefore

A ends with B whatever follow A will follow B

Follow(B) contains Follow(A)

4. $A \rightarrow \alpha B$

Follow(B) contains Follow(A)

13/2/24

$E \rightarrow TE'$
 $E' \rightarrow +TE' | -TE' | \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *T' | /T' | \epsilon$
 $F \rightarrow X \uparrow F | X$
 $X \rightarrow (E) | id$

$F \rightarrow XQ$
 $Q \rightarrow \uparrow F | \epsilon$

Remove left ~~recursion~~ factoring

Left factoring

$A \rightarrow \alpha B \beta | \alpha \gamma$
 takes the next symbol to find the prod rule
 decides the next symbol to find the prod rule

$A \rightarrow \alpha C$
 $C \rightarrow B \beta | \gamma$

Need to delay the X to get the next symbol \uparrow

$F \rightarrow XQ$
 $Q \rightarrow \uparrow F | \epsilon$

Q is following F ,
 F also follows Q
 \downarrow
 recursion

$First(E) = First(T)$
 $= First(F)$
 $= First(X)$
 $= \{C, id\}$

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

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

$First(Q) = \{\uparrow, \epsilon\}$

1. Follow(s) contains \$
2. $A \rightarrow \alpha B \gamma \rightarrow$ Follow(B) contains First(γ) except ϵ
3. $A \rightarrow \alpha B \gamma \rightarrow$ Follow(B) contains Follow(A)
 $\gamma \rightarrow \epsilon$
4. $A \rightarrow \alpha B$ \nearrow

Followed by

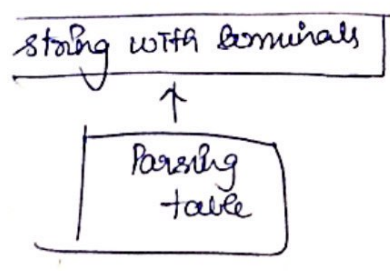
	E	E'	T	T'	F	X	Q
1.	\$		+		*	\$	\$
2.)		+ , -		* , /	↑	
3.			\$,)		+ , - , \$,)	* , / , + , - \$,)	
4.		\$,)		+ , - , \$,)			* , / , + , - , \$,)

$X \rightarrow \alpha$

For all $a \in \text{First}(\alpha)$

$[X, a] \leftarrow X \rightarrow \alpha$

the current input sees parsing table and decide which prod rule to apply



the input start needed to be process \rightarrow

stack \uparrow	id	+	-	*	/	()	\$
E	$E \rightarrow TE'$					$E \rightarrow TE'$		
E'		$E' \rightarrow TE'$	$F' \rightarrow TE'$				$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$					$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$	$T' \rightarrow /FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow XQ$					$F \rightarrow XQ$		
Q		$Q \rightarrow \epsilon$	$Q \rightarrow \epsilon$	$Q \rightarrow \epsilon$	$Q \rightarrow \epsilon$	$Q \rightarrow TF$	$Q \rightarrow \epsilon$	$Q \rightarrow \epsilon$
X	$X \rightarrow id$					$X \rightarrow (E)$		

Passing table

If the stack top is A

and we need to apply $A \rightarrow \epsilon$ to remove it and process follow(A)

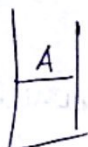
$$X \rightarrow \epsilon$$

for all b in follow(x)

$$[X, b] \leftarrow X \rightarrow G$$

EX $A \rightarrow BCD$

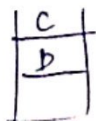
$C \rightarrow aD$



the input needed to be processed is a



$B \rightarrow \epsilon$



→ Many spaces are empty

What happens if the stack top is that ϵ when the input process is that, how to recover from the state

Stack	processing	O/P
$\boxed{\text{stack top}} \ \$ \epsilon$	$\text{id} + \text{id} * \text{id} \$$	$E \rightarrow TE'$ put in rev order in stack to keep T on top
$\$ \epsilon T$	$\text{id} + \text{id} * \text{id} \$$	$T \rightarrow FT'$
$\$ \epsilon T F$	$\text{id} + \text{id} * \text{id} \$$	$F \rightarrow XQ$
$\$ \epsilon T Q X$	$\text{id} + \text{id} * \text{id} \$$	$X \rightarrow \text{id}$
$\$ \epsilon T Q \text{id}$	$\text{id} + \text{id} * \text{id} \$$	pop id, increase read pointer
$\$ \epsilon T Q$	$+ \text{id} * \text{id} \$$	$Q \rightarrow \epsilon$
$\$ \epsilon T $	$+ \text{id} * \text{id} \$$	$T' \rightarrow \epsilon$
$\$ \epsilon $	$+ \text{id} * \text{id} \$$	$E' \rightarrow + TE'$
$\$ \epsilon T $	$+ \text{id} * \text{id} \$$	pop +, inc read pt
$\$ \epsilon T$	$\text{id} * \text{id} \$$	$T \rightarrow FT'$
$\$ \epsilon T F$	$\text{id} * \text{id} \$$	$F \rightarrow XQ$
$\$ \epsilon T Q X$	$\text{id} * \text{id} \$$	$X \rightarrow \text{id}$
$\$ \epsilon T Q \text{id}$	$\text{id} * \text{id} \$$	pop id
$\$ \epsilon T Q$	$* \text{id} \$$	$Q \rightarrow \epsilon$
$\$ \epsilon T $	$* \text{id} \$$	$T' \rightarrow * FT'$
$\$ \epsilon T F *$	$* \text{id} \$$	pop *
$\$ \epsilon T F$	$\text{id} \$$	$F \rightarrow XQ$
$\$ \epsilon T \text{id} Q X$	$\text{id} \$$	pop id, $X \rightarrow \text{id}$

Stack	i/p	o/p
\$ E I T' Q id	id \$	pop id
\$ E I T' Q	\$	Q → ε
\$ E I T'	\$	T' → ε
\$ E I	\$	E' → ε
\$	\$	HALT Accept

One approach:

If an invalid input is processed & if no prod rule for the symbol

then skip the stack symbols ~~from~~ ^{follow} T-ε from the input till we reach the ~~first~~ (stack top) is reached

for

skip stack top

unless follow (stack top) is in the ~~exp~~ i/p under read ptr.

Stack	ip	o/p
E	\$	pop E

Synch - an entry that reports the compilation error

FE error case: id + id * id

$\$ E' T' F *$	$* id \$$	pop *
$\$ E' T' (F)$	$+ id \$$	Synch ←
$\$ E' T'$	$+ id \$$	$T' \rightarrow \epsilon$
$\$ E'$	$+ id \$$	$E' \rightarrow + T E'$
$\$ E' T +$	$+ id \$$	$E' \rightarrow + T E'$ process T
$\$ E' T$	$id \$$	$T \rightarrow F T'$
$\$ E' T' F$	$id \$$	$F \rightarrow X Q$
$\$ E' T' Q X$	$T id \$$	$X \rightarrow id$
$\$ E' T' Q id$	$T id \$$	process id

Bottom up parsing:

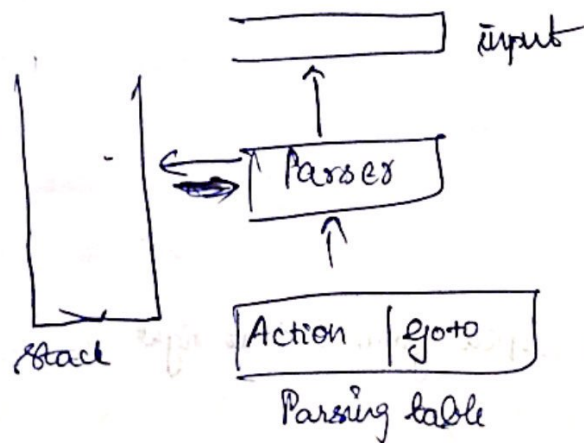
Start with the leaves, goes up one level, to the root

- Shift reduce parsing:

- Operator precedence parsing

Shift reduce parsing:

one type of bottom up parsing where 2 actions are there: shift reduce



→ Terminals / tokens

→ Non terminals

→ State of parser

Parsing table of BUP:

Decide whether which action to take place or whether to go to

Shift:

the input under read pt will be pushed inside the stack.

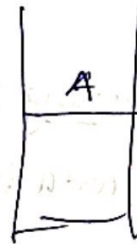
~~A~~ → ABX

Reduce: (doing rightmost derivation in reverse order)
of $A \rightarrow ABX$

and stack:



pop entire ABX, push A



Apply the prod rule on the input (p reduce it to start symbol)
- reduction

Ex:

$S \rightarrow AcBe$

$A \rightarrow Ab|b$

$B \rightarrow d$

Input: abbcd

We always scan the input from left to right

$\begin{matrix} \text{a} & \text{b} & \text{b} & \text{c} & \text{d} & \text{e} \\ & \uparrow & & & & \\ \rightarrow & \text{a} & \text{A} & \text{b} & \text{c} & \text{d} & \text{e} \end{matrix}$

if we put one more $A \rightarrow b$
 \rightarrow aAacde \times we are stuck

\rightarrow aAbcde
 \rightarrow aAcde (1 shift)
 \rightarrow aAcde (2 shift)
 \rightarrow aAcBe
 \rightarrow S

rightmost derivation:

$S \rightarrow$ aAeBc
 \rightarrow aAcde
 \rightarrow aAbcde
 \rightarrow abbcd

We need a marker to specify which part of input we can reduce

Handle:

the string which is reduced will lead one step towards the start symbol

abbcde

\downarrow a handle \rightarrow not a handle

Conditions:

\rightarrow rightmost part of α has to match the righthand side of a production

\rightarrow If α is the derivation of $A \rightarrow \beta$, should lead us one step towards the start. (if the input at rd pt is β)

Process of finding handle of reducing it
 Handle parsing

Read marker : \cdot [Introduced over the input]
 right of \cdot \rightarrow whatever not read

a b b c d e

The dot moves one step at a time and a production rule is applied after passing the symbol

Shift \cdot a b b c d e
 Shift a \cdot b b c d e
 reduce ~~Shift~~ a b b c d e
 handle

What is immediately right of \cdot
 can only be shifted

Shift reduce aA \cdot b c d e
 Reduce Shift a A b \cdot c d e
 Shift Reduce a A \cdot c d e
 Shift Shift a A c \cdot d e
 reduce Shift a A c d \cdot e
 handle
 Shift Reduce a A c B \cdot e
 reduce Shift a A c B e \cdot
 Accept ~~Reduce~~ S

Shift $\alpha \beta$
 Reduce $\alpha A \beta$

$A_i \rightarrow A_j$

$\alpha \Rightarrow \alpha A$

$\alpha \cdot b \beta$

$\alpha b \cdot \beta$

$\alpha A \cdot \beta$

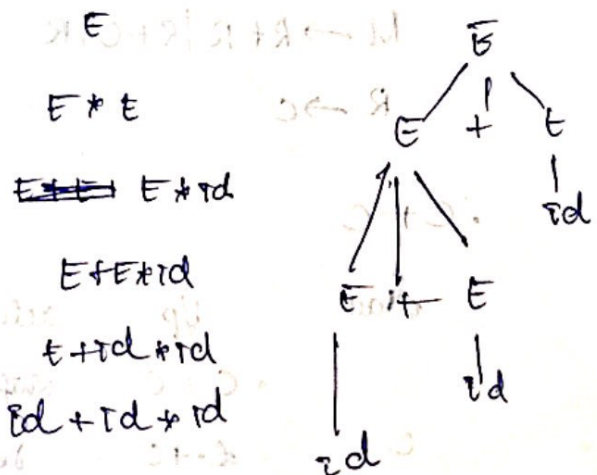
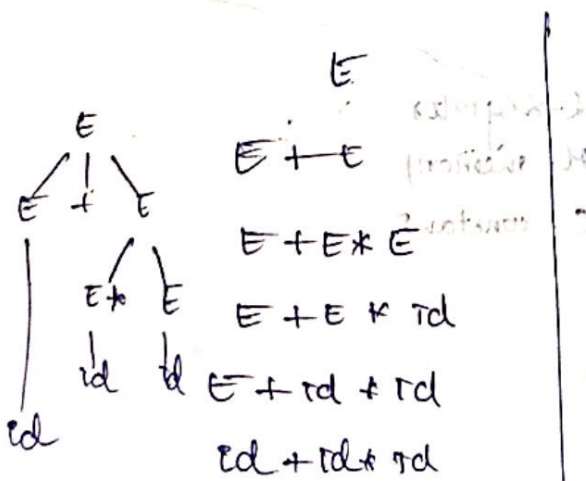
$$E \rightarrow E + E \mid E * E \mid id$$

$$id + id * id$$

Action	$\leq id + id * id$
Shift	$\cdot id + id * id$
Reduce	$id \cdot + id * id$
Shift	$E \cdot + id * id$
Shift	$E + \cdot id * id$
Reduce	$E + id \cdot * id$
Reduce	$E + E \cdot * id$
Shift	$E \cdot * id$
Shift	$E * \cdot id$
Reduce	$E * id \cdot$
Reduce	$E * E \cdot$
Reduce	$E \cdot$
Accept	

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

a situation where we can reduce $E + E$ to E or to shift more of context $id \rightarrow E$. This situation always passes in amb grammar



Another derivation from ambiguous handle

Shift $E+E \cdot Id$

Shift $E+T \cdot * \cdot Id$

Reduce $E+E * Id \cdot$

reduce $E+E * E \cdot$

reduce $E+E \cdot$

Accept $E \cdot$

LR(K) grammar - bottom up parsing

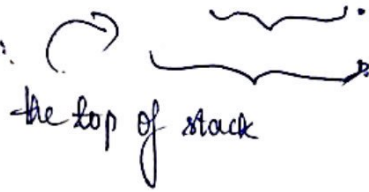
L → left to right

more powerful & expressive

R → rightmost derivation

K → no. of lookahead symbols

A choice of applying production to a set of handles



we can either reduce the small set or the whole from the top of stack

Reduced

$M \rightarrow R+R \mid R+C \mid R$

$R \rightarrow C$

R - register

M - memory

C - constant

$\cdot C + C$

Stack

Up

action

$\cdot C + C$

Shift

C

$\cdot + C$

reduce $R \rightarrow C$

R

$\cdot C$

Shift

R+

$\cdot C$

Shift

R+C
M

reduce $M \rightarrow R+C$
accept

Another method

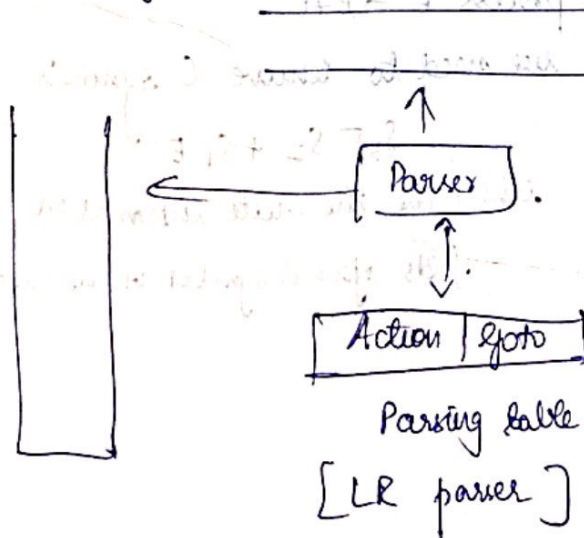
Stack	action / p	action
.	$\cdot C+C$	Reduce shift
.C	$\cdot +C$	Reduce
R	$\cdot +C$	Shift
R+	$\cdot C$	Shift
R+C		Reduce $R \rightarrow C$
R+R		Reduce $R+ \rightarrow R+R$
M		accept

This is called reduced-reduced grammar

where ~~a part of~~ 2 handles and one handle is subset of other

Which prod rule to select?

have to change the grammar to remove ambiguity



$S_0 X_1 S_1 X_2$

$S_{n-1} X_n S_n$

$S \rightarrow$ state symbols

\rightarrow states of the parser

Depending on state, decides when the handle has exhausted

is on top of stack

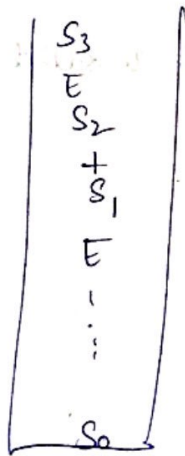
location of read pt.

the part where the \cdot has crossed \rightarrow state

Show the information of whatever has been passed and whatever now in the top of stack is on handle.

