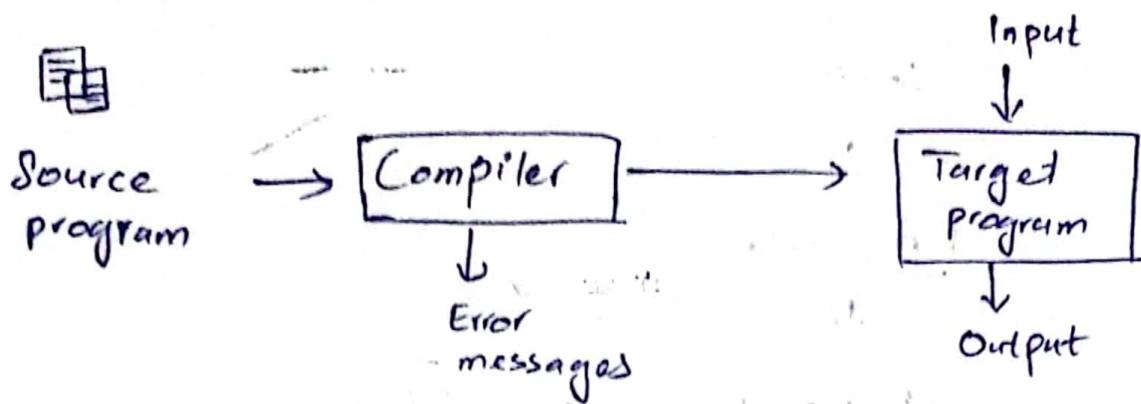


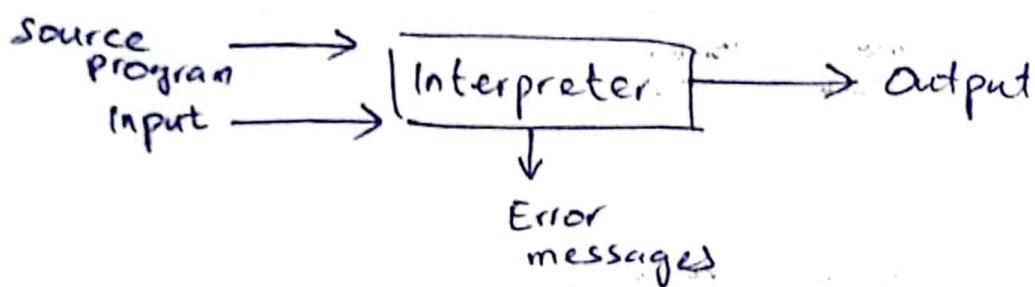
CSC441 - Compiler Design

Prepared by: Mr. Aoun_Haider
fa21-bse-133@cuilahore.edu.pk

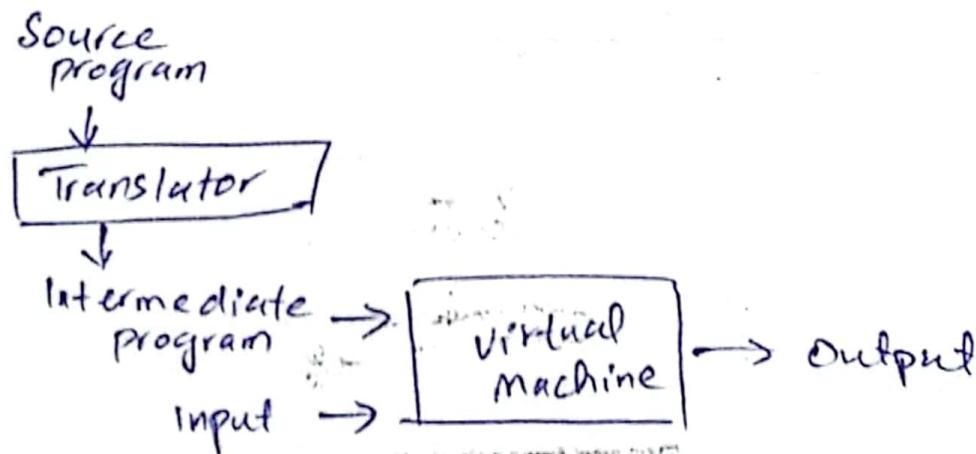
- ① Compiler:- Translation of source program to target program.



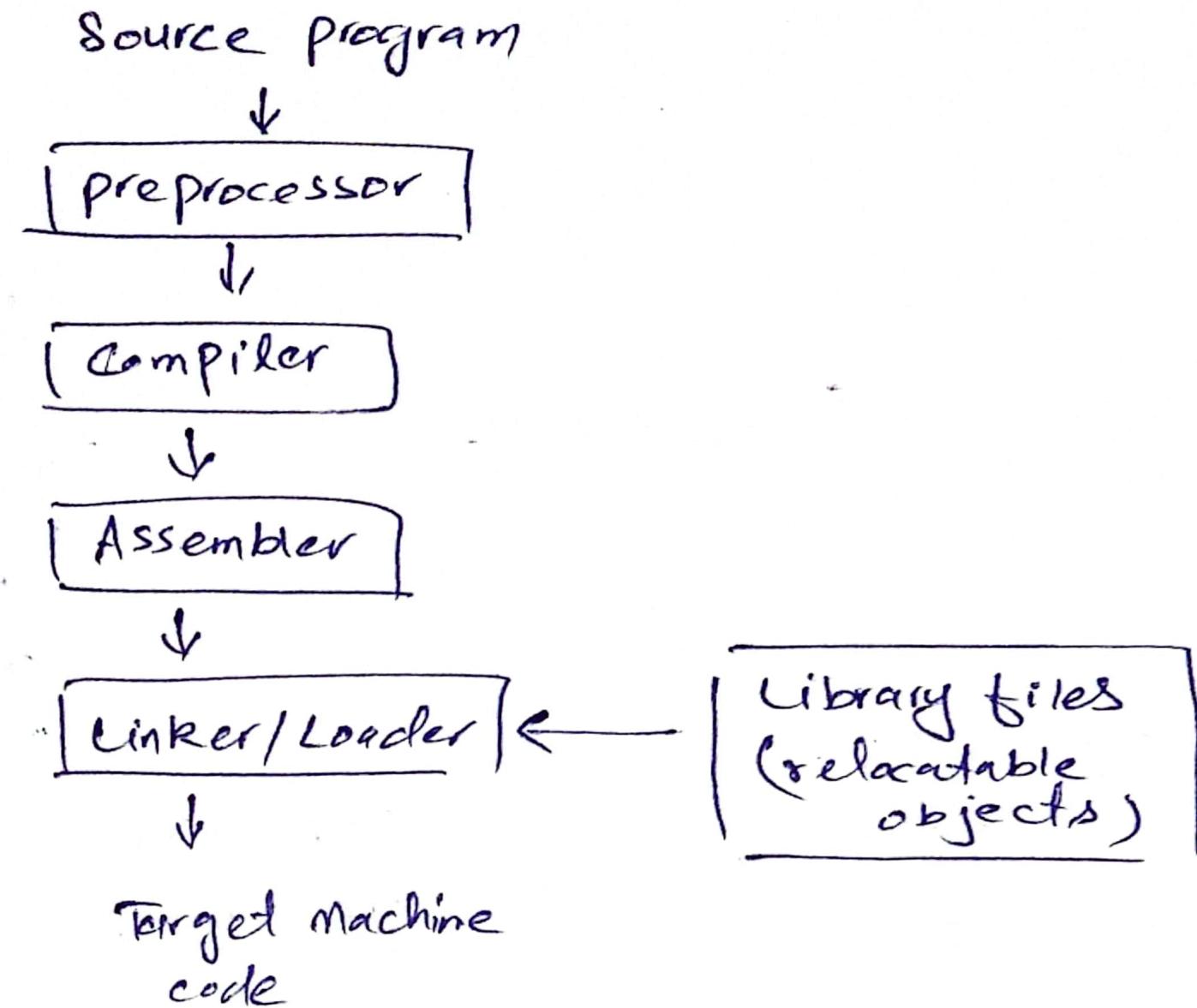
- ② Interpreters:
- line by line compilation
 - performing the operation implied by source program



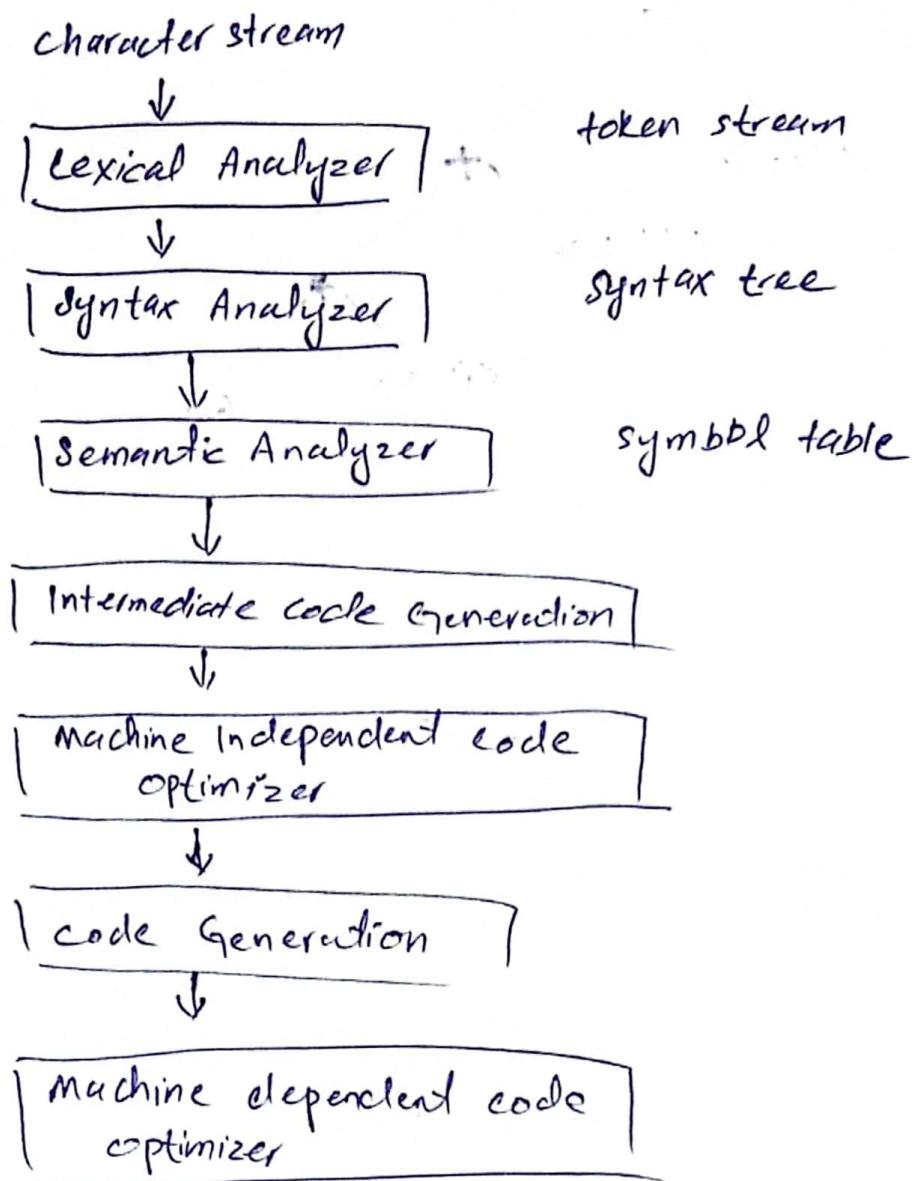
- ③ Hybrid:



Language processing system:



phases of Compiler:



Example:

position = initial + rate * 60 ;

① Lexical Analyzer

- Divide each section into tokens and assign 'id' of symbol table

position = < id, 1 >

initial = < id, 2 >

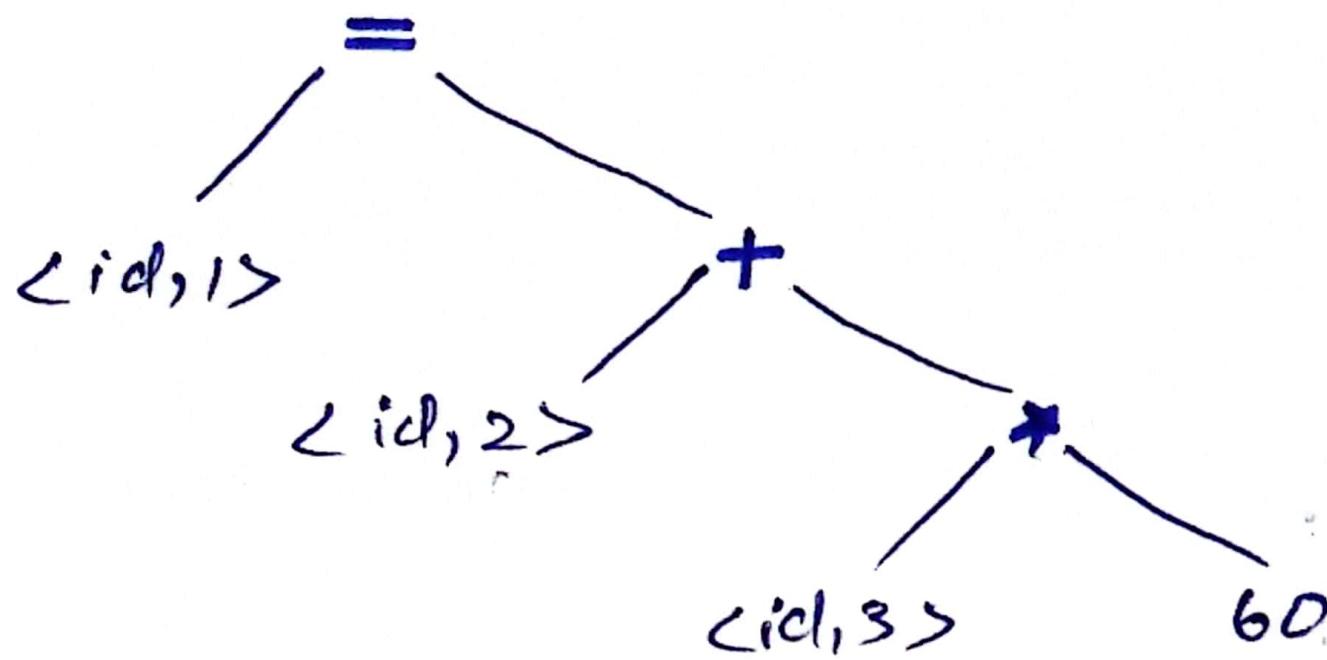
rate = < id, 3 >

< id, 1 > => < id, 2 > < + > < id, 3 >
 < * > < 60 >

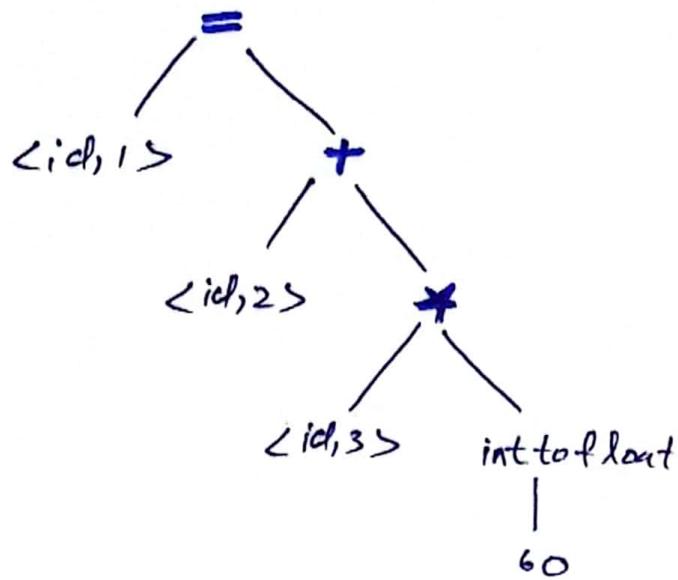
Symbol Table

1	position	...
2	initial	...
3	rate	...

Syntax Analyzer:



3) Semantic Analyzer:



Notes: In this case, semantic analyzer has performed type coercion (implicit type casting) due to the type of position 'float' because int and float are different types and type coercion is required to perform operation.

4) Intermediate code Generation:

$$t1 = \text{inttofloat}(60)$$

$$t2 = id3 \times t1$$

$$t3 = id2 + t2$$

$$id1 = t3$$

5) Machine independent code optimizers:

$$t1 = id3 \times 60.0$$

$$id1 = id2 + t1$$

5) Code Generation:

LDF R2, id3

MULF R2, R2, #160.0

LDF R1, id2

ADDF R1, R1, R2

STF id1, R1

Types of compiler:

Single pass

- 1) Compiler visits the source code ^{not} multiple times
- 2) Require less memory
- 3) Fast
- 4) makes less optimal code
- 5) Detects error in linear fashion and require correction of single error every time

Multipass

- compiler visits the source code multiple times & observe a specific portion each turn
- more memory requirement
- slower
- makes more optimal code

Detect errors back-to-back

Contradiction:

- Compiler is multipass & interpreter is single pass (not always true)
- interpreter is used for dynamically typed languages like python while ~~single~~ compiler is used for **Statically** typed languages.

Grouping

① Front End → analysis (Machine independent)

② Back End → synthesis (Machine dependent)

- Lexical,
- Syntax,
- Semantic,
- Intermediate code

- Code Generation
- Optimization
- (Register allocation, Instruction scheduling)

① Lexical Analysis: (Scanner, Lexer, tokenizer)

- identifier: defined by programmer i.e. variable name
- keyword: language reserved words ('if', 'else', 'for')
- number/constant/literal: 0, 1, 2, ..., 9
- ~~operator~~: Separator {, (,), }, ;
or
punctuators

► Analysis tasks:

① Tokenization

- ② Error messages → Exceeding strings (long variable names)
Illegal character (unending multiline comment)
- ③ Eliminate comments & whitespaces, tabs, \n, \t
NFA or DFA are used in this phase

Example :

④ printf ("i=%d, &i=%x", i, &i);

10 tokens

⑤ main() { /* City Lahore */ → comments are treated as blank characters
a = b ++ + - - - ++ ==; Test product code
printf ("Aocin : %d", a); Inspect code
}

"Look ahead problem"

$i = if \quad || \quad = or \quad ==$

- Ambiguity where a token ends and next token starts

Steps:

- Scan the character stream left to right
 - Identify each token
 - Identify each tokens' category

Regular Expression

Σ = set of character

$$t = \alpha r$$

Φ = zero or more occurrence

$+$ = one or more occurrence

$$A^+ = AA^* \quad \text{set of all alphabets}$$

- ① string of letter or digits starting with letter
~~digit~~ = '0' + '1' + '2' + ... + 'q'

$$\text{digit} = '0' + '1' + '2' + \dots + 'q'$$

Letter = 'a' + 'b' + 'c' + ... + 'z'

= letter (letter + digit)*

- ② phone numbers: (650)-723-3232

$x = \text{digit}^3$

$y = \text{digit}^4$

三

$$y'x'y - x'y - y$$

- ### ③ Unsigned digit

$$\text{digit} = '0' + '1' + \dots + '9'$$

digits = digit⁺

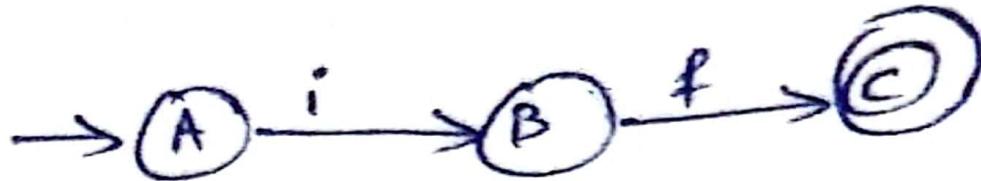
$$\text{fraction} = (\text{' digits}) + \epsilon$$

$$\exp = ('e' (\text{ digits}) + \epsilon)$$

number = digits fraction exp

Tokens' DFA

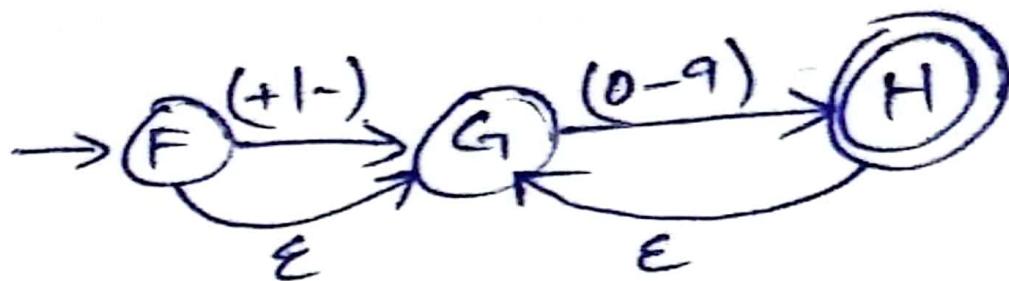
- if :



- identifiers:



- Integer:



ϵ = space, no input

Remaining after 1 page

Symbol Table

- produced by analysis phase (front end) and consumed by ~~synthesis~~ synthesis phase (back end)
- Data structure maintained by compiler to store instances of variable, objects, classes, interface etc.
- Tasks performed by each phase regarding symbol table
 - ① Lexical: create entries for identifier
 - ② Syntax: adds information regarding attributes
 - ③ Semantic: Using available information, updates the symbol table.
 - ④ Intermediate code: Available info helps in adding temporary variables information
 - ⑤ Optimization: Available info is used in machine dependent optimization
 - ⑥ Target code: Generate target code using address info of identifiers
- Sample symbol table:
int varle = 133;

Name	Type	size	Dimension	Line of declaration	Line of usage	Address
varle	int	4	0	5	7 → 9 → 10 → NUL	0X3B076F

- Size of symbol table must be considered
 - too small: no space to store more entries dynamically
 - too large: wastage of space
- So, symbol table must claim dynamic memory

Non-Block Structured language: Fortran

Each variable must be declared once in entire program

Operations: insert(), lookup()

Block Structured:

- Variables can be declared in multiple blocks / ~~scope~~, locality
- Operations:
Insert(), Lookup(), set(), Reset()

10 of 46

All latest languages are block structured

Quiz:

① Which one is not intermediate representation of source program?

