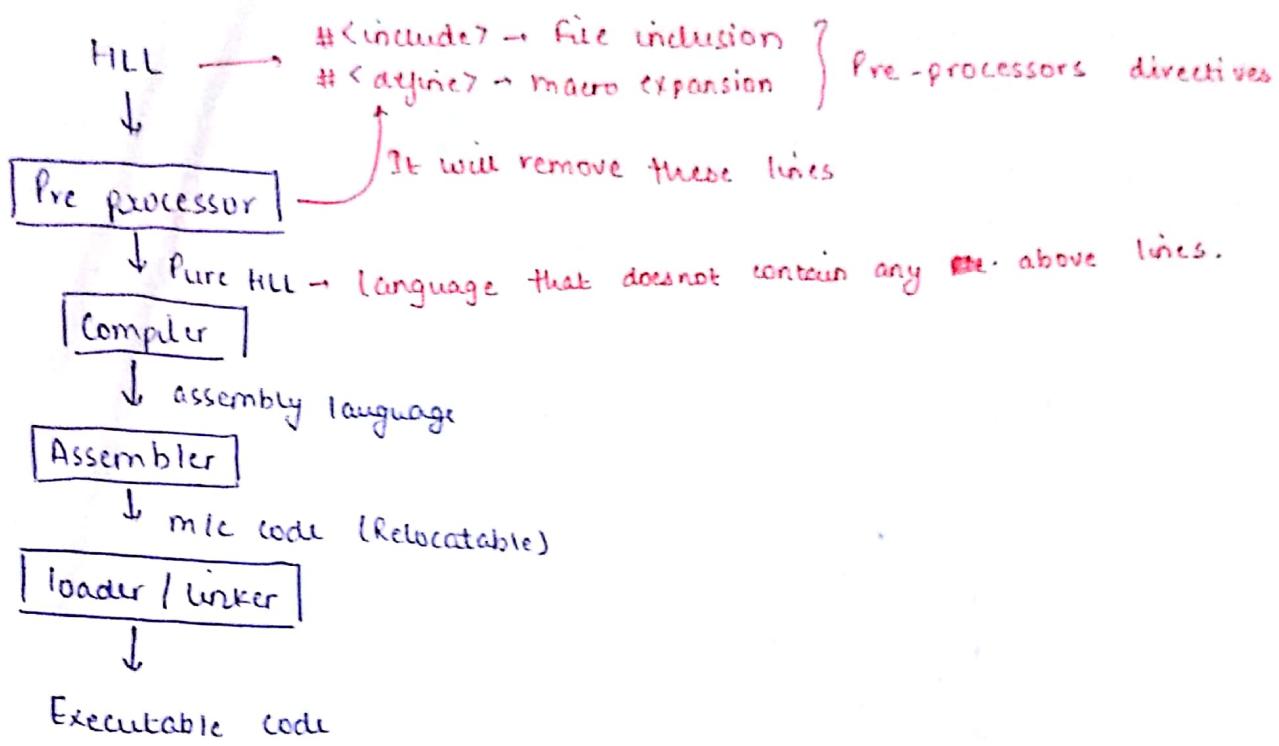


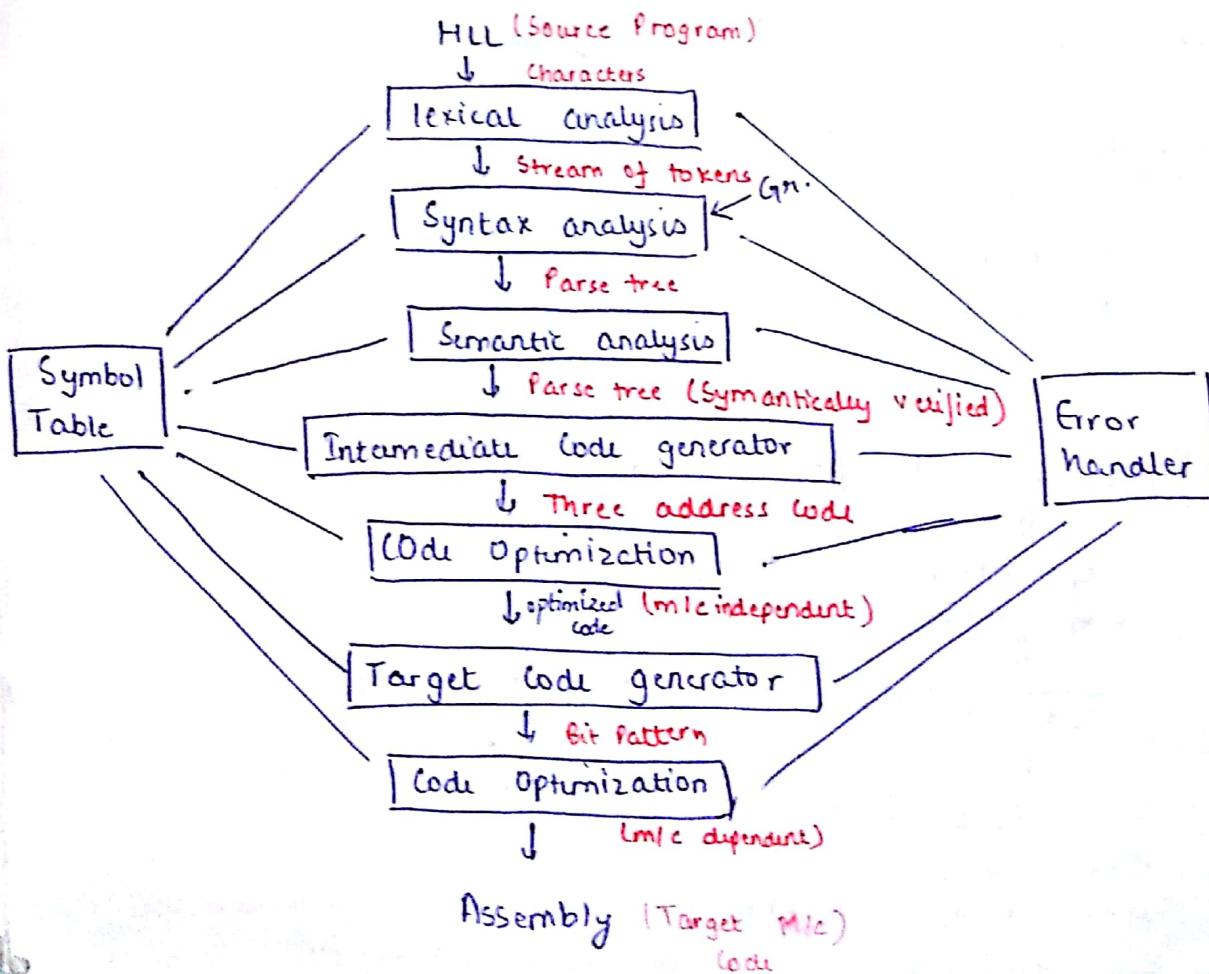
# Compiler Construction

Introduction and various phases of compiler

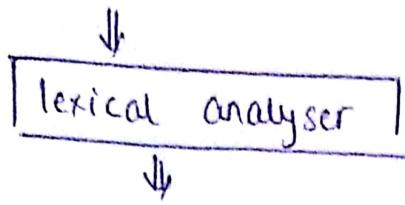
Converts HLL to LLL



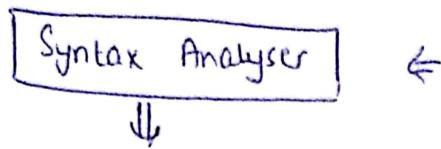
→ Phases of Compiler



IP:  $a = a + b * c$ ; [Character or lexemes]

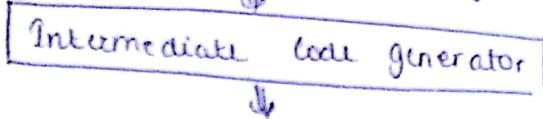
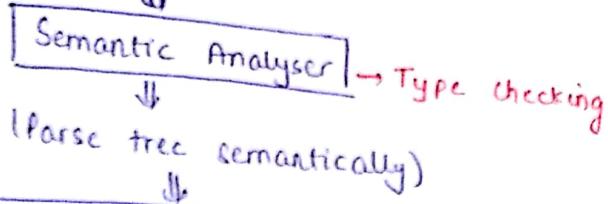
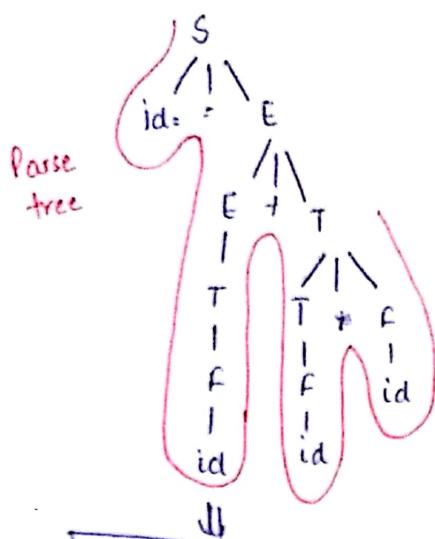


$id = id + id * id$  [Stream of tokens].  
↓



$$\begin{array}{l} S \rightarrow id = E ; \\ E \rightarrow E + T / T \\ T \rightarrow T * F / F \\ F \rightarrow id \end{array}$$

$$\begin{array}{l} E = Expr \\ T = Ter \\ F = Factor \end{array}$$



$t_1 := b * c ;$   
 $t_2 := a + t_1 ;$   
 $n := t_2 ;$

↓

Code Optimization  
↓ → To reduce no. of writes

$t_1 := b * c ;$   
 ~~$n := a + t_1 ;$~~

↓

Target Code Generator  
↓

multiplying  
first 2 storing  
d into R2

Mul	R1, R2		a → R0
Add	R0, R2		b → R1
Mov	R2, X		c → R2

That assembler  
understands.

$$\begin{array}{l}
 \left. \begin{array}{l}
 E \rightarrow E + F \\
 F \rightarrow E \\
 F \rightarrow F
 \end{array} \right\} * \rightarrow * \quad (\text{Defined left associative}) \\
 + \& + \quad (\text{Defined right associative})
 \end{array}$$

### Left Recursion

$$A \rightarrow A\alpha / B$$

$$A \rightarrow BA'$$

$$A' \rightarrow \alpha A' / \beta \in$$

Ambiguous grammar X.

Parser (Kind of syntax analyzer).

Top-down Parser (TDP)

TDP with backtracking

TDP without Backtracking

Recursive descent

considering all the possibilities  
if any poss. goes wrong it will go back to previous poss.

Non-Recursive descent

LL(1)

Bottom-up Parser (BUP)  
(Parser)

Operator Precedence

LR Parser

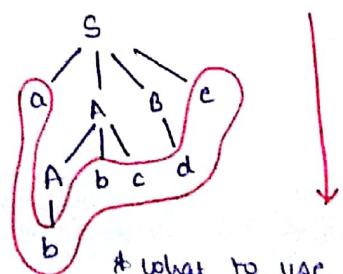
LR(0) SLR(1) LALR(1) CLR(1)

- ① \* Without Backtracking  $\rightarrow$  generate the string exactly what is required.
- ② ↴ Left Recursion X
- ↳ ND X (Non-deterministic).