S_0 start state - signifies \cdot is at beginning

S_n final state - S .



no. of symbols to be removed = twice the

no. of symbols at the top of stack
right of input

S - gives the state information

\therefore to process $E \rightarrow E + E$

we need to remove 6 symbols

$S_3 E S_2 + S_1 E$

which are the state information

which after 6 symbols it needs to be popped

20/10/24
Parsing table

State	Action - terminals						Goto - NT		
	id	+	*	()	\$	E	T	F
0	Sr			S4			1	2	3
1		S6				acc			
2		r2	S7		r2	r2			
3		r4	r4		r4	r4			
4	S5			S4			8	2	3
5		r6	r6		r6	r6			
6	S5			S4				9	3
7	S5			S4					10
8		S6							
9		r1	S7		r1	r1			
10		r3	r3		r3	r3			
11		r5	r5		r5	r5			

- $E \rightarrow E+T$ — ①
 $E \rightarrow T$ — ②
 $T \rightarrow T * F$ — ③
 $T \rightarrow F$ — ④
 $F \rightarrow id$ — ⑤
 $F \rightarrow (E)$ — ⑥
- state
 1, 2, 3 — symbol
 S — shift
 r — ~~reduce~~ reduce
 Whenever reducing goto
 goto part.

id + id * id

Stack	i/p	Action
0	id + id * id \$	Shift 5
0id5	+ id * id \$	Reduce $F \rightarrow id$
0F3	+ id * id \$	Reduce $T \rightarrow F$
0T2	+ id * id \$	Reduce $E \rightarrow T$
0E1	+ id * id \$	Reduce Shift 6

Input	Stack	Action	State
0E1+6	Id * id \$	Shift 5	0
0E1+6id5	* id \$	Reduce F → id	0
E+T * 0E+6 F 3 3	* id \$	Reduce T → F	1
E+T * 0E+6 T 9	* id \$	Shift 7	2
E+T * 0E+6 T 9 * 7	id \$	Shift 5	2
E+T * 0E+6 T 9 * 7 id 5	\$	Reduce F → id	3
0E+6 T 9 * 7 F 10	\$	Reduce T → T * F	3
0E 6 T 9 + 6 T 9	\$	Reduce E → E * T = 6 * 9	3
0E1	\$	accept	3

- $E \rightarrow E + T$
 $\rightarrow E + T * F$
 $\rightarrow E + T * id$
 $\rightarrow E + F + id$
 $\rightarrow E + id * id$
 $\rightarrow T + id * id$
 $\rightarrow F + id * id$
 $\rightarrow id + id * id$

How to build up the state $0 \rightarrow 11$

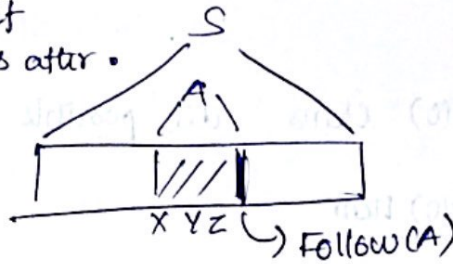
Viable prefix: set of valid symbols for the grammar
 For α to be valid prefix:
 $\alpha.w$: - sentential form
 $\alpha.w$: a configuration of the SR parser
 w is a string of terminals

~~LR(0)~~
 LR(0) Item - simple LR(SLR) parser

$A \rightarrow XYZ$
 $A \rightarrow \cdot XYZ$

$A \rightarrow X \cdot YZ$ - has processed Y, now has to move to process Z
 $A \rightarrow XY \cdot Z$
 $A \rightarrow XYZ \cdot$ Reduction if Follow(A) is after.

You cannot reduce before Follow(A) is after.



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

↓ Augmented Grammar

$G' = (N \cup S', T, S', P \cup S' \rightarrow S)$

Add $S' \rightarrow E$ to the prev grammar

Closure of LR(0) item

$A \rightarrow \alpha \cdot BB$

$B \rightarrow xyz$

$B \rightarrow \cdot xyz$

Ex:

$S' \rightarrow E$

2 possible LR(0) item

$S' \rightarrow \cdot E$
 $S' \rightarrow E \cdot$

This whole derivation is as closure of 1st LR(0) item

2 derivation inside E

$E \rightarrow \cdot E + T$
 $E \rightarrow E \cdot T$

$T \rightarrow \cdot T * F$
 $T \rightarrow T \cdot F$

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

alpha

$S \rightarrow E$
 $E \rightarrow E+T$
 $E \rightarrow T$
 $T \rightarrow T * F$
 $T \rightarrow F$
 $F \rightarrow (E)$
 $F \rightarrow id$

Goto(E, X)

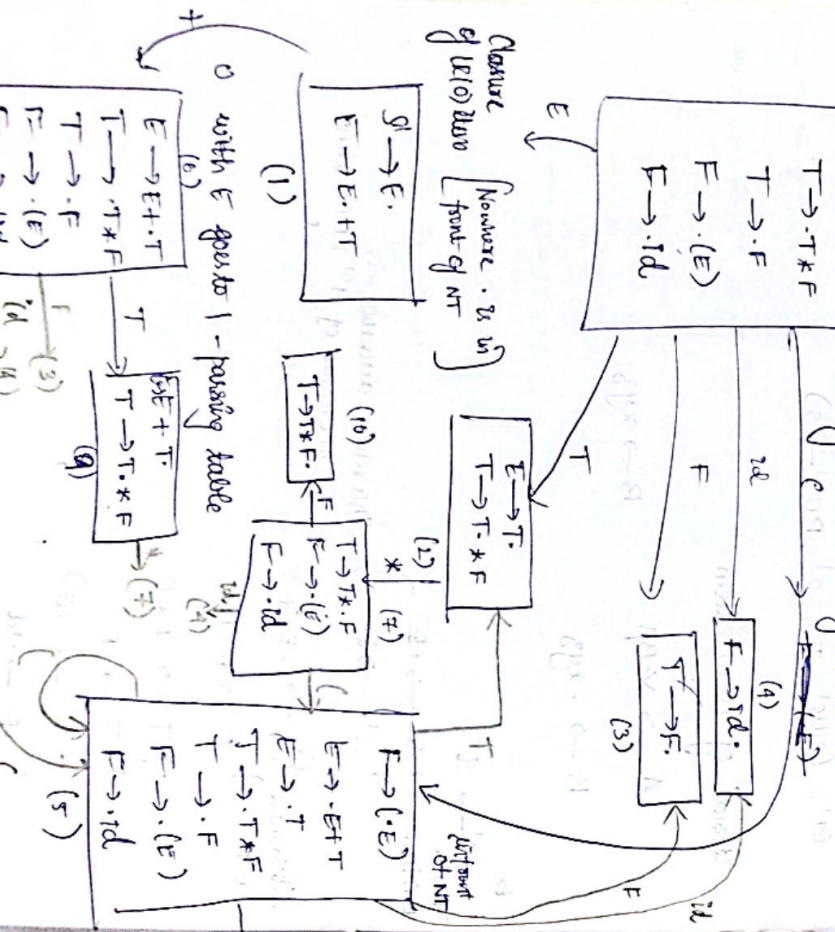
$I: A \rightarrow \cdot X \cdot B \xrightarrow{\text{Goto}(X, X)} A \rightarrow X \cdot B$
 will have to move red pt one step ahead

Set of LR(0) items: all possible LR(0) that can be derived from the above expressions

First LR(0) item

(0) $S \rightarrow \cdot E$
 $E \rightarrow \cdot E+T$
 $E \rightarrow \cdot T$
 $T \rightarrow \cdot T * F$
 $T \rightarrow \cdot F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot id$

Applying goto of each prod rules items and get other LR(0)



closure of LR(0) item

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

(10) $T \rightarrow T * F \cdot$
 $T \rightarrow T * F \cdot$

(7) $F \rightarrow (\cdot E)$
 $F \rightarrow (\cdot E)$

with E goes to 1 - passing kable

(6) $E \rightarrow E + \cdot T$
 $E \rightarrow E + \cdot T$
 $T \rightarrow T * \cdot F$
 $T \rightarrow T * \cdot F$
 $F \rightarrow (\cdot E)$
 $F \rightarrow (\cdot E)$

SLE passing kable

Goto(I, X)

$X = a, a$ is terminal

when is at end of prod rule if we found this follow set of NT

$A \rightarrow XYZ$
 $b \in \text{follow}(A)$

$M[T, b] \leftarrow$ reduce by $A \rightarrow XYZ$
 means we have seen the end of T until

if a set

$A \rightarrow XYZ$
 $B \rightarrow CD$

cannot reduce by SLE power

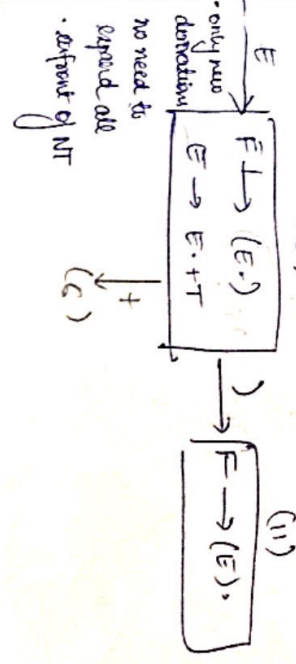
or $A \rightarrow XYZ$
 $B \rightarrow CD$

if it sees a terminal

$M[S_I, a] \leftarrow$ shift S_I

$M[S_I, X] \leftarrow$ goto

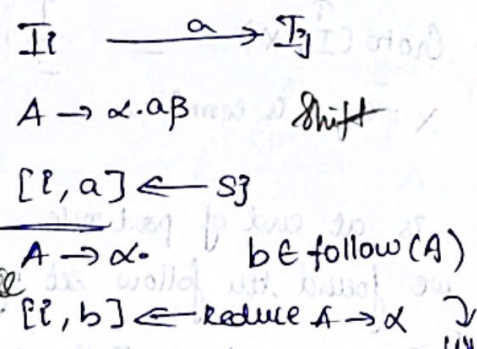
if it sees a terminal



LR parsing table

- $S' \rightarrow E$
- $E \rightarrow E + T \text{ --- (1)}$
- $E \rightarrow T \text{ --- (2)}$
- $T \rightarrow T * F \text{ --- (3)}$
- $T \rightarrow F \text{ --- (4)}$
- $F \rightarrow (E) \text{ --- (5)}$
- $F \rightarrow id \text{ --- (6)}$

all outgoing arrows through terminals are written with shift



Whenever we see a \cdot at the end of a prod rule in any closure we can reduce

State	+	*	(id)	\$	E	T	F
0			S5	S4			1	2	3
1	S6	S7				Accept			
2	r2	S7			r2	r2			
3	r4	r4			r4	r4			
4	r6	r6			r6	r6			
5			S5	S4			8	2	3
6			S5	S4				9	3
7			S5	S4					10
8	S6				S11				
9	r1	S7			r1	r1			
10	r3	r3			r3	r3			
11	r5	r5			r5	r5			

For reducing,

$$\text{Follow}(E) = \{ \$, \rangle, + \}$$

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

If follow of a symbol contains \$,

$$\begin{cases} E \rightarrow T + F & (1) \\ T \rightarrow T * F * F & (2) \\ F \rightarrow (E) \rightarrow (T) \rightarrow (F) \end{cases}$$

$$\begin{cases} E \rightarrow E + T \\ \rightarrow T + T \\ \rightarrow T * F + T \\ F \rightarrow (E) \\ F \rightarrow (T) \\ E \rightarrow T \$ \end{cases}$$

$$\text{Follow}(F) = \{ +, *, \rangle, \$ \}$$

$$\text{Follow}(s') = \{ \$, \rangle \} \rightarrow s' \text{ cannot come any where else}$$

$s' \rightarrow E$. - whenever closure has this make it accept over \$
Accept Follow(s') contains \$

If there is a . at the end of prod rule in any closure
ex $A \rightarrow \alpha \cdot$

write the prod rule no. (ex 1) in cell ~~closure~~ follow set of A

$$I_i \xrightarrow{B} I_j$$

$$A \rightarrow \alpha \cdot B \beta$$

$$[i, B] \leftarrow$$

From the state,

Write the closure number it redirects to on the N.T

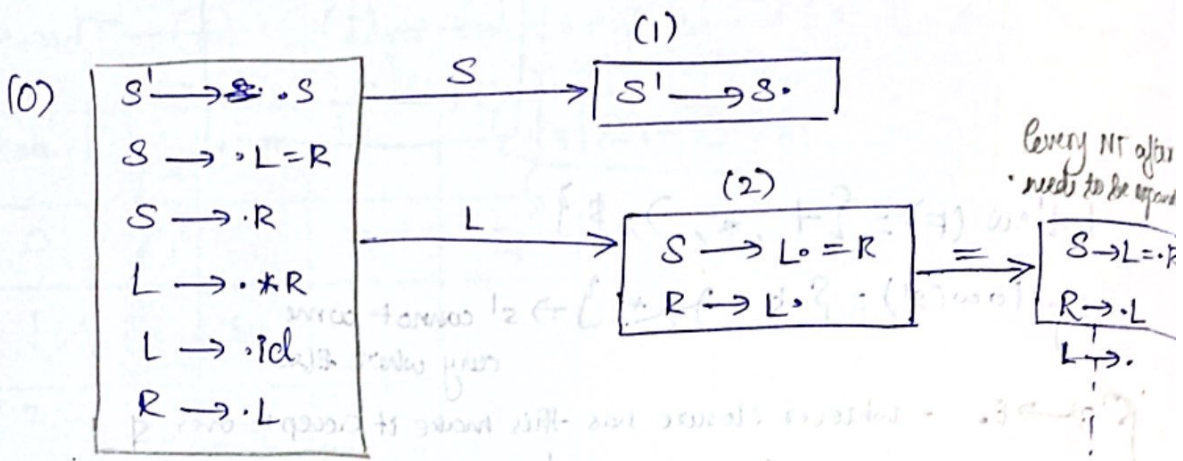
Non terminals

goto part

to p0

Are all unambiguous grammar LR grammar? *(proof is)*

$S \rightarrow S$
 $S \rightarrow L = R$ — (1)
 $S \rightarrow R$ — (2)
 $L \rightarrow *R$ — (3)
 $L \rightarrow id$ — 4
 $R \rightarrow L$ — (5)



	=	
0		$S \rightarrow L = R$
2	S, R	$S \rightarrow *R = R$

$follow(R) = \{ \$, \epsilon \}$
 $follow(R)$ contains =
 (add ϵ to state 2 on =)

cannot be valid in LR

From the parsing table, an example derivation:

- 0 $id = id \$$ Shift S
- 1 $id S = id \$$ Reduce $L \rightarrow id$
- 2 $L 2 = id \$$ Reduce $R \rightarrow L$
- 3 $R 3 = id \$$ No entry for $\$$ over =

Canonical LR parser:

It's LR(1) grammar
 Takes the lookahead whether to reduce or shift

LR(1) item

$[A \rightarrow \alpha \cdot \beta, a]$

Closure

$[A \rightarrow \alpha \cdot \underline{B} \beta, a] \in I^{(1)}$

$B \rightarrow \gamma$, for each terminal $b \in \text{first}(\beta a)$

$[B \rightarrow \gamma, b]$

goto

$I \xrightarrow{B} J$

$[A \rightarrow \alpha \cdot B \beta, a]$

$[A \rightarrow \alpha B \cdot \beta, a]$

↳ the only diff is adding a

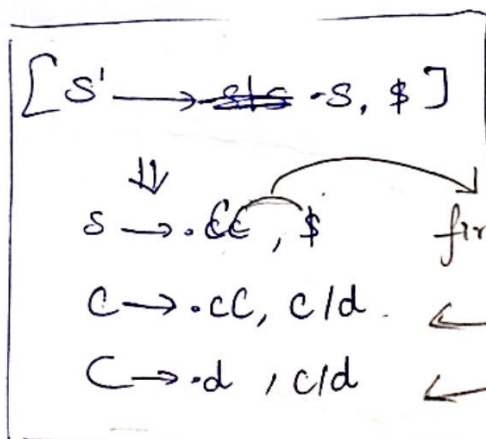
Ex:

$S' \rightarrow S$

$S \rightarrow CC$

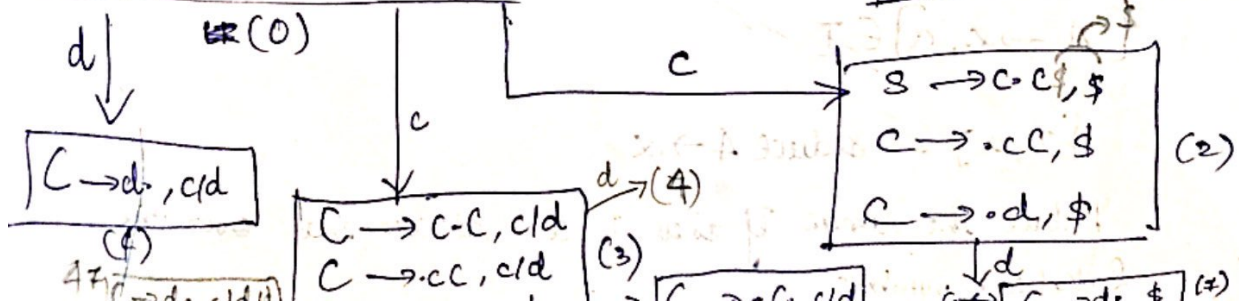
$C \rightarrow cC \mid d$

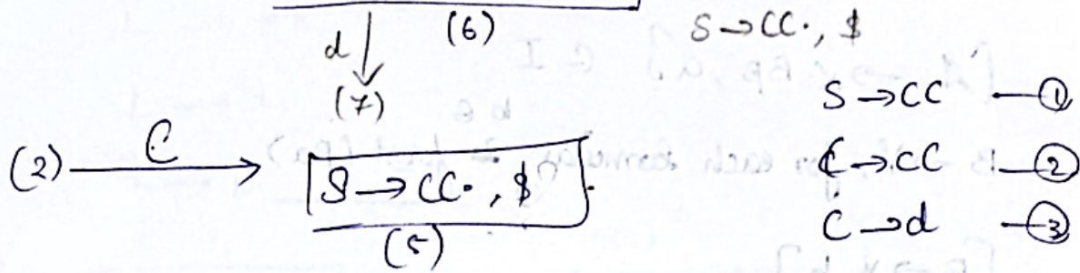
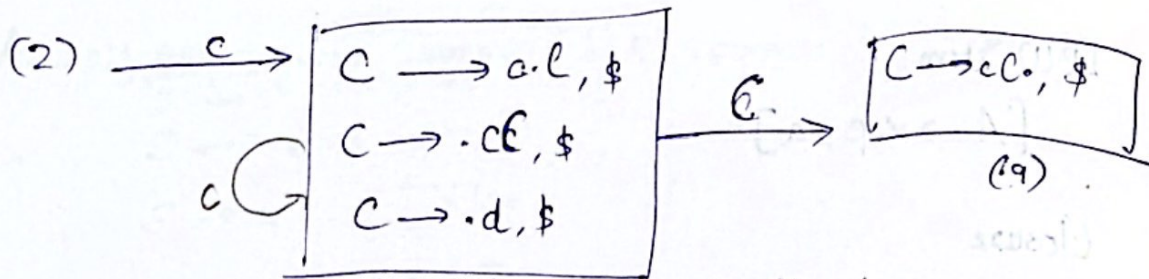
the grammar is $c^* d c^* d$, atleast 2 ds



first LR(1) item, with \$ lookahead symbol

$\text{first}(C\$) = c/d$





State	c	d	\$	S	C
0	S ₃	S ₄		1	2
1			acc		
2	S ₆	S ₇			5
3	S ₃	S ₄			8
4	r ₃	r ₃			
5			r ₁		
6	S ₆	S ₇			9
7			r ₃		
8	r ₂	r ₂			
9			r ₅		

Reduce

$$[A \rightarrow \alpha, a] \in I$$

$$[I, a] \leftarrow \text{reduce } A \rightarrow \alpha$$

Reduce the closure if there is a \cdot at the end over the lookahead symbols

6/3/24

SEMANTIC ANALYSER

Major roles

→ Type checking

Ex: int a, string b;

a = b + 2 → semantic error

↓
no syntax error

→ variable used before declaration

Ex: a = b + c

↑
undefined

→ No. of arguments in function

→ return type

→ array arguments → can be 1D but then assumed as 2D

→ Accessing element / member of array incorrectly

→ Invalid operations Ex: string - string

→ Overloading of operators for code gen phase, we need to distinguish the operator distinctly

→ Inheritance

→ Scope checking

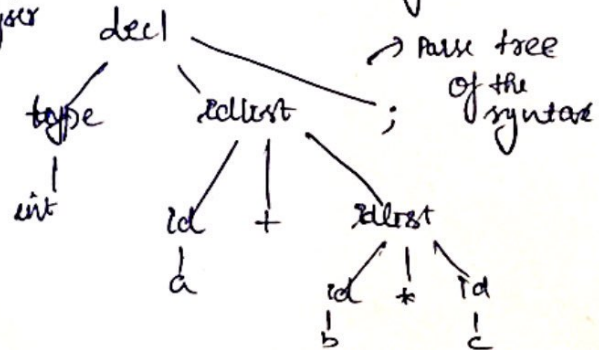
should not have multiple definition of the same variable within the same scope of function

Ex: Addition
Concatenation
sign
depending on the context, meaning changes.

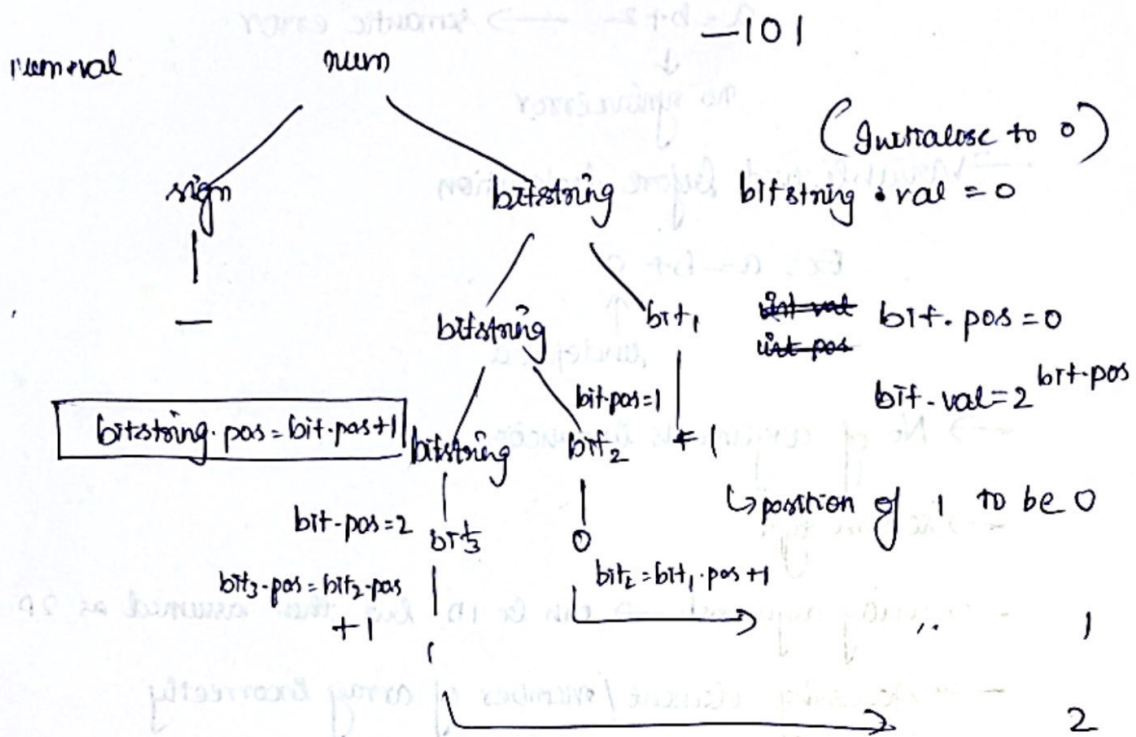
→ Uniqueness of variables / function

for int a, b, c; → output of lex analyser

then:
decl → type idlist ;
idlist → id, idlist | id
type → int | float | char



num \rightarrow sign bitstring
 sign \rightarrow + / -
 bitstring \rightarrow ~~bit~~ bitstring }
 { bitstring } { bit } | bit
 bit \rightarrow 0 | 1



need position of value to convert into decimal by the semantic analyzer.

Attributes here are position of value

bitstring.pos = 0

We are traversing the parse tree again to find the location of variables.

19/3/24

of inherited attribute,

$(\text{got}) \text{ val} \rightarrow \text{xy}$

then we need attribute ready of parents to evaluate attribute of children $(\text{L} \rightarrow \text{got} \text{ val}, \text{L} \text{ got} \text{ val})$

If both inherited & synthesized attribute,

a top down parser is not enough to parse all the inherited attribute.

$L \rightarrow E_n$

$\text{got} = \text{got}$

$E \rightarrow E + T$

$(\text{got}) \text{ val} \text{ val} \text{ val}$

$E \rightarrow T$

$(\text{got}) \text{ val} + (\text{got}) \text{ val} + (\text{got}) \text{ val} \text{ T} + \text{E} \rightarrow$

$T \rightarrow T * F$

$(\text{got}) \text{ val} * (\text{got}) \text{ val} \text{ T} * \text{F}$

$T \rightarrow F$

$F \rightarrow (E)$

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

$F \rightarrow \text{id}$

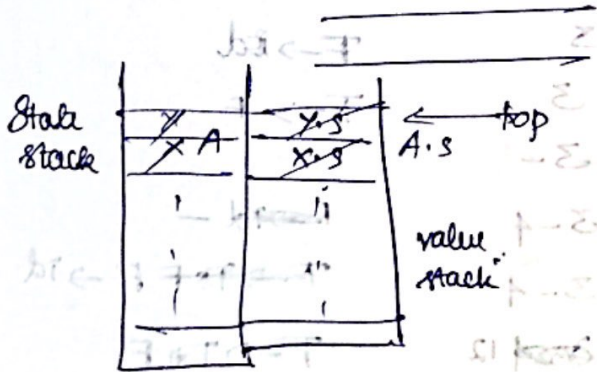
$F \text{ val} = \text{id} \text{ lexval}$

everything root is dependant on children - synthesized attribute

A grammar with all synthesized attribute

\downarrow
S-~~case~~ attribute grammar

Visible prefix: all the seen symbols are there on stack



we need another stack to store var of synthesized it

$A \cdot s = f(x \cdot s, y \cdot s)$

$A \rightarrow x y$

$x \cdot s = a \cdot v$

$y \cdot s = b \cdot v$

$$Y.S \leftarrow \text{value}[\text{top}]$$

$$X.S \leftarrow \text{value}[\text{top}-1]$$

$$n\text{Top} \leftarrow f(\text{value}[\text{top}], \text{value}[\text{top}-1])$$

updated

value of top after reduction $\rightarrow n\text{Top} = \text{top} - \delta + 1$

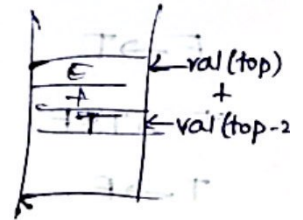
$$\delta = |\text{RHS of production applied for reduction}|$$

$$\text{Top} = n\text{Top}$$

$$L \rightarrow \epsilon \quad \text{print}(\text{val}(\text{top}))$$

$$E \rightarrow E+T \quad \text{val}(n\text{Top}) = \text{val}(\text{top}) + \text{val}(\text{top}-2)$$

$$E \rightarrow T \quad \left[\text{Automatically } T\text{-val is assigned to } E\text{-val} \right]$$



$$T \rightarrow T * F \quad \text{val}(n\text{Top}) = \text{val}(\text{top}) * \text{val}(\text{top}-2)$$

$$T \rightarrow F$$

$$F \rightarrow (E) \quad \text{val}(n\text{Top}) = \text{val}(\text{top}-1)$$

$$F \rightarrow id$$

Value stack helps in getting the new attribute value from old " "

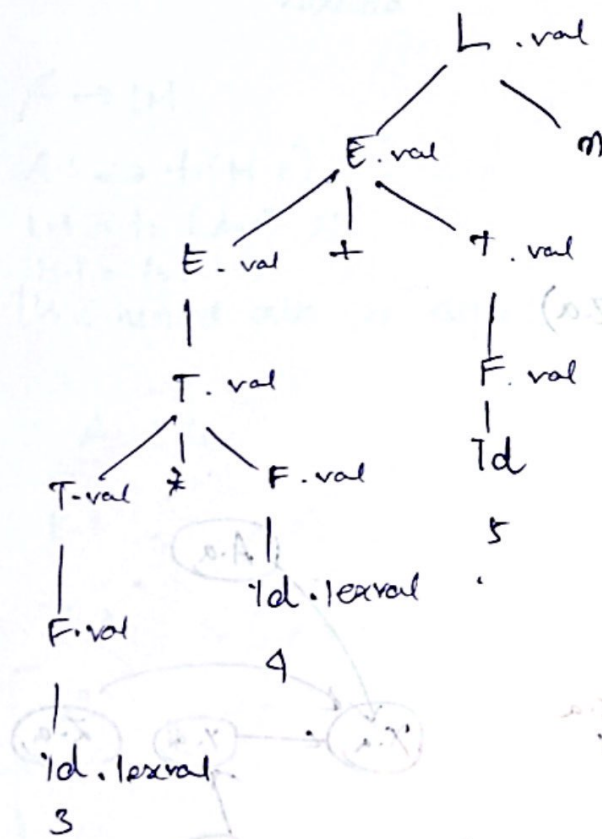
i/p	Stack	Value	Prod
8*4+5n	-	8	
*4+5n	id	3	F → id
*4+5n	F	3	F → id
*4+5n	T	3	T → F
4+5n	T	3-	
+5n	T*id	3-4	F → id
+5n	T*F	3-4	F → id F → id
+5n	T	12	T → T*F
+5n	E	12	E → T

57)	E+	12 -	┌
57)	E+Td	12 - 5	┌
77)	E+F	12 - 5	F → Td
77)	E+T	12 - 5	T → F
77)	E	17	E → E+T
77)	En	17	
	L	17	L → En

Accepting stage,
point the value [top]

↓
print

not using action
" " " "
know of parse tree
know of parse tree
know of parse tree



For all prod rules
B → Y
Val for all variables
→ E.node = newnode(+, E1.node, T.node)
→ F.node = T.node
→ T.node = newnode(*, T1.node, F.node)
→ T.node = F.node
→ F.node = E.node
→ F.node = newnode(7d)

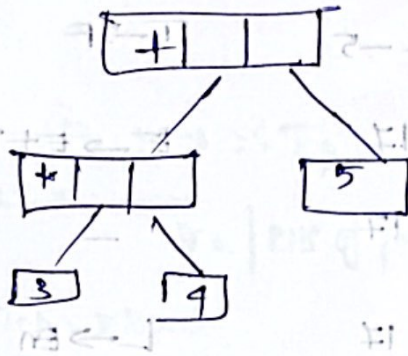
Annotated parse tree - when there are attributes along with the variables

one version

of parse tree

Syntax tree

3 * 4 + 5



Annotated parse trees
 Cop " "

doesn't carry forward
 after semantic analysis
 Only syntax tree will
 carry forward

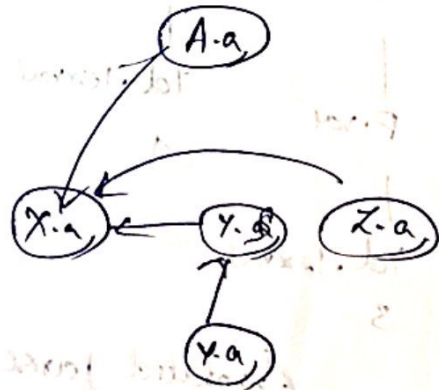
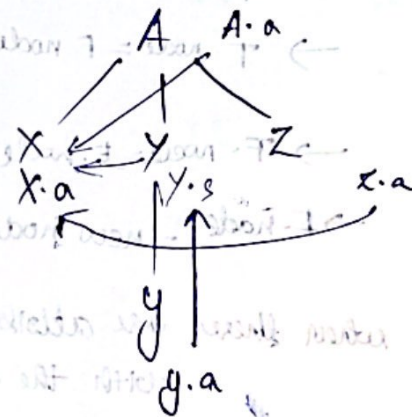
$Y \rightarrow y$

$A \rightarrow XYZ$

$Y.s \Rightarrow f_1(y.a)$

$X.a = f_2(A.a, Y.s, Z.a)$

Dependency Graph



A graph with a node having zero indegree, can start with that
 N_1, N_2, \dots, N_k - nodes
 $\rightarrow N_j$ should be dependant on ~~it~~ any nodes from N_1 to N_{j-1}
 but not on $N_i, i > j$

→ N_1 is zero indegree node

	inh	syn
Top	? ^{L-attr grammar} ✓	??
Bottom	??	✓

→ Depth First search over parse tree

L-attributed grammar

$$A \rightarrow X_1 X_2 \dots X_n$$

$$X_j \cdot a = f(A_i, X_1 \cdot a, X_2 \cdot a, \dots, X_{j-1} \cdot a)$$

↑
inherited

⏟
any attribute

$$A \rightarrow LM$$

$$A \cdot s \Rightarrow f_1(M \cdot s)$$

$$L \cdot l = f_2(A \cdot s) \quad \times$$

$$M \cdot l = f_3(L \cdot s)$$

L's inherited attr can depend only on its parents inh attribute

$$A \rightarrow QR$$

$$R \cdot r = f_4(A \cdot r) \quad \checkmark$$

$$Q \cdot r = f_5(A \cdot r, R \cdot s) \quad \times$$

$$A \cdot s = f_6(Q \cdot s)$$

00/3/14

AGF

- Attribute Grammar

- Translation rules -

$$L \rightarrow \bar{E}n$$

$$L \cdot val = E \cdot val \leftarrow \text{Attribute Grammar}$$

$$\{ \text{print}(E \cdot val) \} \leftarrow \text{translation rule}$$

$$F \rightarrow ?d$$

$$\text{addentry}(?d, \text{id-type})$$

- Syntax directed definition (SDD) - balance b/w A.G. & T-Rule

- " " translation ^{scheme} (SDT) - embedding the action at a particular location where it needs to be executed

Semantic rule $A \rightarrow E \cdot T \rightarrow E \cdot val$

but instead of val , we use $print(E \cdot val)$

↓
This is called side effect.

SDT - to have translation rule with controlled side effect.