- a) 3-address code
- b) Abstract Syntax tree
- c) Symbol table
- d) Control Flow Graph (CFG)

② Access time of symbol table will be $O(\log n)$ if implemented via

- a) Linear list
- b) Search Tree
- c) Hash table
- d) None

Explanation:

→ Linear list

ordered $O(n)$

→ linked

unordered

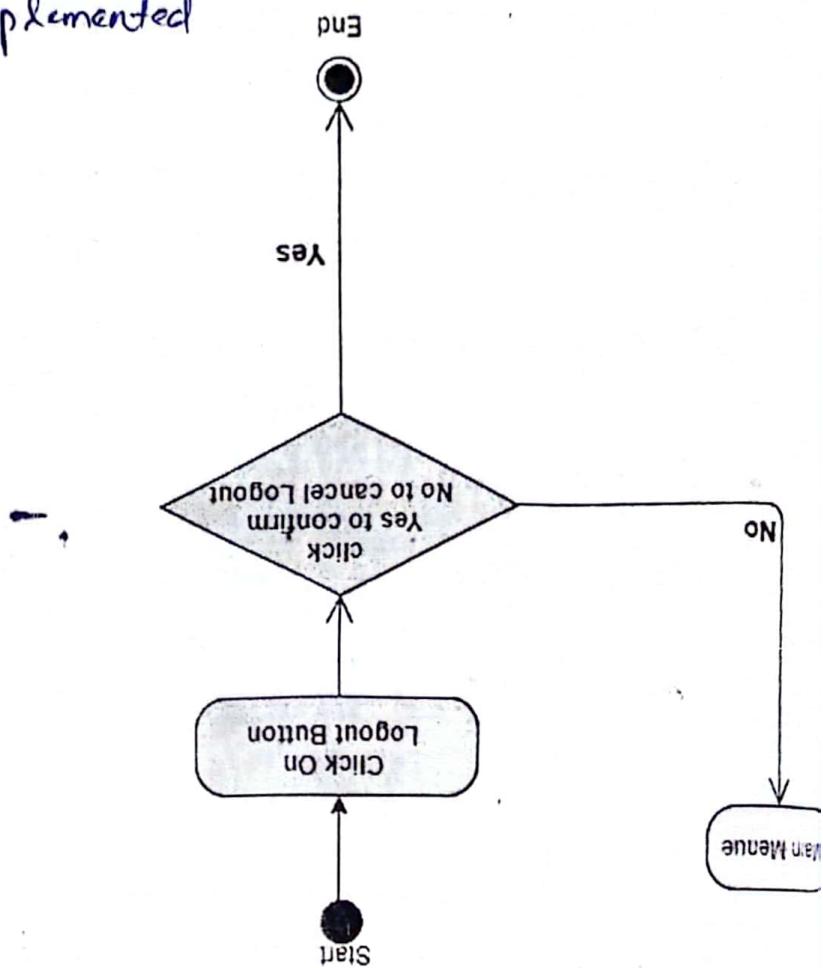
→ Tree

$O(\log n)$

→ Hash table

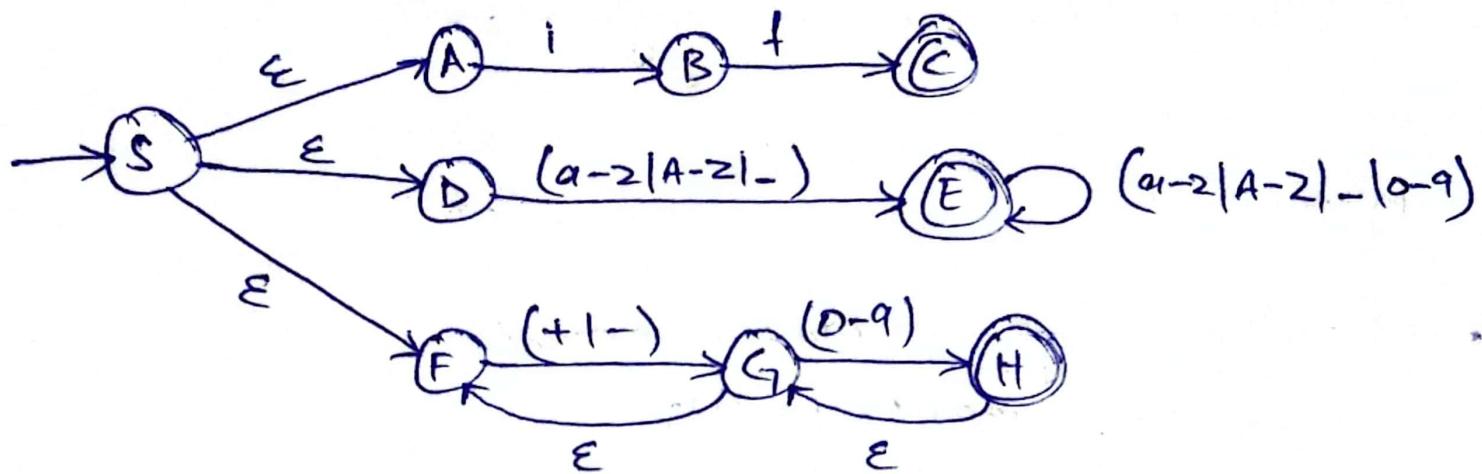
• if no collision $O(1)$

• otherwise $\approx O(n)$

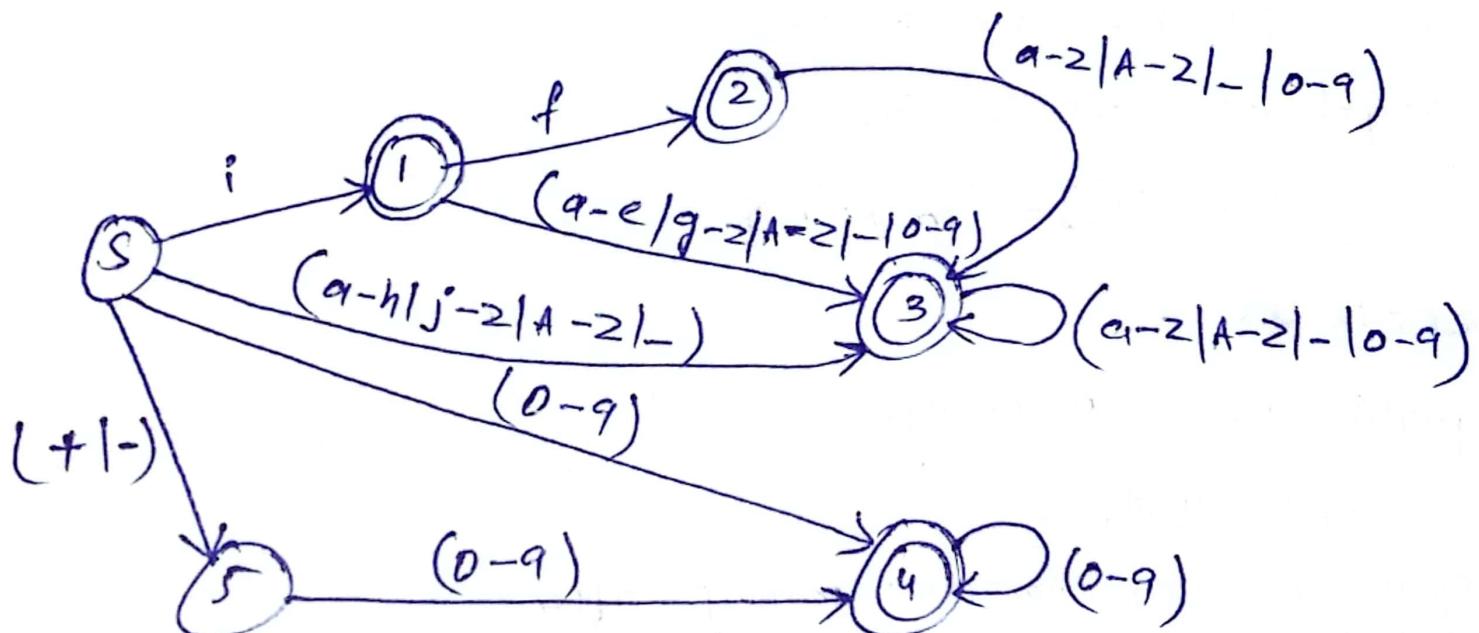


NFA:

- Let say we want to recognize if, identifier 4 digits
we need to combined all the DFAs and introduce
a new state 's'



DFA:



Errors & Error Recovery in Lexical Analysis

• Lexical Errors:

- Identifiers are too long

C allows variable name of length = 247

of 46

C++ " " = 2048
python " " = 79

- Exceeding length of numeric constants

int id = 13345678911;

if 2 bytes allowed, error!

Location Map Table

- Numeric constants are ill-formed

int i = 4567 \$ 913

- Illegal characters addition;

char x[] = "AOUN", \$

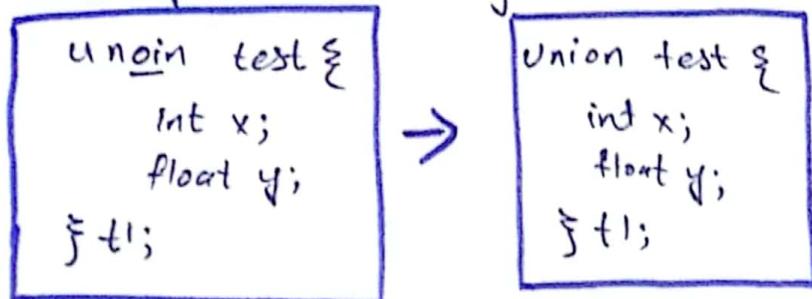
Recovery:

- panic mode recovery:

int yrd;

ignore remaining characters scanning after detecting invalid character

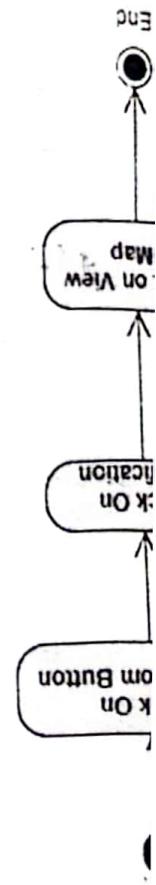
- Transpose of 2 adjacent characters



- Insert missing character



- Delete unknown extra character



lapping

- Replace a character with another

itt id; → int id;

Example :

0 height = (width + 56) * factor(foo);

• Generate symbol table, count # of tokens, create CFG and syntax tree

① Symbol table:

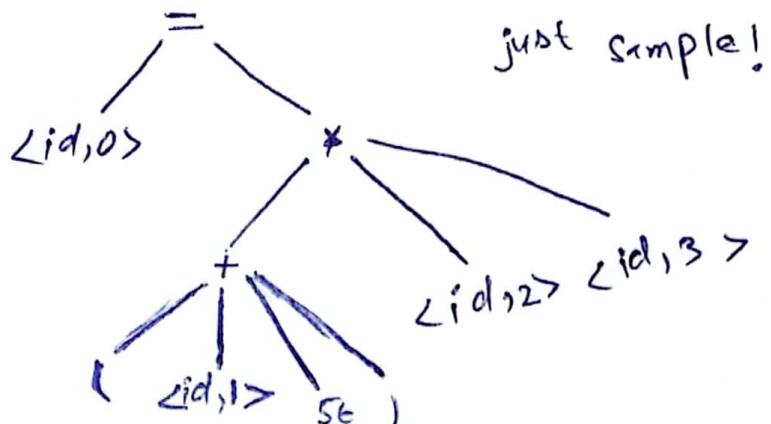
Name	Type	Size	Dimension	Line of declaration	Line of usage	Address
0 height	int	4	0	0 -	0	-
1 width	int	4	0	0 -	0	-
2 factor	int	20	0	-	0	-
3 foo	int	4	0	-	0	-

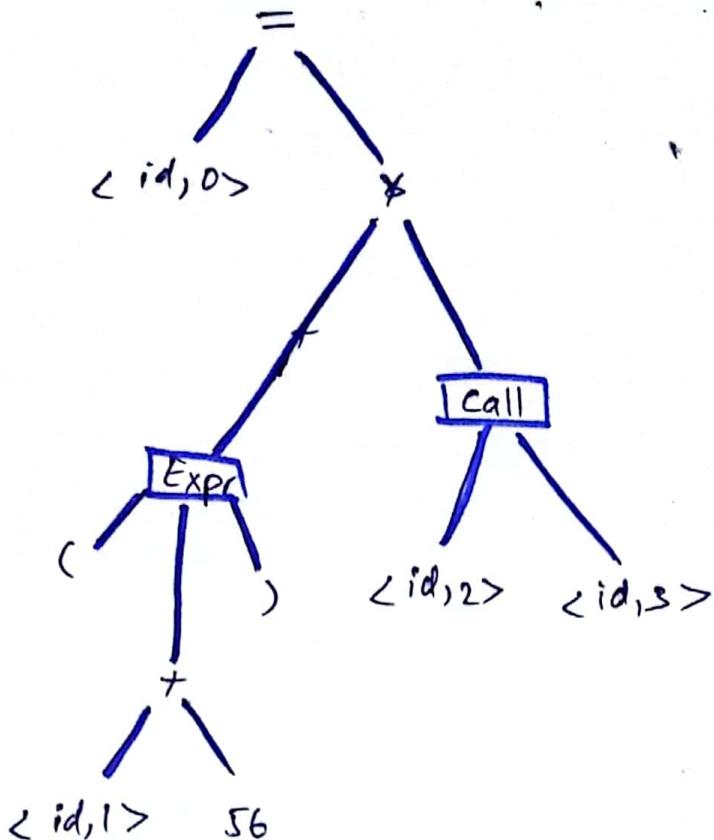
② # of tokens = 13

③ CFG

$\text{expr} \rightarrow (\text{expr}) \mid \text{expr} * \text{expr} \mid \text{int} \mid \text{id}(\text{expr}) \mid \text{id}$
 $\text{expr} = \text{expr}$

④ Syntax tree





$\langle id,0 \rangle \Leftrightarrow \langle (\rangle \langle id,1 \rangle \langle + \rangle \langle 56 \rangle \langle) \rangle \langle * \rangle$
 $\langle id,2 \rangle \langle (\rangle \langle id,3 \rangle \langle) \rangle \langle ; \rangle$

Input buffering:

- Store the source program code in a temporary storage area called buffer for processing
- Use two pointers `lexemBegin` and `forward` in buffer
- Put a sentinel EOF at the end of buffer to indicate end of file

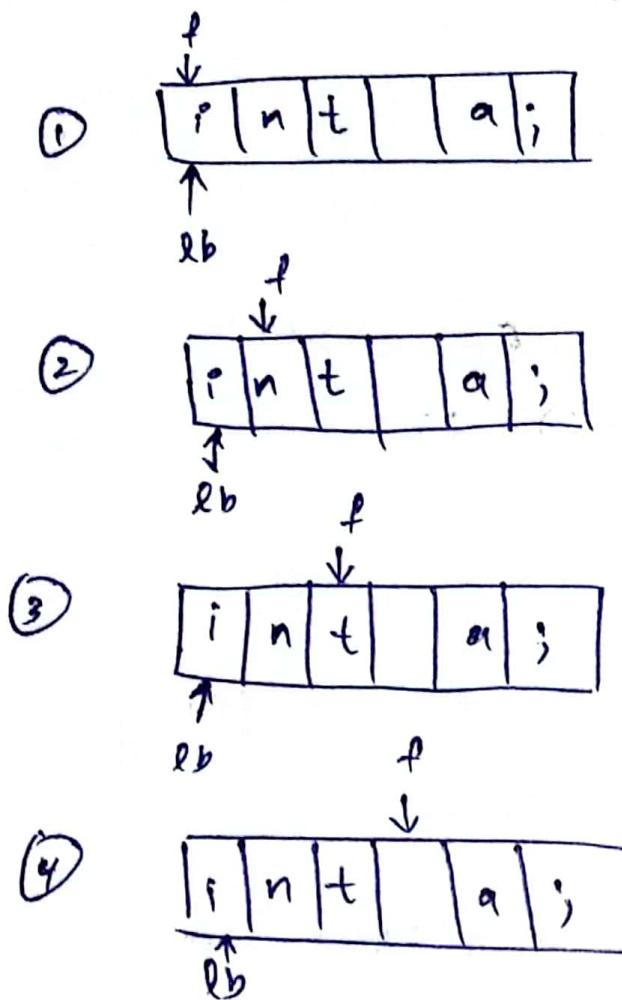
Sometimes, we cannot decide operator by just observing or scanning single character. For example; $=$ vs $==$ or $>=$ or $=$. It means we need some mechanism to keep track of next character as well to tick it as legal character (lexeme). That's why two pointers are used.

M	=	A	*	*	2	;	EOF
---	---	---	---	---	---	---	-----

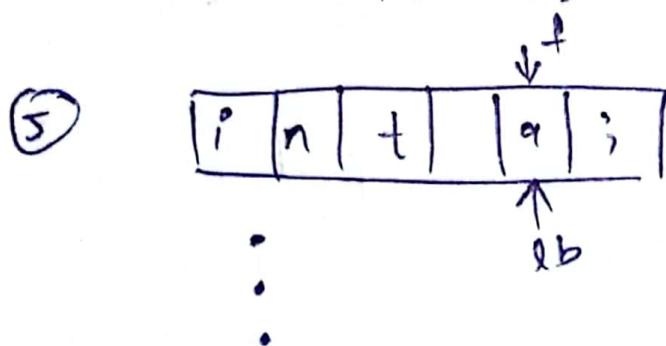
lex
Begin forward

38 of 46

"lexemeBegin" pointer points at the start of token and forward pointer moves one position right and after detecting a token end, identify it as lexeme or not. The lexBegin will move to the start of next token & so on.



token = "int" \Rightarrow lexeme



reading a single character require read() system call. For reading program of 1000 characters, 1000 system calls are needed.

Buffering helps us to reduce this overhead by reading complete program in a single disk block size memory (4K block size = 4K block memory in RAM) requires only one system call.

Schemes:

① One Buffer:

- Use only one buffer to read input stream
- If input stream size is larger than buffer size, buffer will overflow

Example: size of input stream = 500

buffer size = 100

load 1st 100 characters 0-99, after processing 0-99 characters, override the buffer with 100-299 until complete program is ~~not~~ processed.

② Two buffer:

- Use 2 buffers to read input stream
- put sentinel character 'eof' at the end of each buffer which is not part of source code.

i	n	t	i	;	j	;	eof
---	---	---	---	---	---	---	-----

i	=	i	+	1	i	;	eof
---	---	---	---	---	---	---	-----

Double buffering Algorithms

```
switch(*forward++) {
```

case eof:

```
    if (forward is at end of 1st buffer) {
```

```
        reload 2nd buffer;
```

```
        forward = beginning of 2nd buffer;
```

```
}
```

```
else if (forward is at the end of 2nd buffer) {
```

```
    reload 1st buffer;
```

```
    forward = beginning of 1st buffer;
```

```
}
```

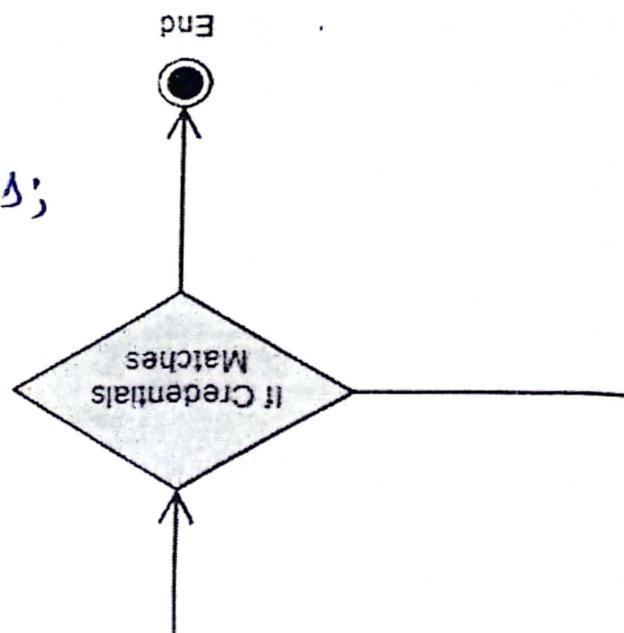
```
else
```

```
    terminate lexical analysis;
```

```
    break;
```

```
// Case of other characters
```

```
}
```



String related terms:

① prefix: Remove zero ~~or~~ more characters from end of string

banana = ban, banan, ε, banana, ba, bana

② Suffix: Remove zero or more characters from start of string

aoun =oun, aoun, un, ε

③ substring: By deleting prefix or suffix from string

aoun = ou, ao

④ subsequence:

Removing zero characters ~~randomly~~ from random position not consecutive

aoun = ou, an, on

Context Free Grammer: CFG

• Way to describe rules of language

Components: $\text{CFG} = (\text{S}, \text{N}, \text{T}, \text{P})$

① Terminals: Non-replaceable

② Non-terminals: replaceable

③ production rules: set of all the rules

④ Starting state 'S'

Ambiguous grammar: having more than one parse tree

Example:

Consider C code:

if($x < 5$) {

b = 2;

} else {

b = 1;

}

CFG:

stmt \rightarrow if expr then stmt |

 if expr then stmt else stmt |
 ε

expr \rightarrow term relop term | term

term \rightarrow id | number

digit \rightarrow [0-9]

digits \rightarrow digit*

number \rightarrow digits (· digits)? (E [+|-] digits)?

relop = relational operator

letter \rightarrow [A-Z a-z]

id \rightarrow letter (letter | letter + '-' | digit)*

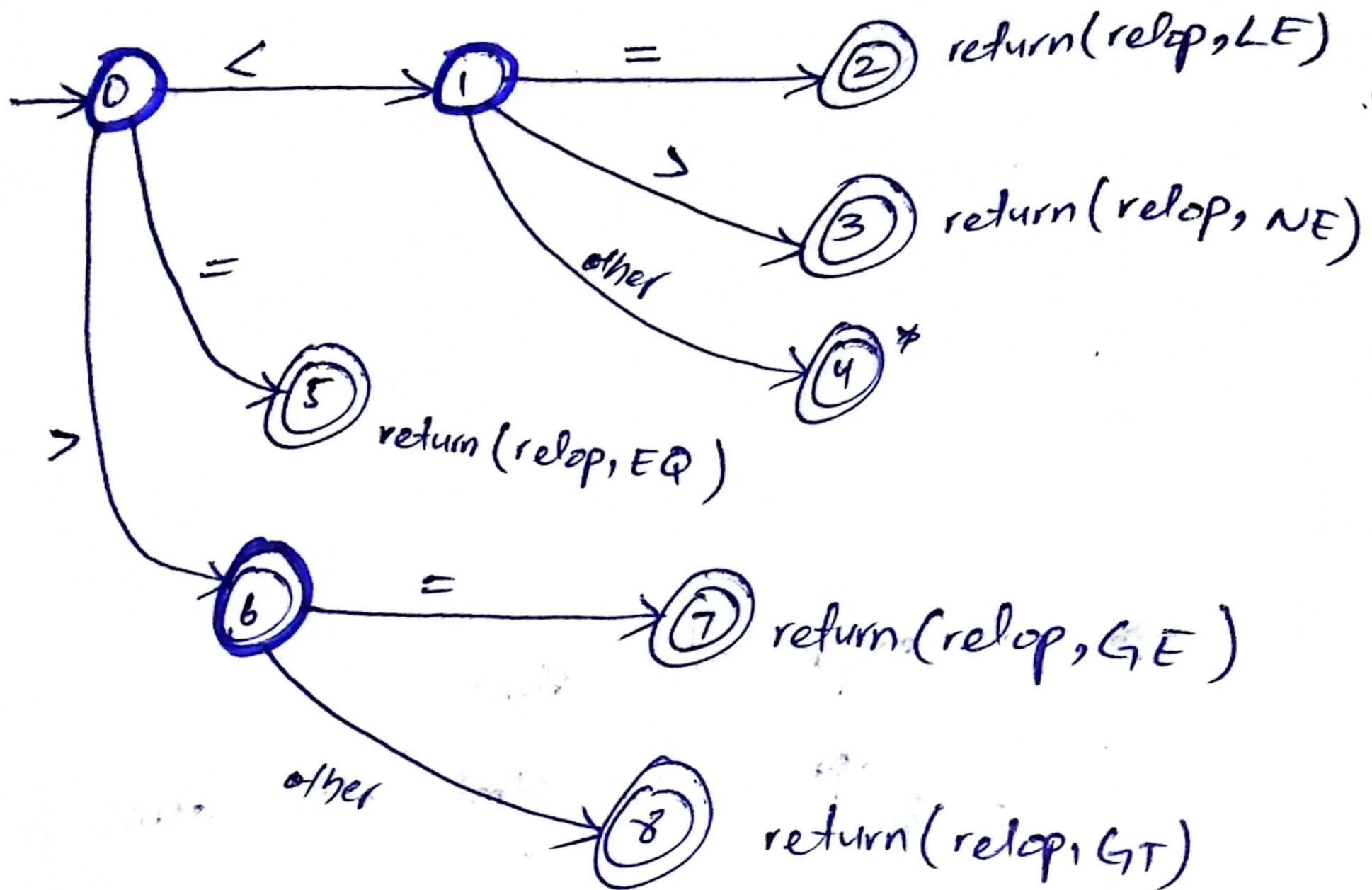
if \rightarrow if

else \rightarrow else

then \rightarrow then

relop \rightarrow < | > | <= | >= | = | <> | !=

Transition Diagram for 'relOp':



How to recognize reserved keywords?

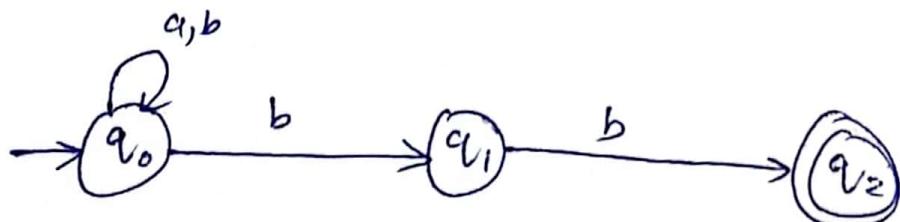
- ① store all reserved keywords in symbol table in advance
- ② create state transition diagram for each keyword separately

For 'then' keyword in pascal:



Conversion of NFA to DFA:

[Subset Construction Algorithm]

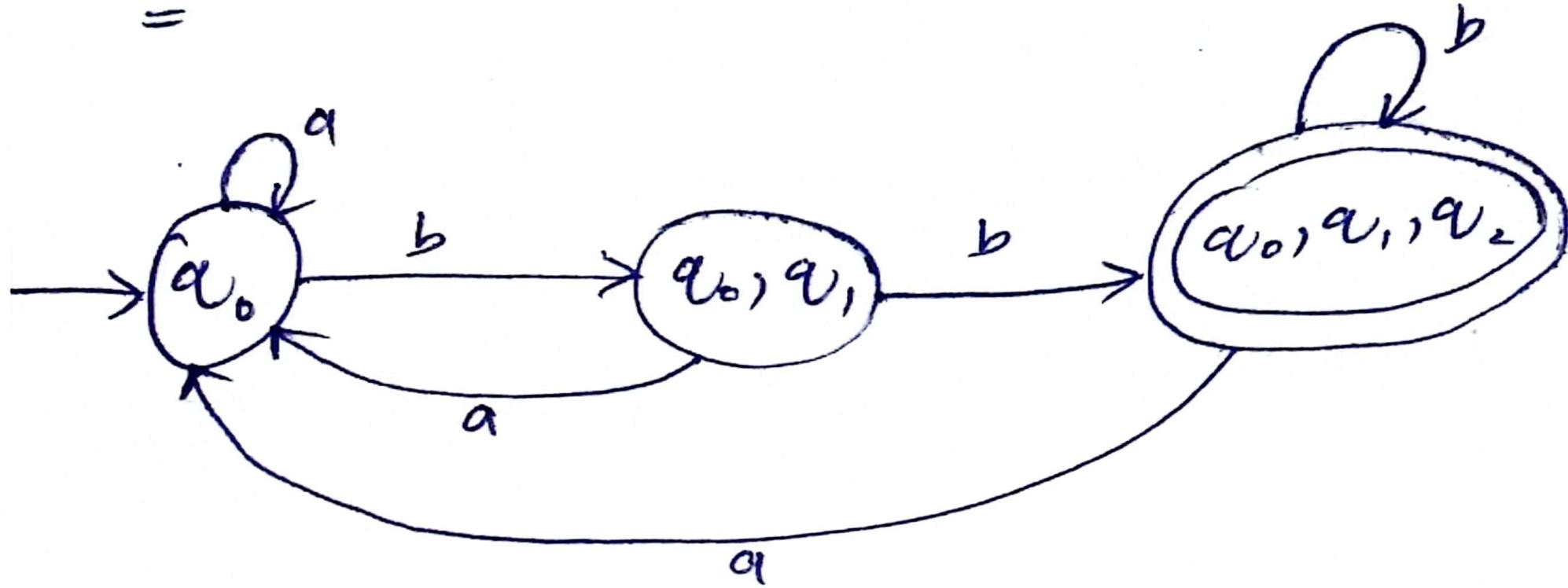


Transition Table:

State	a	b	
q_0	q_0	q_0, q_1	create q_0, q_1 as new state
q_1	0	q_2	
q_2	0	0	

state	a	b	
q_0	q_0	$\{q_0, q_1\}$	select final state having ' q_2 '
$\{q_0, q_1\}$	q_0	$\{q_0, q_1, q_2\}$	
$\{q_0, q_1, q_2\}$	q_0	$\{q_0, q_1, q_2\}$	