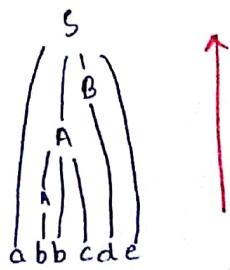
Eg,

$$\begin{array}{l}
 S \rightarrow aABe \\
 A \rightarrow Abc/b \\
 B \rightarrow d
 \end{array}$$

w: abbcde



Eg,



\* whom to read

## \* First & follow

		<u>First</u>	<u>Follow</u>
①	$S \rightarrow ABCDE$	$\{a, \epsilon\} - \epsilon \cup \{b, \epsilon\} - \epsilon \cup \{c\} = \{a, b, c\}$	$\{\$\}$
	$A \rightarrow a/\epsilon$	$\{a, \epsilon\}$	$\{b, \epsilon\} - \epsilon \cup \{\$\} = \{b,$
	$B \rightarrow b/\epsilon$	$\{b, \epsilon\}$	$\{c\}$
	$C \rightarrow c$	$\{c\}$	$\{d, \epsilon\} - \epsilon \cup \{\$\} = \{d, \$\}$
	$D \rightarrow d/\epsilon$	$\{d, \epsilon\}$	$\{e, \epsilon\} - \epsilon \cup \{\$\} = \{e, \$\}$
	$E \rightarrow e/\epsilon$	$\{e, \epsilon\}$	$\{\$\}$

②

$S \rightarrow Bb/Cd$	$\{a, b\} \cup \{c, d\} = \{a, b, c, d\}$	$\{\$\}$
$B \rightarrow aB/\epsilon$	$\{a, \epsilon\}$	$\{b\}$
$C \rightarrow cc/\epsilon$	$\{c, \epsilon\}$	$\{d\}$

③

$E \rightarrow TE'$	$\{id, ()\}$	$\{\$, )\}$
$E' \rightarrow +TE'/\epsilon$	$\{+, \epsilon\}$	$\{\$, )\}$
$T \rightarrow FT'$	$\{id, ()\}$	$\{+, \epsilon\} - \epsilon \cup \{\$, )\} = \{+, \$, )\}$
$T' \rightarrow *FT'/\epsilon$	$\{* , \epsilon\}$	$\{+, \$, )\}$
$F \rightarrow id / (E)$	$\{id, ()\}$	$\{*, \epsilon\} - \epsilon \cup \{+, \$, )\} = \{*, +, \$, )\}$

④

$S \rightarrow ACB \mid CbB \mid Ba$	$\{d, g, h, \epsilon\}$	$\{\$\}$
$A \rightarrow da/Bc$	$\{d, g, h, \epsilon\}$	$\{h, \epsilon\} - \epsilon \cup \{g, \$\} \text{ (crossed out)}, \$, h, g, s$
$B \rightarrow g/\epsilon$	$\{g, \epsilon\}$	$\{\$, a, g, h\}$
$C \rightarrow h/\epsilon$	$\{h, \epsilon\}$	$\{b, \$, g, h\}$

## Construction of LLL(1) Parsing Table

③,

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

## Recursive Descent Parser

$$E \rightarrow iE'$$

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

```

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

```

```

match (char t)
{
    if (x == t)
        x = getchar();
    else
        printf("error");
}

main()
{
    E();
    if (x == '$')
        printf("Parsing successful");
}

```

## → Operator Precedence Parser

### \* Operator grammar

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

No two operands should be adjacent.

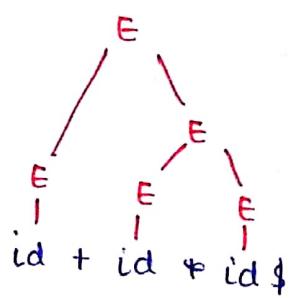
Eg,  $E \rightarrow EAE / id$   
 $A \rightarrow + / *$

X

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

- (4) '\$' & '\$' → Not compared.
- (5) '+' and '\*' always  $\Rightarrow$  because it is left associative.
- (6) '\*' & '\*' always  $\Rightarrow$  —
- (7) '^' & '^' always  $\Leftarrow$  because right associative.

\* less → Push  
 More → Pop.



\$	id	+	id	*	id	\$
----	----	---	----	---	----	----

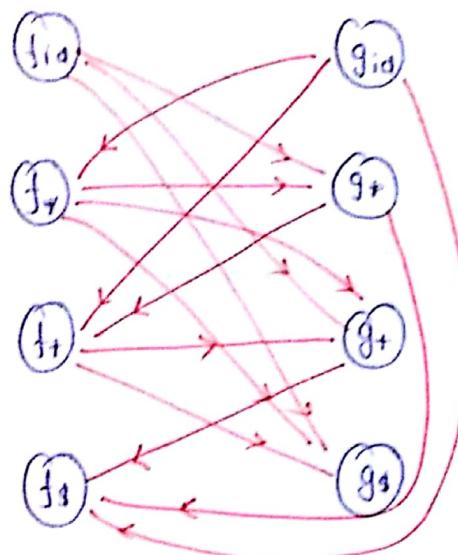
\* We will compare stack with string always. and if precedence is low we will push otherwise pop.

operator precedence	
id	
-	
*	
^	
\$	

## tion precedence Parser

	id	+	*	\$	]g
id	-	>	>	(r)	
+	<	>	<	>	
*	<	>	>	>	
\$	<	<	<	-	

This basically shows the arrow's direction.



$f_id \rightarrow g_+ \rightarrow f_+ \rightarrow g_{\ast} \rightarrow f_{\$}$   
 $g_id \rightarrow f_+ \rightarrow g_{\ast} \rightarrow f_{\$} \rightarrow g_+$   
 $\downarrow$   
 $f_{\$}$

	id	+	*	\$
f	4	2	4	0
g	5	1	3	0

- Function Table.

\* If there is a loop we will stop finding the path.

- ①  $P \rightarrow SR/S \rightarrow ①$
- ②  $R \rightarrow bSR/bS \rightarrow ②$
- $S \rightarrow wbs/w$
- $w \rightarrow L+w/L$
- $L \rightarrow id$

The above grammar is not an operator grammar as "SR" is adjacent. So, now we will make it operator grammar.

$$P \rightarrow SbSR/Sbs/S \quad [\text{from } ②]$$

Again adjacent then,

$$\left. \begin{array}{l} P \rightarrow SbP/Sbs/S \\ S \rightarrow wbs/w \\ w \rightarrow L+w/L \\ L \rightarrow id \end{array} \right\}$$

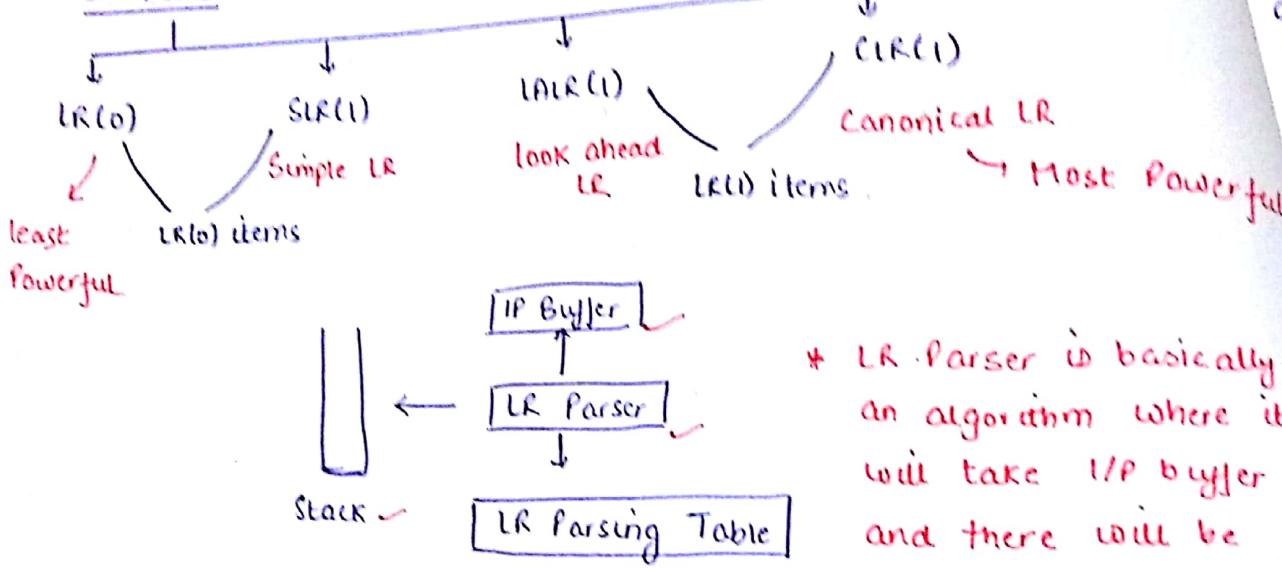
Now, it is an operator grammar.

Now, constructing operator precedence table.

	id	*	b	\$
id	→	→	→	→
*	↓	↓	↓	↓
+	↓	↓	↓	↓
b	↓	↓	↓	-
\$	↓	↓	↓	-



→ LR Parser



\* LR Parser is basically an algorithm where it will take I/P buffer and there will be stack for all L.

But the difference will be only in the construction of "LR Parsing Table".