S-attributed SDD \Rightarrow bottom up parser

\Rightarrow Post order traversal of parse tree

$A \rightarrow X_1 X_2 X_3 \rightarrow actions$

$AS = f(X_1.S, X_2.S, X_3.S)$

SDT

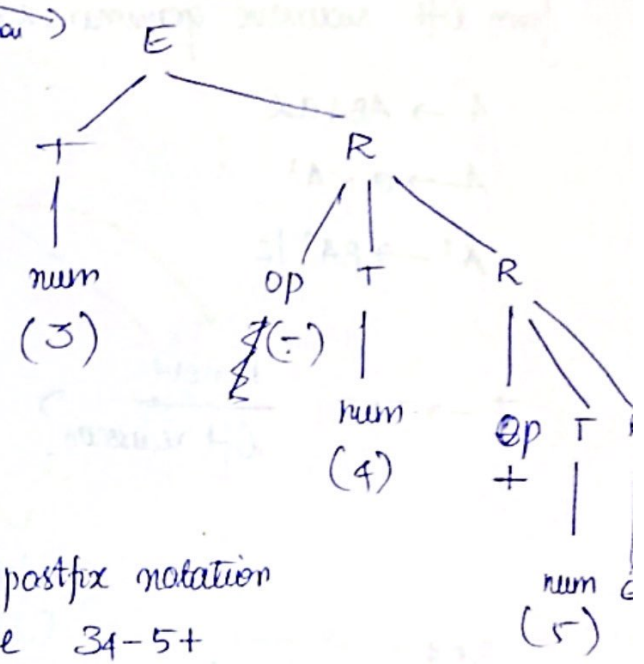
The action is put at the rightmost end of the production rule in a synthesized attribute grammar.

SDT

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

$E \rightarrow TR$
 $R \rightarrow opTR | \epsilon$
 $T \rightarrow num$
 $op \rightarrow + / -$

parse tree for this grammar

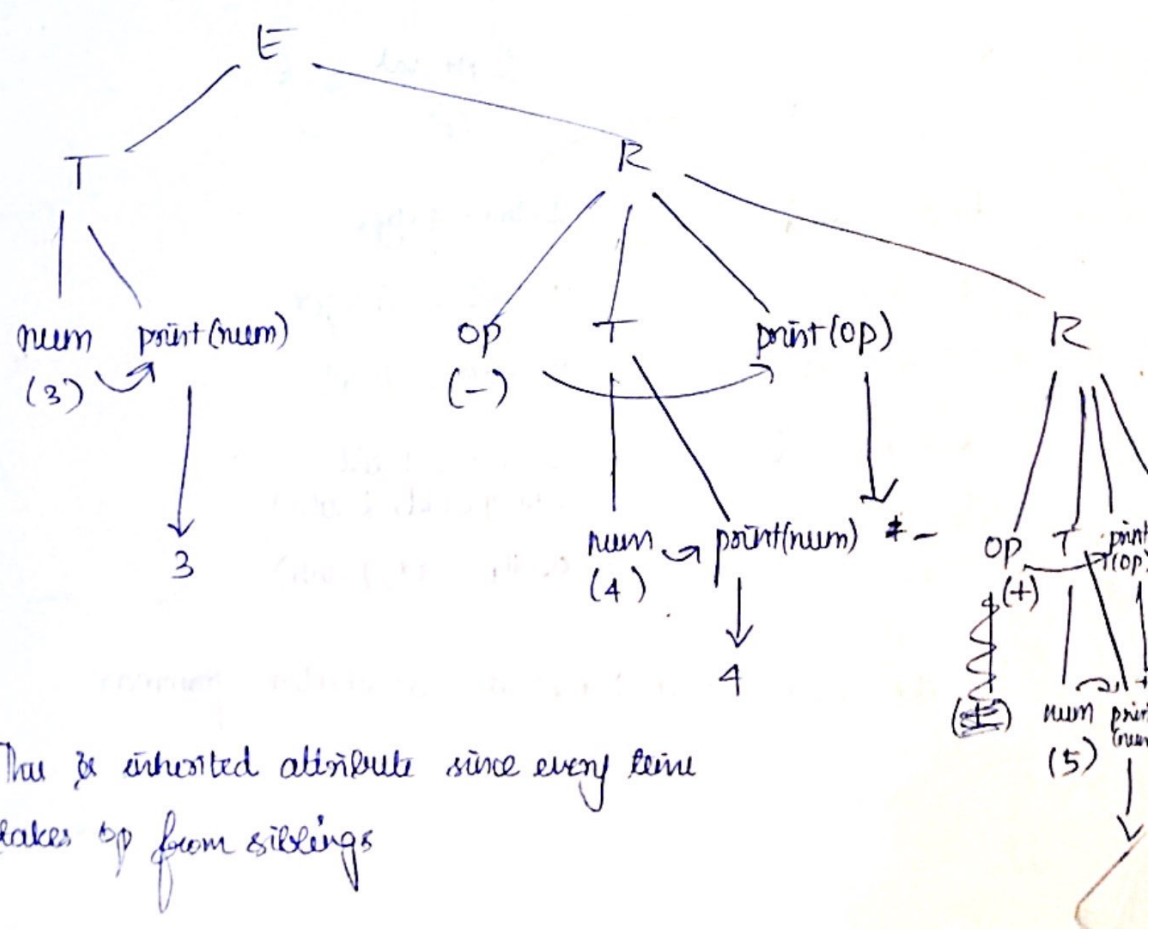


To get the full expression,
 $3-4+5$

Add the action s.t it reads the postfix notation
 i.e $3-5+$

$E \rightarrow TR$
 $R \rightarrow opT \{print(op)\} R | \epsilon$
 $T \rightarrow num \{print(num)\}$
 $op \rightarrow + / -$

\therefore after every T inverted to num in R, we need to print op symbol



This is inherited attribute since every time it takes op from siblings

from left recursive grammar, we cannot construct top down parser

$$A \rightarrow AB | aA$$

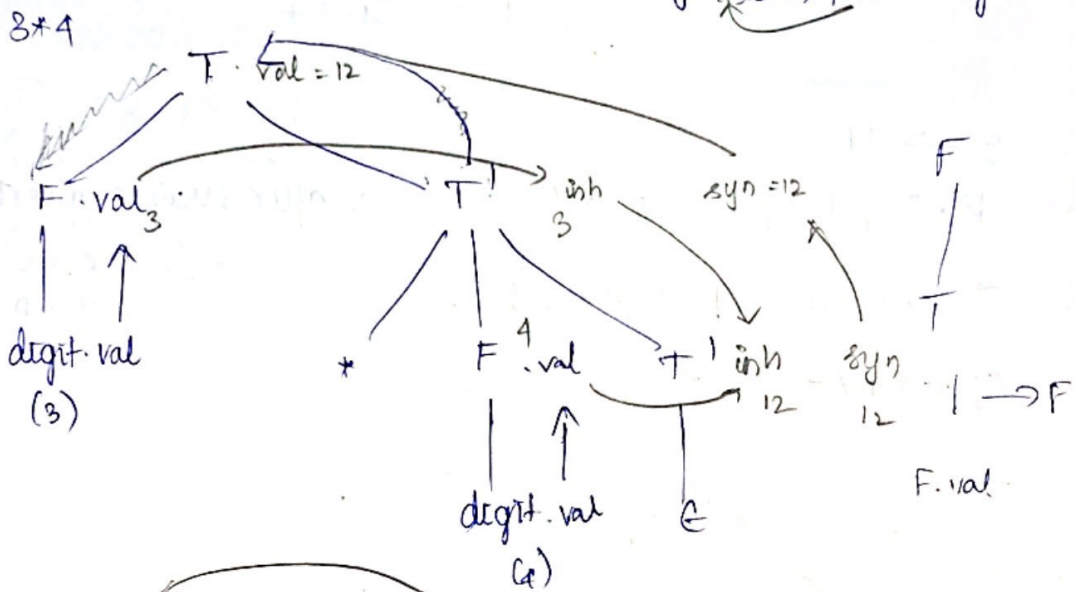
$$A \rightarrow a \times A'$$

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

$$T \rightarrow T * F \xrightarrow[\text{left recursion}]{\text{Remove}}$$

sem rule

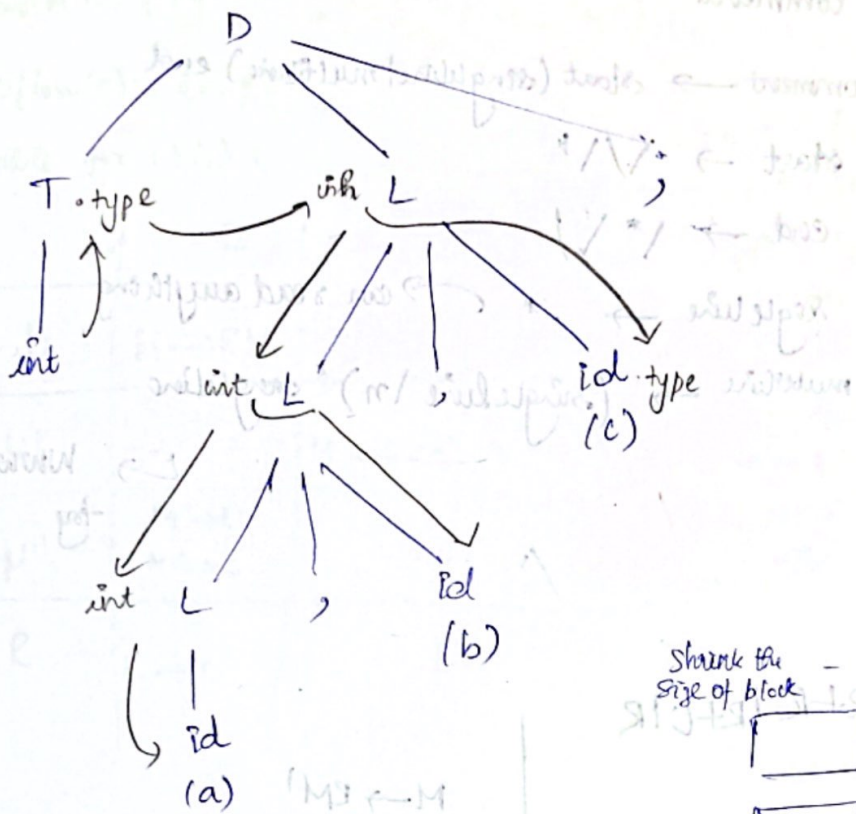
$$\begin{aligned}
 T &\rightarrow FT \{ \} \rightarrow T'.inh = F.val \\
 &\quad T.val = T'.syn \\
 T' &\rightarrow + FT \{ \} \rightarrow T'.inh = T'.inh + F.val \\
 &\quad T'.syn = T'.syn \\
 T' &\rightarrow \epsilon \{ \} \rightarrow T'.syn = T'.inh \\
 T &\rightarrow digit \{ \} \rightarrow T.val = digit.lexval
 \end{aligned}$$



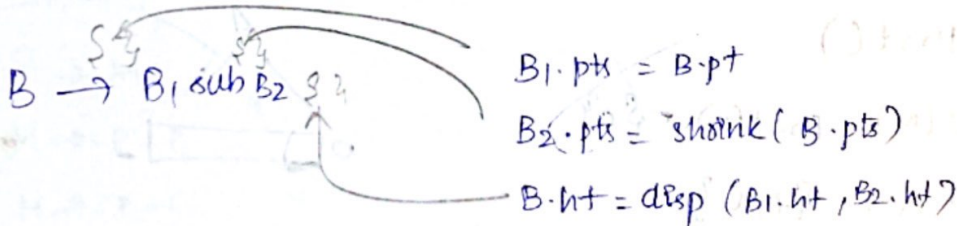
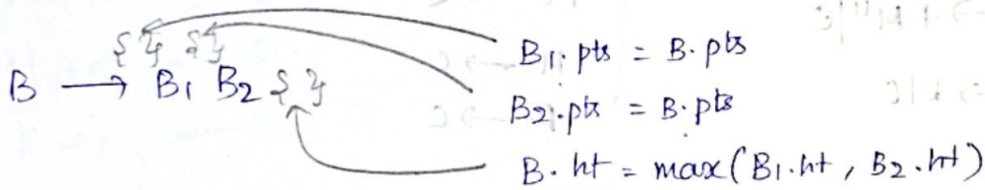
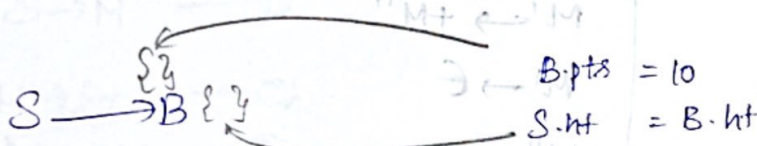
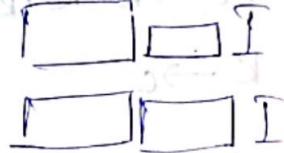
$$\begin{aligned}
 D &\rightarrow TL; && L.inh = T.type \\
 T &\rightarrow \text{int } \{ \} && T.type = \text{integer} \\
 T &\rightarrow \text{float } \{ \} && T.type = \text{float} \\
 L &\rightarrow \{ \} id && L.inh = L.inh \\
 &&& addtype(id, L.inh) \\
 L &\rightarrow id && addtype(id, L.inh)
 \end{aligned}$$

Put the action just before the attributed grammar

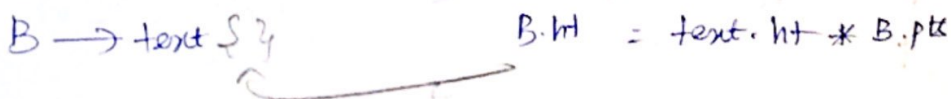
alt a, b, c;



shrink the size of block - Subscript



↓ displacement ! calculating subscript
block size along with main block size



C9

RE for comments

comment \rightarrow start (singulare/multiline) end

start $\rightarrow \backslash \backslash ^*$

end $\rightarrow \backslash ^* \backslash$

singulare $\rightarrow .*$ \rightarrow can read anything

multiline \rightarrow (singulare $\backslash n$)^{*} singulare

\rightarrow www.regex101.com
-toy

" / + "

$M \rightarrow R+R | R+CR$
 $R \rightarrow c$
 $M \rightarrow RM'$
 $M' \rightarrow +R | +C | \epsilon$
 $M' \rightarrow +M'' | \epsilon$
 $M'' \rightarrow RC$

$M \rightarrow RM'$
 $M' \rightarrow +M''$
 $M' \rightarrow \epsilon$
 $M'' \rightarrow R$
 $M'' \rightarrow c$
 $R \rightarrow c$

First()

$First(M) = First(R) = \{c\}$

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

$First(M'') = First(R) = \{c\}$



	follow set			
	M	M'	M''	R
1	\$			
2				+
3				\$
4		\$	\$	

$$\text{Follow}(M) = \{\$, \gamma\} = \text{Follow}(M') = \text{Follow}(M'')$$

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

Parse table for LL(1):

	C	+	\$
M	$M \rightarrow RM'$		
M'		$M' \rightarrow +M''$	$M' \rightarrow c$
M''	$M'' \rightarrow R$ $M'' \rightarrow c$		
R	$R \rightarrow c$		

On the first set of the respective symbols fill the prod that made that symbol first set