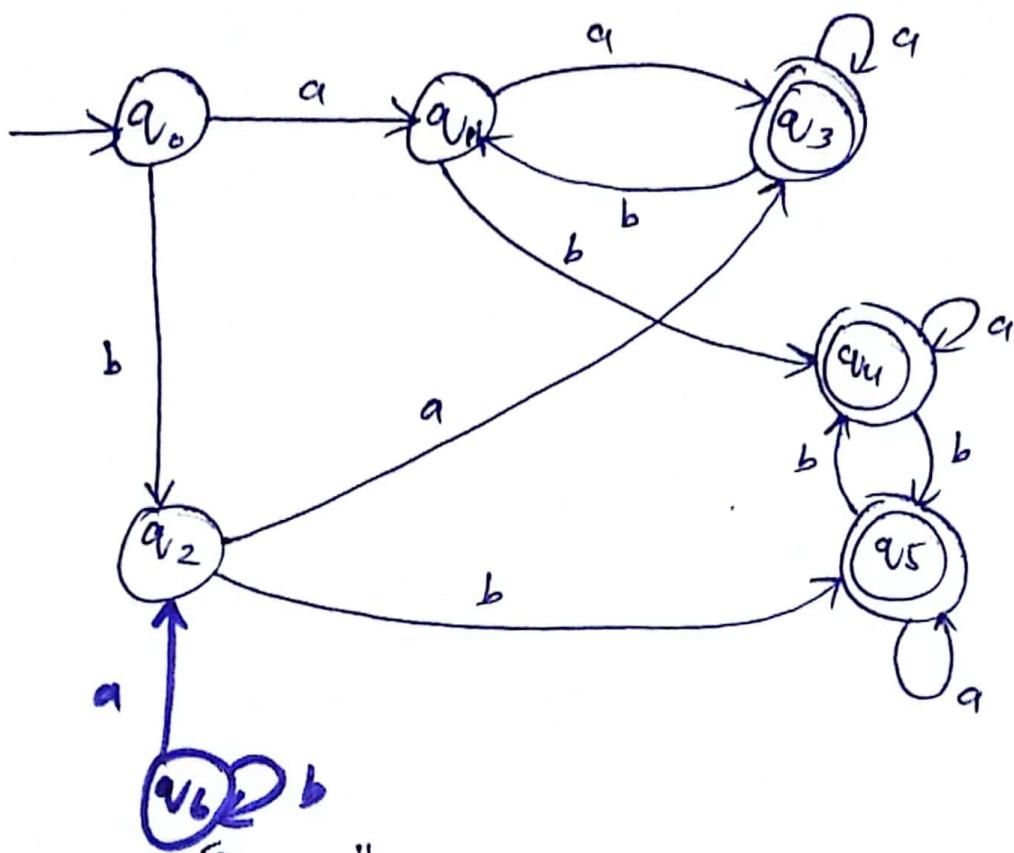
DFA:



DFA minimization:

① Remove

- Dead states \rightarrow termination in q_6 state will
- Unreachable states
 - ↓
unable to go from initial state q_0 to that specific state



$\rightarrow q_6$ is eliminated

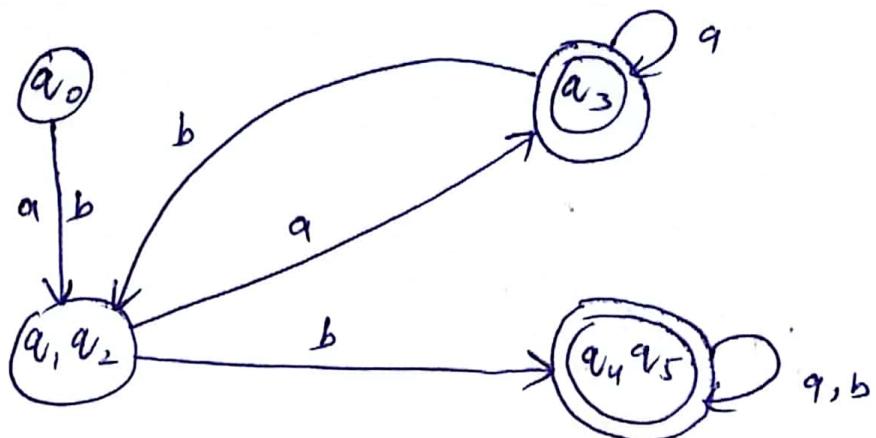
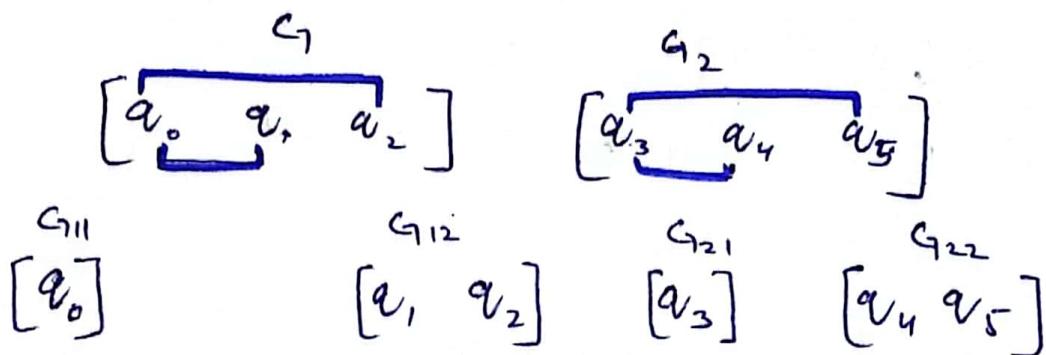
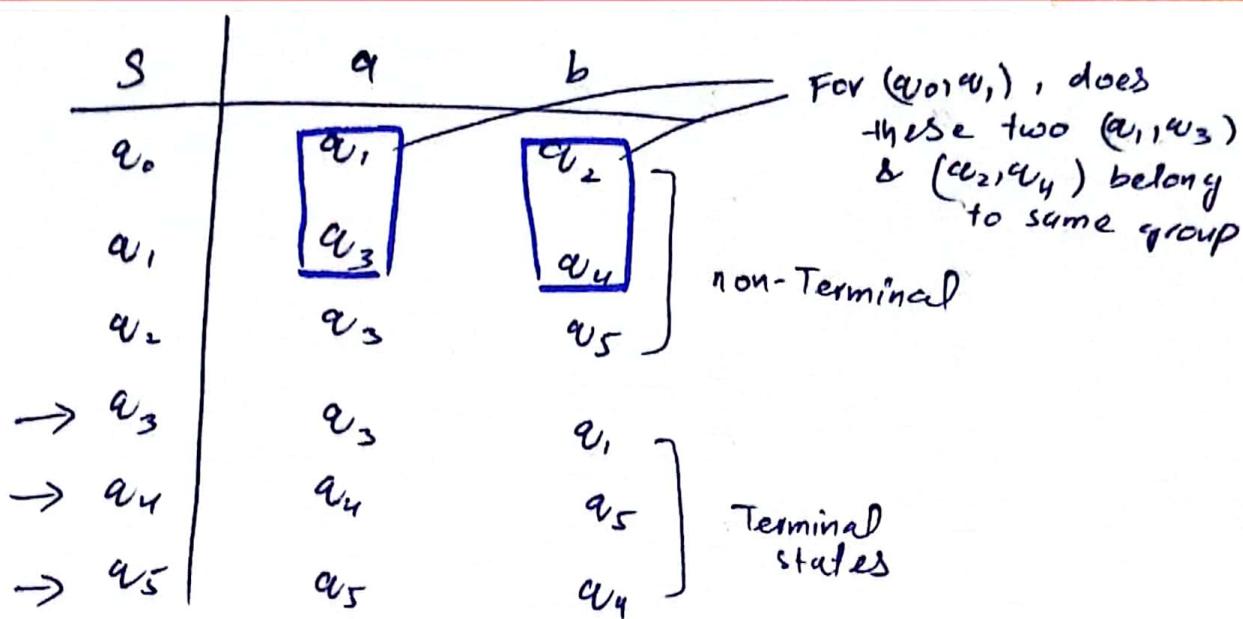
② Create group of terminal and non-terminal states

$$\{q_0, q_1, q_2\}, \{q_3, q_4, q_5\}$$

\rightarrow create transition table

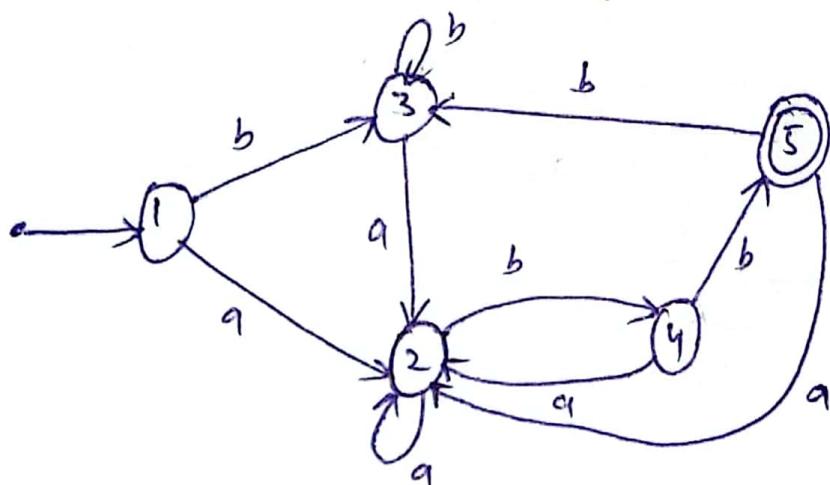
\rightarrow pick two states from one group & check their destination belong to same group, then pick next two otherwise split/divide into two group

③ Pick q_0, q_1



practise:

Minimize -the states of given DFA:



Transition Table

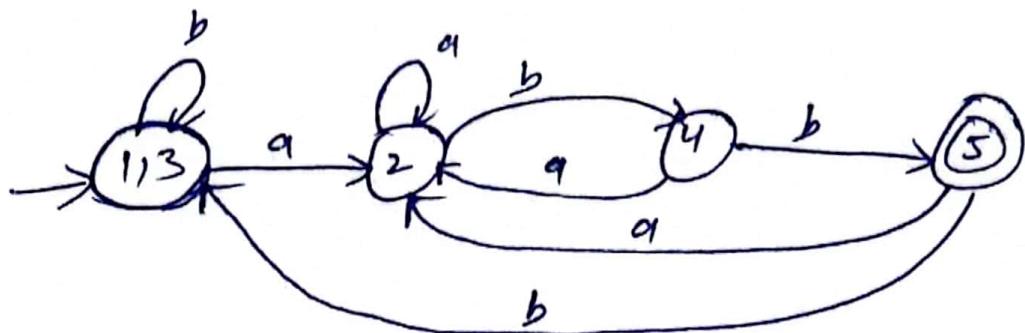
States	a	b
1	2	3
2	2	4
3	2	3
4	2	5
5	2	3

$$[\boxed{2 \ 3 \ 4}] \quad [5]$$

1 4 2

$\rightarrow (1, 2, 3) (4) (5)$

$\rightarrow (1, 3) (2) (4) (5)$



Lex:

- Tool used to create scanner

Syntax:

%{

// declarations &
preprocessor directives

%}

// regular expressions

-/-/

// methods

to %

int main() {

}

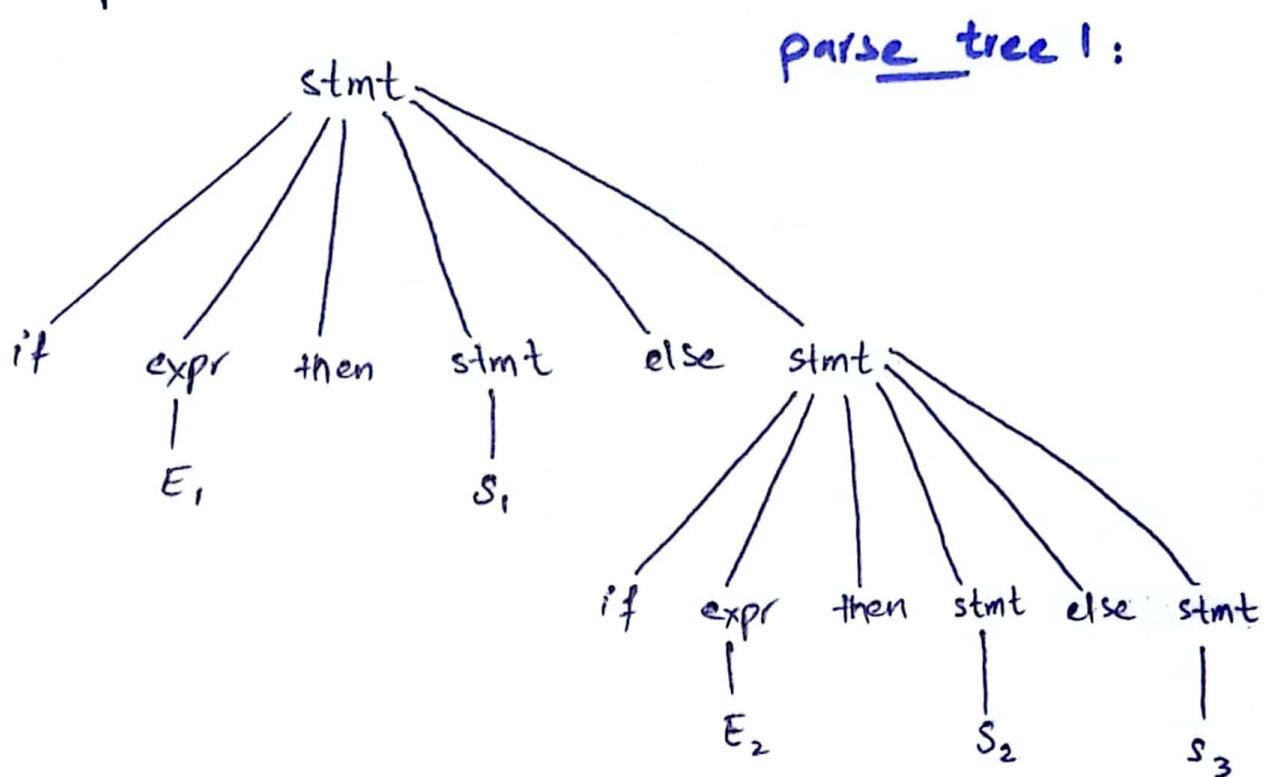
- We just need to tell the rules & states, tokens will be generated automatically

Syntax Analysis (parsing):

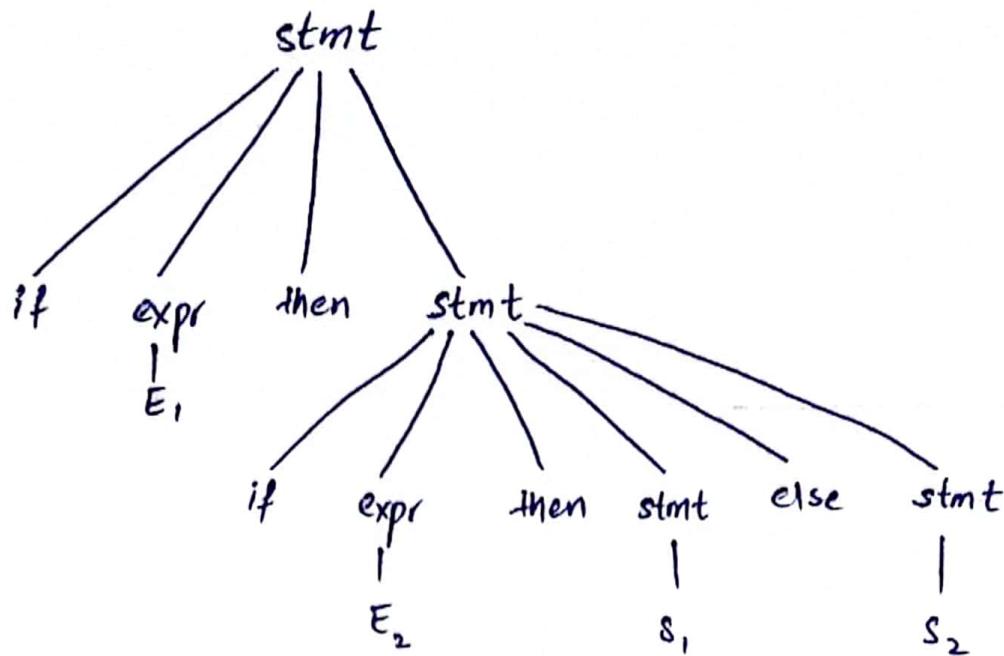
- This phase takes stream of tokens from lexical analyzer and generate parse tree via CFG
- The main issue in defining language structure is that CFG must be unambiguous (There must be only one possible parse tree for each input string)

For example:

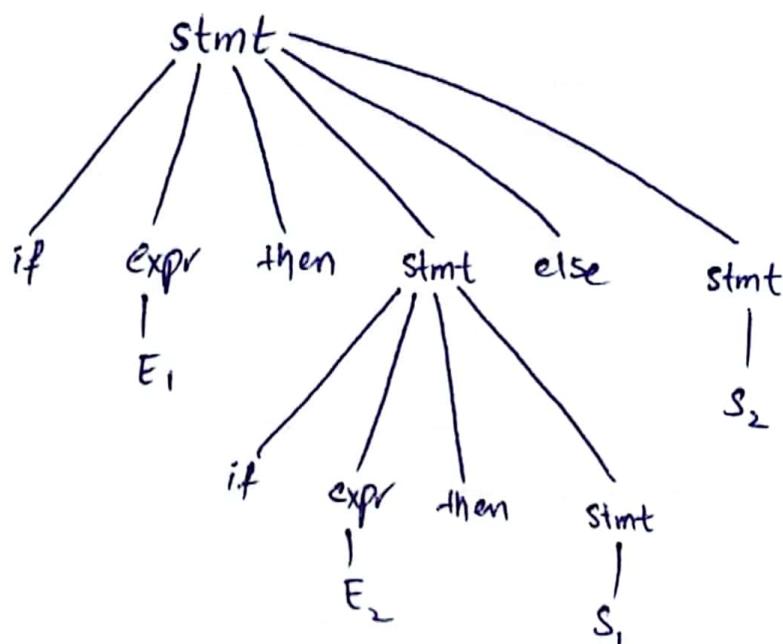
- $\text{stmtl} \rightarrow \text{if expr then stmt} \mid \text{if expr then stmt else stmt} \mid \text{other}$
- input = if E_1 then S_1 else if E_2 then S_2 else S_3
- possible parse trees:



parse tree 1: if E_1 then if E_2 then S_1 , else S_2



parse Tree 2:



"Ambiguous Grammar"

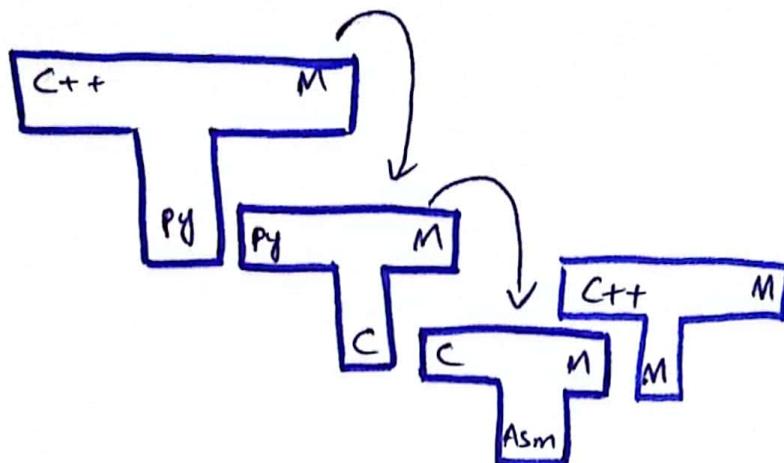
Types of parsing

→ Top ~~Down~~ Down Construct parse tree from root to leaf

→ Bottom up " leaf to root or ~ Top Down

Bootstrapping :

Language in which compiler is written is not understandable by target machine, we need to write another compiler which translates its instructions in target machine



Left & Right Recursion:

- ① Left Recursion: Nonterminal before ' \rightarrow ' is present after ' \rightarrow ' in 1st position.

For example:

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

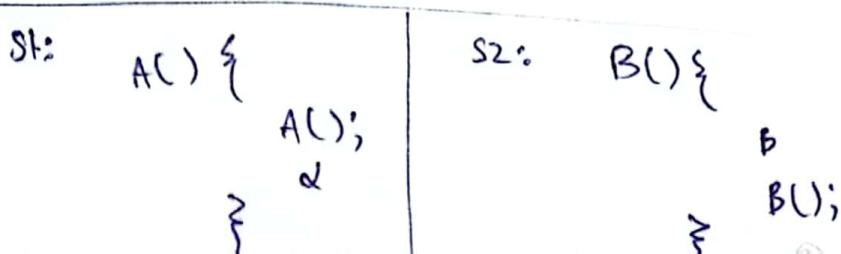
- ② Right Recursion: opposite of left recursion

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

\rightarrow Left Recursion need to remove from CFG

Why not Right recursion?

Consider two scenarios of recursion in programming



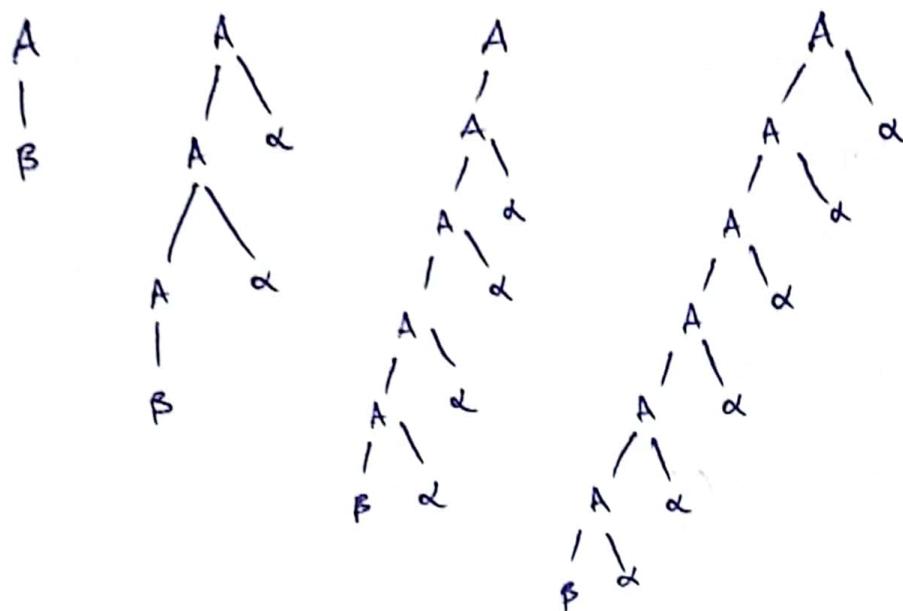
gl will have issue that no pre-condition will be checked due to left recursion while other case is perfect and has no major issues.

That's why left recursion must be removed from CFG

Example:

$$A \rightarrow A\alpha | \beta$$

parse Trees!



possible values = { β , $\beta\alpha$, $\beta\alpha\kappa$, $\beta\alpha\kappa\alpha$, $\beta\alpha\kappa\alpha\dots$, }

$$RE = \beta \alpha^*$$

Rewriting CFG:

$$A \rightarrow B A'$$

$$A' \rightarrow \alpha A')/\epsilon$$

We can map this example on all other cases by just recognizing α, β, A & A'

Example:

① Remove left recursion from given CFG

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

Sol:

$$\frac{E}{A} \rightarrow \frac{E + T}{A} \mid \frac{T}{B}$$

$$E = A, +T = \alpha, T = B$$

$$\text{Applying } A \rightarrow A\alpha \mid B$$

$$E \rightarrow TE'$$

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

② Remove LR from CFG

$$\frac{S}{A} \rightarrow \frac{SOSIS}{A} \mid \frac{O}{B}$$

$$S \rightarrow OS'$$

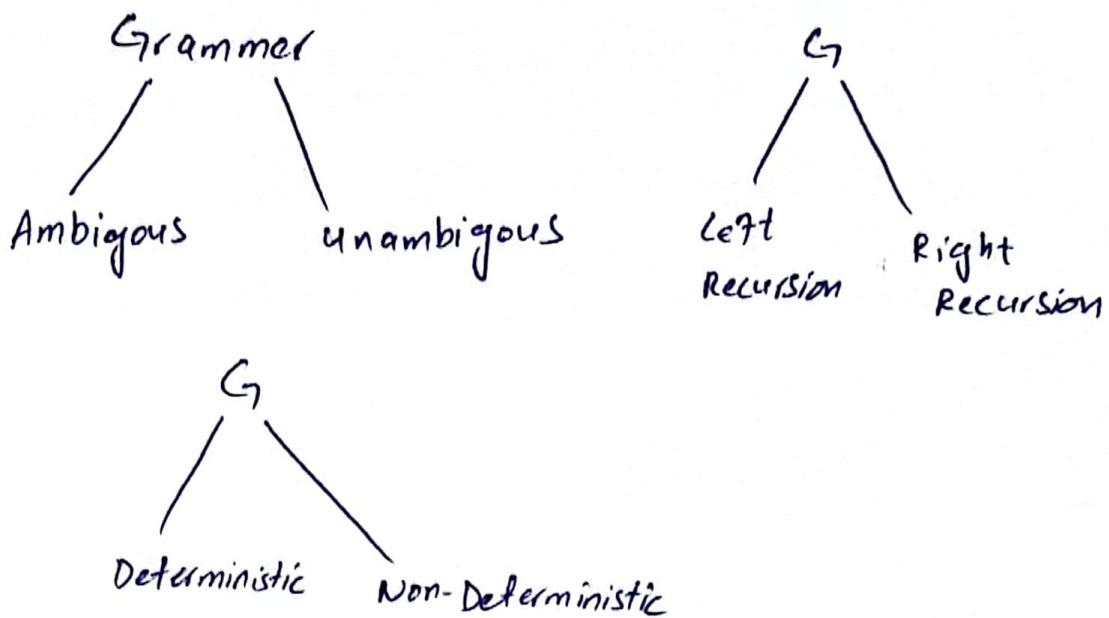
$$S' \rightarrow OSSI S' \mid \epsilon$$

③ Identify LR exist?