\* Construction of LR(0) construction table

The given grammar is,

$$S \rightarrow AA$$

$$A \rightarrow aA/b.$$

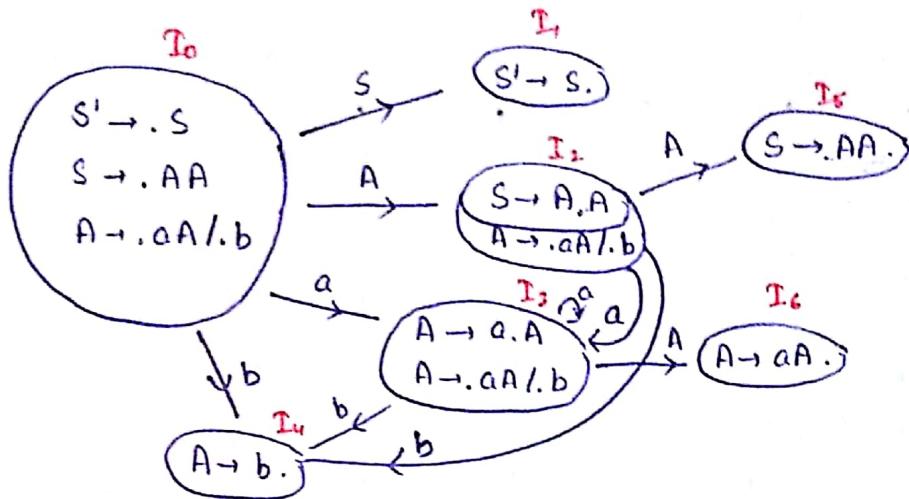
Now, we will add augmented state,

$$S' \rightarrow S$$

$$S \rightarrow AA$$

$$A \rightarrow aA/b.$$

finding the goto of productions,



goto (I<sub>0</sub>, s)  
 {  
   S' → S. } I<sub>1</sub>

goto (I<sub>2</sub>, a)  
 {  
   A → a. A  
   A → .aA/.b } I<sub>3</sub>

goto (I<sub>0</sub>, A)  
 {  
   S' → A. A  
   A → .aA/.b } I<sub>2</sub>

goto (I<sub>2</sub>, b)  
 {  
   A → b. } I<sub>4</sub>

goto (I<sub>0</sub>, a)  
 {  
   A → a. A  
   A → .aA/.b } I<sub>3</sub>

goto (I<sub>2</sub>, A)  
 {  
   S → AA. } I<sub>5</sub>

S → AA ①  
 A → aA/b ② ③

goto (I<sub>0</sub>, b)  
 {  
   A → b. } I<sub>4</sub>

goto (I<sub>3</sub>, A)  
 {  
   A → aA. } I<sub>6</sub>

## ② LR(0) Parsing Table

shift.  
 Terminal going to col with State will be the S

Shift.	Action Table			Goto Table	
	a	b	\$	A	S
0	S <sub>3</sub>	S <sub>4</sub>		2	1
1			Accept		
2	S <sub>3</sub>	S <sub>4</sub>		5	
3	S <sub>3</sub>	S <sub>4</sub>		6	
4	R <sub>3</sub>	R <sub>3</sub>	R <sub>3</sub>		
5	R <sub>1</sub>	R <sub>1</sub>	R <sub>1</sub>		
6	R <sub>2</sub>	R <sub>2</sub>	R <sub>2</sub>		

Reduce

We will find this when '.' (dot) is at the end of production (A → b.) then we will notice that in which state it is and what is the no. of that production. R →

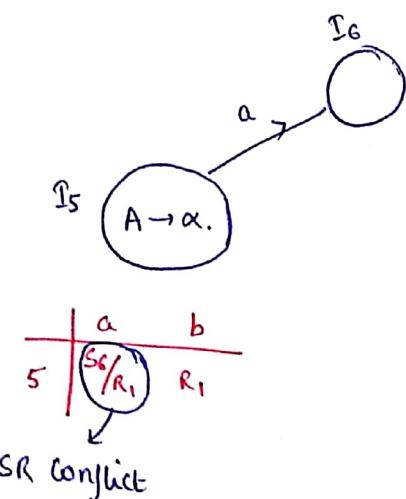
\* SLR(1) Parsing Table for the previous one.

	Action			Goto	
	a	b	\$	A	S
0	$s_3$	$s_4$		2	1
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	$r_2$	$r_2$	$r_2$		

Only where that production's left side non-terminal follows.

- \* Shift is found same procedure
- \* for Reduce in we will find where the ":" (dot) at the end then will check the no. production which for  $R_-$ . But we'll fill  $R_-$  in that block ①

- \* SLR(1) is more powerful than LR(0) because this has less no. of Reduce nodes so we this will find the error quickly and at earlier stage
- \* LR(0) is invalid when there is SR or RR conflict. (Shift-Reduce or Reduce-Reduce conflict).

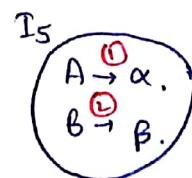


	a	b
5	$s_3/R_1$	$R_1$

SR conflict

This happens when a state is going to another state i.e. goto  $(I_5, a)$  and  $I_5$  contains the final state  $(A \rightarrow \alpha.)$

dot at the end.



	a	b	c
5	$r_1/r_2$	$r_1/r_2$	$r_1/r_2$

RR conflict

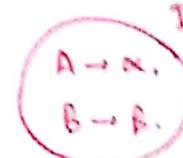
This happens when there is two final production in one state.

is invalid when the follow(A) = {a}  $\Rightarrow$  SR conflict.

↳ in case of  $A \rightarrow \alpha.$   $T_S \xrightarrow{a} I_6.$

follow(A) = follow(B)  $\neq \emptyset$ .

↳ in case of  $A \rightarrow \alpha.$   $T_S \xrightarrow{a} I_6.$

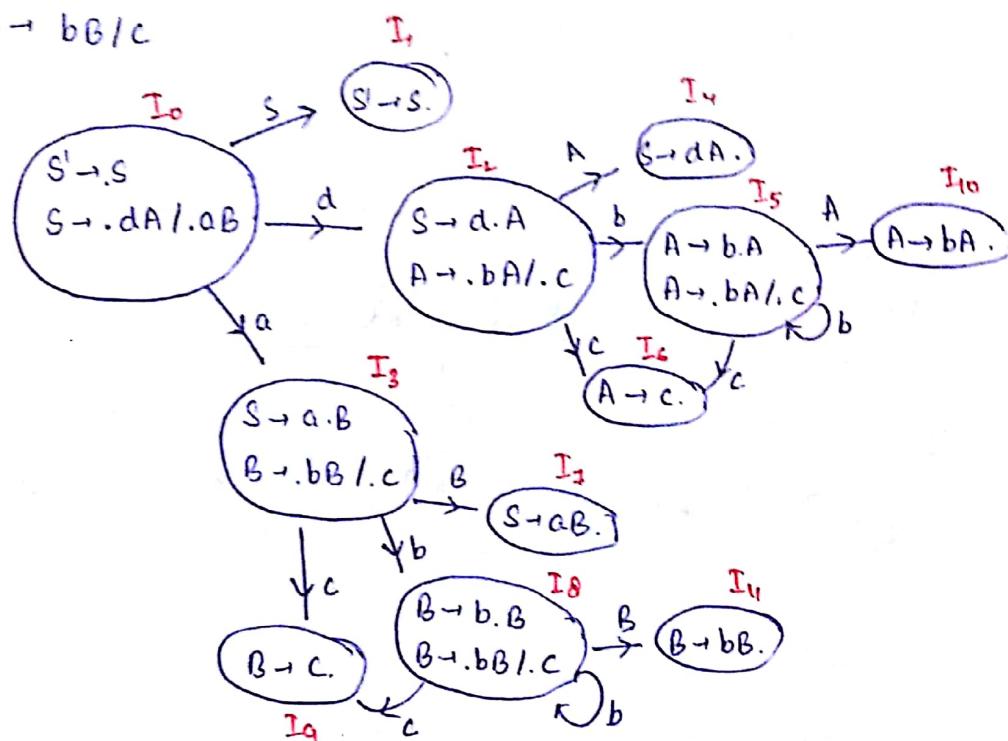


①

$$S \rightarrow dA / aB$$

$$A \rightarrow bA / c$$

$$B \rightarrow bB / c$$



\* Whether this grammar is LL(1) or not.

↳ Yes, it is because  $\left. \begin{array}{l} \text{first}(S) = \{d, a\} \\ \text{first}(A) = \{b, c\} \\ \text{first}(B) = \{b, c\} \end{array} \right\}$  There no repetition of first's.

\* LR(0) or not

↳ It's a LR(0) because there is no SR or RR conflict in all the final states.

\* SLR(1) or not

↳ It is SLR(1) because it is LR(0) & a grammar is LR(0) so it will not have any conflicts by itself.

## LR(1) and CLR(1)

LR(1) item = LR(0) items + look ahead

Eg,

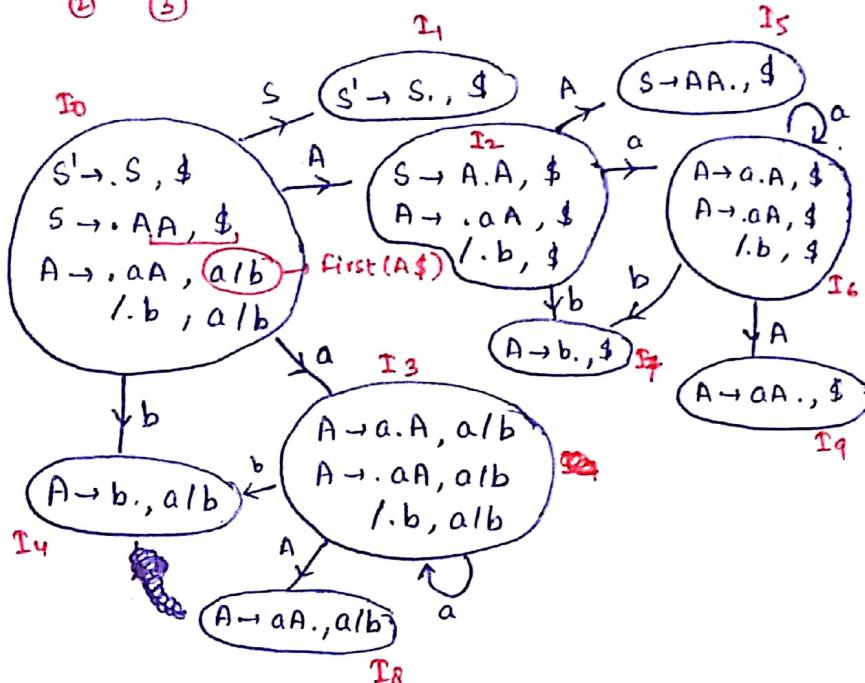
$$S \rightarrow .aA, a/b$$

### NOTE

\* For lookahead,

Difference  $A \rightarrow \alpha.BB, a/b$   
 $B \rightarrow .V, \underline{\text{First of the}} \underline{\text{remaining}}$  after B.

$$\begin{array}{l} S \rightarrow AA \quad ① \\ A \rightarrow aA/b \quad ② \quad ③ \end{array}$$



- ① By which production it comes from.
- ② After that, the first of the remaining will be the lookahead.

- \* While applying the transition lookahead will not change
- \* While applying closure LA may change

→ Shifting the dot

## CLR(1) Parsing Table

	a	b	\$	S	A
0	$s_3$	$s_4$		1	2
1			accept		3
2	$s_6$	$s_7$			5
3	$s_3$	$s_4$			3
4	$r_3$	$r_3$			
5			$R_1$		
6	$s_6$	$s_7$			9
7			$R_3$		
8	$r_2$	$r_2$			
9			$R_2$		

- ①  $I_3, I_6 \xrightarrow{I_{36}}$
  - ②  $I_4, I_7 \xrightarrow{I_{47}}$
  - ③  $I_8, I_9 \xrightarrow{I_{89}}$
- Same productions, but diff. LA.

$$CLR(1) \geq LR(0) \Rightarrow CLR(1) = SLR(1)$$

} No. of states.