If there is ϵ , then see the follow set symbol of the prod there

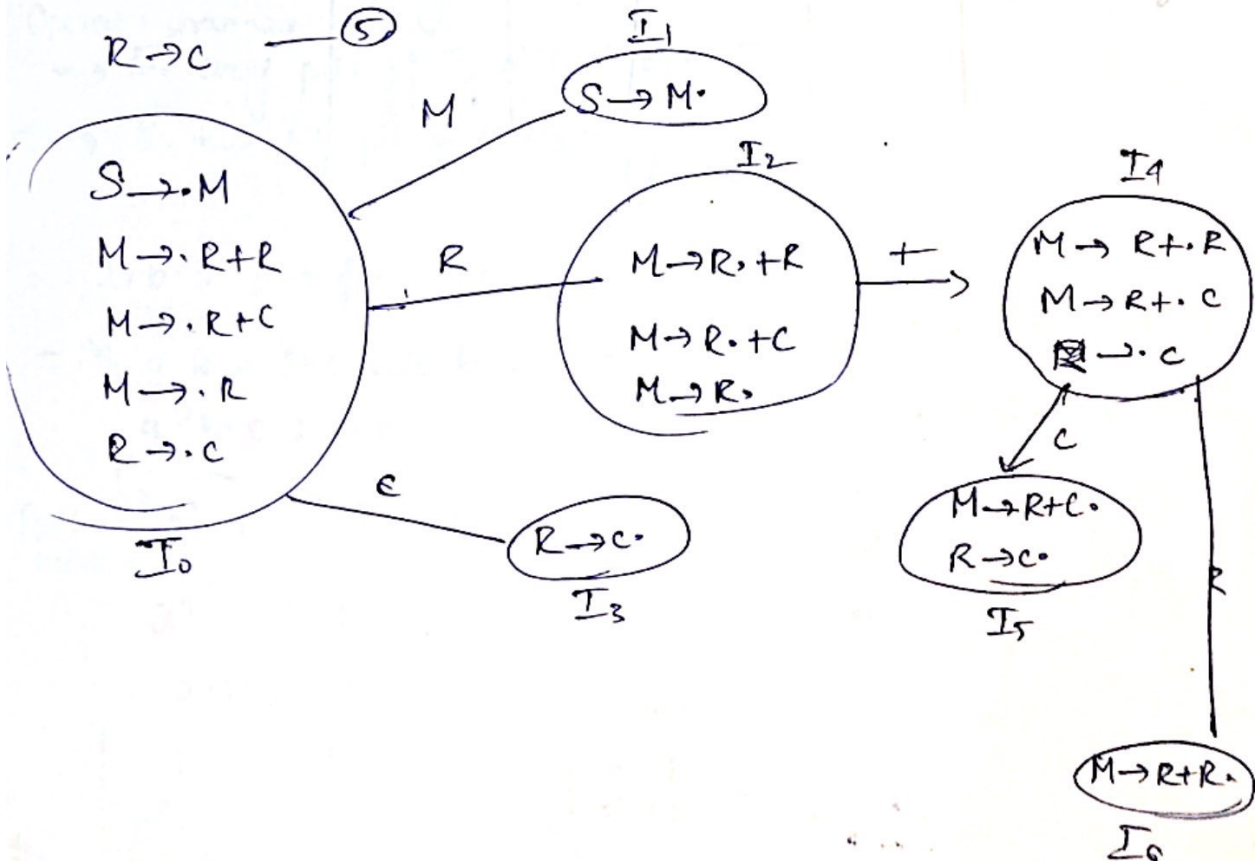
$$S \rightarrow M \quad \text{--- ①}$$

$$M \rightarrow R+R \quad \text{--- ②}$$

$$M \rightarrow R+c \quad \text{--- ③}$$

$$M \rightarrow R \quad \text{--- ④}$$

$$R \rightarrow c \quad \text{--- ⑤}$$



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

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

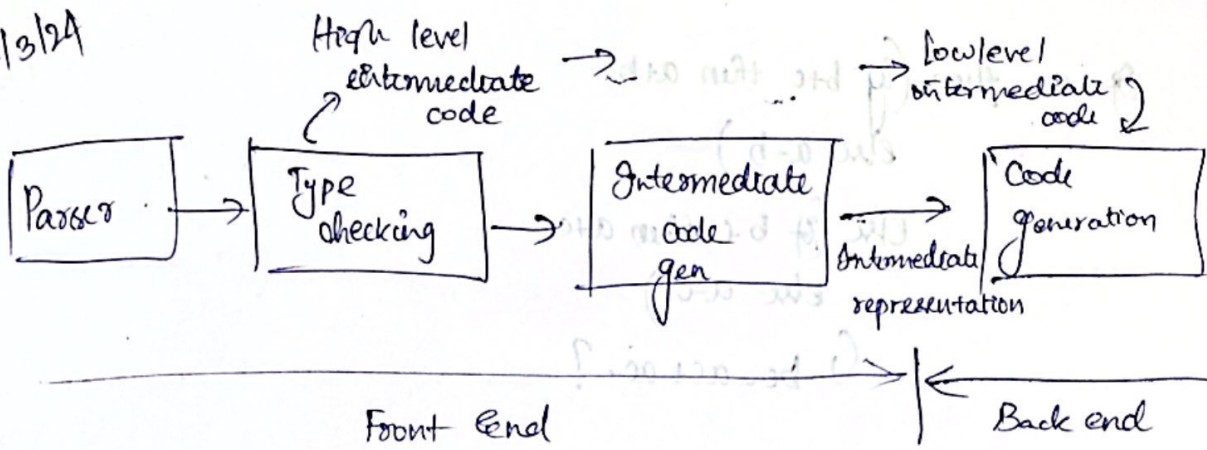
$$\text{Follow}(R) = \{ +, \$ \}$$

LR parse table

M is the start symbol $\therefore S \rightarrow M$ is there $\therefore S1$
 goto $\text{Follow}(S) = \{ \$ \}$

States	Action					
	C	+	\$	S	M	R
0	S ₃				1	2
1			Accept			
2		S ₄	r ₄			
3		r ₅	r ₅			
4	S ₅					6
5		r ₅	r ₃ r ₅			
6			r ₂			

27/3/24



- Post fix notations
- syntax tree
- Three address Code

$$A \rightarrow \alpha \beta \gamma$$

α, β, γ can be ϵ or a single non-terminal.

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

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

not an operator grammar

Operator grammar:

- For every pair of T_i , atleast one NT
- No two NT placed consequently

$$E \rightarrow ETE \mid E-E \mid E * E \mid$$

$$E \uparrow E \mid (E) \mid id$$

operator grammar

$a+b$ in post fix: $ab+$

- If a then $b+c$ else $b-c$
- If a then (if $b+c$ then $a+b$) else if $b-c$ then $a+c$

$$a ? b+c ; b-c$$

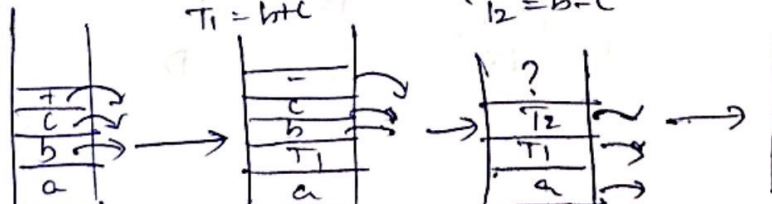
Postfix notation $abc+bc- : ?$

$$a ? b+c \quad b-c$$

$$abc+bc- ?$$

$$T_1 = b+c$$

$$T_2 = b-c$$



If a then (if b+c then a+b
 else a-b)
 else (if b-c then a+c
 else a-c)
 ↳ bc-ac+a-c-?

a b+c a+b a-b-? bc-ac+a-c-??

Quadruples

OP	ARG1	ARG2	RESULT
ADD	B	C	T1
EQ	T1	A	

(1) ADD B C
 (2) EQ (1) A
 Triplets/Triples

A = B * C + D

MUL	B	C	T1
ADD	T1	D	T2
EQ	T2	A	

(1) MUL B C
 (2) ADD (1) D
 (3) EQ (2) A

Easy implementation

Space efficient
 Changing an expression affects entirely

Indirect Triples

order If any changes needed, the order table changing would do.

key	addr
(1)	10
(2)	11
(3)	12
(4)	13

	OP	ARG1	ARG2
10	ADD	B	C
11	EQ	10	A
12	SUB	B	C
13	EQ	12	D

1/19/24

Jumping
 unconditional conditional

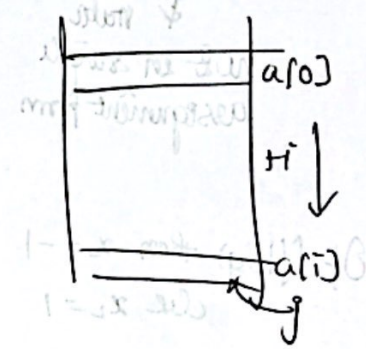
Reverse engineering notes

In an if/else block, these exist common code

call stack
 param arg 1
 " " " 2
 ...
 call (<function>, n)

<function> (arg 1, arg 2, ..., arg n)

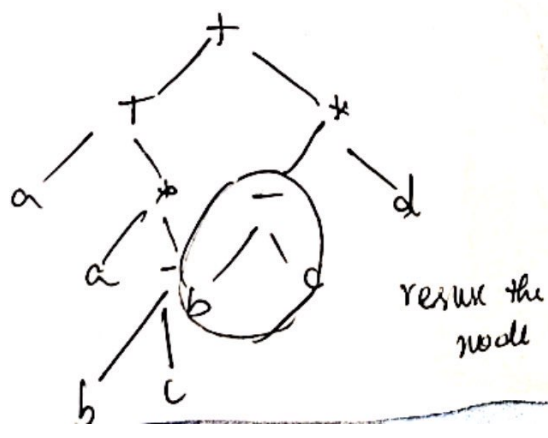
a[i] = j
 j = a[i]



TAC for

if $x > y$ then $x = y + 10$
 else $y = x + 10$
 if $x > y$ goto L1 → conditional
 t1 = $x + 10$
 y = t1
 goto L2 → uncond
 L1 t1 = $y + 10$
 x = t1
 L2 ...

⇒ $a + a * (b - c) + (b - c) * d$



Table

1.	id	a	
2.	id	b	
3.	id	c	
4.	-	2	3
5.	*	1	4
6.	+	1	5
7.	id	d	
8.	+	4	7
9.	+	6	8

Static single assignment form

$$p = a + b$$

$$q = p - c$$

$$p = q + d$$

$$p = e - p$$

$$q = p + q$$

↓ note
not in single
assignment form

$$p_1 = a + b$$

$$q_1 = p_1 - c$$

$$p_2 = q_1 + d$$

$$p_3 = e - p_2$$

$$q_2 = p_3 + q_1$$

↓
SSA → can't reuse same variable
without defining

If (flag) then $x = -1$
else $x = 1$
 $y = x + a$

∴ we can't predict the val of x in
 $y = x + a$, we take $x_3 = f(x_1, x_2)$
which outputs the correct value

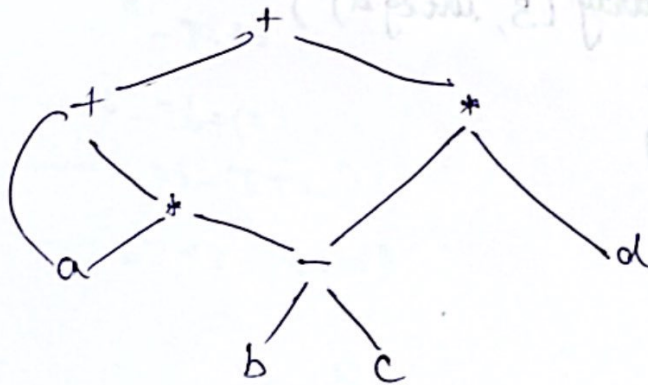


10/4/24

DAG:

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

Start →



Three address codes:

1. Assignment : $x = y \text{ op } z$, op is arithmetic or logical operation
2. " : $x = \text{op } y$, op is unary operation
3. Copy instructions : $x = y$
4. Unconditional jump → goto L
5. Conditional jumps
 - a. If x goto L
 - b. If false x goto L
 - c. If x relop y goto L
6. Procedure calls & returns:

param x1
 param x2
 ⋮
 param xm

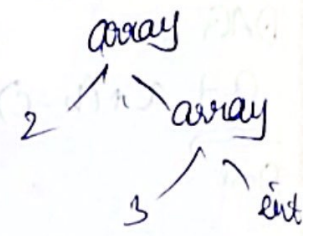
7. Indexed copy instructions
8. Address of pointer assignments
 - $x = \&y$ sets the r-value of x to be location of y
 - $x = *y$ y is pointing to a value which is a address
 - $*x = y$ sets r-value of the object pointed to by x to the r-value of y

⇒ SSA :

⇒ Type expressions

array = int[2][3]

array(2, array(3, integer))



$x = a + ix$
 $x = y[i]$

Ex: int[2][3]

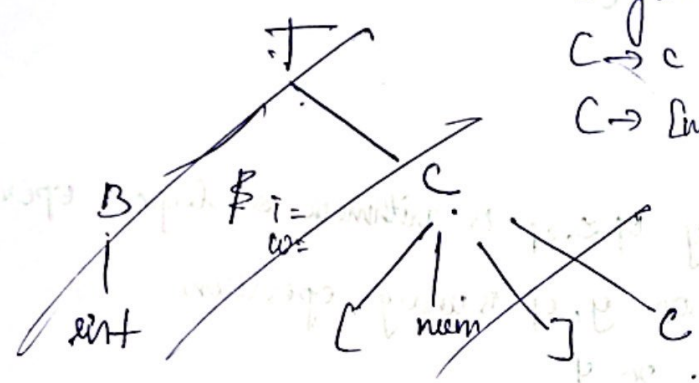
$T \rightarrow B \{ t = B \text{ type}, w = B \text{ width} \}$
 B C

$B \rightarrow \text{int} \{ B \text{ type} = \text{int}, B \text{ width} = 4 \}$

$B \rightarrow \text{float} \{ B \text{ type} = \text{float}, B \text{ width} = 8 \}$

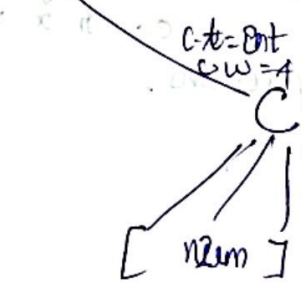
$C \rightarrow c \{ c \text{ type} = t, c \text{ width} = w \}$

$C \rightarrow [\text{num}] c_1 \{ \text{array}(\text{num-value}, c_1 \text{ type}), c_1 \text{ width} = \text{num-value} \times c_1 \text{ width} \}$

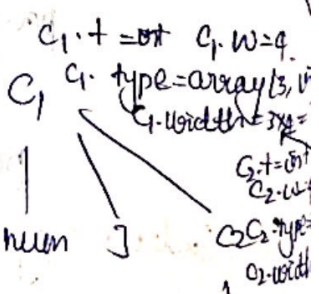


$t = \text{int}$
 $w = 4$
 $T \text{ type} = \text{array}(2, \text{array}(3, \text{int}))$
 $T \text{ width} = 24$

$B \text{ type} = \text{int}$
 $B \text{ width} = 4$



$t \text{ type} = \text{array}(2, \text{array}(3, \text{int}))$
 $c \text{ width} = 2 \times 2 \times 4 = 24$



Seq of declarations

computing the relative addr of declared name

$D \rightarrow \text{offset} = 0$

$D \rightarrow T \text{ id}; E \text{ top} = \text{put id (lexeme)}$

16/11/24

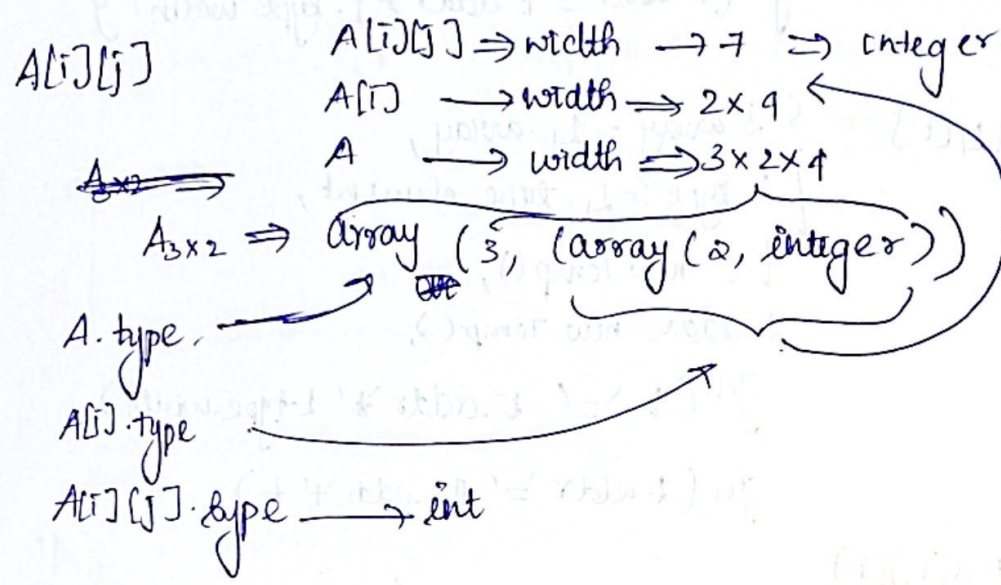
$$A[i][j] = \text{base} + i \times w_1 + j \times w_2$$

using no. of elements

K-dimension \Rightarrow $\boxed{\text{base} = i_1 \times w_1 + i_2 \times w_2 + \dots + i_k \times w_k}$

$\text{base} + (i \times m_2 + j) \times w \Rightarrow$

$\text{base} + ((i_1 \times m_2 + i_2) \times m_3 + i_3) \dots \times m_k + i_k \times w_k$



$L \rightarrow L[E] \mid \text{id}[E]$

$E \rightarrow E + E \mid \text{id} \mid (E)$

$S \rightarrow \text{id} = E, \{ \text{gen}(\text{top.get}(\text{id.lexeme}) = 'E.addr'); \}$

Creating S addr code

$L = E \{ \text{gen}(l.array.base[l.addr] = 'E.addr'); \}$

base addr

$E \rightarrow E_1 + E_2 \{ \text{gen}(E.addr = \text{new Temp}(), \text{gen}(E_1.addr '+' E_2.addr)); \}$