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

b) $S \rightarrow SL \mid \alpha$

c) $L \rightarrow L, S/S$



Deterministic Grammars: (Left factoring)

- If multiple opportunities not available for a single input in CFG, it is deterministic otherwise non-deterministic

Example:

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \alpha \beta_3$$

lets say compiler got $\alpha \beta_3$, it will ask for production rule related to $\alpha \beta_3$. CFG will return $\alpha \beta_1$ then $\alpha \beta_2$ next time. The main problem is not that we have to traverse the rules, but we have to compare α' everytime when we got rule from CFG which is redundant.

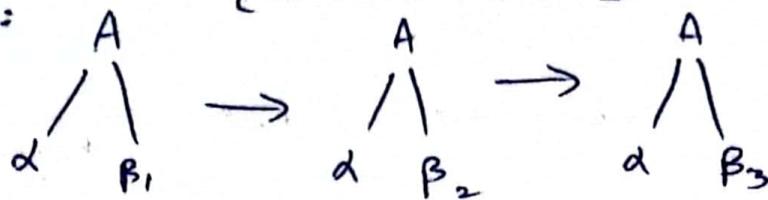
we need to put determinism in CFG by taking common as a separate rule:

$$\begin{aligned} A &\rightarrow \alpha A' \\ A' &\rightarrow \beta_1 \mid \beta_2 \mid \beta_3 \end{aligned}$$

In this case, we introduced new rule A' by taking ' α' as common.

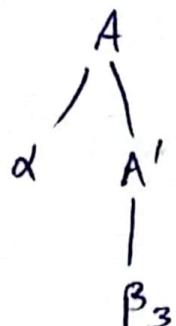
Comparison of parse tree for both cases:

Case#1:



Case#102:

[Deterministic]



Example 2:

Convert CFG to deterministic one:

$$S \rightarrow iEtS / iEtSeS / a$$

Sol

$$\text{Common} = 'EtS' \Rightarrow S'$$

$$S \rightarrow iEtSS' / a$$

$$S' \rightarrow \epsilon / eS$$

Example 3:

$$S \rightarrow aSSbS / aSaqSb / abb / b$$

Sol

$$S \rightarrow aS' / b$$

$$S' \rightarrow SSbS / SaSb / bb$$

still require determinism

$$S \rightarrow aS' / b$$

$$S' \rightarrow SS'' / bb$$

$$S'' \rightarrow SbS / aSb$$

Why searching is not used as compare to parsing?

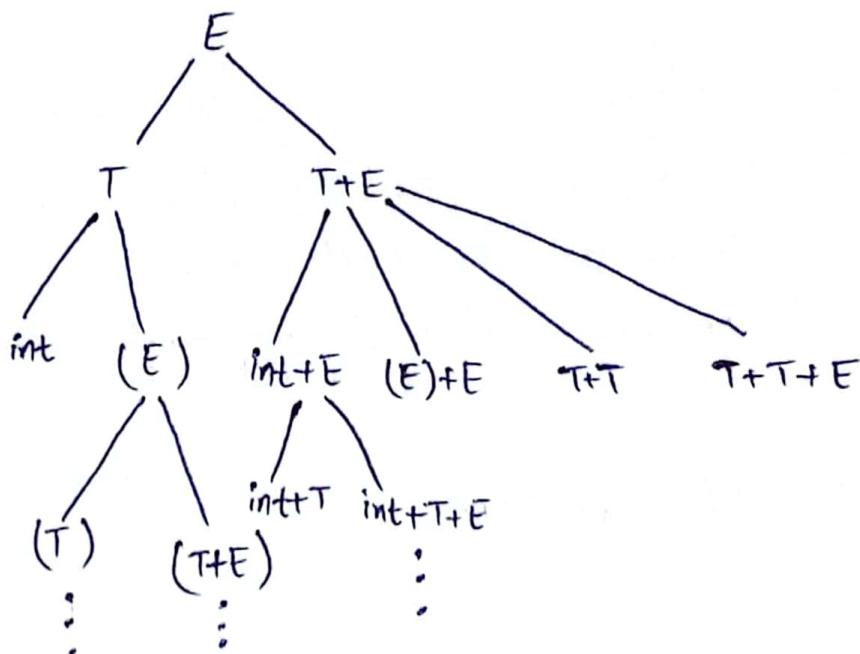
Searching a string from CFG requires finding the appropriate rule each time. As # of production rules increases, # of child increases. What if string not exist in CFG rules, wastage of time.

Example:

Let's use Breadth first Search (**BFS**) for searching:

CFG:

$$\begin{aligned} E &\rightarrow T \mid T+E \\ T &\rightarrow \text{int} \mid (E) \end{aligned}$$



BFS = Level wise from left to right traversal

lets say input string is 'int+int' , we need to find all paths from tree till int+int found (if exist according to rule).

That's why searching via BFS is not used.

Branching factor = max # of child allowed for tree

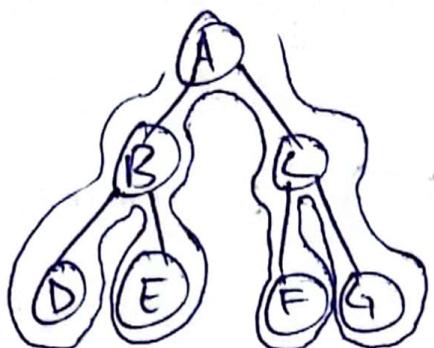
Left most BFS:

Only expand left most child and compare single character with input string whether to expand more or pick another node.

Page 25 of 46

Still require lot of search to reach valid sequence.

DFS



Activity Diagram Generate Reports

It is in favour of leftmost but it can cause problem in case of graphs and leads to infinite search.

We need algorithm which guide exact rules which need to apply to get final result.

Predictive Parsing is used to get exact rules and provides efficient approach towards parsing.

It uses LL(1) parser.

LL(1)

- left to right
- left most derivation
- one token of lookahead



Generate Report

CFG :

$$E \rightarrow \text{int} \mid (E \text{ op } E)$$

$$\text{op} \rightarrow + \mid *$$

parse Table

		terminals			
		int	()	+
		int	$(E \text{ op } E)$ <th></th> <th>*</th>		*
E					
op				+	*

rules
if 'int'
matched if '('
 matched

Guideline:

- pick up tokens and check corresponding rule from parse table, if not found →
- if '\$' matched at end, accept
- remove matched character one by one

Example:

$(\text{int} + (\text{int} * \text{int}))$

put '\$' at end of rules & input string

Rules

$E \$$

'(' $(E \text{ op } E) \$$

'int' $\text{int op } E \$$

'+' ~~$\text{int} + E \$$~~

'(' $(E \text{ op } E) \$$

'int' $\text{int op } E \$$

'*' $\text{op } E \$$

$* E \$$

$E \$$

String $(\text{int} + (\text{int} * \text{int})) \$$

$\uparrow \text{int} + (\text{int} * \text{int}) \$$

$\uparrow + (\text{int} * \text{int}) \$$

$\uparrow (\text{int} * \text{int}) \$$

$\uparrow \text{int} * \text{int}) \$$

$\uparrow \text{int} * \text{int}) \$$

$\uparrow * \text{int}) \$$

$\uparrow * \text{int}) \$$

$\uparrow \text{int}) \$$

'int' int)) \$ int)) \$

"matched"

How to generate "parse table"? 94.

put terminals on top row and non-terminals on left most column.

Issue in L(1) parser is that we are looking only first token, what if 1st token is itself non-terminal which leads toward tracing until terminal is found.

For example:

stmt \rightarrow if expr then stmt ①

| while expr do stmt ②

| expr; ③ ④

expr \rightarrow Term \rightarrow id | zero? term ⑤

| not expr ⑥ + id ⑦ - id ⑧

term \rightarrow id ⑨ | constant ⑩

Rule (3) and (4) have 'expr' &

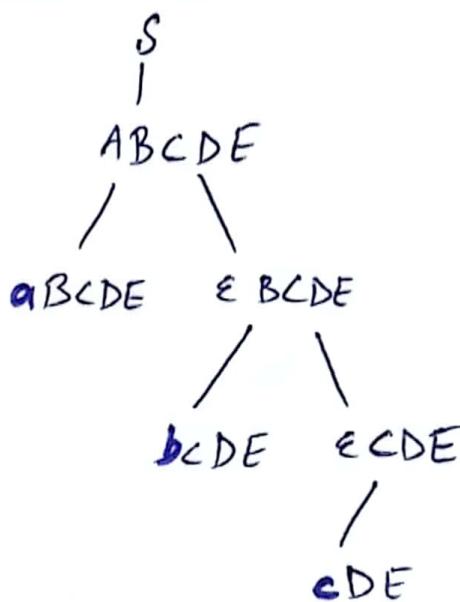
'term' as 1st token which are not terminal.

For creating 'parse table', these entries can be easily applied using `first()` and `follow()` set.

first() set:

- place all the first non-terminals in a set

Rules	First set	follow()
$A \rightarrow a/\epsilon$	$\{a, \epsilon\}$	$\{b, c\}$
$S \rightarrow ABCDE$	$\{a, b, c\}$	$\{\$ \}$
$B \rightarrow b/\epsilon$	$\{b, \epsilon\}$	$\{c\}$
$C \rightarrow c$	$\{c\}$	$\{d, e, \$\}$
$D \rightarrow d/\epsilon$	$\{d, \epsilon\}$	$\{b\}$ $\{e, \$\}$
$E \rightarrow e/\epsilon$	$\{e, \epsilon\}$	$\{\$\}$



Example 2:

Rules	first() set	follow()
$S \rightarrow Bb/cd$	$\{a, b, c, d\}$	$\{\$\}$
$B \rightarrow aB/\epsilon$	$\{a, \epsilon\}$	$\{b\}$
$C \rightarrow CC/\epsilon$	$\{c, \epsilon\}$	$\{d\}$

For follow set:

→ for starting ~~rule~~ rule, it will be always $\{\$\}$

→ if no follower, then put all the follow set elements of leftmost (before arrow)

For first(): put all the first terminals in a separate set
it non-terminal when trace till terminal found

Example 3:

Find first() and follow() set of given grammar:

23 of 46

Activity Diagram Login

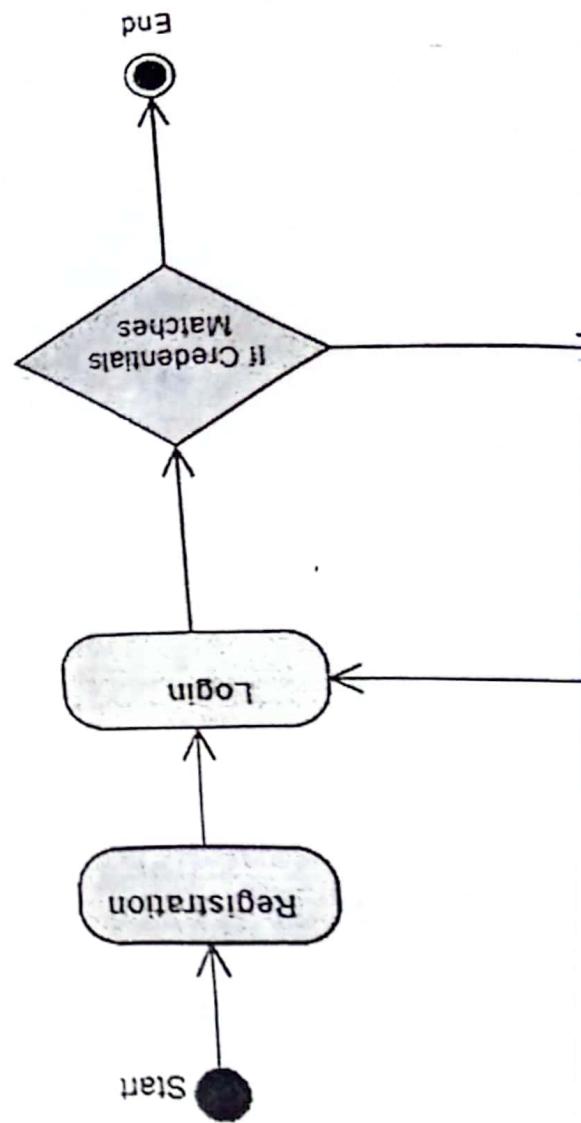
- 1) $E \rightarrow TE'$
- 2) $E' \rightarrow +TE'/\epsilon$
- 3) $T \rightarrow FT'$
- 4) $T' \rightarrow *FT'/\epsilon$
- 5) $F \rightarrow id \mid (E)$

→ First():

- 1) {id, ()}
- 2) {+, ε}
- 3) {id, ()}
- 4) {* , ε}
- 5) {id, ()}

→ Follow():

- 1) { \$, }¹
- 2) { \$,) }²
- 3) { +, \$,) }³
- 4) { +, \$,) }⁴
- 5) { *, +, \$,) }⁵



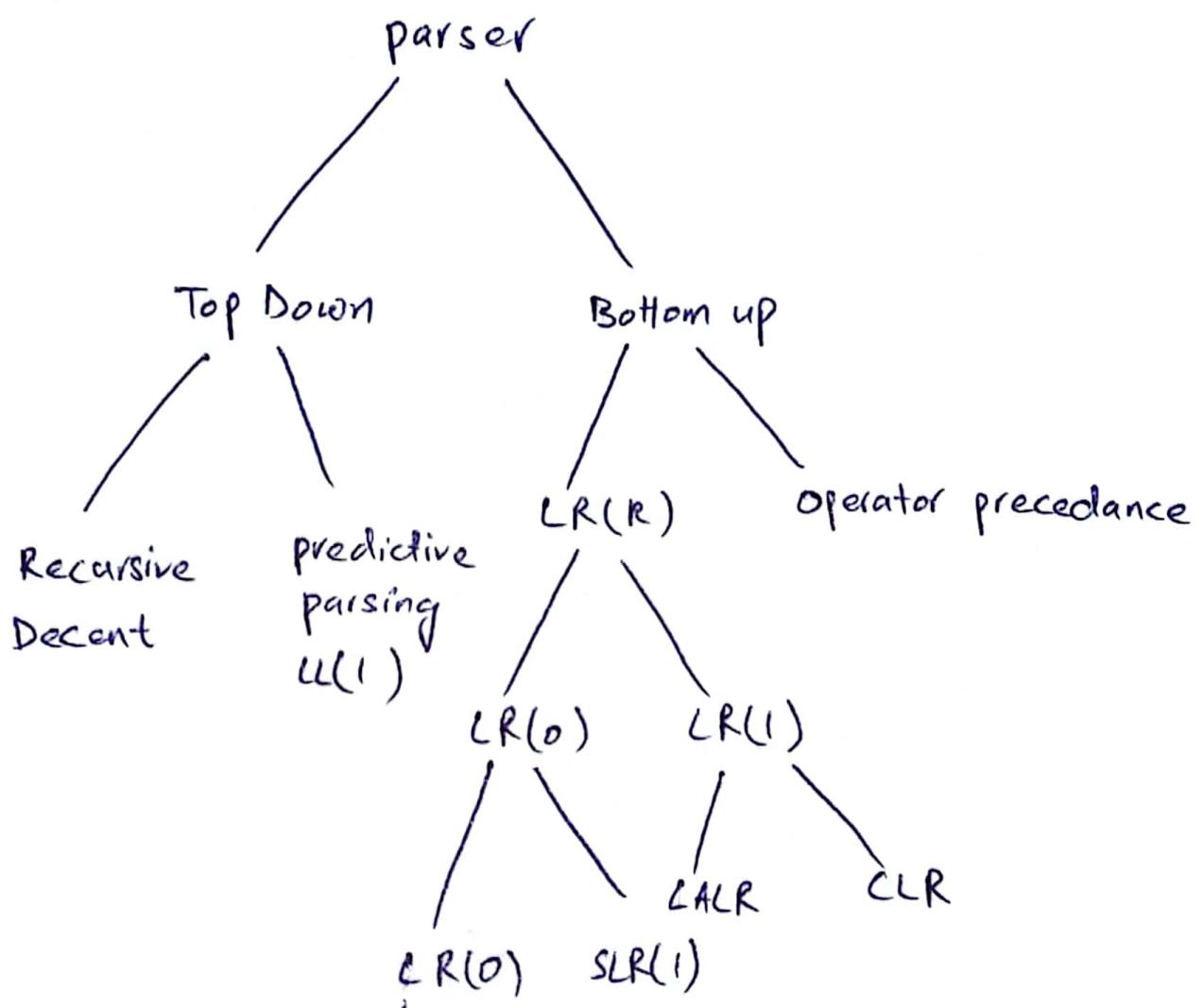
$S \rightarrow ABC \mid ghi \mid jkl$	$\text{first}(S) = \{a, g, j\}$	$\text{follow}(S) = \{\$\}$
$A \rightarrow a \mid b \mid c$	$\{a, b, c\}$	$\{b\}$
$B \rightarrow b$	$\{b\}$	$\{c\}$
$D \rightarrow d$	$\{d\}$	$\{\}$
$C \rightarrow c$	$\{c\}$	$\{\$\}$

② $S \rightarrow ABC \quad \{\$\}$

$A \rightarrow DEF \quad \{\$\}$

$B, C, D, E, F \rightarrow \epsilon \quad \{\$\}$

Types of parsing:



How to check Grammer is LL(1)?

- 1) If epsilon ' ϵ ' not in a rule, find ~~intersection~~ of first() set, if equal to ϕ , LL(1) otherwise not
- 2) if epsilon ' ϵ ' exist in rule, also find follow() set of that non-terminal and take ~~intersection~~ with first() of others

Example:

→ Case 1:

$$S \rightarrow \alpha_1 \quad | \quad \alpha_2 \quad | \quad \alpha_3$$

$$\text{first}(\alpha_1) \cap \text{first}(\alpha_2) \cap \text{first}(\alpha_3) = \emptyset$$

→ Case 2:

$$S \rightarrow \alpha_1 | \alpha_2 | \alpha_3 | \varepsilon$$

$$\text{first}(\alpha_1) \cap \text{first}(\alpha_2) \cap \text{first}(\alpha_3) \cap \text{follow}(S) = \emptyset$$

$s \rightarrow asa \mid bsb \mid c$

is LCL(1) or not?

`first(asq) n first(bs) n first(c)`

$$a \wedge b \wedge c = \emptyset$$

Yes, g1s LL(1) Grammer

$s \rightarrow ictss_i | a \quad \text{if } \text{first}(ictss_i) \cap \text{first}(a) \\ i \neq a = \emptyset$

$$S_1 \rightarrow eS | \epsilon \quad \text{② } \text{first}(eS) \cap \text{follow}(S_1) \\ C \rightarrow b \quad e \cap \$e \neq \emptyset$$

gives not LL(1) parsing Grammer!

Apply LL(1) parsing:

① $S \rightarrow aABb$

② $A \rightarrow c|\epsilon$

③ $B \rightarrow d|\epsilon$

	First()	follow()
1)	{a}	{\$}
2)	{c, ε}	{a, b}
3)	{d, ε}	{b}

parsing Table:

	a	b	c	d	\$
S	$S \rightarrow aABb$				
A		$A \rightarrow \epsilon$	$A \rightarrow C$	$A \rightarrow \epsilon$	
B		$B \rightarrow \epsilon$		$B \rightarrow d$	

if ' ϵ ' in rule, then put follow() set entries the rule

For example: in 2nd rule $A \rightarrow C\epsilon$, epsilon exist, so check follow set which are $\{d, b\}$. put the rule $A \rightarrow \epsilon$ in 'b' and 'd' column in parsing table

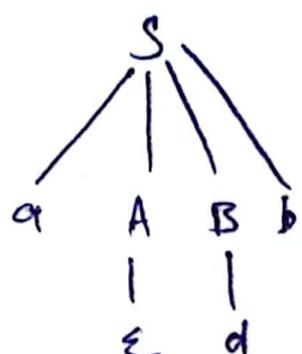
i/p buffer	a	o	b	\$
	\downarrow	\downarrow	\downarrow	

o
ϵ
a
C
B
ϵ
g
\$

if \$ symbol at the end in the stack, accept the string.

parse tree will be generated

side by side upon pushing in stack.



Example:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow id \end{aligned}$$

Is LL(1) parsing possible?

Try to create parse table, if more than one entry in a table \Rightarrow box exist, then LL(1) parsing not possible.

As in the given example, there is no ' ϵ ', so just find first() set.

$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow id \end{aligned}$	first() $\begin{cases} \{id\} \\ \{id\} \\ \{id\} \end{cases}$
---	--

Since Left Recursion exist, LL(1)^{not} exist.

(2)

<u>S</u>	first()	follow()
$A \rightarrow Bb \mid Cd$	$\{\underline{a}, b, c, d\}$	$\{\$ \}$
$B \rightarrow aB \mid \epsilon$	$\{\underline{a}, \epsilon\}$	$\{b\}$
$C \rightarrow CC \mid \epsilon$	$\{\underline{c}, \epsilon\}$	$\{a\}$

Check intersection of first() if no ' ϵ ' otherwise also take intersection with follow() excluding ' ϵ '

For previous example

Rule 1 $\{\underline{a}, b\} \cap \{\underline{c}, d\} = \emptyset$ because no ' ϵ ', just pick first() set

Rule 2 $\{\underline{a}\} \cap \{b\} = \emptyset$ because ' ϵ ' exist, pick first & follow

Rule 3 $\{\underline{c}\} \cap \{d\} = \emptyset$ same as rule 2

So, LL(1) possible!

Operator precedence parser:

- Bottom up parser
- CFG must have operator b/w non-terminals (operator grammar)

Example:

① $E \rightarrow E + E \mid E * E \mid id$

It's operator grammar because non-terminals are separated by +, * operator

② $E \rightarrow E op E \mid id$

$op \rightarrow + \mid - \mid *$

It's not operator grammar due to 'op' non-terminal between two non-terminals.

Example:

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

Apply operator precedence parsing. (operator relational table)

	id	+	*	\$	
id	-	>	>	\geq	precedence of 'id' is greater than '*'
+	<	>	<	>	
*	<	>	>	>	identifier has highest precedence
\$	<	\leq	<	accept	$id > * \mid \div > + \mid -$

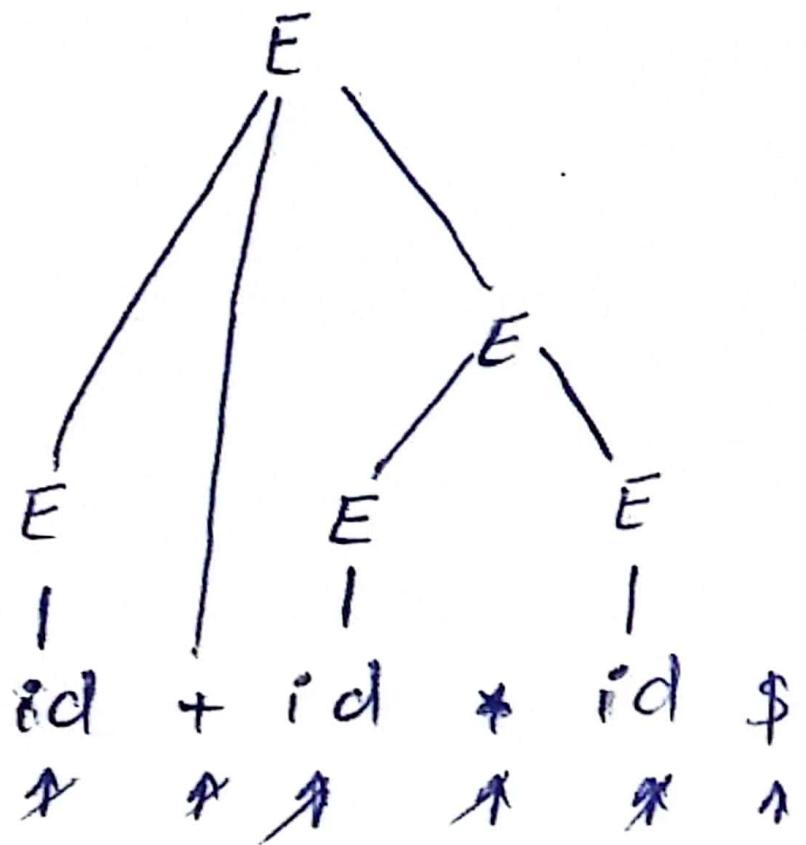
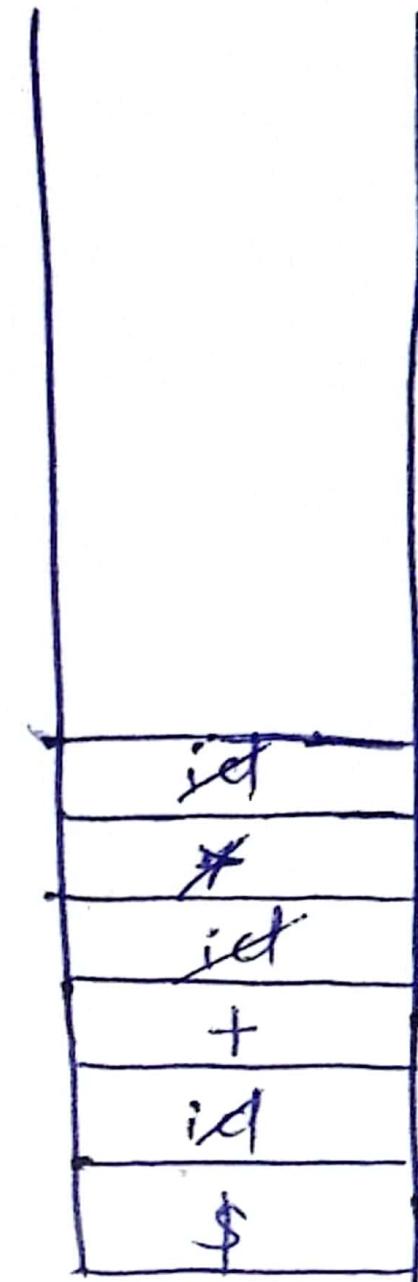
precedence of leftmost '*' is > than rightmost '*'.