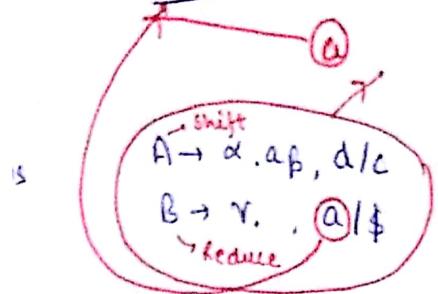
→ For LALR(1) Parsing Table,

	a	b	\$	S	A
0	$s_{36}$	$s_{47}$		1	2
1			accept		
2	$s_{36}$	$s_{47}$			5
36	$s_{36}$	$s_{47}$			89
47	$R_3$	$R_3$	$R_3$		
5			$R_1$		
89	$s_{36}$	$s_{47}$		89	(Merged)
47			$R_3$		(Merged)
89	$R_2$	$R_2$	$R_2$		(Merged)
89			$R_2$		

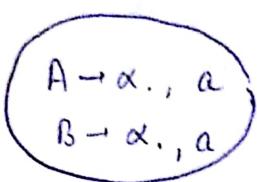
LALR(1) Parsing Table,

	a	b	\$	S	A
0	$s_{36}$	$s_{47}$		1	2
1			accept		
2	$s_{36}$	$s_{47}$			
36	$s_{36}$	$s_{47}$			5
47	$R_3$	$R_3$	$R_3$		89
5			$R_1$		
89	$R_2$	$R_2$	$R_2$		

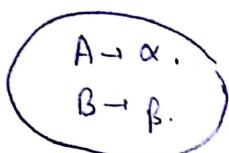
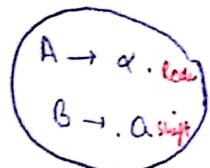
SR conflict



RR conflict



erns

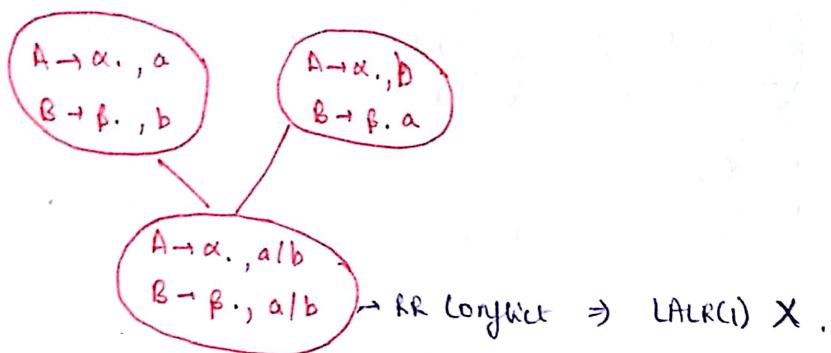


CLRL(1) ✓ then LALR(1) may or not contain list,

\* If CLRL(1) X then

LALR(1) X  
for sure

e.g.,



$$S \rightarrow AaAb / BbBa$$

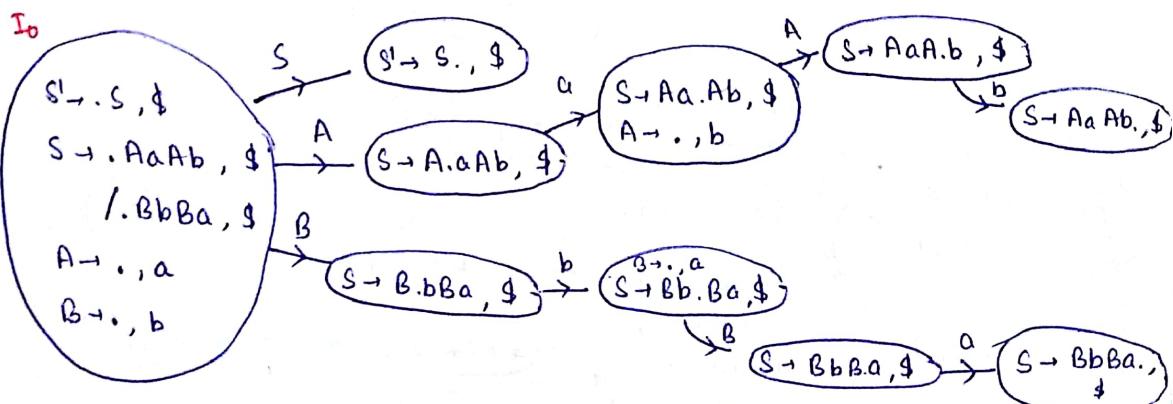
$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

LL(1) ✓  
LR(0) X  
SLR(1) X

CLRL(1) ✓

LALR(1) ✓



## Syntax Directed Translation (SDT)

Grammar + Semantic action/Rules = SDT

Transforming the tokens into structure/tree.

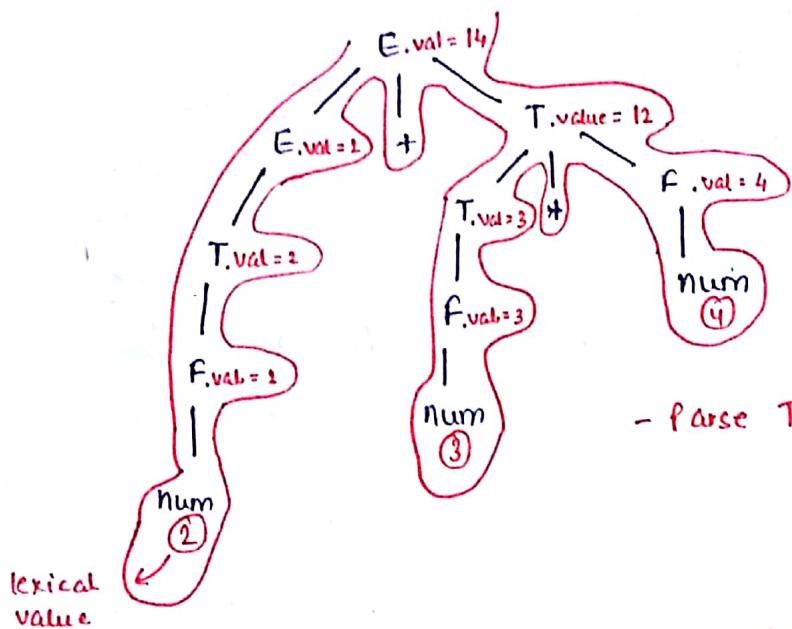
SDT for evaluation of exp.

①  $E \rightarrow E + T \quad \{ E.\text{value} = E.\text{value} + T.\text{value} \}$   
 $/ T \quad \{ E.\text{value} = T.\text{value} \}$

$T \rightarrow T * F \quad \{ T.\text{value} = T.\text{value} + F.\text{value} \}$   
 $/ F \quad \{ T.\text{value} = F.\text{value} \}$   
 $F \rightarrow \text{num} \quad \{ F.\text{value} = \text{num}. \text{value} \}$

↳ lexical value.

$$W = 2 + 3 * 4$$



Top - bottom - left - right .

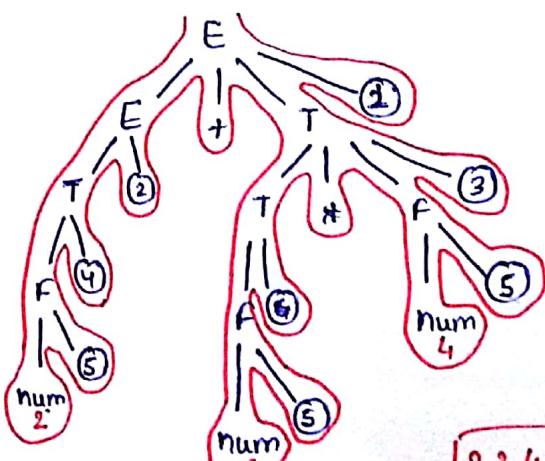
- Parse Tree.

② SDT, <sup>①</sup> semantic action.

$E \rightarrow E + T \quad \{ \text{printf}( "+ "); \}$   
 $/ T \quad \{ \} \quad \text{②}$

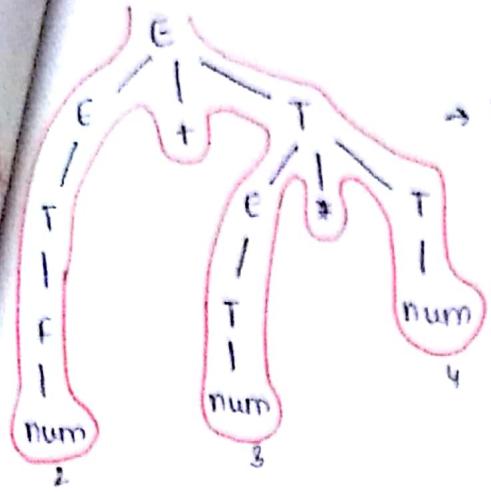
$T \rightarrow T * F \quad \{ \text{printf}( "\times "); \}$   
 $/ F \quad \{ \} \quad \text{③}$

$F \rightarrow \text{num} \quad \{ \text{printf}( \text{num}. \text{val}); \}$



Top - down

2 3 4  $\times +$   
 $2 + 3 * 4$  → Postfix exp.  
from infix.



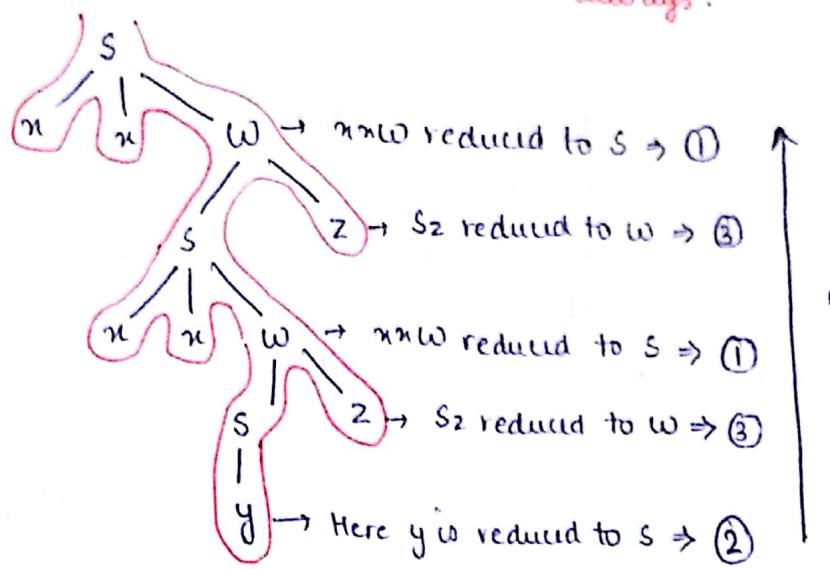
→ This is how bottom-up works.