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

L $\{ E.addr = new Temp();$
 $gen(E.addr := L.array.base['L.addr']);$

$L \rightarrow id[E]$

$\{ L.array = top.get(id.lexeme); \}$

$L.type = L.array.type.element,$

$L.addr = new Temp();$

$gen(L.addr := E.addr * L.type.width) \}$

$L_1[E]$

$\{ L.array = L_1.array,$

$L.type = L_1.type.element,$

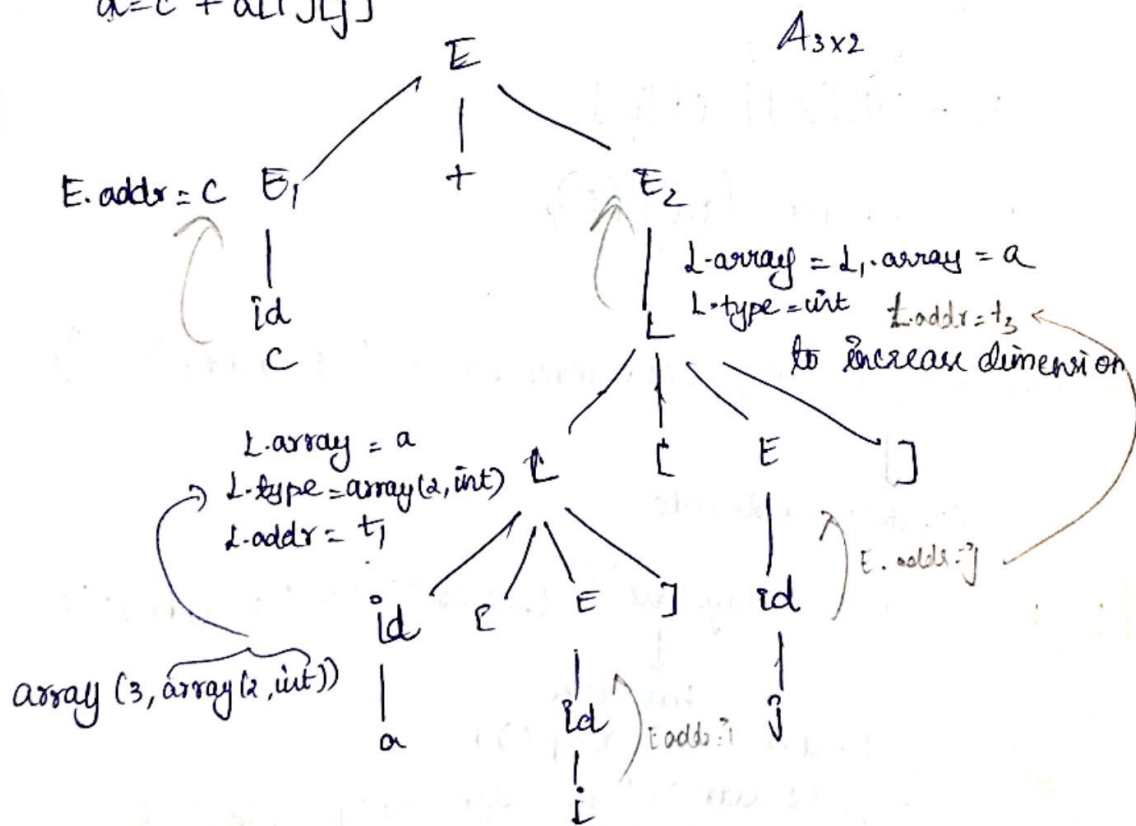
$t = new Temp(),$

$L.addr = new Temp(),$

$gen(t := E.addr * L.type.width),$

$gen(L.addr := t_1.addr + t);$

$d = c + a[i][j]$



$$1. t_1 = i * (2 * 4)$$

↑
type.width

$$2. t_2 = j * 4$$

↑
E.addr L-type.width (int)

$$3. t_3 = t_1 + t_2$$

↑
L.addr

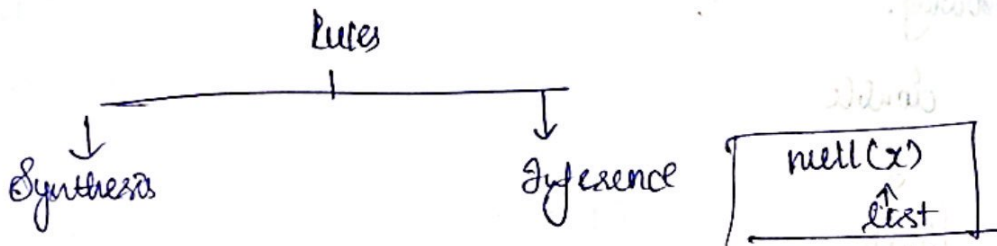
$$4. t_4 = a[t_3]$$

↑
L.array.base ↑
L.addr

17/11/24

Type Checking

- sound type system - compile time error free
- Strongly typed - language that guarantees that there won't be any type error



'ML - doesn't require var to be declared before using -
depending on context decides the type
inference type

If f has type $s \rightarrow t$ and x has type s then
expression $f(x)$ has type t

$$a = b + c$$

If $f(x)$ is an expression

then for some α and β , f has type $\alpha \rightarrow \beta$ as x has type α

finding length (x):

if null(x) then 0
else length

Type widening - default type conversion in most prog lang

double

↑

float

↑

long

↑

int

char ← byte

short

not using any info,
giving more space than reqd.

Type narrowing:

double

↓

float

↓

long

↓

int

→ short → char → byte

Casting: explicit
programmer has to define

Coercion: implicit

$$E \rightarrow E_1 + E_2 \quad \left. \begin{array}{l} E\text{-type} = \max(E_1\text{-type}, E_2\text{-type}), \\ a_1 = \text{wider}(E_1.\text{addr}, E_1\text{-type}, E\text{-type}), \\ a_2 = \text{wider}(E_2.\text{addr}, E_2\text{-type}, E\text{-type}), \\ E.\text{addr} = \text{new temp}(), \\ \text{gen}(E.\text{addr} = 'a_1' + a_2), \end{array} \right\}$$

Control flow

- Boolean variable
- Relational operator
- if
- if then else
- while loop

$$B \rightarrow B | B | B \& B | \neg B | (B) | E \text{ rel } E | \text{true} | \text{false}$$

④
③
②
①

rel : — allowed along with boolean operations

<, <=

>, >=

=, !=

more precedence
↑

if (x < 100 || x > 200 && x != y) then x = 0;

↓

if x < 100 goto L2 (some level)

if ~~not~~ x > 200 goto L1

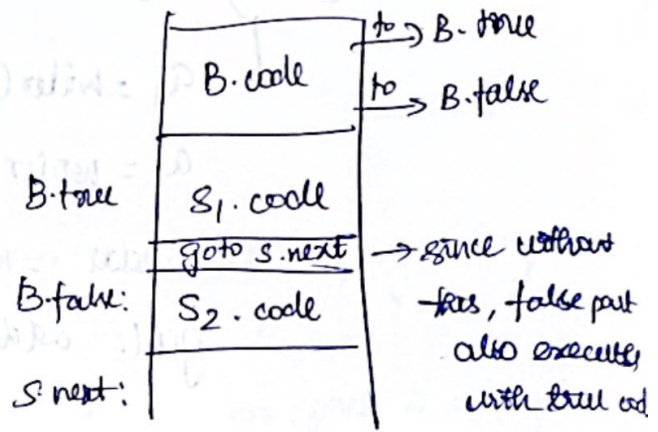
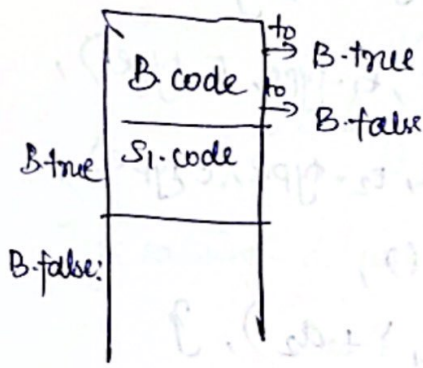
if ~~False~~ x != y goto L1

L2: x = 0

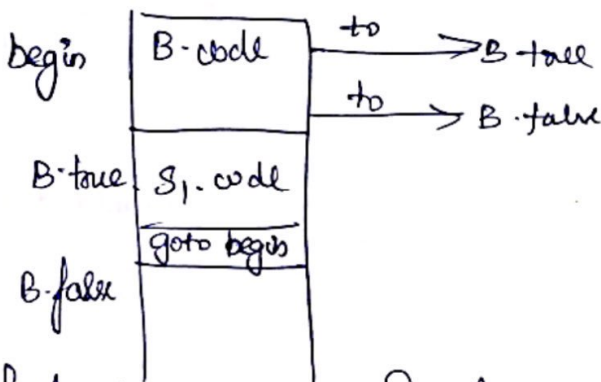
L1

short correct code

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



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



Proced rules

$P \rightarrow S$

Semantic rules

$S.\text{next} = \text{newlabel}() \Sigma,$

$P.\text{code} = S.\text{code} \parallel \text{label}(S.\text{next})$

$S \rightarrow \text{if } (B) \text{ then } S_1$

$B.\text{true} = \text{newlabel}()$

$B.\text{false} = S_1.\text{next} = S.\text{next}$

$S.\text{code} = B.\text{code} \parallel \text{label}(B.\text{true}) \parallel S_1.\text{code}$

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

$B.\text{true} = \text{newlabel}()$

$B.\text{false} = \text{newlabel}()$

$S_1.\text{next} = S_2.\text{next} = S.\text{next}$

$S.\text{code} = B.\text{code} \parallel \text{label}(B.\text{true})$

$\parallel S_1.\text{code} \parallel \text{gen}(\text{goto } S.\text{next})$

$\parallel \text{label}(B.\text{false}) \parallel S_2.\text{code}$

$S \rightarrow \text{while } (B) S_1$ $\text{begin newlabel}()$
 $B.\text{true} = \text{newlabel}()$
 $B.\text{false} = S.\text{next}$
 $S_1.\text{next} = \text{begin}$
 $S.\text{code} = \text{label}(\text{begin}) \parallel B.\text{code} \parallel$
 $\text{label}(B.\text{true}) \parallel S_1.\text{code} \parallel$
 $\text{gen}(\text{goto begin})$

For comparison operators:

Prod rules { Semantic rules (Precedence)

$B \rightarrow B_1 \parallel B_2 \text{ with } B_1:\text{true}$

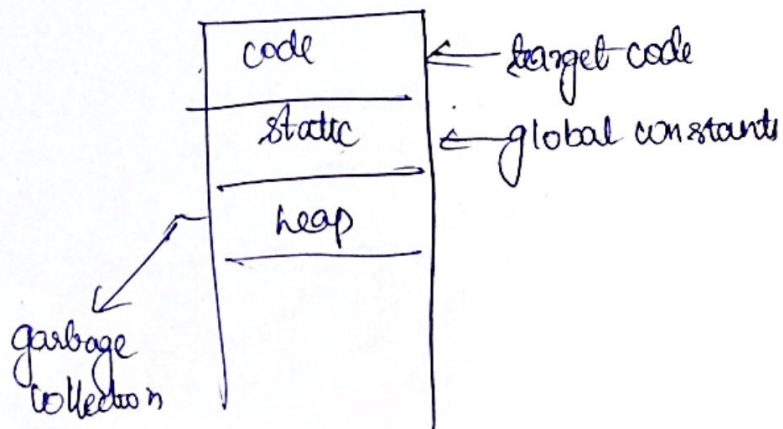
24/4/24

Run time environment

Code optimization



code generation

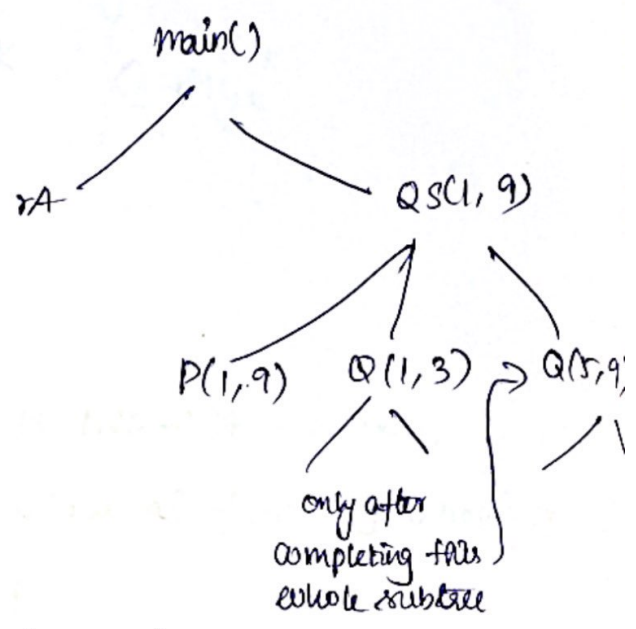


23/9/21
afternoon-extra

Activation - executing procedure

- main
- readArray
- Qsort(1, 9)
- partition(1, 9)
- QS(1, 3)
- ⋮
- QS(5, 9)

Activation tree



pre-order traversal: the seq of procedure call which find out seq of procedure are executed

post-order "

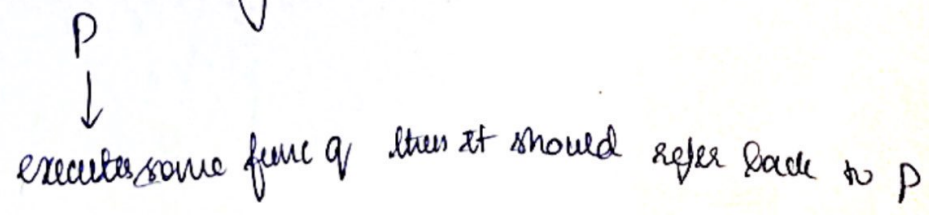
- normal termination
- termination con of execution
 - if p is able to handle the exception returned by q
 - notable → then p also terminates

→ if q is aborted ↔ p also aborted

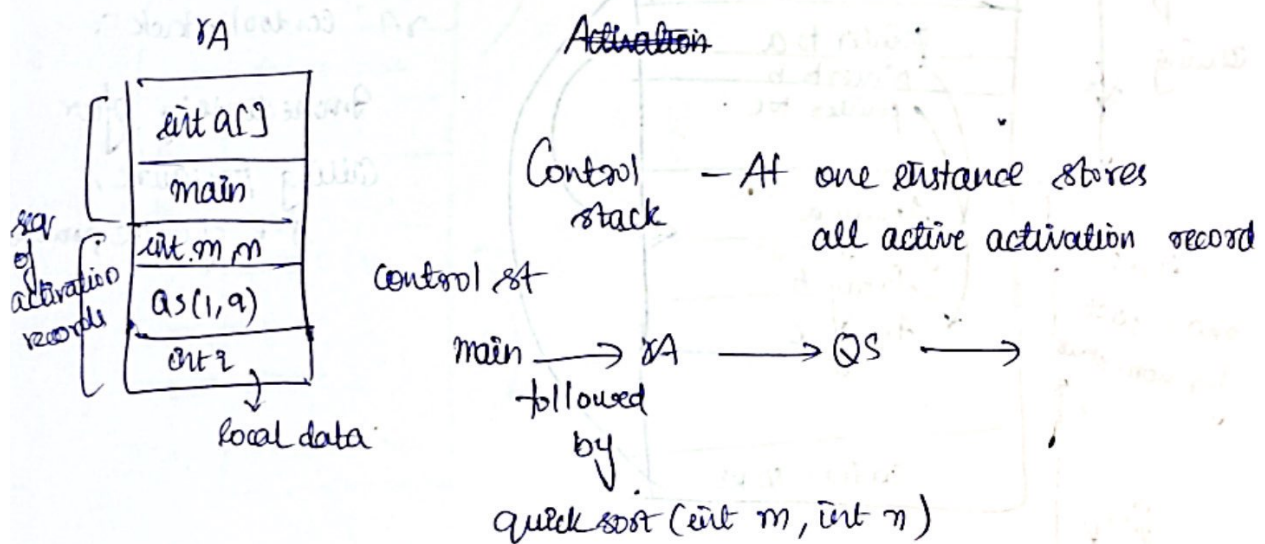
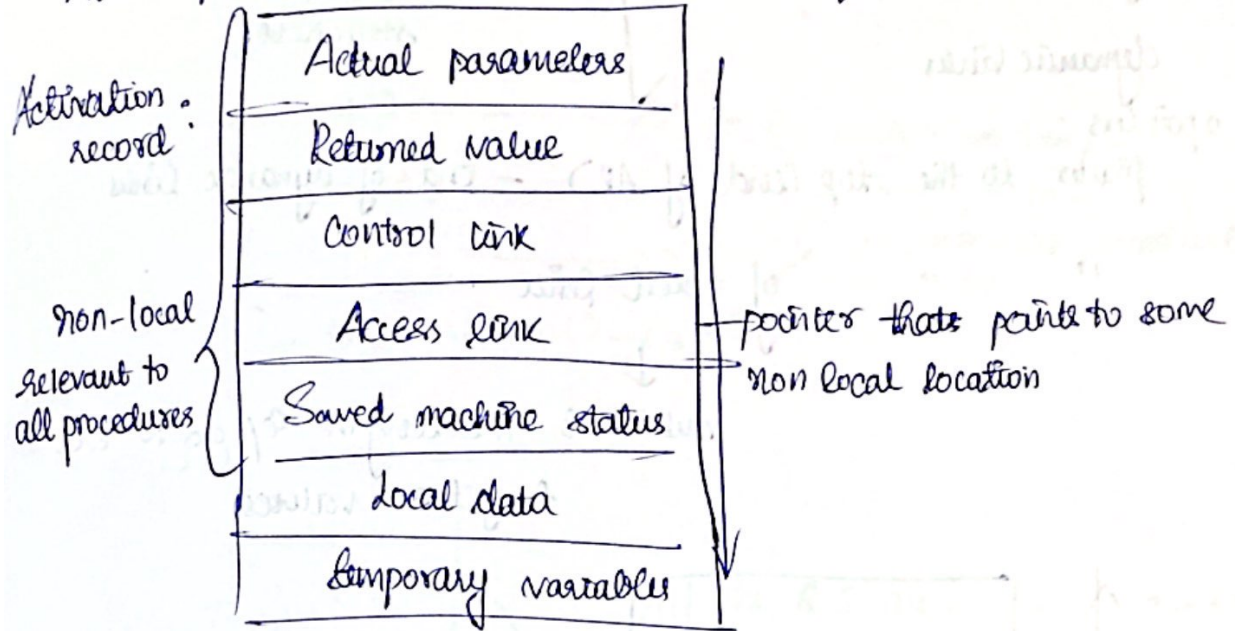
A process being terminated then the main also terminates