② Create stack and

push = element of current pointer has \geq precedence than stack element

pop = element in the stack has high precedence

Stack



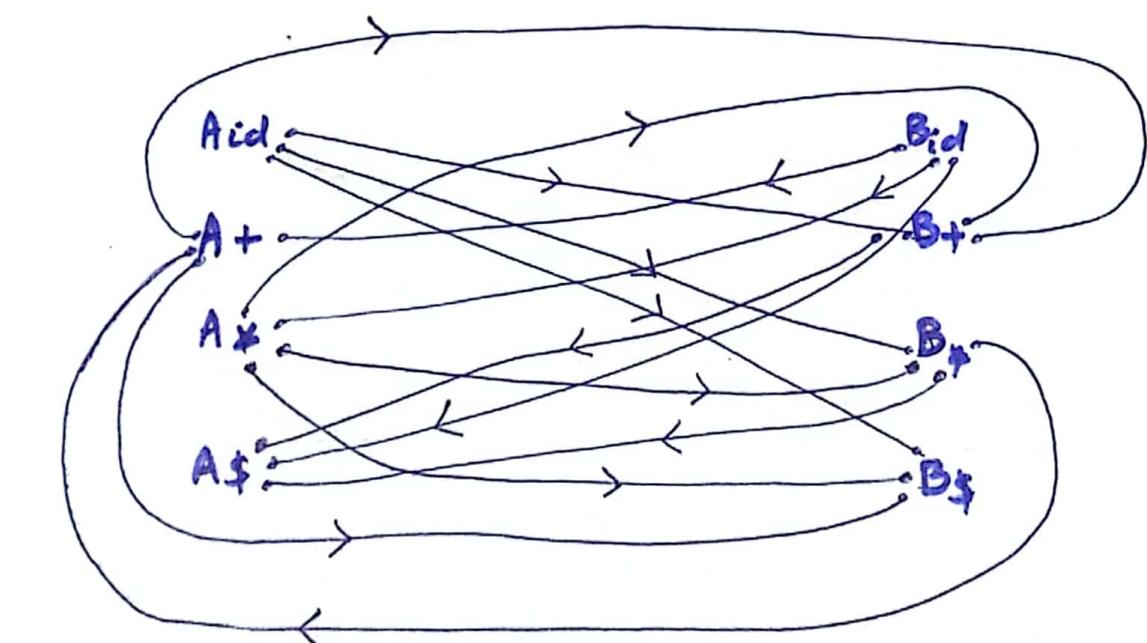
→ Compare element in the pointer
with element in the stack
& decide to push or
pop

In operator precedence parser, for 'n' terminals table size will be $n \times n$ which will consume more space if grammar size is large.

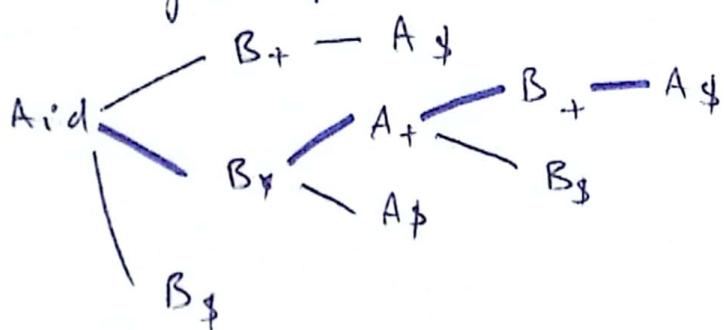
We need some mechanism to reduce table size.

This parser also accepts ambiguous grammar.

	id	$\leftarrow B$	$\rightarrow B$	$+$	$*$	$\$$	
id	-	>	>	>			
$+$	<	>	<	>			
$*$	<	>	>	>			
$\$$	<	<	<			Accept	



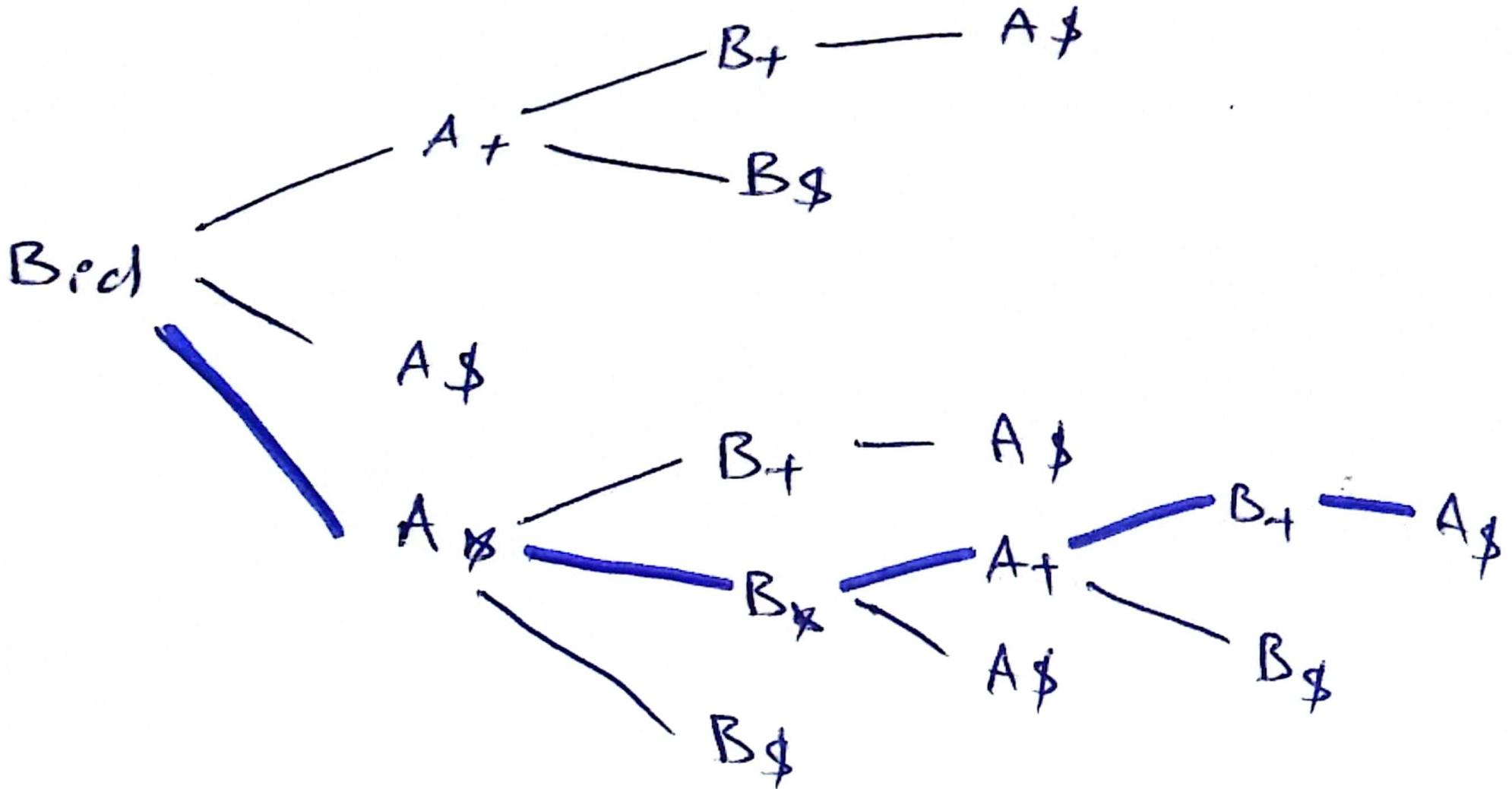
Find longest path



put $A_{\text{id}} \rightarrow B_{*} \rightarrow A_{+} \rightarrow B_{+} \rightarrow A_{\$}$ is

missing 3 terminals $A_{*}, B_{*}, B_{\text{id}}$

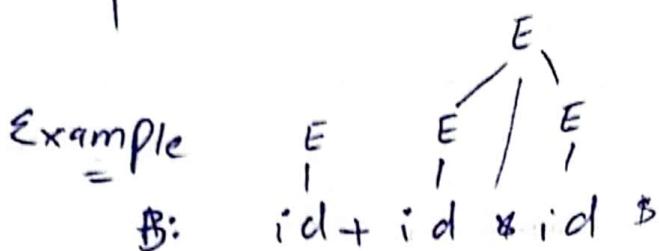
= B_{id} now!



$A_{id} \rightarrow B_x \rightarrow A_+ \rightarrow B_+ \rightarrow A_\$$

$B_{id} \rightarrow A_x \rightarrow B_x \rightarrow A_+ \rightarrow B_+ \rightarrow A_\$$

	id	*	+	\$	
A	4	4	2	0	"operator functional Table"
B	5	3	1	0	



A	\$	id	+	id	*	id	\$
Stack							

- | | |
|---|---|
| 1) $A_\# > or < B_{id}$
$4 < 5$ "push" | 7) $A_x < > B_{id}$
$4 < 5$ "push" |
| 2) $A_{id} > or < B_+$
$4 > 1$ "pop" | 8) $A_{id} < > B_\$$
$4 > 0$ "pop" |
| 3) $A\$ < or > B_+$
$0 < 1$ "push" | 9) $A_x < > B_\$$
$4 > 0$ "pop" |
| 4) $A_+ < > B_{id}$
$2 < 5$ "push" | 10) $A_+ < > B_\$$
$2 > 0$ "pop" |
| 5) $A_{id} < > B_x$
$4 > 3$ "pop" | 11) $A_\# = B\$$
"accept" |
| 6) $A_+ < > B_x$
$2 < 3$ "push" | |

precedence & associativity in CFG:

Some times, CFG violates precedence or associativity property which leads toward invalid results.

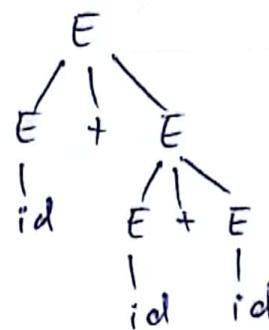
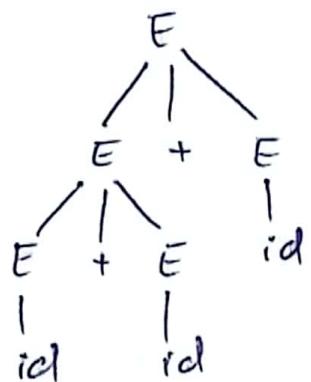
For example :

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

Let say we have input string

$id + id + id$

parse Tree:



Ambiguity exist in this case, one tree favours right associativity $(id+id)+id$ and other one left associativity $id+(id+id)$.

To make right associative, we can update grammar by making it right recursive. In the same way, for left associative, make grammar left recursive. Higher/deeper the level of operator, highest the precedence will be.

As in this case, we need ~~right~~ left associative, convert grammar to left recursive one.

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

$b \rightarrow b \text{ and } b$

| $b \text{ or } b$

| $\text{not } b$

| true

| false

$b \rightarrow b \text{ and } c \mid c$

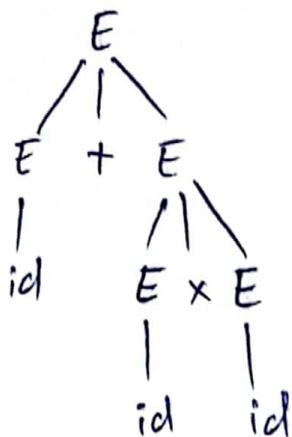
$c \rightarrow c \text{ or } d \mid d$

$d \rightarrow \text{true} \mid \text{false} \mid \text{not } d$

Example 2: (precedence violation)

$$E \rightarrow E+E \mid E \times E \mid id$$

input string = id + id * id

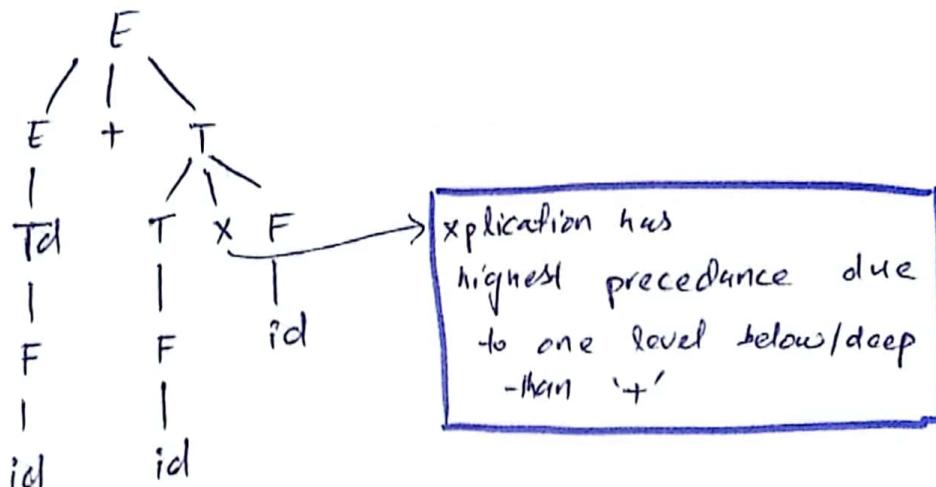


Making grammar left recursive

$$\begin{aligned} E &\rightarrow E+T \mid T \\ T &\rightarrow T\times F \mid F \\ F &\rightarrow id \end{aligned}$$

"right associative"

Again making parse tree on new grammar



Example 3:

Determine precedence of operators in given grammar

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

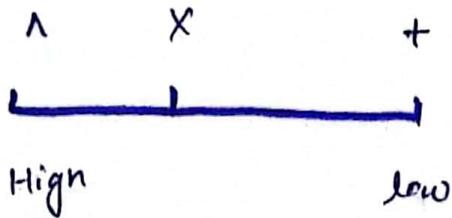
$$T \rightarrow T\times F \mid F$$

$$F \rightarrow G \wedge F \mid G$$

$$G \rightarrow id$$

Try to parse the target operator from start symbol, operator having longest depth will have highest precedence

precedence



Example 4:

Consider given ambiguous grammar:

$$\begin{aligned} bExp \rightarrow & bExp \text{ AND } bExp \mid bExp \text{ OR } bExp \mid \text{NOT } bExp \\ & \mid \text{True} \mid \text{False} \end{aligned}$$

precedence of operators NOT > AND > OR

Solution

$$bExp \rightarrow bExp \text{ OR } bExp^1 \mid bExp^1$$

$$bExp^1 \rightarrow bExp^1 \text{ AND } bExp^2 \mid bExp^2$$

$$bExp^2 \rightarrow \text{NOT } bExp^2 \mid \text{True} \mid \text{False}$$

Example 5:

$$R \rightarrow R + R \mid R \cdot R \mid R^* \mid a \mid b \mid c$$

precedence = $\star > \cdot$ (and) $> +$ (or)

(left to right associativity)

Solution

$$R \rightarrow R + R^1 \mid R^1$$

$$R^1 \rightarrow R^1 \cdot R^2 \mid R^2$$

$$R^2 \rightarrow R^2^* \mid a \mid b \mid c$$

Example 6:

$$S \rightarrow TxP$$

$$T \rightarrow U \mid TxU$$

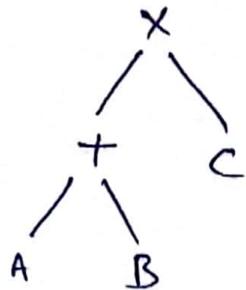
$$P \rightarrow Q + P \mid Q$$

$$Q \rightarrow id$$

$$U \rightarrow id$$

In this example,
'+' has highest precedence
than 'x'
+ = right associative
due to right recursion
x = left associative

Which of the following expression is presented by parse tree?



- a) $(A+B) \times C$
- b) $A + \times BC$
- c) $A + B \times C$
- d) $A \times C + B$

Example :

$BE \rightarrow BE \text{ and } BF \mid BE \text{ or } BF \mid (BE) \mid B$
 $BF \rightarrow \text{true} \mid \text{false}$

⇒ Rewrite CFG by applying rule
() > and > or => precedence

Solution :

$BE \rightarrow BE \text{ or } BF \mid B1$
 $B1 \rightarrow B1 \text{ and } B2 \mid B2$
 $B2 \rightarrow (B2) \text{ or } BF \mid BF$
 $BF \rightarrow \text{true} \mid \text{false}$

Example :

$E \rightarrow E+E \mid E-E \mid E \times E \mid E/E \mid \text{id} \mid \text{num}$

+ has highest precedence while other have same.

right to left associativity

Sol :
 $E \rightarrow E1-E \mid E1 \times E \mid E1/E \mid E1$
 $E1 \rightarrow E2+E \mid E2$
 $E2 \rightarrow \text{id} \mid \text{num}$
 $E2 \rightarrow E2+E \mid E2 \times E \mid E2/E \mid E2$
 $E2 \rightarrow \text{id} \mid \text{num}$

① $E \rightarrow E < E \mid E > E \mid E == E \mid E != E \mid \text{num}$

~~precedence~~ precedence of equality & inequality is $>$ than other

Sol:

$E \rightarrow E < E \mid E > E \mid E != E$

$E != \rightarrow E == E \mid E != E \mid E_2$

$E_2 \rightarrow \text{num}$

② $BE \rightarrow BE \text{ and } BE \mid BE \text{ or } BE \mid \text{not}(BE) \mid \text{true} \mid \text{false}$

not > and > or \Rightarrow precedence

a) $BE \rightarrow BE \text{ or } BE \nmid BE_1$

$BE_1 \rightarrow BE_1 \text{ and } BE_1 \mid BE_2$

$BE_2 \rightarrow \text{not } BE_2 \mid BE_3$

$BE_3 \rightarrow \text{true} \mid \text{false} \mid (BE)$

b) make ' $\&$ ' and ' $\#$ ' right associative

$BE \rightarrow BE \text{ or } BE_1 \mid BE_1$

$BE_1 \rightarrow BE_1 \text{ and } BE_2 \mid BE_2$

$BE_2 \rightarrow \text{not } BE_2 \mid BE_3$

$BE_3 \rightarrow \text{true} \mid \text{false} \mid (BE)$

③ $S \rightarrow S \& S \mid S = S \mid S \# S \mid \text{id} \mid \text{num}$

$\&, \# > \text{other} \Rightarrow$ precedence (left assoc)

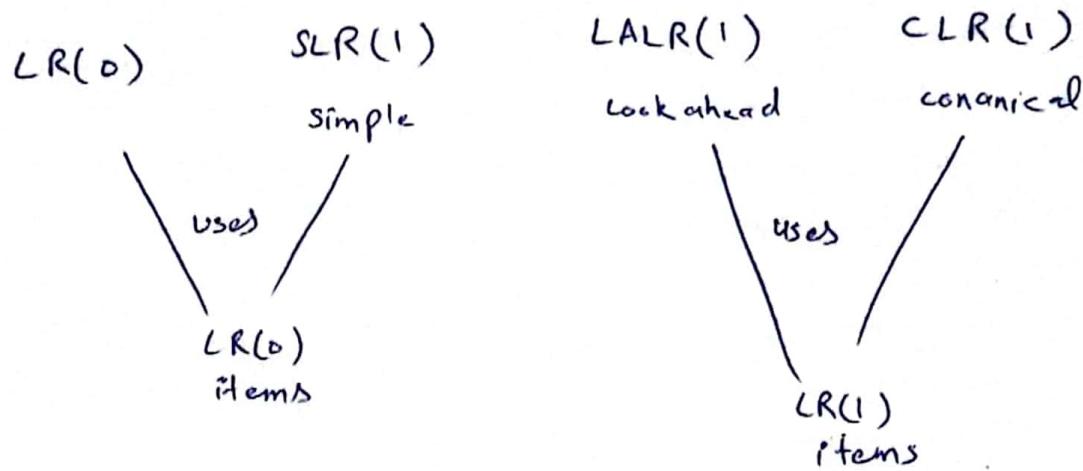
Sol

$S \rightarrow S = S_1 \mid S_1$

$S_1 \rightarrow S_1 \& S_2 \mid S_1 \# S_2 \mid S_2$

$S_2 \rightarrow \text{id} \mid \text{num}$

Bottom up parser:



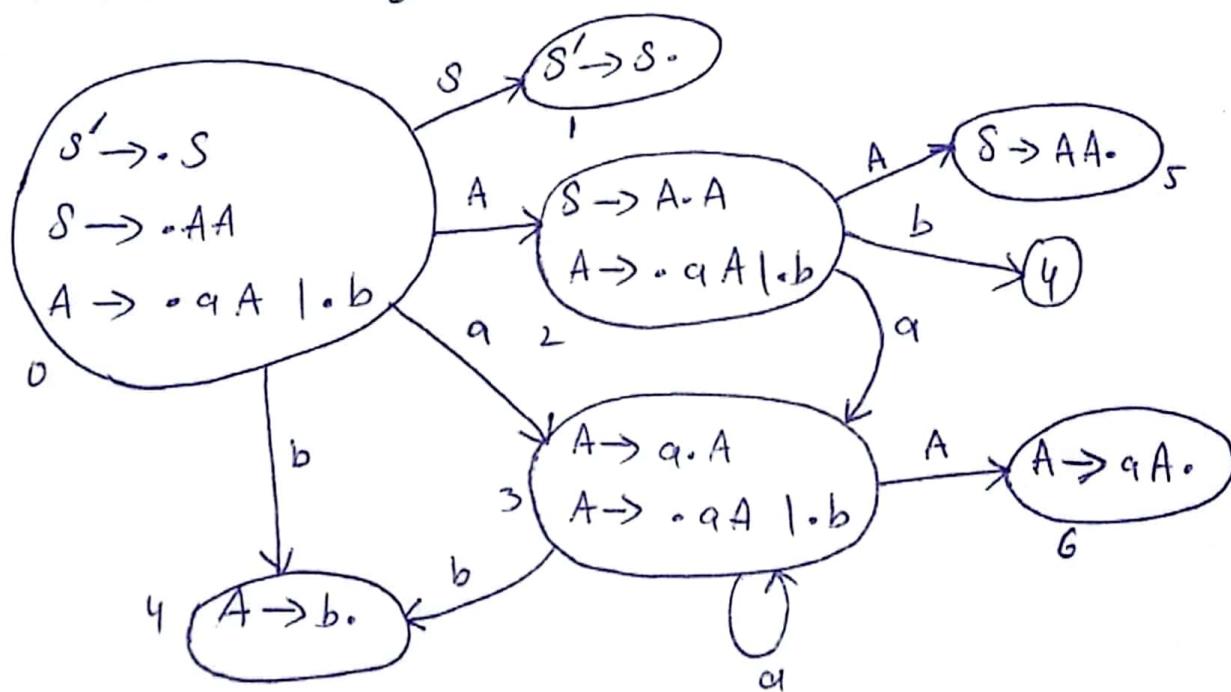
$LR(0)$ items:

-Used in $LR(0)$ & $SLR(1)$ with little modification

Example:

$$\begin{array}{ll}
 S' \rightarrow S & \text{rule 0} \\
 S \rightarrow A A & \text{rule 1} \\
 A \rightarrow a A / b & \text{rule 2 \& 3}
 \end{array}$$

put dot '.' at start of each right side symbol
 shift '.' one place - If non-terminal exist after dot '.', expand it
 Continue till dot shift at the end



Canonical collection of $LR(0)$ items

States	Action			Go to	
	a	b	\$	S	A
0	s_3	s_4		1	2
1			Accept		
2	s_3	s_4			5
3	s_3	s_4			6
4					
5					
6					
7					

Till this point, LR(0) & SLR(1) are same

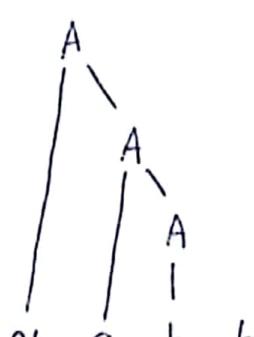
Now LR(0):

States	Action			Go to	
	a	b	\$	S	A
0	s_3	s_4		1	2
1			Accept		
2	s_3	s_4			5
3	s_3	s_4			6
4	r_3	r_3	r_3		
5	r_1	r_1	r_1		
6	r_2	r_2	r_2		

Let's parse "aabbb\$"

aabb\$
↑↑↑

0	a	3	a	3	b	1	A	6	A	6
---	---	---	---	---	---	---	---	---	---	---



for 1st character 'a', state=0,

Action = $s_3 \rightarrow$ push 'a' & '3' on stack

(There must be a number on top of stack
if rule is encountered, do "push")

- ① 2 times pop
- ② push LMS
- ③ Tree
- ④ put goto num

If r_n is encountered

- 1) pop 2 ~~elements~~ times # of element on right side
of rule
characters
- 2) push corresponding rule left most side non-terminal

$$A \rightarrow b$$

Pick number before 'A' which is 3 in this case,

push element having Action state = 3, $q_0 + 0 = A$
which is 6, so push '6'

Pick rule of state 6

$$A \rightarrow aA$$

SLR(1):

- find follow() set of grammar

$$\begin{aligned} S' &\rightarrow S & \{S\} \\ S &\rightarrow AA & \{a, b, \$\} \\ A &\rightarrow aA \mid b \end{aligned}$$

State	a	Action	b	\$	S	A
0	S_3		S_4		1	2
1				Accept		
2	S_3		S_4		5	
3	S_3		S_4		6	
4	r_3		r_3	r_3		
5				r_1		
6	r_2		r_2	r_2		

- ① $A \rightarrow b$
state = 4
rule = 3
put ' r_3 ' under
follow() set
terminals of A
 $\{a, b, \$\}$
it means
 $'r_3'$ will come
under $a, b, \$$
because these
are follower of
'A'
Same for other
- ② $A \rightarrow aA$
 $s = 6, r = 2$
follow = $\{a, b, \$\}$

Remaining part
of parsing is same
as LR(0).