284 \* 4  
Postfix.

$$\begin{aligned} \textcircled{3} \quad S &\rightarrow n w w \quad \{ \text{Printf}(1); \} \\ / y &\quad \{ \text{Printf}(2); \} \\ w \rightarrow S_2 &\quad \{ \text{Printf}(3); \} \end{aligned}$$

w: unnnnyyzz

\* SDT's semantic actions performed on reduction always.



$$\begin{aligned} \textcircled{2} \quad \textcircled{4} \quad E &\rightarrow E + T \quad \{ E.\text{val} = E.\text{val} + T.\text{val} \} \\ / T &\quad \{ E.\text{val} = T.\text{val} \} \end{aligned}$$

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

$$\textcircled{1} \quad / F \quad \{ T.\text{val} = F.\text{val} \}$$

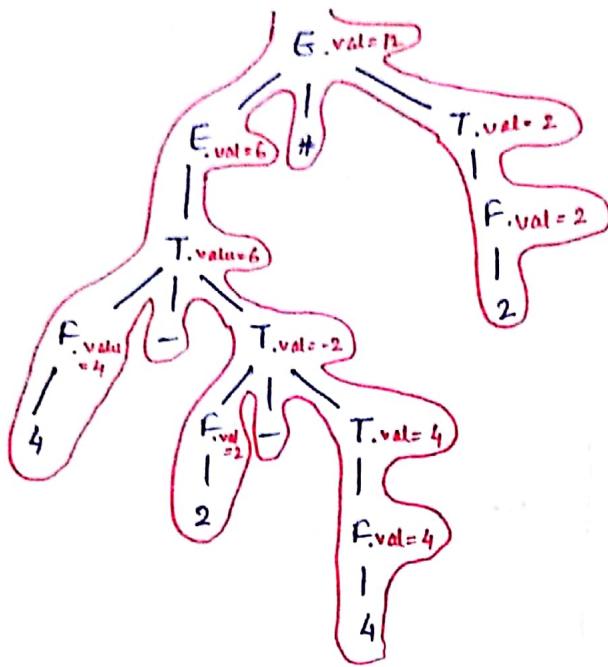
$$F \rightarrow 2 \quad \{ F.\text{val} = 2 \}$$

$$/ 4 \quad \{ F.\text{val} = 4 \}$$

$$w = 4 - 2 - 8 + 2$$

$$\Rightarrow ((4 - (2 - 4)) * 2)$$

6



$$(4) \quad E \rightarrow E \# T \quad \{ E.val = E.val \# T.val \}$$

$$/ T \quad \{ E.val = T.val \}$$

$$T \rightarrow T \& F \quad \{ T.val = T.val + F.val \}$$

$$/ F \quad \{ T.val = F.val \}$$

$$F \rightarrow \text{num} \quad \{ F.val = \text{num.lval} \}$$

$$w = 2 \# 3 \& 5 \# 6 \& 4$$

↓

$$w = 2 * (3 + 5) * (6 + 4)$$

↓

$$(2 * 8) * 10$$

$$16 * 10$$

$$= 160$$

## to build Syntax Tree

$E \rightarrow T \quad \{ \text{E.nptr} = \text{mknode} (E.\text{nptr}, '+' , T.\text{nptr}); \}$

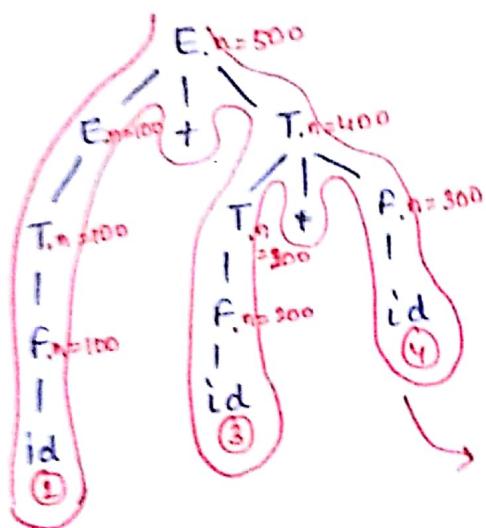
$/T \quad \{ \text{E.nptr} = \text{T.nptr}; \}$

$\rightarrow T * F \quad \{ \text{T.nptr} = \text{mknode} (T.\text{nptr}, '*' , F.\text{nptr}); \}$

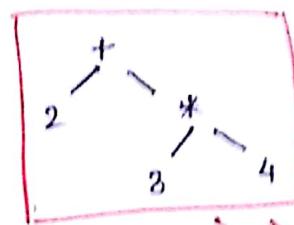
$/F \quad \{ \text{T.nptr} = \text{F.nptr}; \}$

$F \rightarrow \text{id} \quad \{ \text{F.nptr} = \text{mknode} (\text{null}, \text{id.name}, \text{null}); \}$

$$W = 2 + 3 * 4$$

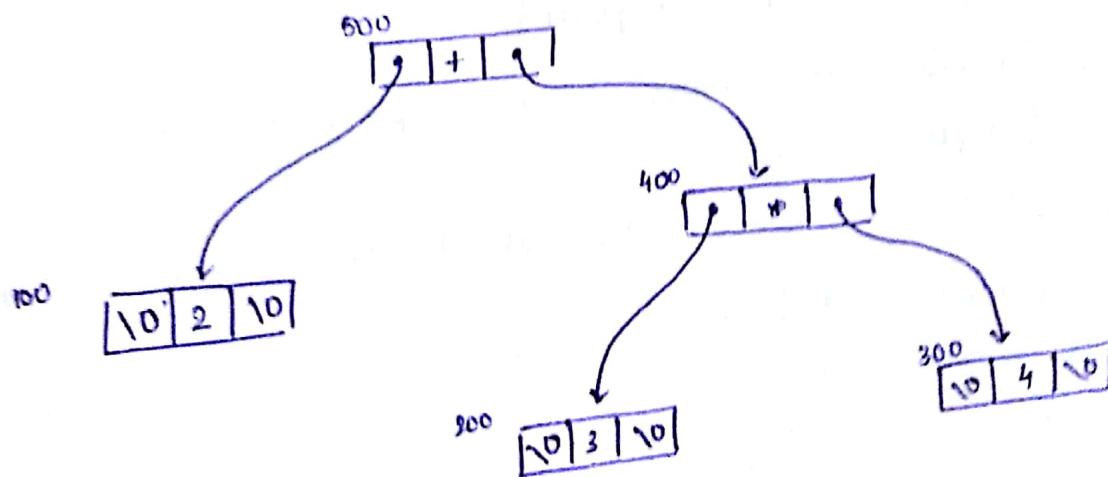


Concrete parse tree.  
(Syntax tree)



Abstract parse tree  
(Syntax tree).

Parse tree where  
variables are not  
shown.



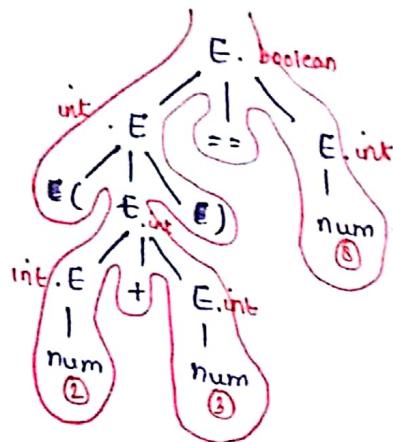
Here both are E  
but to distinguish we are writing  $E_1$  &  $E_2$

$E \rightarrow E_1 + E_2 \quad \{ \begin{array}{l} y ((E_1.type == E_2.type) \& (E_1.type = int)) \text{ then } E.type = int \\ y ((E_1.type == E_2.type) \& (E_1.type = int / bool)) \text{ then } E.type = int / bool \end{array}$

$/ E_1 == E_2 \quad \{ \begin{array}{l} y ((E_1.type == E_2.type) \& (E_1.type = int / bool)) \text{ then } E.type = int / bool \\ y ((E_1.type == E_2.type) \& (E_1.type = boolean)) \text{ then } E.type = boolean \end{array}$

$/ (E_1) \quad \{ E.type = E_1.type \}$   
 $/ \text{num} \quad \{ E.type = int; \}$   
 $/ \text{True} \quad \{ E.type = boolean; \}$   
 $/ \text{False} \quad \{ E.type = boolean; \}$

$w = (2+3) = 8$



→ SOT to generate three address code

$S \rightarrow id = E \quad \{ \text{gen } (\text{id.name} = E.place); \}$

$E \rightarrow E_1 + T \quad \{ E.place = \text{new Temp}(); \text{ gen } (E.place = E_1.place + T.place); \}$

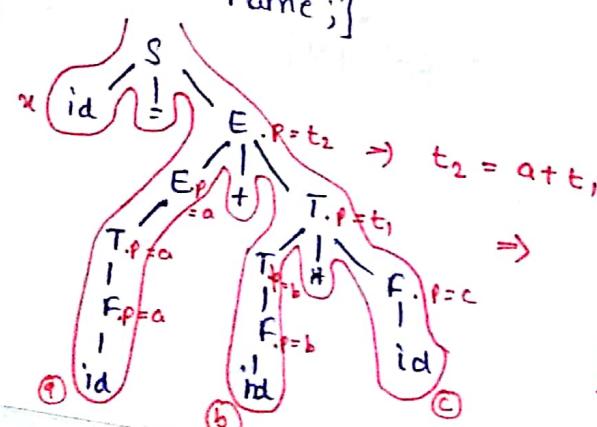
$/ T \quad \{ E.place = T.place; \text{ gen } (E.place = E_1.place + T.place); \}$

$T \rightarrow T_1 * F \quad \{ E_1.place = \text{new Temp}(); \text{ gen } (T.place = E_1.place * F.place); \}$

$/ F \quad \{ T.place = \text{new Temp}(); \text{ gen } (T.place = T.place * F.place); \}$

$F \rightarrow id \quad \{ F.place = id.name; \}$

$w: \boxed{x = a + b * c}$



Type checking

in a node take synthesized attributes  
then a node takes all inherited attributes  
when it is synthesized.  
eg

$$\begin{cases} t_1 = b * c \\ t_2 = a + t_1 \\ n = t_2 \end{cases}$$

Three address code for

$$\boxed{n = a + b * c}$$

## Synthesized attributes

When a node takes the value from its children then it is known as synthesized.

## Inherited attributes

When a node takes the value from its parent or from its siblings then it is inherited attribute.

Eg,  $A \rightarrow BCD$

$$A.s = f(B.s, C.s, D.s)$$

↳ Synthesized

$$\begin{aligned} C.i &= A.i && \text{from parent.} \\ C.i &= B.i && \text{Inherited.} \\ C.i &= D.i && \text{from siblings} \end{aligned}$$

### S-attributed SDT

- ① Uses only synthesized attr.

- ② Semantic actions are placed at right end of production

$$A \rightarrow BCC\{$$

\* Postfix SDT

- ③ Attributed are evaluated during Bottom-up parsing.  
\* Post order traversing.

### L-attributed SDT

- ① Uses both synthesized as well as inherited attr. But each inherited attr. is restricted to inherit from either parent or from left sibling only.

Eg,  $A \rightarrow XYZ$  ✓ ✗  
 $4.s = A.s$ ; ✓  $4.s = X.s$ ,  $4.s = Z.s$

- ② Semantic actions are placed anywhere on RHS.

$$\begin{aligned} A &\rightarrow \{BC \\ &\quad / D\} \\ &\quad / FG\{ \end{aligned}$$

- ③ Attributes are evaluated by traversing parse tree depth first, left, to right.  
# Pre order traversing

## Syntax Directed Translation (SDT)

As we know that after syntax analysis, the semantic actions of input program is needed.

"Syntax Directed Translation Scheme" is the process of attaching semantic actions to the context free grammar productions.

Two notations are used to associate the semantic rules to the production.