Activation record:

A DS that stores lot of info which are relevant to executing p



Actual parameters :- actual values that are passed



partition activation record comes after RA

Activation record - store the communicating parameters

↓
top of activation record

↓
One to be shared between callee procedure and calling procedure

callee q → P calls q
calling P

parameters to be shared

fixed links

dynamic links

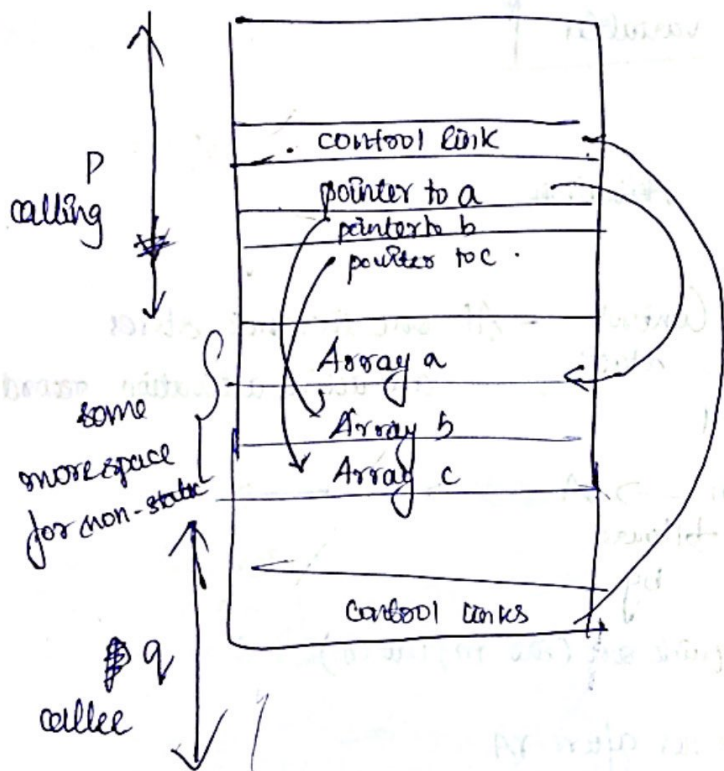
2 pointers:

pointer to the top (end of AR) - end of dynamic links

" " of fixed links

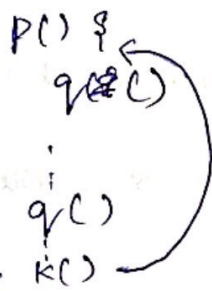
subtract the length of go to ad to get values

activation record structure



RA control stack:

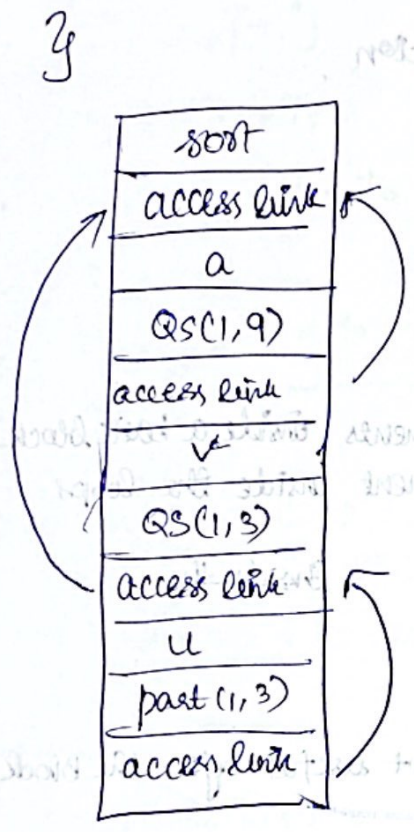
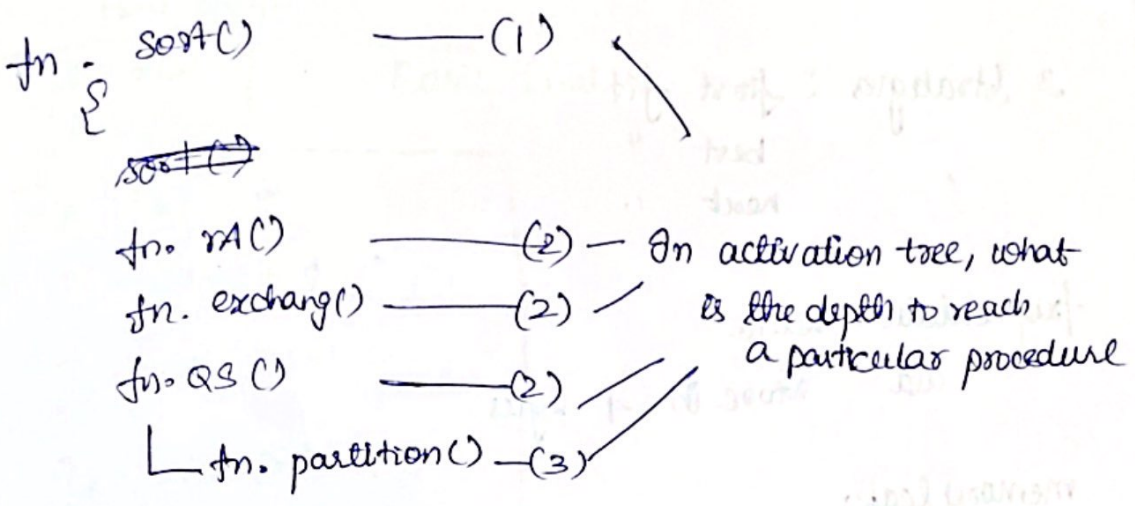
Immediately after calling procedure, A.R of callee proc



Let say, K - was array of procedure P
a proc K

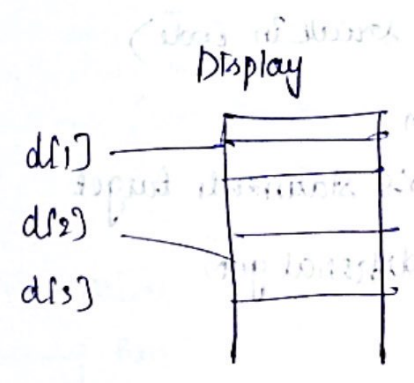
goes to the access link
pointer to a non-locals array

Using access link,
can go to any non local procedure



All the Q-S are under sort function

access link will point to the immediate nesting procedure.



temporal locality: stores the recently accessed items

spacial " : stores the items nearest to recently accessed

has higher chance of being called next

3 Strategies : first fit
 best "
 next "

fragmentation reduce:
 also those in 1 bytes

memory leak:
 not handling garbage collection.
 so fragmented :- cannot use

Dangling pointer
 excessive freeing of memory

4/14/24

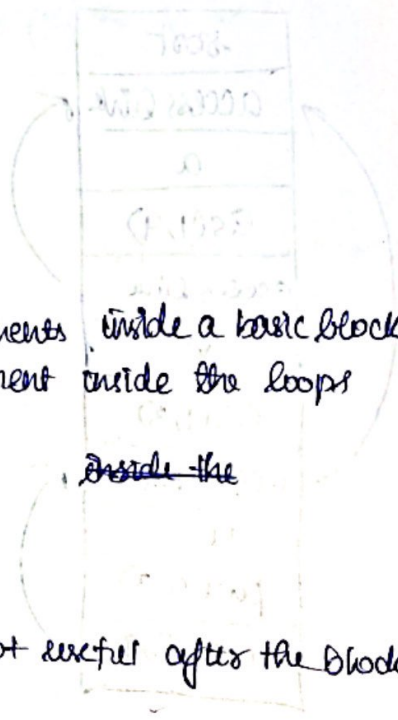
Basic blocks of

Code optimization :-> Minimize the no. of statements inside a basic block
 Improve the no. of statement inside the loops

code improvement
 no feasible approach

-> Dataflow analysis:

temp variables not useful after the block



Leader : (wherever there is break in code)

1. First statement of the program
2. Conditional or uncond goto's statements target
3. Immediate next stmt after conditional goto

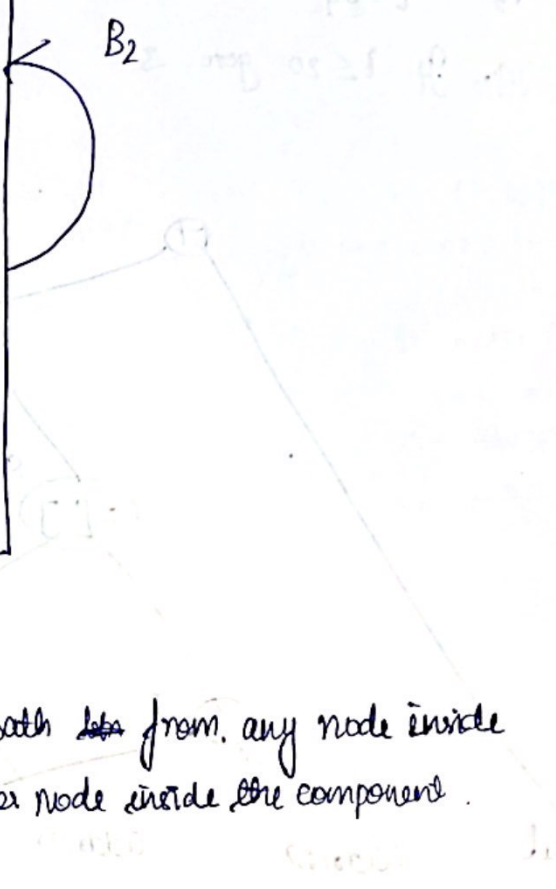
```

prod = 0
i = 1
do
  begin
    prod = prod + a[i] * b[i];
    i++;
  end
while i <= 20
  
```

1. $prod = 0$
2. $i = 1$

Basic Block B₁

3. $T_1 = 4 * i$
4. $T_2 = \text{addr}(a) - 4$
5. $T_3 = T_2 [T_1]$
6. $T_4 = \text{addr}(b) - 4$
7. $T_5 = T_4 [T_1]$
8. $T_6 = T_3 * T_5$
9. $prod = prod + T_6$
10. $i = i + 1$
11. $i \leq 20$ go to 3



Strongly connected comp

There should be path from any node inside the component to any other node inside the component.

Flow graph:

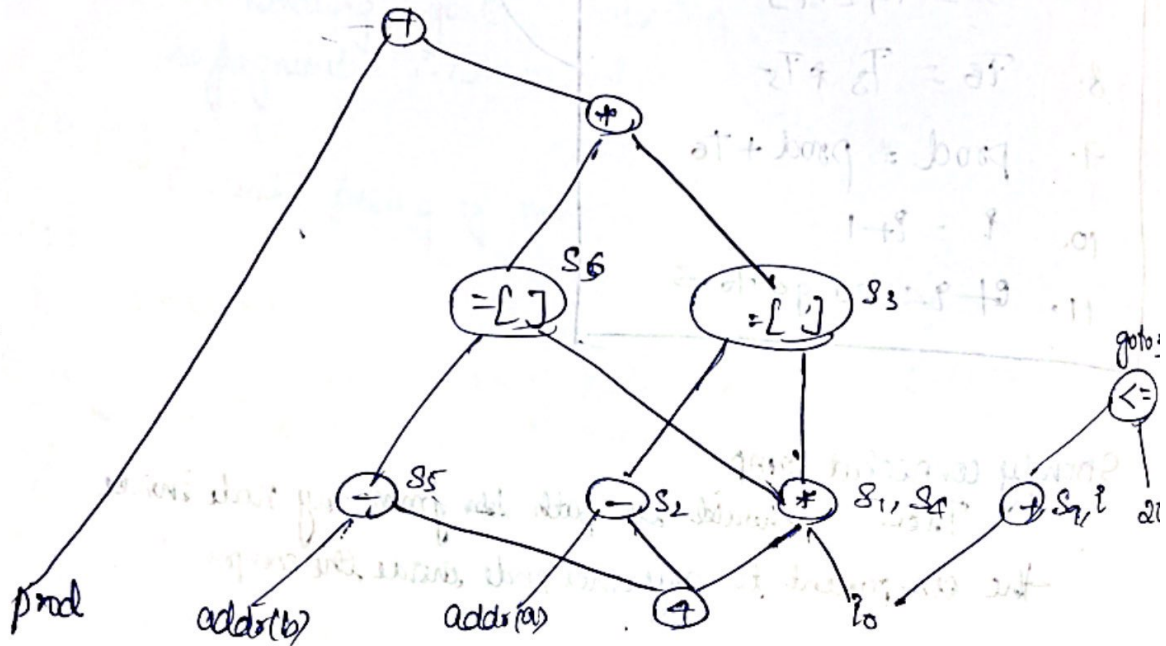
- Only one entry to go inside the component
- Strongly connected component

These conditions determines whether there is a loop in the flow graph.

- $4 * i$ need to calculated twice to write TAC but use optimized to one
- escaped using another temp var for argument, at $prod$ ②

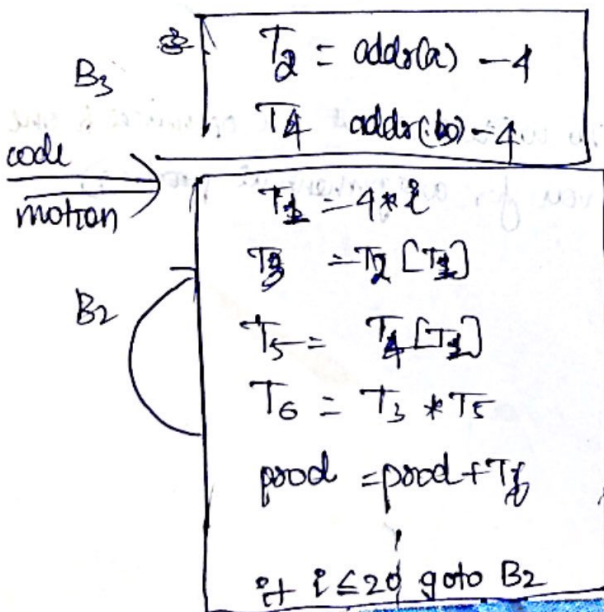
3. $S_1 = 4 * i$
4. $S_2 = \text{addr}(a) - 4$
5. $S_3 = S_2 [S_1]$
6. $S_4 = 4 * i$
7. $S_5 = \text{addr}(b) - 4$
8. $S_6 = S_3 * S_5$

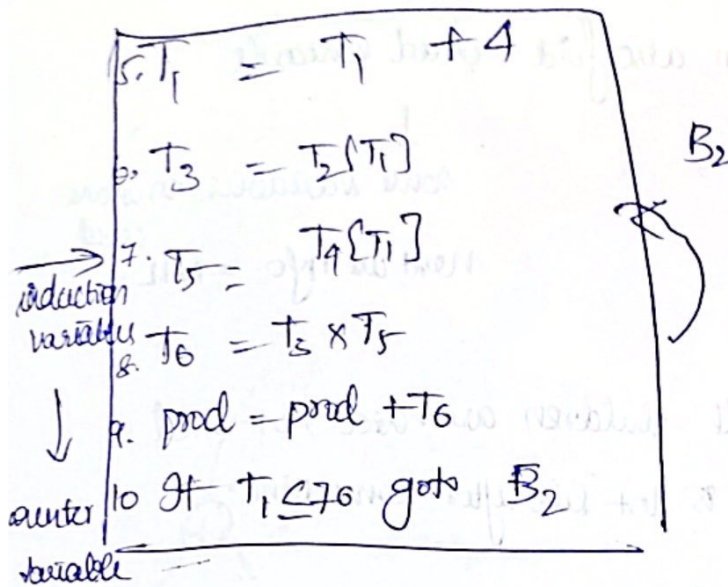
9. $S_7 \leftarrow S_8 * S_8$
10. $S_8 \leftarrow \text{prod} - S_1$
11. $\text{prod} = S_8$
12. $S_9 \leftarrow i + 1$
13. $i = S_9$
14. If $i \leq 20$ goto 3



Code motion

- B1 1. $\text{prod} = 1$
2. $i = 1$





24
 $i = 1, 2, \dots, 20$
 $T_1 = 4, 8, \dots, 80$
 \times
 T_6 only till

$\text{if } T_1 = 76 \text{ it will go one more loop for } 80$
 but $80 \text{ it gives } T_1$
 then it will run one more iteration