③ $S \rightarrow AA$

$$r=1, s=5$$

$$\text{follow} = \{\$\}$$

We can see LR(0) has

more entries than SLR(1). The less the
entries in table, more chance of detecting error.

Example:

Applying LR(0) & SLR(1) parsing:

$$S \rightarrow dA \mid aB$$

$$A \rightarrow bA \mid c$$

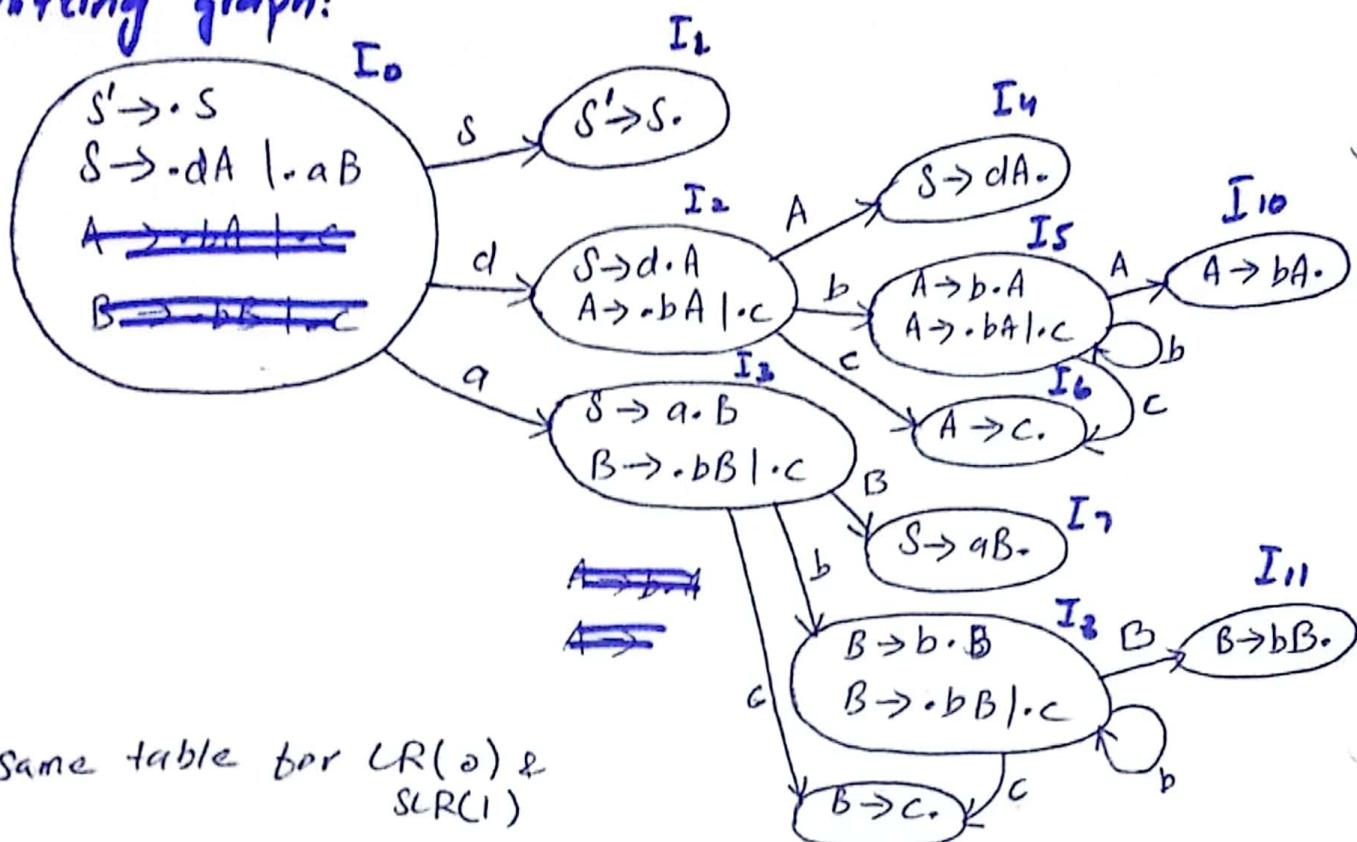
$$B \rightarrow bB \mid c$$

Solution:

adding additional state
 $s' \rightarrow s$

(~~initial~~)

Shifting graph:



Same table for LR(0) & SLR(1)

States	Actions					Goto	I_9
	a	b	c	d	$\$$		
0	δ_3			δ_2		S	
1							
2		δ_5	δ_6				4
3		δ_8	δ_9				7
4							
5		δ_5	δ_6				10
6							
7							
8		δ_8	δ_9				11
9							
10							
11							

LR(0):

$s' \rightarrow s$

$s \rightarrow dA$

$| dB$

$A \rightarrow bA$

$| c$

$B \rightarrow bB$

$| c$

rule # & state #

follow()

— (1) s_1, s_4

{\\$}

— (2) s_2, s_7

{\\$}

— (3) s_3, s_{10}

{\\$}

— (4) s_4, s_6

— (5) s_5, s_{11}

{\\$}

— (6) s_6, s_9

Putting rules in table:

Action

s_n = state #

r_n = rule #

State	Action							
	a	b	c	d	\$	s	A	B
0	s_3				s_2	1		
1								
2		s_5	s_6			4		
3		s_8	s_9				7	
4	r_1	r_1	r_1	r_1	r_1			
5		s_5	s_6			10		
6	r_4	r_4	r_4	r_4	r_4			
7	r_2	r_2	r_2	r_2	r_2			
8		s_3	s_9				11	
9	r_6	r_6	r_6	r_6	r_6			
10	r_3	r_3	r_3	r_3	r_3			
11	r_5	r_5	r_5	r_5	r_5			

Just remove box entries for SLR(1)

Main difference in LR(0) and SLR(1) is table generation for each rules

In LR(0), we fill complete row of corresponding state with rule #, but in SLR(1), just put rule # under their 'follow set' and 'Action' common terminals

LALR(1) & CLR(1)

- Similar to LR(0) & SLR(1) but add "look ahead" as well
- Lookahead is like follow set
For starting state, it will always be '\$'

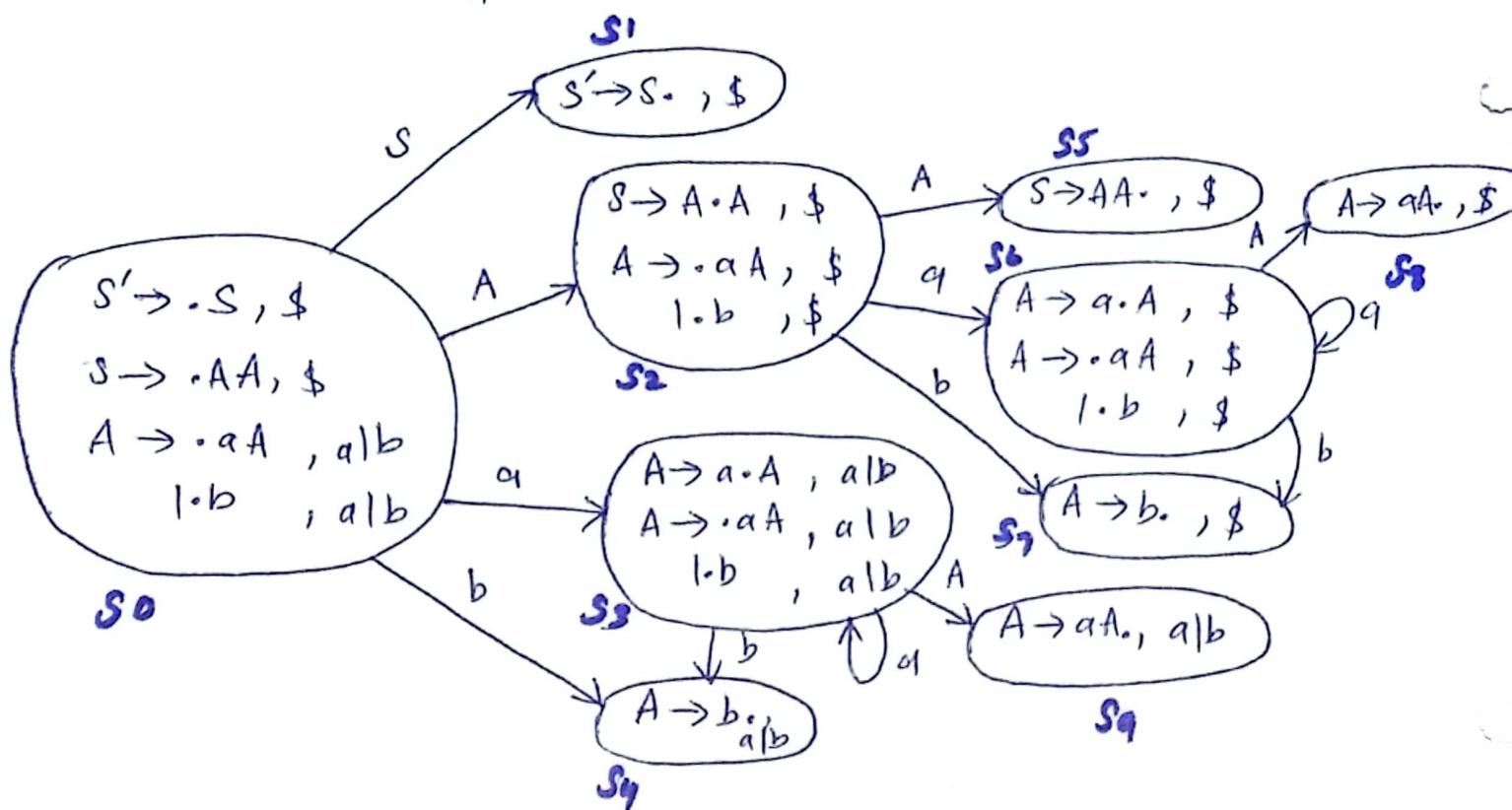
Example:

$$S \rightarrow AA$$

$$A \rightarrow aA \mid b$$

Sol: we always add additional state $s' \xrightarrow{\text{lookahead}} s$, in all bottom up parsers.

$s' \xrightarrow{} s$	$\$$	→ copy lookahead of non-terminal state as it is
$s \xrightarrow{} AA$	$\$$	
$A \xrightarrow{} aA$	$a \mid b$	



LALR(1)

State	Action			S	A
	a	b	\$		
0	s_3	s_4		1	2
1					
2	s_6	s_7			5
3	s_3	s_4			9
4	γ_3	γ_3			
5			γ_1		
6	s_6	s_7			8
7			γ_3		
8			γ_2		
9	γ_2	γ_2			

$$S \rightarrow A A$$

rule # 'Y_n', state #, S_n, look-ahead
 s_5, γ_1 { \$ }

$$A \rightarrow aA$$

$$S_9, S_2, \gamma_2$$

$$\{ \$ \}, \{ a \mid b \}$$

$$| b$$

$$S_7, S_4, \gamma_3$$

$$\{ a \mid b \}, \{ \$ \}$$

CLR(1)

Looking at transition diagram, below states are similar

$$S_8, S_9 \Rightarrow S_{89} \text{ (new state)}$$

$$S_3, S_6 \Rightarrow S_{36}$$

$$S_4, S_7 \Rightarrow S_{47}$$

Combine these states in the table computed via LALR(1)

Table must be in LALR(1) to find CLR(1)

CLR(1) reduces # of states but increases table entries

LALR(1) & CLR(1) uses LR(1) items

State	Action			Go to	
	a	b	\$	s	A
0	s_{36}	s_{47}		1	2
1					
2	s_{36}	s_{47}			5
36	s_{36}	s_{47}			89
47	r_3	r_3	r_3		
5				r_1	
89	r_2	r_2	r_2		

problem in LALR(1) & CLR(1)

if same lookahead exist in a state , where multiple rules exist then there is a conflict

For example:

$$A \rightarrow ab$$

$$B \rightarrow cd$$

$A \rightarrow ab \cdot , \$$
 $B \rightarrow cd \cdot , \$$

so

Table

State	Action				\$
	a	b	c	d	
0					1/2

↓
problem/
conflict

Example:

$$S \rightarrow Aa \mid bAc \mid Bc \mid bB^a$$

$$A \rightarrow d$$

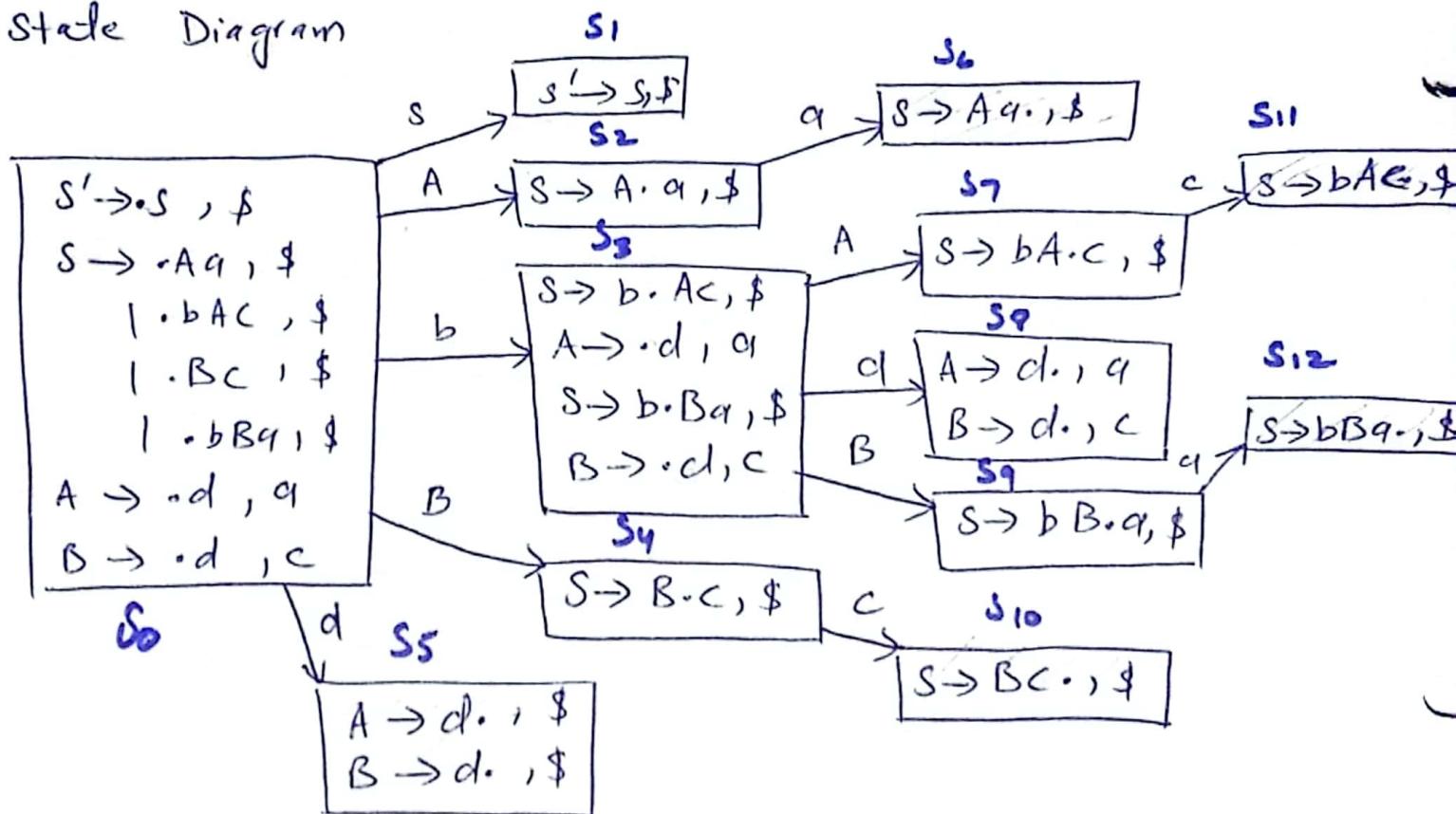
$$B \rightarrow d$$

Create LALR(1) & SLR(1) parser table.

Sol:

	rule #1, state #1, lookahead
$S' \rightarrow S$	" \$
$S \rightarrow Aa$	$\gamma_1, S_6 \mid \$$
$\mid bAc$	$\gamma_2, S_{11} \mid \$$
$\mid Bc$	$\gamma_3, S_{10} \mid \$$
$\mid bB^a$	$\gamma_4, S_{12} \mid \$$
$A \rightarrow d$	$\gamma_5, S_5 \mid S_2 a$
$B \rightarrow d$	$\gamma_6, S_5 \mid S_1 C$

State Diagram



State	Action							
	a	b	c	d	\$	S	A	B
0			s_3		s_5	s_1	s_2	s_4
1								
2		s_6						
3					s_7		s_7	s_9
4				s_{10}				
5	r_5		r_6					
6					r_1			
7			s_{11}					
8	r_5		r_6					
9		s_{12}						
10					r_3			
11					r_2			
12					r_4			

f CLR(1) is not possible due to conflict,
 LALR(1) is also not possible.

Recursive Descent parser

- top down parser
- implement separate method/routine/function for each non-terminal

Example:

$$E \rightarrow iE'$$

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

code:

E() {

```
if( look_ahead == 'i' )  
{  
    match('i');  
    E'();  
}
```

}

E'() {

```
if (look_ahead == '+') {  
    match('+');  
    match('i');  
    E'();  
}  
else  
    return;
```

}

match(char c) {

```
if( look_ahead == c ) {  
    look_ahead = getnextchar();  
} else  
    print("ERROR");
```

}

main() {

```
E();  
if( look_ahead == '$' ) print("parsing successful!");
```

}

Phase 5:

Semantic analysis

- Lexical analyzer gives stream of tokens to Syntax analyzer/parser which generates parse tree
- Semantic analyzer adds additional information in the parse tree for type checking.

→ Syntax directed definition (SDD):

- Attributes are associated with grammar symbol
- & semantic rules are associated with productions

Attributes may be number, string, reference, datatype etc

Grammar + Semantic rules

Semantic Rule

$$E \rightarrow E + T$$

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

$$E \rightarrow T$$

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

$$T \rightarrow T * F$$

$$T.\text{val} = T.\text{val} * F.\text{val}$$

$$T \rightarrow F$$

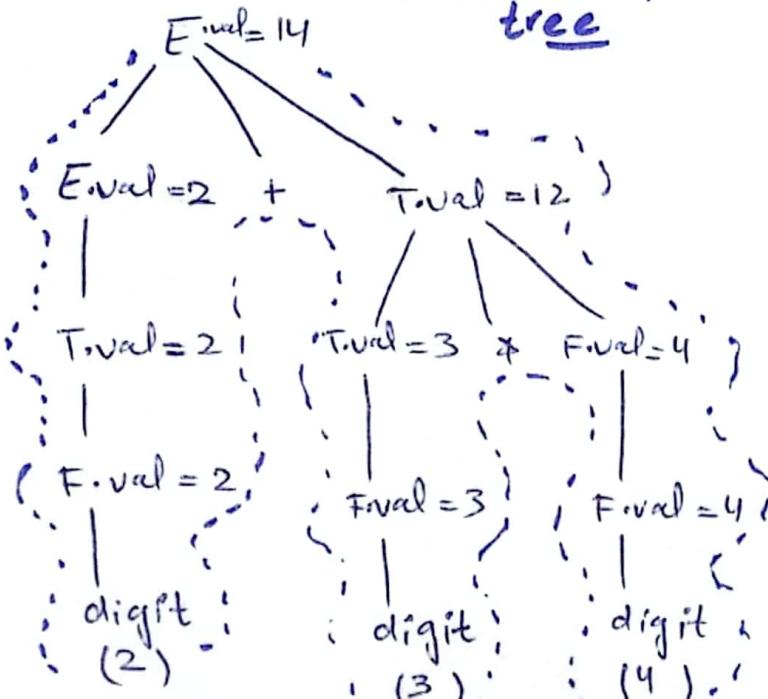
$$T.\text{val} = F.\text{val}$$

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

$$F.\text{val} = \text{digit-val} \\ (\text{lexical value})$$

$$2 + 3 * 4$$

Annotated parse tree



→ we can observe associativity and precedence by just looking at grammar
 → left recursive = left associative
 → symbol farthest from start symbol has highest precedence

cout << "A"
cout <<

Example:

production

$$E \rightarrow E \# T$$

$$E \rightarrow T$$

$$T \rightarrow T \& F$$

$$T \rightarrow F$$

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

Semantic Rule

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

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

$$T.\text{val} = T.\text{val} - F.\text{val}$$

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

$$F.\text{val} = \text{digit}.\text{val}$$

If the expression "8#12&4#16&12#4&2" is evaluated to 512 then which is best replacement to achieve this result

- a) $T.\text{val} = T.\text{val} \times F.\text{val}$
- b) $T.\text{val} = T.\text{val} + F.\text{val}$
- c) $T.\text{val} = T.\text{val} - F.\text{val}$
- d) None

Solution:

Instead of replacing each in parse tree, we can observe that Grammer is left recursive for '#' and '&', so, if it is left associative. & has highest precedence than # due to farthest from start symbol.

→ Trying option (a)

$$8 \# (12 \& 4) \# (16 \& 12) \# (4 \& 2)$$

$$\# = \times$$

$$\& = \times$$

$$8 \times (12 \times 4) \times (16 \times 12) \times (4 \times 2) \neq 512$$

→ Trying option (b)

$$\# = *, \& = +$$

$$8 \times (12+4) \times (16+12) * (4+2) \neq 512$$

→ Trying option (c)

$$\# = *, \& = -$$

$$8 \times (12-4) \times (16-12) \times (4-2) = 512$$

So (c) is correct!

Example:

Semantic Rules

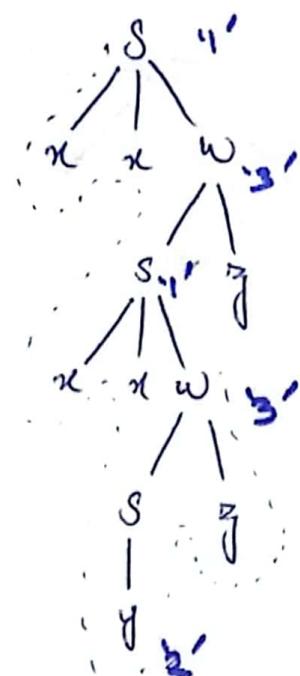
$$S \rightarrow nxw \mid y \quad \{ \text{print('1')} \} \mid \{ \text{print('2')} \}$$

$$w \rightarrow xy \quad \{ \text{print('3')} \}$$

Perform semantic analysis on "nxxyyjj".

Sol:

Parse Tree:

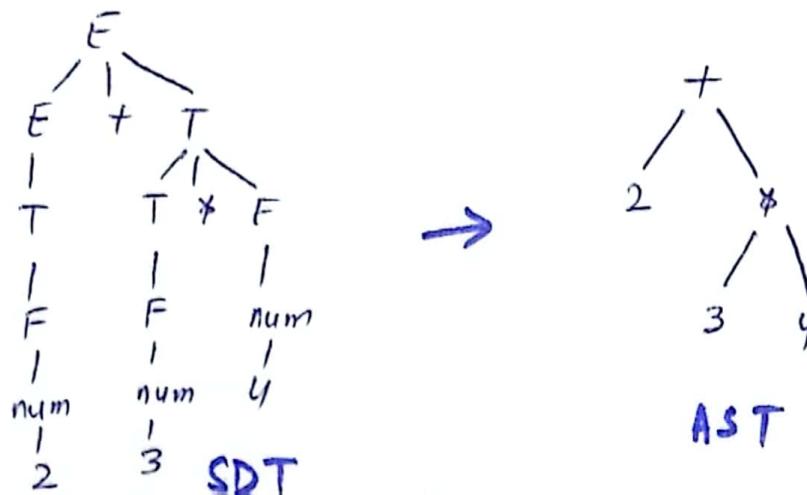


= 23131

Abstract Syntax Tree (AST):

- After removing irrelevant information from parse tree, AST is generated.