- ① Syntax Directed Definition (SDD)
- ② Translation Scheme.

① SDD is generalization of the context free grammar in which each grammar symbol either terminal or non-terminal has associated set of attributes, i.e. string, number, type, etc.

associated with the set of semantic rules of the form.

$$n := f(a_1, a_2, \dots, a_k)$$

attributes are associated with the grammar symbol i.e. A

\* Rules are associated with the production i.e.  $\alpha$ .

### \* Annotated Parse Tree

A parse tree containing the values of attributes at each node of the tree.

\* SDT can be implemented by using extra fields in the parser stack entries corresponding to grammar symbol. So, Stack is implemented by two pair of arrays. STACK & VAL.

STCK	VAL
2	2.00
4	4.00
3	3.00
1	1.00

Before

Stack Val

	Stack	Val
z		z.VAL
y		y.VAL
x		x.VAL
:		
:		
:		

Before reduction

Stack Val

	Stack	Val
A		A.VAL
:		
:		

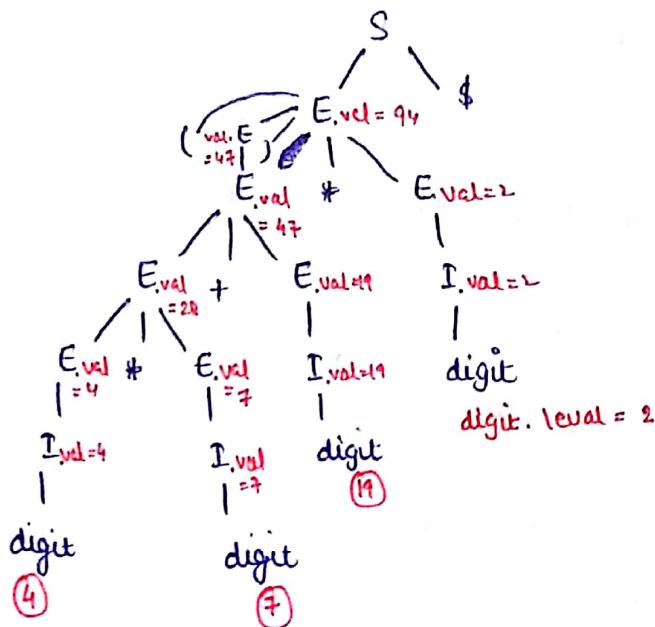
TOP →

After reduction

Eg,

$S \rightarrow E\$$	{ printf E.Val }
$E \rightarrow E^{(1)} + E^{(2)}$	{ E.Val = $E^{(1)}$ .Val + $E^{(2)}$ .Val }
$E \rightarrow E^{(1)} * E^{(2)}$	{ E.Val = $E^{(1)}$ .Val * $E^{(2)}$ .Val }
$E \rightarrow (E^{(1)})$	{ E.Val = $E^{(1)}$ .Val }
$E \rightarrow I$	{ E.Val = I.Val }
$I \rightarrow I^{(1)} \text{ digit}$	{ I.Val = $10 \times I^{(1)}$ .Val + lex val }
$I \rightarrow \text{ digit}$	{ I.Val = lex val }

(4 \* 7 + 19) \* 2 \$



## Intermediate code generation

Intermediate code  $\rightarrow$  In most of the compilers the source code is translated to a language which is in complexity between HLL & machine code such as called intermediate code.

Most common representation of source code are

- ① Postfix notation
- ② Syntax trees
- ③ Quad triples
- ④ Triples
- ⑤ Indirect triple

### ① Postfix Notation

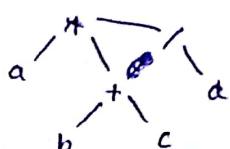
- a)  $e_1 \Theta e_2 \Rightarrow e_1 e_2 \Theta$
- b)  $(a+b) * c \Rightarrow ab+c*$
- c) If 'e' then 'n' else 'y' by using ?  $\Rightarrow eny?$
- d) If a then y c-d the a+c else a\*c else a+b

$acd-ac+ac?ab+?$  This ternary operator will come after every else statement.

### ② Syntax Tree

It is a variant to parse tree in which each leaf represents an operand and each interior node represents an operator.

Eg,  $a * (b + c) / d$ ,

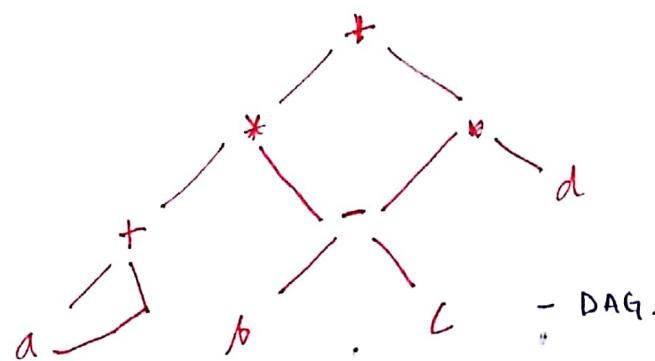
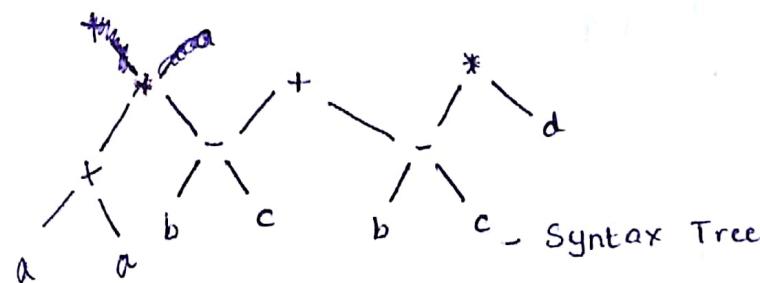


## Directed Acyclic Graph (DAG)

similar to Syntax tree having leaves corresponding to constants and interior node representing the operators.

Difference b/w DAG & Syntax tree is that a node N in a DAG has more than one parent if N represents a common subexpression in a syntax tree. means DAG represents no generation of efficient code to evaluate the expression.

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



→ Three address code (at most 3 operand per instruction and one operator in RHS). It is a linearized representation of a syntax tree or a DAG in which explicit names corresponding to the interior nodes of the graph are used.

It is a sequence of statements of the form,

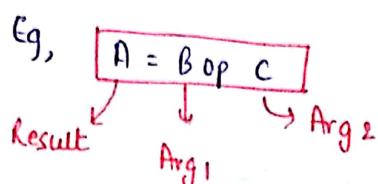
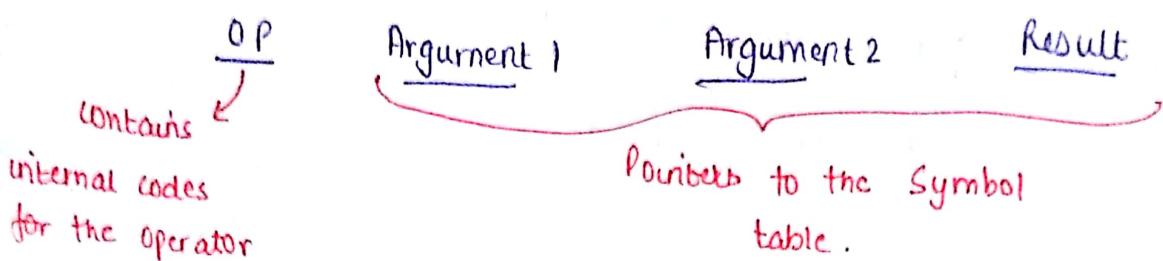
$A = B \text{ op } C$   
↓  
2 operands  
Result  
Three address code.

## Implementation of three-address code

3 address code is an abstract form of intermediate. These statements can be implemented as record for the operator & operands.

### ① Quadruples

It is a record structure using 4 fields.



Eg,

①  $A := -B * (C + D)$

$$\left. \begin{array}{l} T_1 := -B \\ T_2 := C + D \\ T_3 := T_1 * T_2 \\ A := T_3 \end{array} \right\} \text{Three - address code}$$

Now, Quadruple table,

	OP	Arg 1	Arg 2	Result
(0)	unary minus	B	-	$T_1$
(1)	+	C	D	$T_2$
(2)	*	$T_1$	$T_2$	$T_3$
(3)	$:=$	$T_3$		A

Temporary

entering

avoid temporary names into the symbol table, temporary  
we can be detected by the position of the statement  
at computes it.

o, in this, the three-address statements can be represented  
by record structure with only three fields.

OP      Arg1      Arg2

These are either pointers to the symbol table  
or into the structure itself.

Eg,  $A := -B * (C+D)$

$T_1 := -B$

$T_2 := C+D$

$T_3 = T_1 * T_2$

$A = T_3$

	OP	Arg1	Arg2
(0)	un minus	$T_1$	-
(1)	+	C	D
(2)	*	(0)	(1)
(3)	:=	A	(2)

### → Indirect Triples

consist of listing of pointers  
triples themselves in desired

order to triples rather than a listing of

Eg,  $A := -B * (C+D)$

	Statement
(0)	(9)
(1)	(10)
(2)	(11)
(3)	(12)

	OP	Arg1	Arg2
(9)	un minus	B	-
(10)	+	C	D
(11)	*	(9)	(10)
(12)	:=	A	(11)

→ Some examples of generating three-address code

### ① Assignment statement of integer type

Eg,  $a = b + c + d$

$$\Rightarrow t_1 := b + c$$

$$t_2 := t_1 + d$$

$$a := t_2$$

### ② Assignment Statement with mixed type

Eg, ①  $x = y + I * J$

Where,

$y$  and  $I$  → Real

$I$  and  $J$  → Integer

$$\Rightarrow T_1 := I * J$$

$$T_2 := \text{int to Real}(T_1)$$

$$T_3 := y + T_2$$

$$x := T_3$$

b)  $x := y * z + y * r$

Where,

$r$  and  $y$  are integers

$x$  and  $z$  are float

$$T_1 := \text{int to float } y$$

$$T_2 := \text{int to float } r$$

$$T_3 := T_1 + T_2$$

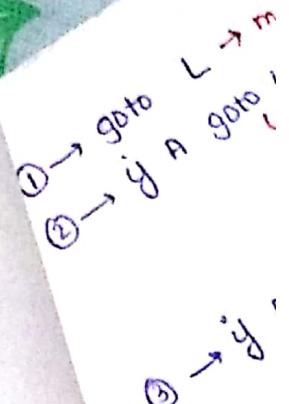
$$T_4 := T_1 * z$$

$$T_5 := T_3 + T_4$$

$$x := T_5$$

### ③ Boolean expression

For generating the three address code for boolean exp.,  
Some branching statements of the following form  
are used



(Unconditional Jump statement)

①  $\rightarrow \text{goto } L$   $\rightarrow$  means three address code statement with label  $L$  is next to be executed.

②  $\rightarrow \text{if } A \text{ goto } L$

↳ Conditional Jump Statement

(means if  $A$  is non-zero/true then statement with label  $L$  is executed otherwise next statement)

③  $\rightarrow \text{if } A \text{ relop } B \text{ goto } L$

↳ means if  $A$  relop  $B$  is true statement with label  $L$  is executed else the next statement in the normal flow sequence is executed.