$* - \text{costly} \xrightarrow{\text{converted to } +}$

20/4

Optimization: (Reductions)

- code motion
- induction variables
- Reduction in Strength
- live variable analysis:

find the next use possible for the variable
 ↓
 next use information

$i : x = y \text{ op } z$
 line number i
 is the next use loop

$T_3 \text{ or } T_5 - \text{next use info} \dots 8$
 $\text{prod}, T_6 - \dots 9$
 $T_1 - \dots 10$

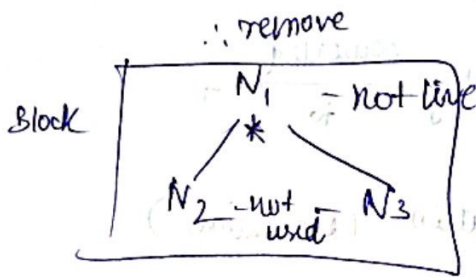
Dead code elimination:

live var analysis can also find dead variable

↓
 some variable no where used
 next use info = NIL

A block of code where all children are ~~not~~ not used elsewhere. If the parent is not live after some time ∴

remove the variable



Use of Algebraic Identities

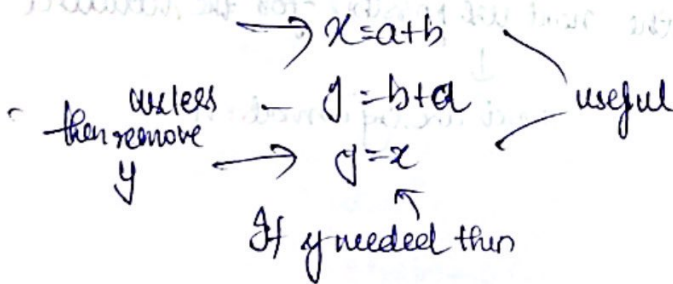
$$i = i \times 1$$

$$i = 1 \times i$$

$$x = x + 0$$

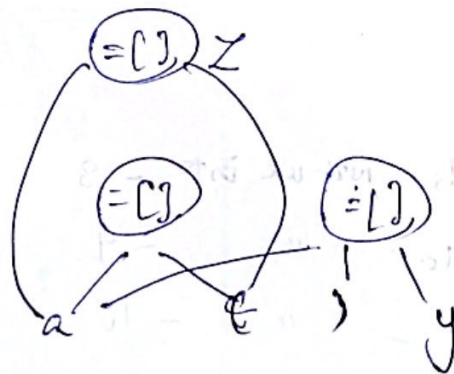
$$x = x - 0$$

useless computation



$$a = b + c$$

$$e = c + d + b$$

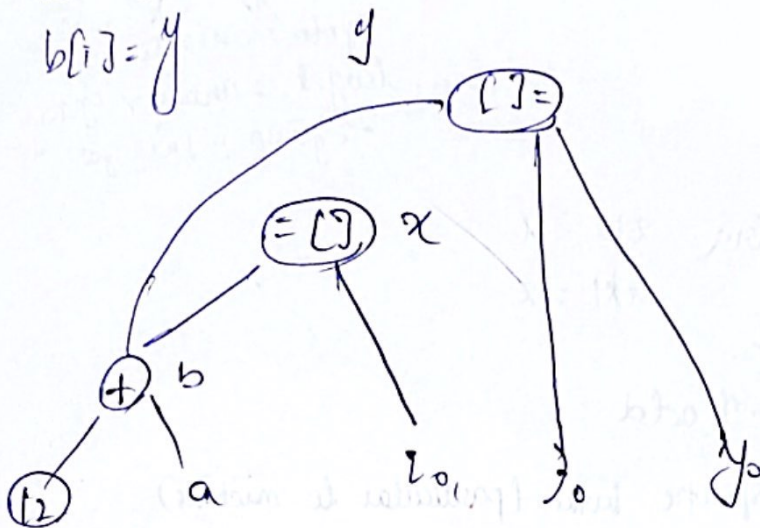


10/01/14

$$b = 12 + a$$

$$x = b[1]$$

$$b[i] = y$$



30/4/24

Loop hole optimisation:

- Eliminating redundant loads & stores (1)
- Eliminating unreachable code (dead code) (2)
- Flow of control optimisation (3)
- Algebraic simplification (4)
- Reduction on strengths (5)
- Use of machine idioms (6)

(1) LD a, R0
ST R0, a

↳ continuous statements

↳ either of the statements may be removed

(2) If debug == 1 goto L1
goto L2
L1 print debugging info

↓ L2
If 0! = 1 goto L2
↓ L2

- L1 will never be reached if debug is set to 0
- If L1 is not used anywhere, it can completely be removed

(3) goto L1

L1 goto L2

↓

goto L2

peephole optimization eliminates jumps
replace L1 with L2

over jumps

goto statements
target is another goto
→ jump is final goto

(4) Apply rules of algebra

$$x+0 = x$$

$$x*1 = x$$

(5) pow(x, 2) → x * x

x + 2 → replace with add

(6) Embed func to do specific tasks (particular to machine)

INC of increment/decrement
DEC

instead of $x = x + 1$ write INC x

↓
underlying arch

ISA

CODE GENERATION

Generated code should have the following characteristics

- correctness
- easy to implement
- easy to debug
- efficient

Input to code generator : optimised intermediate representation

Assumptions

- 1- lexical, syntactic or semantic error \times
- 2- code has been optimised (intermediate repr is efficient)

Output : target language (usually assembly)

Design Issues

→ input to code generator

→ target ~~issue~~ program

↳ underlying architecture : CISC/RISC
↳ resources of architecture

→ special case : JVM (runs here on a stack-based arch)

↓
intermediate : byte code
repr

→ uses push/pop

considered absolute (almost)

↓
directly convert to target (mc) code
no intermediate (assembly) code created.

Instr selection

LD R0 b } a = b + c
ADD R0 R0 c

ST	a	R0
LD	R0	a

→ extra block as R0 contains a & can be used directly

ADD R0 R0 e } d = a + e
ST d R0

limited no. of registers
need to know which registers are "safe"
values don't change

Register access has min cost

$$R_1 \leftarrow R_1 + R_2$$

$$R_1 \leftarrow \text{content}(100 + \text{content}(R_2)) - 2$$

$$R_1 \leftarrow \text{ " (content}(100 + \text{content}(R_2)) - 3$$

↓ 3 'content'
access max 3x

$$\text{cost} = 3$$

Exercises

1. $x = y - 5$

2. $a[i] = c$

3. $x = *p$

4. $*p = y$

5. $\text{if } x < y \text{ goto L}$

(3) LD R1, p
LD R2, 0(R1)
ST x, R2

1/5/24

LD R1, a

ADD R1, R1, R2

LD R1, 100(R2)

LD R1, *100(R2)

LD R1, #100

$R_1 \leftarrow R_1 + R_2$

$R_1 \leftarrow \text{content}(100 + \text{content}(R_2))$

$R_1 \leftarrow \text{content}(\text{content}(100 + \text{content}(R_2)))$

$R_1 \leftarrow 100$

R:

$a[x] = c$

↓

LD R1, c

LD R2, j

MUL R2, R2, 8

ST a(R2), R1

$x_p = y$

↓

LD R1, P

LD R2, y

ST 0(R1), R2

if $(x < y)$ goto L

LD R1, x

LD R2, y

SUB R1, R2, R3

BLTZ R1, L

A Simple Code Generation

2 Data Structures used

→ Register descriptor

→ Address "

- list of curr var in every reg

- for every var, its location info

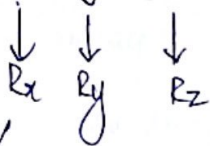
Register access - min cost - unit 1

memory access > 1

\therefore 2 register contents - 2 units

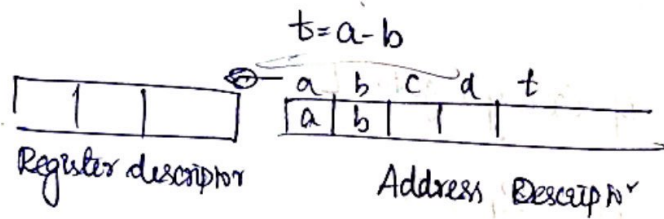
" " - 3 units

getReg(x = y op z)



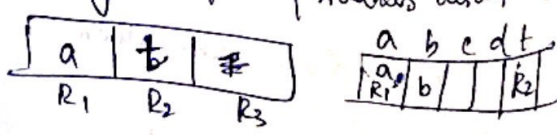
One of R_x for any other var should be nil \therefore R_x should now be the address of x which is updated.

- 1) update register descriptor of R_x to hold only x & remove R_x from the address descriptor of every other variable.
- 2)

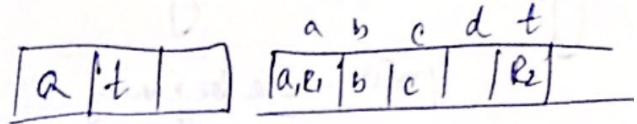


```
LD R1 a
LD R2 b
SUB R3 R1 R2
```

update Reg descriptor & Address descriptor

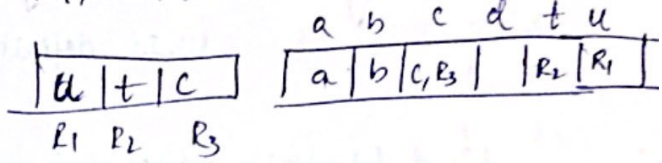


$$u = a - c$$



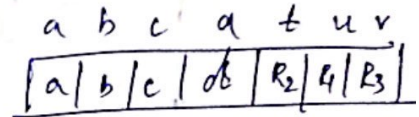
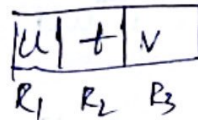
LD R3, c

SUB R1, R1, R3



$$v = t + u$$

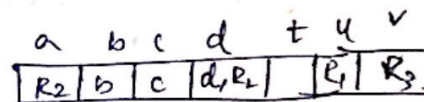
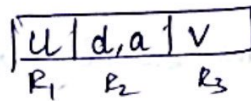
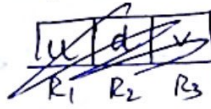
~~LD~~
ADD R3, R1, R2



$$a = d$$

LD R2, d

not doing any op for assignment, just changing descriptors



meaning
both has same value

exit

while exiting store R2 to a

ST a R2

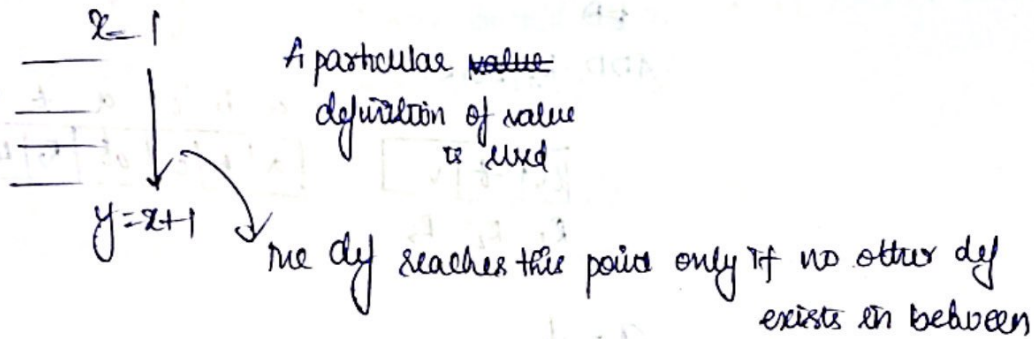
t, u, v are temp variables - no need of string
the defined variables are need to be stored

Global Data Flow Analysis

Point - the beginning of any statements or just after the statement

Used definition - when usage of a value which definition of value is active

Used definition (ud) chaining - find where the value is used which definitions are active



Reaching definitions:

For every block the variables that are active at the start of the block ~~at~~ at the end of block.

IN[B] → definitions reaching before start of block

OUT[B] → " " end of block.

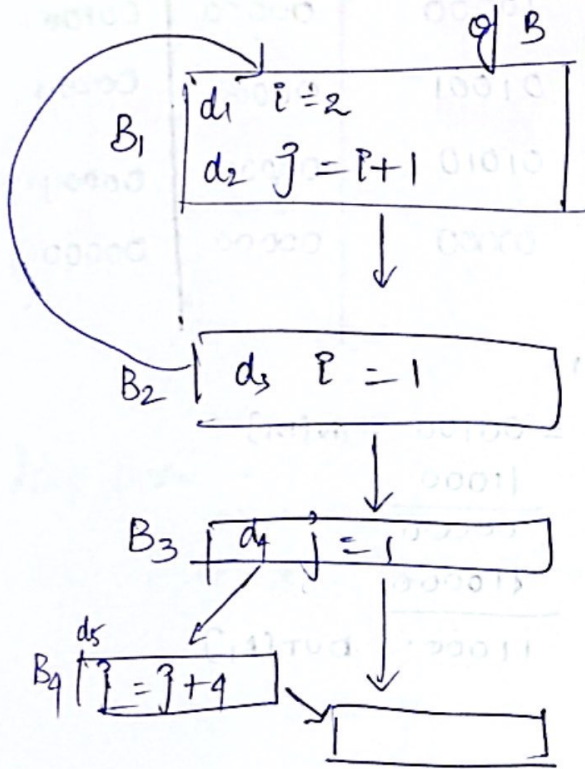
GEN[B] → what are the definitions generated inside that block

KILL[B] → what are the definitions that are killed due to re-definitions

$$OUT[B] = IN[B] - KILL[B] \cup GEN[B]$$

$$IN[B] = \cup OUT[p]$$

p is predecessor blocks



A vector of all definitions active

$$GEN[B_1] = \begin{matrix} d_1 & d_2 & d_3 & d_4 & d_5 \\ 1 & 1 & 0 & 0 & 0 \end{matrix}$$

$$KILL[B_1] = \begin{matrix} 0 & 0 & 1 & 1 & 1 \\ \text{Kill all other def of } i, j \end{matrix}$$

$$GEN[B_2] = \begin{matrix} 0 & 0 & 1 & 0 & 0 \end{matrix}$$

$$KILL[B_2] = \begin{matrix} 1 & 0 & 0 & 0 & 0 \end{matrix}$$

$$GEN[B_3] = \begin{matrix} 0 & 0 & 0 & 1 & 0 \end{matrix}$$

$$KILL[B_3] = \begin{matrix} 0 & 1 & 0 & 0 & 0 \\ \text{Kill all other def of } j \end{matrix}$$

Kill all other def of j

$$OUT[B_1] = 00000 - 00111 \cup 11000$$

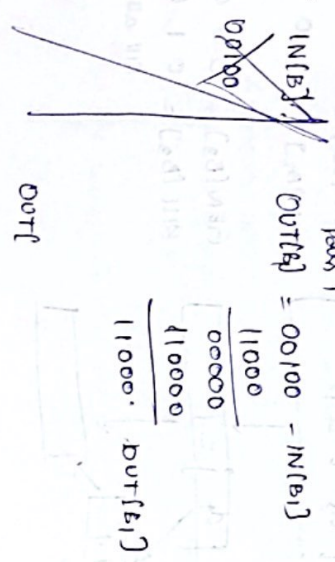
$$A - B = A \wedge \neg B$$

$$\begin{array}{r} A \quad 00000 \\ \quad 11000 \\ \hline \quad 00000 \\ \vee \quad 11000 \\ \hline 11000 \end{array} \quad \begin{array}{l} \text{AND} \\ \text{OR} \end{array}$$

$$\begin{cases} OUT[B_1] = \cancel{11000} 11000 \\ IN[B_1] = 00100 - \text{loop} \end{cases}$$

$$OUT[B_2] = 00100$$

Block	GEN[B _i]	KILL[B _i]	IN[B _i]	OUT[B _i]	IN[B _i]	OUT[B _i]	IN[B _i]	OUT[B _i]
B ₁	11000	00111	00000	11000	01100	11000	01111	11000
B ₂	00100	10000	00000	00100	01111	01111	01111	01111
B ₃	00010	01001	00000	00010	01111	01111	01111	00110
B ₄	00001	01010	00000	00001	00110	00101	00110	00101
B ₅	00000	00000	00000	00000	00111	00111	00111	00111



IN[B₁] = USE[OUT[P]]

OUT[B₁] = IN[B₂] - KILL[B₁] ∪ USE[B₁]

Pass 1	Pass 2	Pass 3			
IN[B ₁]	OUT[B ₁]	IN[B ₁]	OUT[B ₁]	IN[B ₁]	OUT[B ₁]
00100	11000	01100	11000	01111	11000
11000	01100	11111	01111	11111	01111
01100	00110	01111	00110	01111	00110
00110	00101	00110	00101	00110	00101
00111	00111	00111	00111	00111	00111

Loop unrolling:

Increase the body of loop or decrease the no. of comparison

for (i=1 to 100)

Σ A[i] = 0

→

for (i=1 to 100)

Σ A[i] = 0

i = i + 1

A[i] = 0

Loop forwarding:

save loop indices but diff body : combine

for i = 1 to 10

for j = 1 to 10

A[i, j] = 0

→

for i = 1 to 10

for j = 1 to 10

A[i, j] = 0