For example:

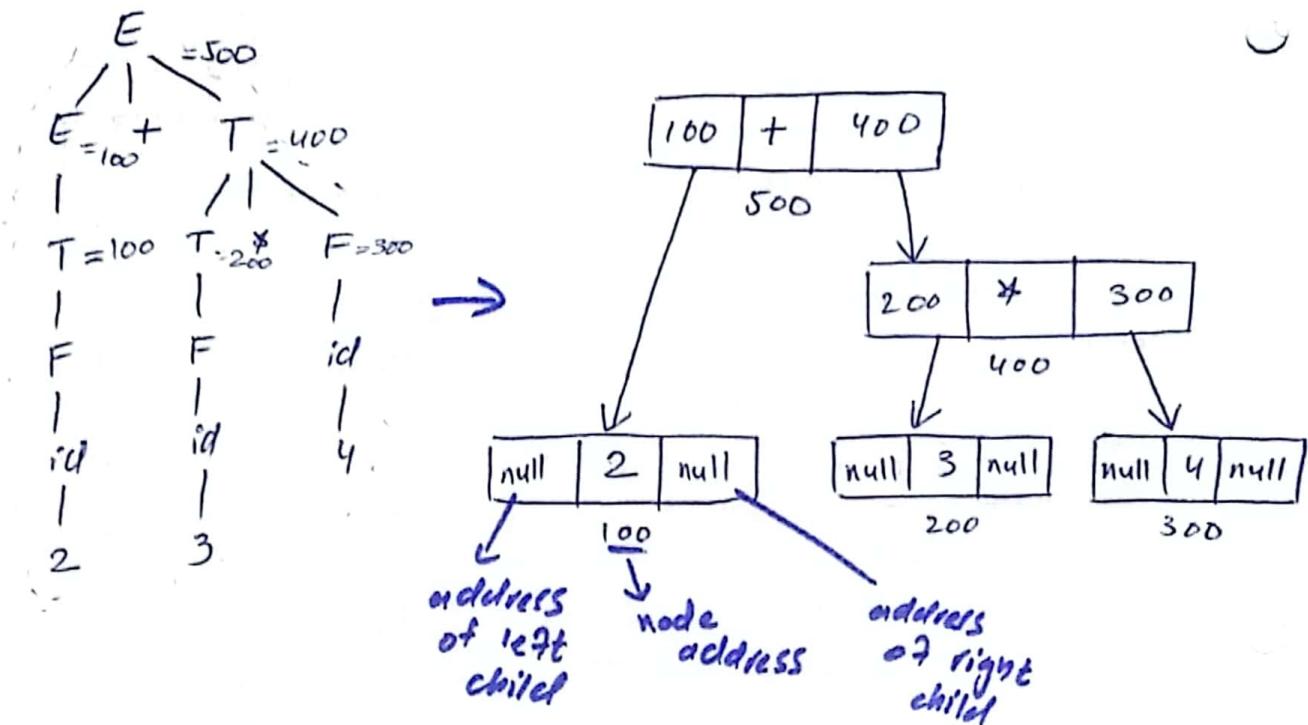


Example:

$E \rightarrow E + T \quad \text{if } T$
 $T \rightarrow T * F \quad \text{if } F$
 $F \rightarrow id$

Semantic rule
 $\{ E \cdot nptr = \text{newNode}(E \cdot nptr, '+', T \cdot nptr) \}$
 $\{ E \cdot nptr = T \cdot nptr \}$
 $\{ T \cdot nptr = \text{newNode}(T \cdot nptr, '*', F \cdot nptr) \}$
 $\{ T \cdot nptr = F \cdot nptr \}$
 $\{ F \cdot nptr = \text{newNode}(\text{null}, \text{id}.val, \text{null}) \}$

Perform semantic analysis & generate result on "2+3*4"



This semantic rules generate Abstract syntax tree (AST) and syntax directed tree (SDT).

Example:

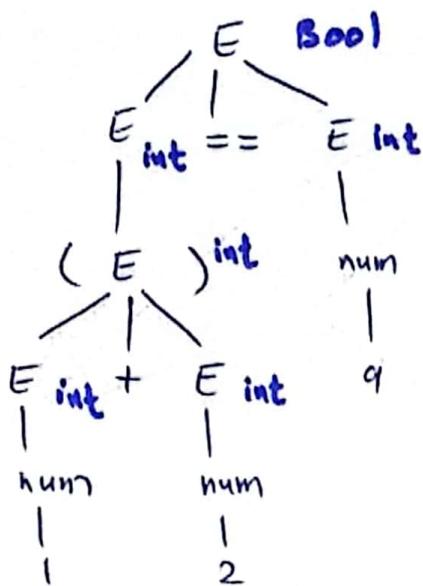
Semantic Rules

$E \rightarrow E + E \quad \{ \text{if}(E \cdot type == E \cdot type) \& (E \cdot type = \text{int}) \text{ then } E \cdot type = \text{int} \text{ else error} \}$
 | $E == E \quad \{ \text{if}(E \cdot type == E \cdot type) \& (E \cdot type = \text{int/Bool}) \text{ then } E \cdot type = \text{Bool} \text{ else error} \}$
 | $(E) \quad \{ E \cdot type = E \cdot type \}$
 | num $\{ E \cdot type = \text{int} \}$
 | True $\{ E \cdot type = \text{Bool} \}$
 | False $\{ E \cdot type = \text{Bool} \}$

what will be type of $(5+7) == (11+2)$?

$$(1+2) == 9$$

what will be type returned by semantic rules?



'Bool' will be final type.

Types of SDD

① S-Attributed

- A SDD that uses only synthesized attribute is called S-Attributed SDD.

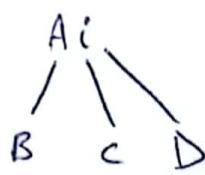
Example:

$$A \rightarrow BCD$$

$$Ai = Bi$$

$$Ai = Ci$$

$$Ai = Di$$



Synthesized attributed:- child node value is consumed/used by parent node

- Semantic actions are placed on right most side of CFG
- Attributes are evaluated bottom up because child value is accessed by parent node.

② L-Attributed

- A SDD that uses both inherited and synthesized attribute is L-attributed. But each inherited attribute is restricted to inherit from parent or left sibling because DFS traversal is used.



'C' can only inherit attribute from A & B not D.

Inherited attributed:- child inherit attributes from parent node

- Semantic rules are placed anywhere on RHS of production
- Depth first search from left to right is used.

Example:

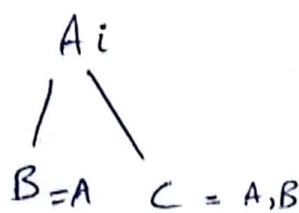
Identify productions are S-attributed, L-attributed or both

1) $A \rightarrow BC \quad \{ B_i = A_i, C_i = B_i, A_i = C_i \}$

2) $A \rightarrow QR \quad \{ R_i = A_i, Q_i = R_i, A_i = Q_i \}$

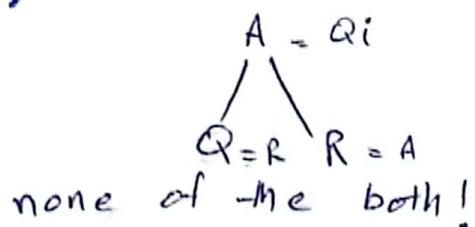
3) $A \rightarrow PQ \quad \{ A_i = P_i, A_i = Q_i \}$

(1)



As 'C' can inherit from A & B, parent & left sibling, it is L-attributed

(2)



(3)



If attribute is 'S', it will also be 'L' attributed

Example:

Given Grammar converts ~~binary~~ binary to decimal:

$$S \rightarrow L \quad \{ S.dv = L.dv \}$$

$$L \rightarrow LB \quad \{ L.dv = 2 \times L.dv + B.dv \}$$

$$L \rightarrow B \quad \{ L.dv = B.dv \}$$

$$B \rightarrow 0 \quad \{ B.dv = 0 \}$$

$$B \rightarrow 1 \quad \{ B.dv = 1 \}$$

create annotated parse tree
for "1010"

write CFG & Semantic rule for given C statement

int [2][3];

Semantic Rules

$T \rightarrow BC;$

$B \rightarrow \text{int}$

$B \rightarrow \text{float}$

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

$C \rightarrow \epsilon$

$\{T.t = C.t, C.b = B.t\}$

$\{B.t = \text{int}\}$

$\{B.t = \text{float}\}$

$\{C.t = \text{array}(\text{num}.val, C.t)\}$

$\{C.b = C.b\}$

$\{C.t = C.b\}$

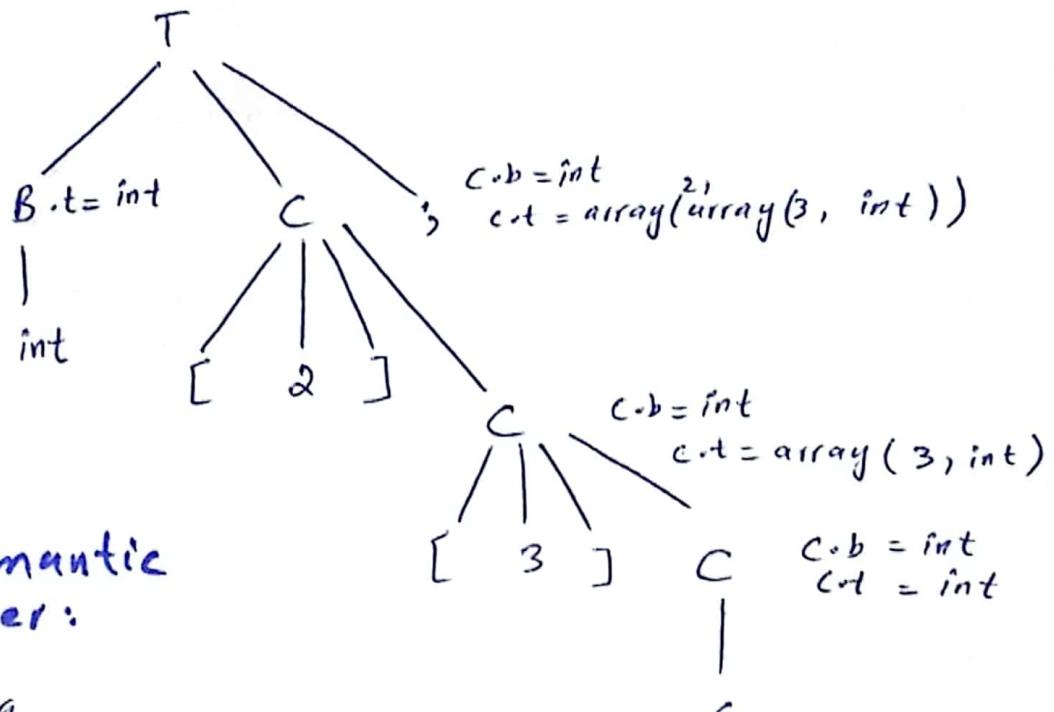


Table of Semantic analyzer:

- 1) Type checking
- 2) store variable type in symbol table
- 3) Create syntax tree
- 4) Evaluate arithmetic expression
- 5) Infix to postfix or prefix
- 6) Conversion from binary to decimal
- 7) Counting # of reductions

Example:

Given grammar converts floating binary to decimal

$$S \rightarrow L \cdot L \quad \{ S \cdot dV = L \cdot dV + L \cdot dV / 2^{L \cdot nb} \}$$

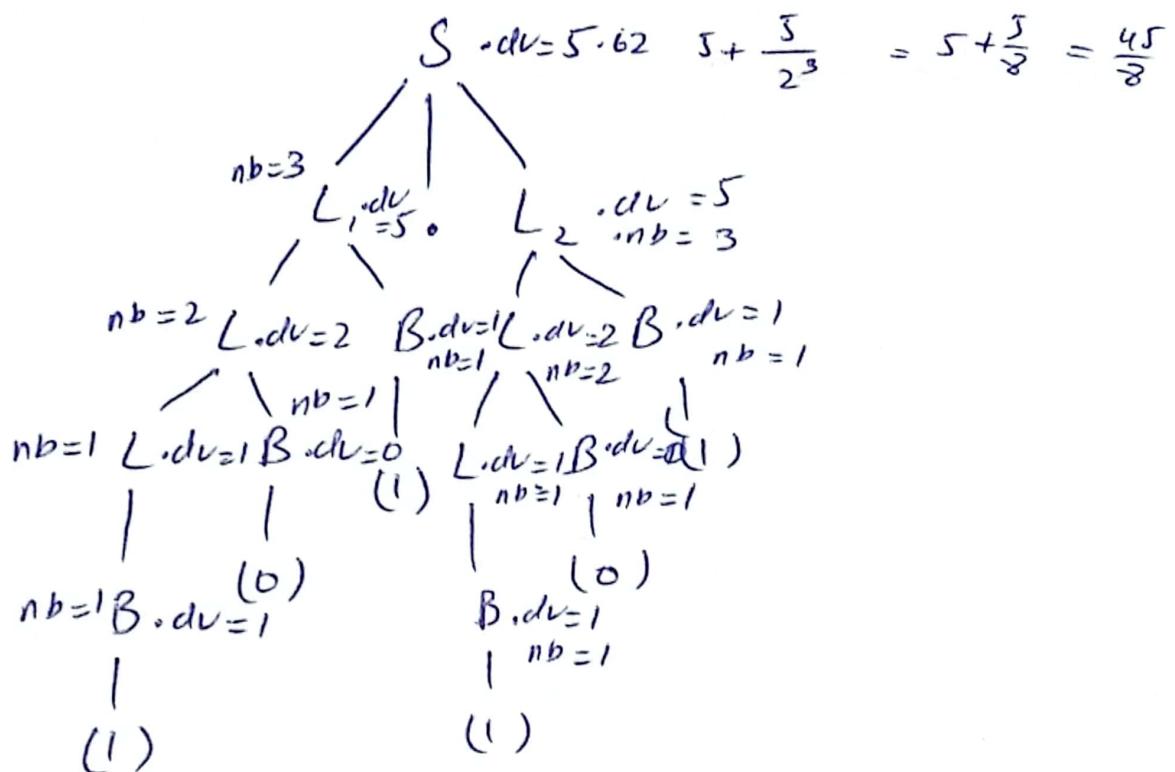
$$L \rightarrow L \cdot B \quad \{ L \cdot dV = 2 \times L \cdot dV + B \cdot dV, L \cdot nb = L \cdot nb + B \cdot nb \}$$

$$L \rightarrow B \quad \{ L \cdot dV = B \cdot dV, L \cdot nb = B \cdot nb \}$$

$$B \rightarrow 0 \quad \{ B \cdot dV = 0, B \cdot nb = 1 \}$$

$$B \rightarrow 1 \quad \{ B \cdot dV = 1, B \cdot nb = 1 \}$$

Create annotated parse tree for 101.101



SDT to Generate 3 address code (3AC):

6

$S \rightarrow id = E$

$\{ id.name = E.place \}$

$E \rightarrow E_1 + T$

$\{ E.place = new temp();$

$gen(E.place = E_1.place + T.place) \}$

$E \rightarrow T$

$\{ E.place = T.place \}$

$T \rightarrow T * F$

$\{ T.place = new temp(); gen(T.place = E.place + T.place) \}$

$T \rightarrow F$

$\{ T.place = F.place \}$

$F \rightarrow id$

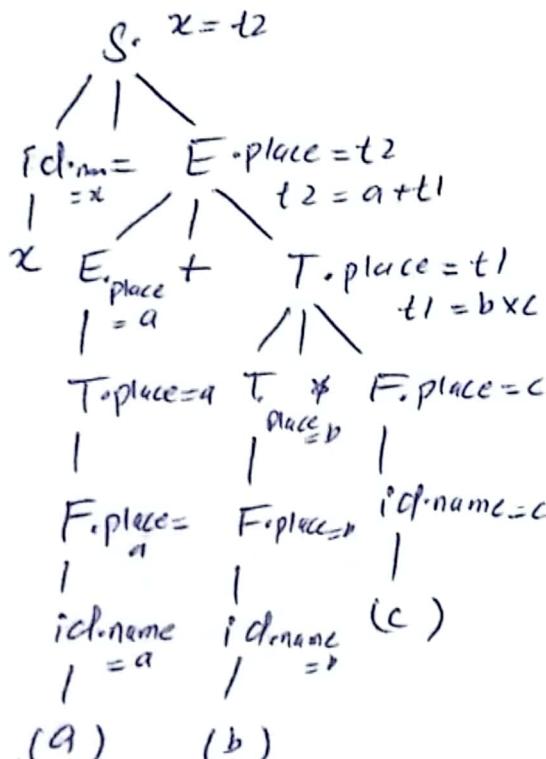
$\{ F.place = id.name \}$

In 3AC, maximum 3 addresses are allowed like in this expression $a = b + c * d$, addresses are 4 and need to create a temporary variable to convert it to 3AC.

3AC is used in intermediate code generation.

Lets parse and generate 3AC of

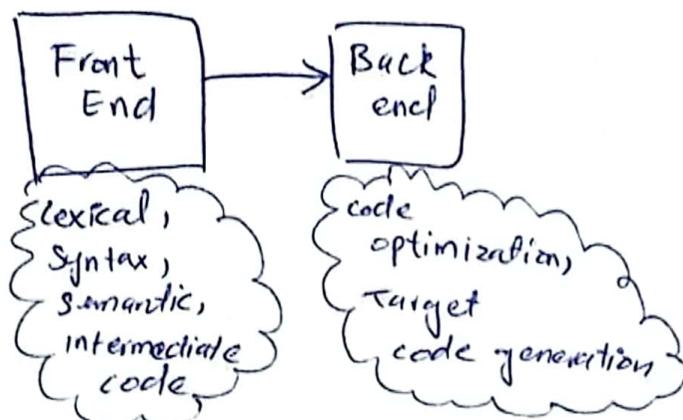
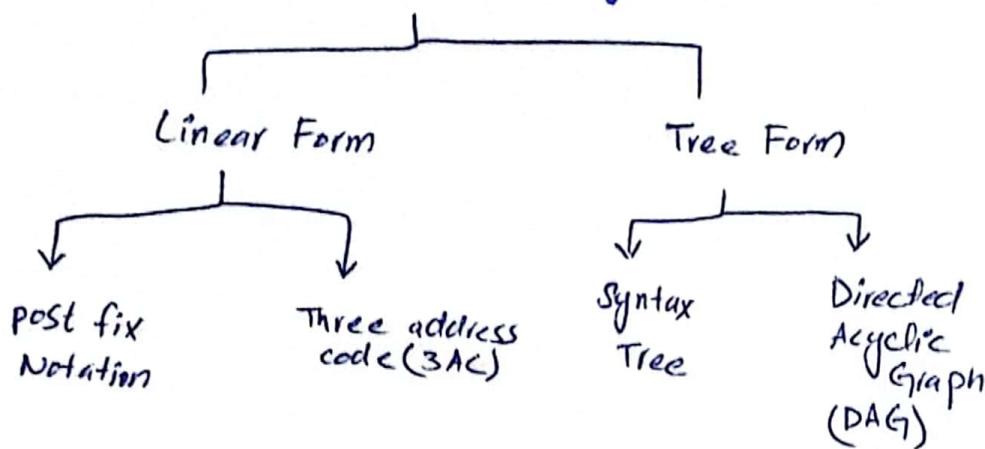
$$x = a + b * c;$$



3AC

$t1 = b * c$
$t2 = a + t1$
$x = t2$

Phase 4: "Intermediate code generation"

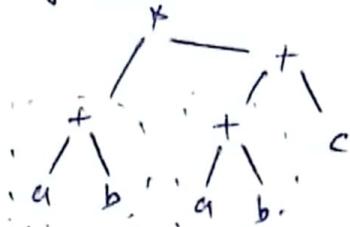


➤ If I want to design compiler for windows & Linux OS, then I just need to manipulate Backend part while front end will be same for all machines. It means we do not have to write compiler from scratch. Compiler construction tools like lex, yacc, bison are also used to design front end, we just need to decide the structure & rules of our language.

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

$$\text{post fix} = a\ b + a\ b + c + * \quad |$$

Syntax Tree:



$$\begin{aligned} 3AC = \quad t1 &= a+b \\ &t2 = a+b \\ &t3 = t2+c \\ &t4 = t1*t3 \end{aligned}$$

DAG: Remove duplicate parent-child subtree



→ Types of 3AC statements

Try to put 3 address per instruction

1) $x = y \text{ op } z$

op = binary operator $\{+, -, *, /, \dots\}$

Example :

$$n = a + b * c$$

$$\boxed{\begin{array}{l} t1 = b * c \\ t2 = a + t1 \\ n = t2 \end{array}}$$

2) $x = \text{op } z$

op = unary operator $\{++, --\}$

3) $x = y$

4) if ($x \text{ relop } y$) goto L (conditional jump)

relop = relational operator $\{<=, ==, >=, !=, ..\}$

5) goto L (unconditional jump)

6) $A[i] = x$

7) $x = *p$

$$y = \&x$$

→ Representation of 3AC :

- 1) Quadruple 2) Triples 3) Indirect triples

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

3AC: $t1 = a * b$

$$t2 = -t1$$

$$t3 = c * d$$

$$t4 = t3 * e$$

$$t5 = t2 + t4$$

1) Quadruple:

[Consider SAC of previous example]

	operator	op1	op2	result
0	*	a	b	t1
1	-	t1		t2
2	*	c	d	t3
3	+	t3	e	t4
4	+	t2	t4	t5

- Require extra space
- Re-ordering of instructions is possible if there is no inconsistency in result occurs.

2) Triples:

- Remove last column from Quadruple
- put index of table in place of temporary variables
 $\{t1, t2, \dots, t5\}$

	operator	op1	op2
0	*	a	b
1	-	(0)	
2	*	c	d
3	+	(2)	e
4	+	(1)	(3)

- Requires less space than Quadruple
- Reordering is not possible due to table index as address

t2 t4, entry no. 3 in
 entry quadruple last
 no. 1 column latest
 in table t5
 of Quadruple

3) Indirect Triples

100	(0)	→ index no. 0 of table of "Triples"
101	(1)	
102	(2)	
103	(3)	
104	(4)	

- Requires extra memory address
- Re-order of instruction is possible

Backpatching

Leaving table as empty and filling them latter is backpatching.

We have seen 7 types of instruction in programming

For if-else, SAC will be:

Example:

$\text{if}(a < b) \text{ then } t = 1 \text{ else } t = 0$

SAC:

Line#

1. $\text{if}(a < b) \text{ goto 4}$
2. $t = 0$
3. goto 5
4. $t = 1$

Example 2:

$\text{if}(a < b) \& \& (c < d) \text{ then } t = 1 \text{ else } t = 0$

1 $\text{if}(a < b) \text{ goto 4}$

2 $t = 0$

3 ~~$\text{if}(c < d)$~~ goto 7

4 $\text{if}(c < d) \text{ goto 6}$

5 goto 2

6 $t = 1$

7

Loop : (3AC)

for(i=1 ; i<n ; i++) {

 x = a + b * c ;

}

3AC

1	i = 1
2	if(i < n) goto 4
3	goto 9 // false case
4	t1 = b * c
5	t2 = a + t1
6	x = t2
7	i = i + 1
8	goto 2

OR

1	i = 1
2	if(i > n) goto 8
3	t1 = b * c
4	t2 = a + t1
5	x = t2
6	x = i + 1
7	goto 2
8	

switch Case:

switch(i) {

case 1:

 x1 = a1 + b1 * c1;
 break;

case 2:

 x2 = a2 + b2 * c2; break

default:

 x3 = a3 + b3 * c3;
 break;

}

3AC

1	if(i == 1) goto 7
2	if(i == 2) goto 11
3	t1 = b3 * c3;
4	t2 = a3 + t1;
5	x3 = t2
6	t1 goto 14
7	t1 = b1 * c1
8	t2 = a1 + t1
9	x = t2
10	goto 14
11	t1 = b2 * c2
12	t2 = a2 + t1
13	x = -12
14	

Arrays: (1D)

100	102	104	106	108	110	112
4	5	9	3	2	0	1

size of int = 2

$$\begin{aligned} A[3] &= \text{base_address} + (i - \text{lower bound}) \times \text{size of single element} \\ &= 100 + (3 - 0) \times 2 \\ &= 106 \end{aligned}$$

Example:

int A[10], B[10];

int x = 0; i = 0;

for (i = 0; i < 10; i++) {

x = x + A[i] * B[i];

}

3AC:

- 1 x = 0
- 2 i = 0
- 3 if ($i \geq 10$) goto 15
- 4 t1 = base address of A
- 5 t2 = $i * 4$
- 6 t3 = t1[t2]
- 7 t4 = base address of B
- 8 t5 = $i * 4$
- 9 t6 = t4[t5]
- 10 t7 = t6 * t5
- 11 t8 = x + t7
- 12 x = t8
- 13 i = i + 1
- 14 goto 3
- 15

2D Array :

$$A[3][3] = \begin{bmatrix} 2 & 3 & 4 \\ 5 & 0 & 1 \\ 1 & 2 & 5 \end{bmatrix}$$

$$A[i][j] = \text{Base_address} + (i \times \cancel{\text{size of datatype}} + j) \times \underset{\substack{\# \text{ of} \\ \text{columns}}}{\cancel{\text{size of datatype}}}$$

$$A[1][3] = 0 + (1 \times 3 + 3) \times 4 \\ = 24$$

int i=10, j=15;

$$X = A[i][j]$$

or

$$t1 = i * 15$$

$$t2 = t1 + j$$

$$t3 = t2 * 2$$

$$t4 = \text{Base address of } A$$

$$t5 = t4[t3]$$

$$X = t5$$

$$t1 = \text{base addr of } A$$

$$t2 = i \times 4$$

$$t3 = t1[t2]$$

$$t4 = j \times 4$$

$$t5 = t3[t4]$$

$$X = t5$$

3D Array :

int A[a][b][c];

$$A[i][j][k] = A[i \times b \times c + j \times c + k] \times \text{size of datatype}$$

Example

Give 3AC code, identify correct option. assume array of char is 8 bits & integer 32 bits ≈ 4 bytes

$$t0 = i \times 1024$$

$$a) \text{ int } X[32][32][8];$$

$$t1 = j \times 32$$

$$b) \text{ int } X[4][1024][32];$$

$$t2 = k \times 4$$

$$c) \text{ char } X[4][32][8];$$

$$t3 = t1 + t0$$

$$d) \text{ char } X[32][16][2];$$

$$t4 = t3 + t2$$

$$*5 = X[t4]$$

Try option (a)

$$\begin{aligned} X[i][j][k] &= X[i \times 32 \times 8 + j \times 8 + k] \times 4 \\ &= \cancel{X[i \times 32 \times 8 + j \times 8 + k]} - \\ &= X[i \times 32 \times 8 + j \times 8 + k \times 4] \\ &= X[i \times 1024 + j \times 32 + k \times 4] \end{aligned}$$

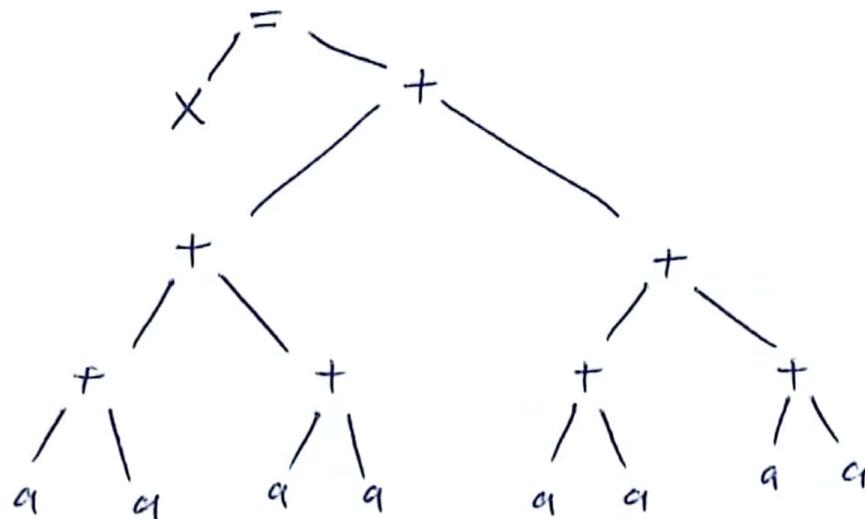
So option (a) is right one!