3,

① If  $A < B$  then 1 else 0

Three address code is given as,

or

If  $A < B$  then  $L_1$ ,

$T = 0$

goto  $L_2$

$L_1 : T_1 = 1$

$L_2 :$

(1)  $\text{if } A < B \text{ goto } (4)$

(2)  $T = 0$

(3)  $\text{goto } (5)$

(4)  $T = 1$ .

② If  $A < B$  or  $C$  then 1 else 0.

(1)  $\text{if } A < B \text{ goto } (4)$

(2)  $T_1 = 0$

(3)  $\text{goto } (5)$

(4)  $T_1 = 1$

(5)  $T_2 := T_1 \text{ or } C$

③ If  $a > b$  then  $x = a+b$ .

(1)  $\text{if } a > b \text{ goto } (3)$

(2)  $\text{goto } (5)$

(3)  $T_1 := a+b$

(4)  $x := T_1$

(5)

$y \text{ } a > b \text{ goto } L_1$

goto  $L_2$

$L_1 : T_1 := a+b$

$x := T_1$

$L_2 :$

## → Code Optimization

This <sup>refers to</sup> techniques which is used to improve the quality of the generated object code by compiler.

It includes various transformations on intermediate representation of the program. These transformations must preserve the semantics of the program.

### \* Features

- ① Attempt from compiler to generate better object code.
- ② Code improvement
- ③ Replacing pattern by more equivalent & more efficient constraints.
  - ↳ may be local or global.
  - ↳ may be machine dependent / independent.
- ④ Inner loop optimization (90%-10% Rule)
- ⑤ Constant folding
- ⑥ Elimination of common sub-expression
- ⑦ Elimination of loop-invariant
- ⑧ Elimination of induction variables.

There are mainly 3 interrelated areas of code optimization,

- ① local optimization
  - ② loop optimization
  - ③ Data flow analysis.
- ① Which is performed in a straight line (basic block of code means a block with no jump except at the beginning & no jumps except at the end).
- ③ Transmission of useful information from all parts of the program to the places where the information can be of use.

## Blocks and flow graph

Graph representation of three address statement is called flow graph where the nodes in flow graph represents the computation and the edges represents the flow of control. Flow graph can be used to collect information about the intermediate code and to find the inner loop where a program is expected to spend most of its time.

### Basic Block

It is a sequence of consecutive statements in which

- flow of control can only enter the block in the beginning. No jumps are allowed in the middle of the block.
- flow of control can leave the block at the end without halt on possibilities of branching ex

### Steps to construct basic block

① We first determine the set of leaders. Rules are

- The first statement is a leader.
- Any statement which is a target statement of conditional or unconditional goto statement is a leader.
- Any statement which immediately follows the conditional goto is a leader.

② From one leader statement to the just prior statement of the next leader (if no such leader statement is present then upto the last statement) will form a basic block.

\* If there is any statement which has not been included to any basic block can be removed & discarded.

Eg,  
①

```
Begin
    PROD := 0;
    I := 1;
    do begin
        PROD := PROD + A[I] * B[I];
        I := I + 1;
    end
    while I ≤ 20.
end
```

- Program Segment

*m = 0;  
= 10;  
while i <= m  
{ sum =  
i = i;  
}  
out:  
TAC:*

(1) PROD := 0

(2) I := 1

B<sub>1</sub>

(3) T<sub>1</sub> := 4 \* I

(4) T<sub>2</sub> := addr(A) - 4

(5) T<sub>3</sub> := T<sub>2</sub>[T<sub>1</sub>]

(6) T<sub>4</sub> := addr(B) - 4

(7) T<sub>5</sub> := T<sub>4</sub>[T<sub>1</sub>]

(8) T<sub>6</sub> := T<sub>3</sub> \* T<sub>5</sub>

(9) PROD := PROD + T<sub>6</sub>

(10) I := I + 1

(11) if I ≤ 20 goto (3)

B<sub>2</sub>

- Three address code for a machine with 4 bytes / word. (Basic Block).

PROD := 0

I := 1

B<sub>1</sub>

T<sub>1</sub> := 4 \* I

T<sub>2</sub> := addr(A) - 4

T<sub>3</sub> := T<sub>2</sub>[T<sub>1</sub>]

T<sub>4</sub> := addr(B) - 4

T<sub>5</sub> := T<sub>4</sub>[T<sub>1</sub>]

T<sub>6</sub> := T<sub>3</sub> \* T<sub>5</sub>

PROD := PROD + T<sub>6</sub>

I := I + 1

if I ≤ 20 goto (3)

B<sub>2</sub>

\* The relationship among basic block is represented

- FLOW GRAPH

To block beginning with the statement following Statement (11)

```

Sum = 0 ;
i = 10 ;
while i <= 10 do
{
    Sum = Sum + a[2*i];
    i = i + 1;
}
average = sum / i;

```

TAC id  $\Rightarrow$

- (1) Sum := 0 ]  $B_1$ ,
- (2) i := 10 ]  $B_1$ ,
- (3) if  $i \leq 10$  then goto (11) ]  $B_2$ ,
- (4)  $T_1 := 2 * i$
- (5)  $T_2 := a[T_1]$
- (6)  $T_3 := \text{Sum} + T_2$
- (7)  $\text{Sum} := T_3$
- (8)  $T_4 := i + 1$
- (9)  $i := T_4$
- (10) goto (3)
- (11)  $T_5 := \text{Sum} / i$
- (12) average :=  $T_5$ . ]  $B_4$

### local Optimization

A basic block computes a set of expressions. Two basic blocks are said to be equivalent if they compute the same set of expressions.

Optimization can be achieved by applying no. of transformation within each blocks. Two classes of local transformation can be applied to basic blocks.

- ① Structure preserving transformation (common subexpression elimination)
- ② Algebraic transformation .

## ① Structure Preserving Optimization

### A) Common sub-expression elimination

Eg, consider the following basic blocks.

$$\begin{array}{ll} 1) a = b + c & 1) a = b + c \\ 2) b = a + d & \Rightarrow \quad 2) b = a + d \\ 3) c = b + c & \quad \quad \quad 3) c = b + c \\ 4) d = a + d & \quad \quad \quad 4) d = b \end{array}$$

Transformed block as 2 statements with the same exp. exp.

Now,  $a =$

\* But 1<sup>st</sup> and 3<sup>rd</sup> statements are same with the exp. takes the different value of b.

### B) Dead Code Elimination

A variable is said to dead if its value is not used anywhere in the program. Subsequently an assignment statement made to the dead variable can be safely removed.

$x = t_3 \Rightarrow \text{dead code.}$

$a[t_2] = t_5$

$a[t_4] = t_2$

Print  $a[t_4]$

### C) Copy Propagation

Let we have following expression,

$$\begin{array}{ccc} a = d + e & & b = d + e \\ & \searrow & \swarrow \\ & \text{new variable } c = d + e & \end{array}$$

Here the common subexpression is  $c = d + e$  is eliminated using new variable to hold the value of  $d + e$ .

## Constant folding

during at compile time that the value of an expression is a constant and using the constant instead is known as constant folding.

Eg,  $a = 36 * 2$   
print a

Here,  $a = 72$  will be computed at compile time.

## ② Algebraic Transformation

Many algebraic transformations can be used to change the set of expressions computed by the basic block into simplified expressions or to replace the expensive operation into cheaper operations.

Eg, a)  $n = n + 0$  and  $x = x * 1$   
 $\downarrow$   
 $n = x$                              $\downarrow$   
 $n = x$

b)  $n = y^{\text{power}}$  can be replaced by  $n = y * y$ .

c)  $2 * n$   
 $\downarrow$   
 $n + n$

d)  $x / 2$   
 $\downarrow$   
 $x + 0.5$

e)  $x^2$   
 $\downarrow$   
 $x * x$ .

## → loop Optimization

- |   |                                      |
|---|--------------------------------------|
| ① loop invariant code $\Rightarrow$ Elimination | } Structure preserved transformation |
| ② Induction variable elimination.               |                                      |

Reduction in Strength  $\rightarrow$  Algebraic transformation

## b) Induction Variable elimination

Induction variables are those variables whose values in the loop are of the form  $I = I + c$  where  $c$  is a constant. Elimination of induction variables also performs strength reduction.

Eg,

L:  $T_1 := 4 + I$   
 $T_3 := T_2 [T_1]$   
 $I := I + 1$   
if  $I \leq 20$  goto L

$\Rightarrow$

L:  $T_1 = 0$   
 $T_1 := T_1 + 4$   
 $T_3 := T_2 [T_1]$   
if  $T_1 \leq 76$  goto L

Note,

$T_1$  and  $I$  are in the lock stage. When we increase the value of  $I$  like 1, 2, ..., 20 then the value of  $T_1$  also increases from 4, 8, ..., 80. So,  $I$  and  $T_1$  are the induction variables.

Eg,

## \* Reduction in Strength

Replacement of an expensive operation by a cheaper one.

Eg,  
 $n=1$   
for ( $i=1; i \leq 10; i++$ )  
{      $n=i+2;$   
}

Can be optimised as,

$n=0$   
for ( $i=1; i \leq 10; i++$ )  
{      $n = i + \frac{i}{2};$   
}  
//

## Controlled Acyclic Graph

useful data structure for analysing basic block.

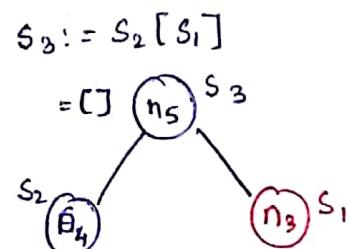
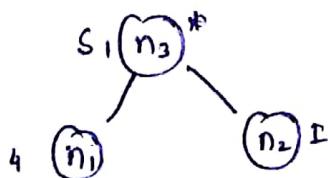
It is a directed graph with no cycle.

Determines common sub-expressions in a basic block.

Determines which names are used inside but evaluated outside of the block.

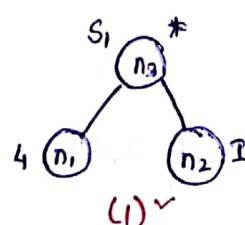
5) Determines which statement of the block could have their values outside the block.