Syntax tree & DAG:

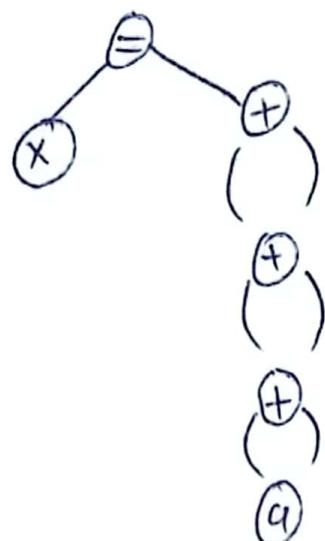
Example: (directed acyclic graph)

$$X = ((a+a)+(a+a)) + ((a+a)+(a+a))$$

Syntax tree

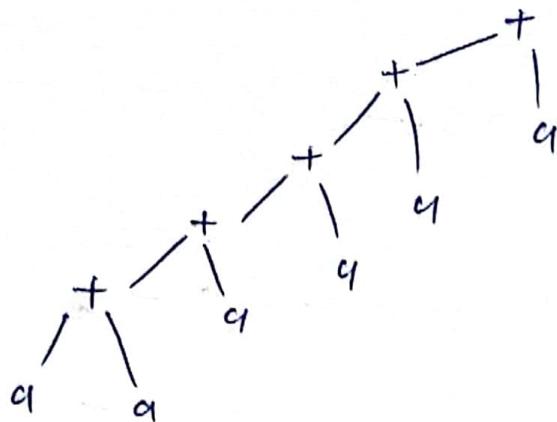


DAG: Remove same subtree & create a loop

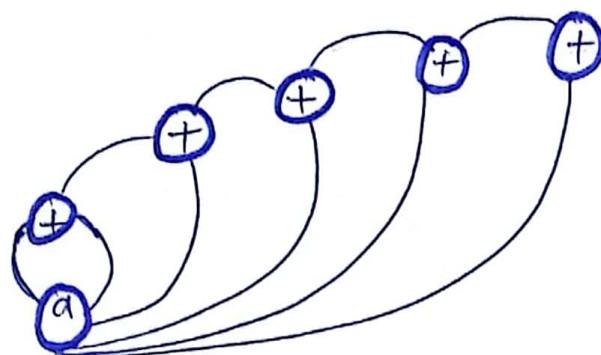


2) $X = a + a + a + a + a + a$

Syntax tree:

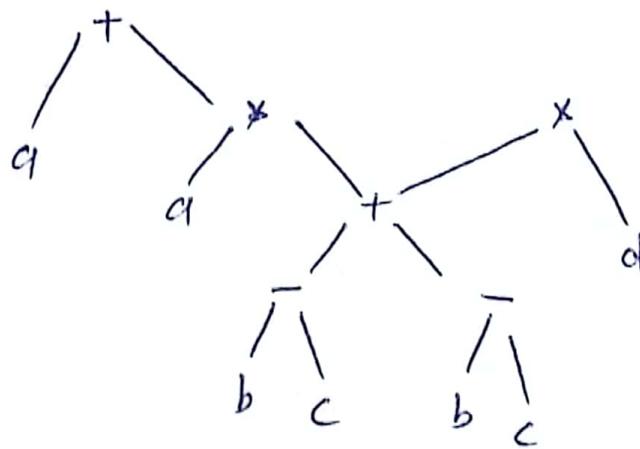


DAG:

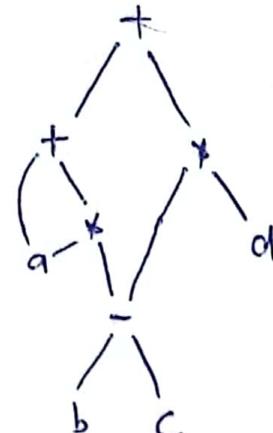
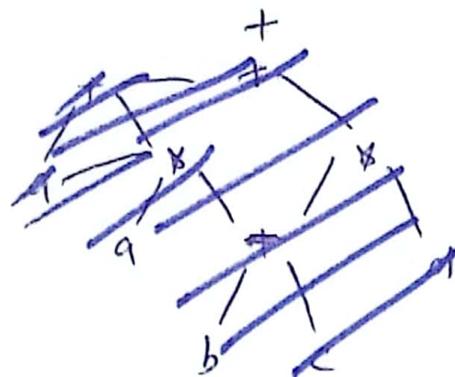


3) $X = a + a * (b - c) + (b - c) * d$

Syntax Tree



DAG:



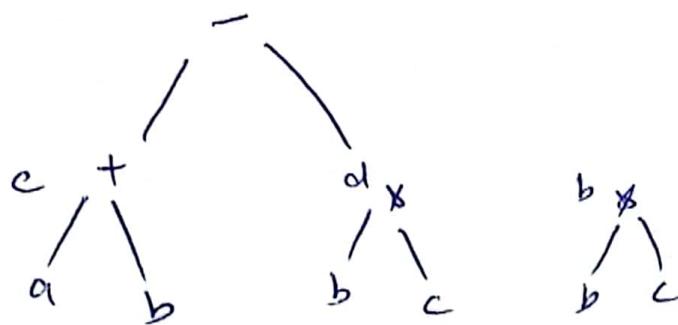
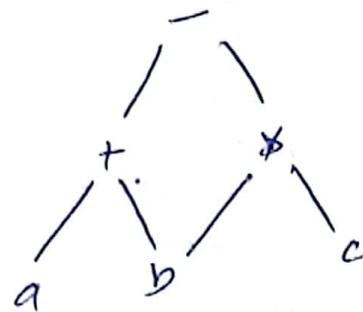
4)

$$d = b * c$$

$$e = a + b$$

$$b = b * c$$

$$a = e - d$$

Syntax Tree:DAG:

5)

~~$d = b * c$~~

$$a = b + c$$

~~$a = a - d$~~

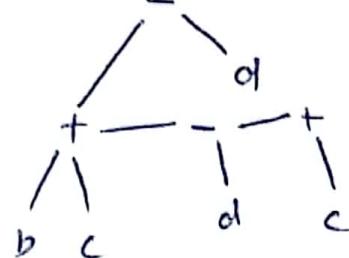
$$b = a - d$$

~~$b = b * c$~~

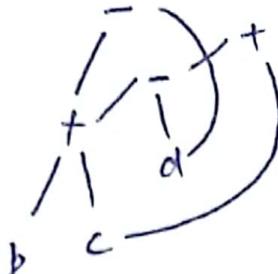
$$c = b + c$$

~~$d = a - d$~~

$$d = a - d$$

Syntax Tree:

DAG



Example : 6

$$a = b + c$$

$$c = a + d$$

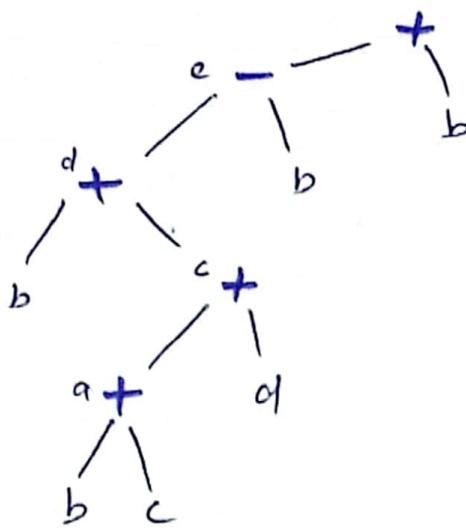
$$d = b + c$$

$$e = d - b$$

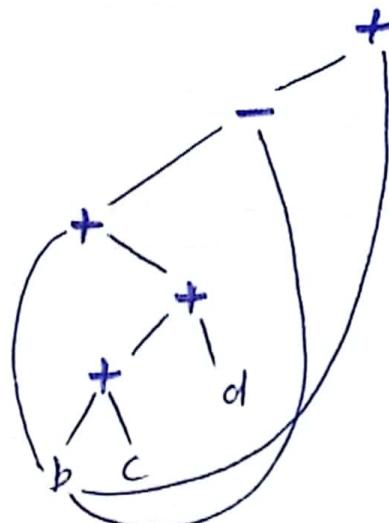
$$a_1 = e + b$$

find min # of nodes in DAG.

Syntax
tree



DAG:



Reducing by putting values in equation (substitution) :

$$a_1 = e + b$$

$$= d - b + b$$

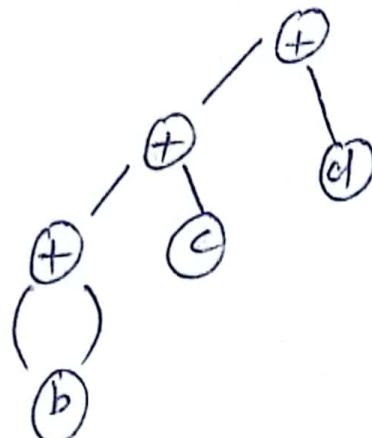
$$= b + c$$

$$= b + a + d$$

$$= b + b + c + d$$

$$\text{min # of nodes} = 6$$

$$\text{# edges} = 6$$



3AC for procedures:

$n = f(a[i]);$

$t1 = \text{Base address of } a$

$t2 = i \times 4$

$t3 = t1[t2]$

param $t3$

$t3 = \text{call } f, 1$

$n = t3$

SDT to generate 3AC for Control Flow statements

$P \rightarrow S$

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

$P.\text{code} = S.\text{code}$

$S \rightarrow \text{assign}$

$S.\text{code} = \text{assign}.\text{code}$

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

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

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

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

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

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

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

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

$S.\text{code} = B.\text{code} \parallel S_1.\text{code} \parallel S_2.\text{code}$
 $\parallel \text{gen(goto } S.\text{next}) \parallel$
 $\text{label}(B.\text{false})$

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

$\text{begin} = \text{newlabel}()$

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

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

$S_1.\text{next} = \text{begin}$

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

|| label(B.true) || gen('goto begin')

$S \rightarrow S_1 S_2$ $S_1.\text{next} = \text{newlabel}()$

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

$S.\text{code} = S_1.\text{code} || S_2.\text{code} || \text{label}(S.\text{next})$

$B \rightarrow B_1 || B_2$ $B_1.\text{true} = B.\text{true}$

$B_1.\text{false} = \text{newlabel}();$

$B_2.\text{true} = B.\text{true}$

$B_2.\text{false} = B.\text{false}$

$B.\text{code} = B_1.\text{code} || B_2.\text{code} ||$
 $\text{label}(B_1.\text{false})$

$B \rightarrow B_1 \& B_2$ $B_1.\text{true} = \text{newlabel}()$

$B_1.\text{false} = B.\text{false}$

$B_2.\text{true} = B.\text{true}$

$B_2.\text{false} = B.\text{false}$

$B.\text{code} = B_1.\text{code} || B_2.\text{code} ||$
 $\text{label}(B_1.\text{true})$

$B \rightarrow !B_1$ $B_1.\text{true} = B.\text{false}$

$B_1.\text{false} = B.\text{true} -$

$B.\text{code} = B_1.\text{code}$

$B \rightarrow E_1 \text{ relop } E_2$ $B.\text{code} = E_1.\text{code} || E_2.\text{code} /$
 $/ \text{gen}('if' E_1.\text{addr} \text{ relop } E_2.\text{addr}$
 $\text{'goto' } B.\text{true})$
 $/ \text{gen}('goto' B.\text{false})$

$B \rightarrow \text{true}$ $B.\text{code} = \text{gen}('goto' B.\text{true})$

$B \rightarrow \text{false}$ $B.\text{code} = \text{gen}('goto' B.\text{false})$

Example:

Write SDT to compute data type and size of array. Consider int occupies 4 bytes & float 8 bytes

$T \rightarrow BC$

$\{ t = B.type ; w = B.width ;$
 $T.type = C.type ; T.width = C.width \}$

$T \rightarrow \text{int}$

$\{ T.type = \text{integer} ; T.width = 4$

$| \text{float}$

$T.type = \text{float} ; T.width = 8 \}$

$C \rightarrow \epsilon$

$\{ C.type = t ; C.width = w \}$

$C \rightarrow [\text{num}] C_1$

$\{ C.type = \text{array}(\text{num.value}, C.type) ;$
 $C.width = \text{num.value} \times C_1.width \}$

Phase 5: Code Optimization

→ Machine Independent

- 1) Loop optimization
 - a) Code motion or frequency reduction
 - b) Loop unrolling
 - c) Loop scanning
- 2) Folding
- 3) Redundancy elimination
- 4) Strength reduction
- 5) Algebraic Simplification

→ Machine dependent

- 1) Register allocation
- 2) Instruction selection
- 3) Use of Addressing modes
- 4) Peephole optimization
 - a) Redundant LOAD/STORE
 - b) Flow of control optimization
 - c) Use of machine idioms

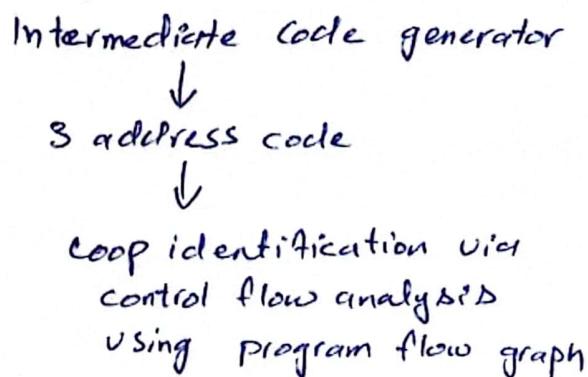
Loop Optimization:

- 1) To apply loop optimization, we must first detect loops.
- 2) For detecting loops, control flow analysis via control flow graph (CFG) is used.
- 3) To build program flow graph, we need to find basic blocks
- 4) A basic block is a sequence of SAC statements where control enters at the beginning & leaves only at the end without any jump or halt.

~~cont~~

Construction of program flow graph:

1) Basic block:



→ In order to find basic blocks, we need to find leaders in 3AC

→ A basic block will start from one leader to next leader but not including next leader

* Finding leaders in a basic block:

① 1st 3AC instruction in intermediate code is leader

② Any instruction that is the target of conditional goto or unconditional goto is a leader

$\text{if}(x < 3) \text{ goto } 3$

// '3' is leader

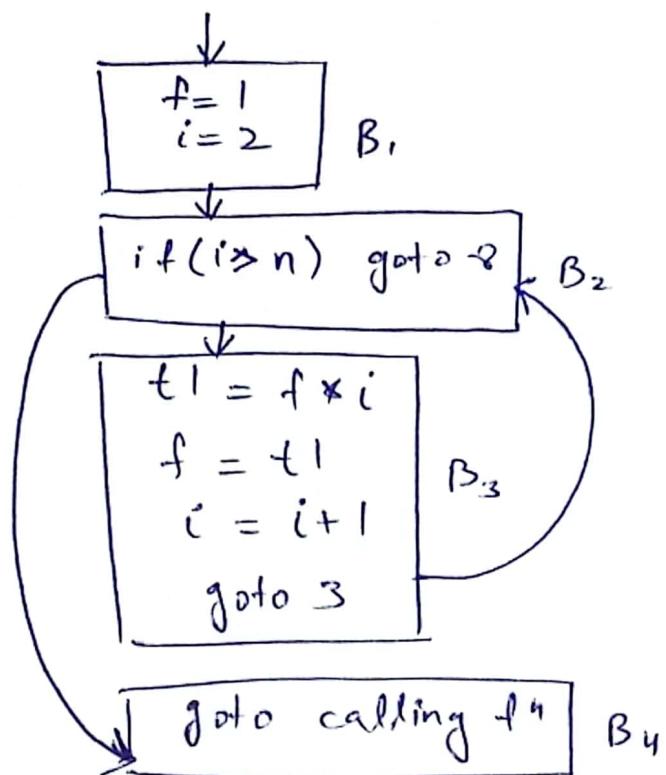
③ Any instruction that immediately follows a conditional or unconditional jump is a leader.

$x = y + z$ // leader

$\text{if}(x == 0) \text{ goto } -$

Example:

- * 1 $f = 1$
- 2 $i = 2$
- * 3 $\text{if } (i > n) \text{ goto } 8$
- * 4 $t1 = f * i$
- 5 $f = t1$
- 6 $i = i + 1$
- 7 $\text{goto } 3$
- * 8 $\text{goto calling } f^n$



→ If loop exist, there will be cycle in the graph.
→ optimization can be performed at each block basic.

Example 2:

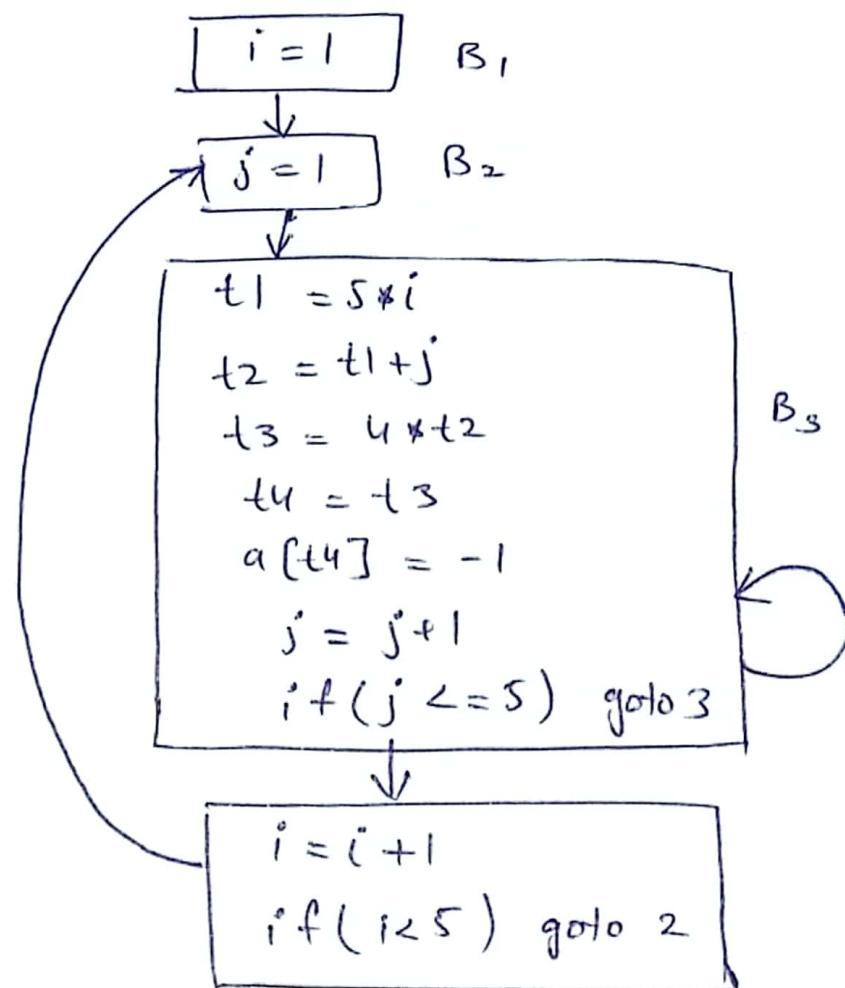
```
* 1      i = 1  
* 2      j = 1  
* 3      t1 = 5 * i  
4      t2 = t1 + j  
5      t3 = 4 * t2  
6      t4 = t3  
7      a[t4] = -1  
8      j = j + 1  
# 9      if (j <= 5) goto 3  
* 10     i = i + 1  
* 10     if (i < 5) goto 2
```

Construct program flow graph.

Identify leader

" basic block

" ~~construct~~ construct Graph



Loop optimization

① Frequency Reduction: (code motion)

A statement or expression which can be moved outside the loop without affecting the semantics of program

```
int P;  
char *name = "Aoun";  
for( int i=0; i< len(name); i++ ) {  
    print( name[i] );  
    P=D;  
}
```

problem

'len' will be invoked 5 times, $P=0$ is constant expression & can be kicked out.

② Loop unrolling:

Reducing # of times comparisons are made in the loop.

```
for( i=0; i<10 ; i++ ) {  
    print( "Mr. Aoun" );  
}
```



```
for( i=0; i<10; i+=2 ) {  
    print("Mr. Aoun");  
    print("Mr. Aoun");  
}
```

③ Loop jamming:

Combining bodies of multiple loops

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



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

Dead code:

- Unreachable code
- code never executed or used or participate in final outcome

For example, in if-else condition, 'if' is always executed then 'else' block is dead block.

① Unused variable, functions

② Redundant code

③ Constant folding

int $x = 10 * 5 + 3;$

here final outcome can be stored in 'x'

Induction variables & strength reduction

- Variable increasing or decreasing by a constant factor
 $t^4 = i \times 4$
 is induction variable
- Strength reduction is to replace ~~strong~~ strong operator like ' \wedge ' with cheap operator ' \wedge'

Liveness analysis:

Register allocation:

variables are stored in registers. It makes problem when no program variables are greater than available registers. We need to synchronize some variables in some way.

Purpose of register allocation is to assign a single register to multiple variables without changing program behaviour/semantics

X is live at statement s_i if

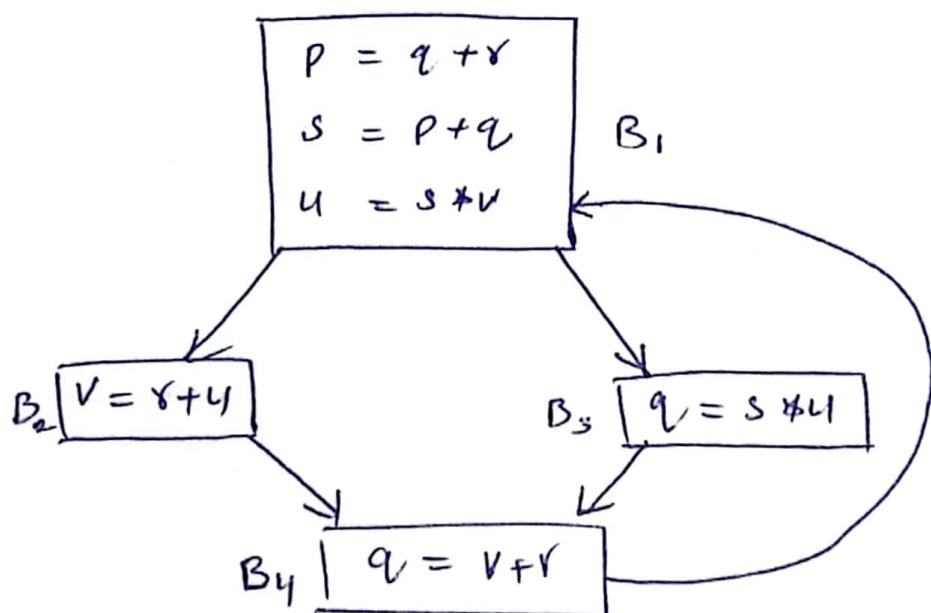
- 1) There is a statement s_j using X
- 2) There is a path from s_i to s_j
- 3) There is no new definition to X before s_j

1. $X = a+b$
2. $Y = d+c$
3. $X = X+b$
4. $a = b+d$
5. $X = a+b+c$

	a	b	c	d	X	Y
1	L	L	L	L	D	D
2	D	L	L	L	L	D
3	D	L	L	L	L	D
4	D	L	L	L	D	D
5	L	L	L	D	D	D

Example 2:

Perform liveness analysis.



→ check 3 conditions of liveness

	P	q	r	S	U	V
2	D	D	L	D	L	D
3	D	D	L	L	L	L

Example 3:

Assume that all operations take their operands from register, what is the minimum # of registers needed to execute this code without spilling?

$$\begin{aligned}a &= 1 \\b &= 10 \\c &= 20 \\d &= a+b\end{aligned}$$

↓
Storing variables in
main memory when
of register are
allocated to variables
too much!

$$\begin{aligned}e &= c+d \\f &= c+e \\b &= c+e \\e &= b+f \\d &= 5+e\end{aligned}$$

- a) 2 b) 3 c) 4 d) 6

Hint: perform liveness analysis.

Sol:

$$\begin{array}{ll} \cancel{a=1} & r_1 \\ \cancel{b=10} & r_2 \\ \cancel{c=20} & r_3 \\ \cancel{d=a+b} & r_1 \quad // 'a' is dead, reclaim 'r_1' \\ \cancel{e=c+d} & r_1 \quad // 'd' is dead, reclaim 'r_1' \\ \cancel{f=c+e} & r_2 \quad // 'b' is dead \\ \cancel{b=c+r_3} & r_1 \\ \cancel{e=b+f} & r_1 \\ \cancel{d=5+e} & r_1 \end{array}$$

Registers used = r_1, r_2, r_3

so, minimum 3 are required

Peephole optimization (machine dependent)

performed during code generation

a) Redundant LOAD/STORE elimination:

$$\begin{aligned} a &= b + c \\ d &= a + e \end{aligned}$$

LOAD R₀, b

LOAD R₀, c

STORE a, R₀

LOAD R₀, b

ADD R₀, c

STORE a, R₀] useless instructions

LOAD R₀, a

STORE d, R₀

b) Flow of control optimization:

① Avoid jumps on jumps

L₁: jump L₂

≡

L₂: jump L₃

≡

L₃: jump L₄

≡

L₄: x = a + b/c

② Eliminate dead code

L₁: jump L₄

≡

L₄: x = a + b/c

int i = 0;

if(i == 1)

print("COMSATS")

c) Use of M/C idioms:

$$i = i + 1;$$

LOAD $R_{o,i}$

ADD R₀, #1 →

STORE i, R.

INC i