Eg,  $S_1 := 4 * I$

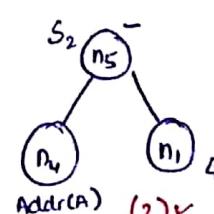


## Constructing a DAG example

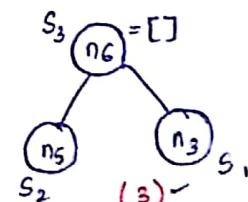
(1)  $S_1 := 4 * I$



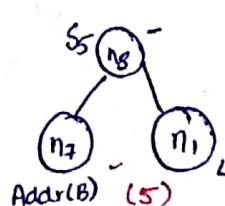
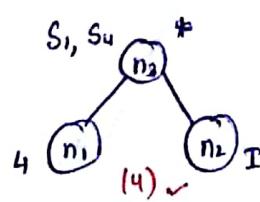
(2)  $S_2 := \text{Addr}(A) - 4$



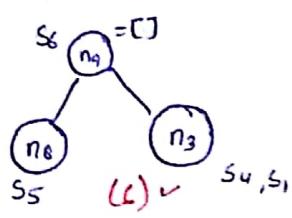
(3)  $S_3 := S_2 [S_1]$



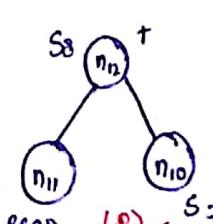
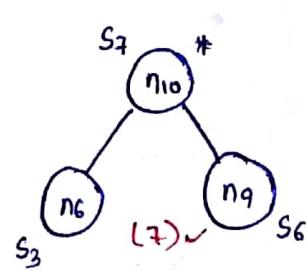
(4)  $S_4 := 4 * I$



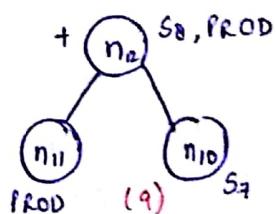
(5)  $S_5 := \text{Addr}(B) - 4$



(6)  $S_6 := S_5 [S_4]$

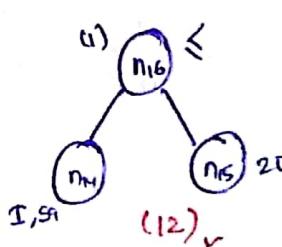
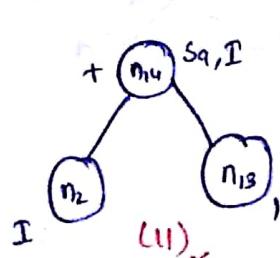


(7)  $S_7 := S_3 + S_6$



(8)  $S_8 := \text{PROD} + S_7$

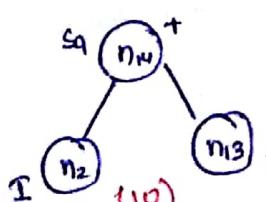
(9)  $\text{PROD} := S_8$



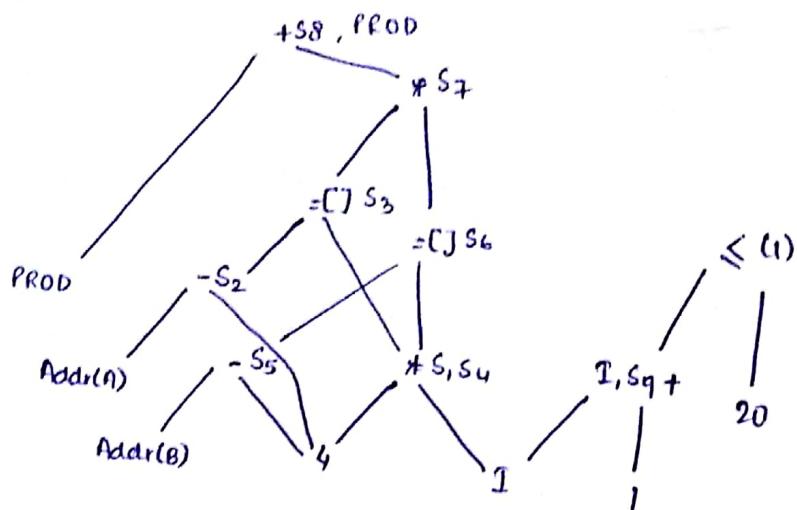
(10)  $S_9 := I + 1$

(11)  $I := S_9$

(12) if  $I \leq 20$  goto (1)



Now,  
Constructing DAG,



→ Basics

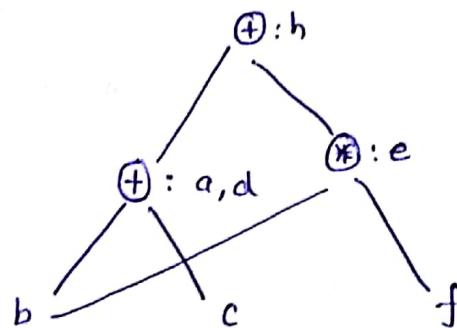
(1)

$$a = b + c$$

$$d = b + c$$

$$e = b * f$$

$$h = a + e$$

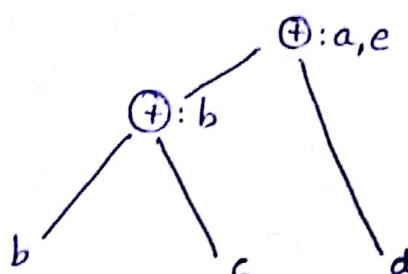


(2)

$$b = b + c$$

$$a = b + d$$

$$e = a$$



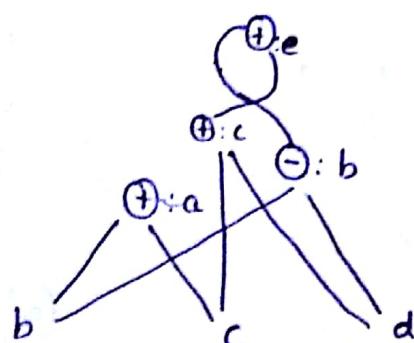
(3)

$$a = b + c$$

$$b = b - d$$

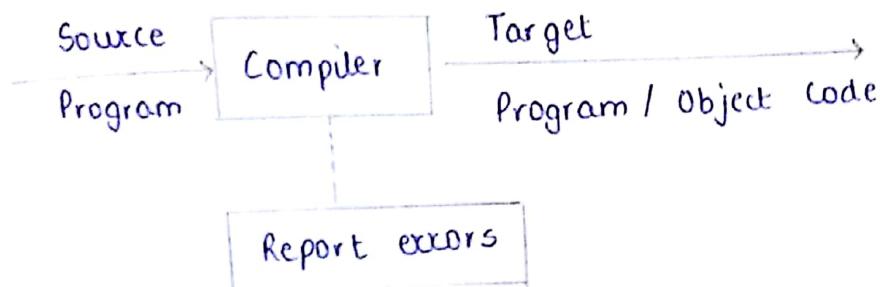
$$c = c + d$$

$$e = b + c$$



What do you mean by compiler? Write different types of compilers.

A compiler is a program that reads a program written in one language known as source language and translates into an equivalent program into another language called target language compiler also reports errors to the user.



→ Types of Compilers.

① Native code compiler: The compiler used to compile a source code for the same type of platform only is known as native code compiler.

Eg, ocamplc

② Cross compiler: A compiler that runs on one machine and produces code for another machine is called a cross compiler. It compiles for multiple platforms from one machine.

Eg, GCC collections of compiler can be set up to cross compile.

to source compiler: The compilers that take a high-level language as input and output source code of another high-level language only. It is also known as a transpiler.



Eg, script sharp, Cfront.

④ One pass compiler: It is a type of compiler that compiles the whole process in only one pass.

Eg, early PASCAL compilers.

⑤ Threaded code compilers: The compilers which simply replace a string by an appropriate binary code.

Eg, BASIC compiler.

⑥ Incremental compiler: The compiler which compiles only the changed lines from the source code and update the object code.

Eg, IBM visual Age C++ compiler 4.0.

⑦ Source compiler: The compiler which converts the source code high-level language code into assembly language only.

plain analysis Synthesis model of a compiler. Also plain analysis tools.

In analysis synthesis model of a compiler, compilation process is divided into two parts:

1. Analysis
2. Synthesis

### ① Analysis

It breaks up the source program into constituent pieces and creates an intermediate representation of the source program. During analysis, the operations implied by the source program are determined and recorded in a syntax tree. The analysis part also collects the important information about the source program and stores it in a data structure called symbol table, which is passed to synthesis part.

### ② Synthesis

It constructs desired target program by using intermediate code and information stored in the symbol table provided by the analysis part.

Various analysis tools are:

#### a) lexical analysis: analyzer,

It reads the streams of characters from the source program and groups it into meaningful sequences called lexemes which produces tokens of the form,  
< token name, attribute value >

Lexical Analyzer,  
in, characters or tokens are grouped hierarchically  
noted collections. The parser uses the tokens  
produced by lexical analyzer to create syntax tree.

Semantic Analyzer,  
It uses the syntax tree and information in the symbol  
table to check the source program for semantic  
consistency. It also gathers type information, performs  
type checking and some implicit conversions called  
coersions.

the following,

table,

symbol table is a data structure containing a record for each identifier with attribute details of the identifier. When an identifier in a source program is detected by lexical analyzer, the identifier is entered into the symbol table. The remaining phase enter the information about identifier into the ST and then use their information in various ways.

## ② Error detection and Reporting,

Each phase of the compiler encounters errors. The lexical analyzer phase can detect errors only where the characters remaining in the input do not form any token of the language. Syntax error detects errors when token streams violates the structural rules. Semantic analyzer tries to detect the errors for the programmatical constructs that have correct syntactic structure but no meaning.

## ③ Text formatter

A text formatter takes input that is a stream of characters, most of which is text to be typeset but some of which includes commands to indicate paragraphs, figures or mathematical structures like subscripts.

Preprocessor,

It is a program that processes its input data to produce output that is used as input to another program. They may perform the following functions.

1. Macro processing
2. File inclusion
3. Rational preprocessor
4. Language extensions.

b) Assembler,

It writes object code by translating assembly instruction mnemonics into machine code. They may be one-pass assembler or two-pass compiler.

c) Linkers,

A linker is a program that takes one or more objects generated by a compiler and combines them into a single executable program.

d) Loader

A loader is a part of the OS that is responsible for loading programs in memory.

## Compiler Construction Tools

Development tools that are available for implementation of one or more compiler phases are known compiler construction tools. These are as follows;

### a) Parsing Generator,

they produce syntax analysis normally from context free grammar input.

Eg, PIC, EGM.

### b) Scanner Generator,

These generators automatically generate lexical analyzer, normally based on a regular expression.

### c) Syntax-directed Translation Engines,

They generate intermediate code with three add format based on the parse tree.

### d) Automatic Code Generator,

It takes a collection of rules that defines the translation of each operation of the intermediate language for the target machine.

### e) Data flow engines,

It gathers informations about how values are transmitted from one part of the program to each other part.

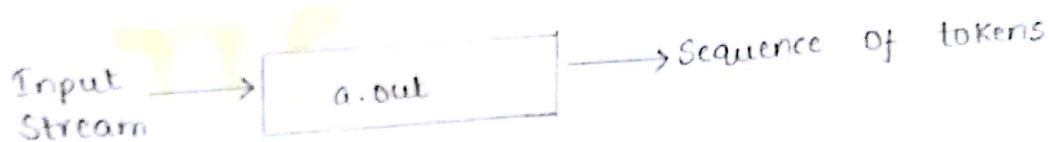
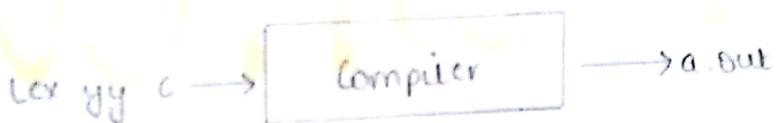
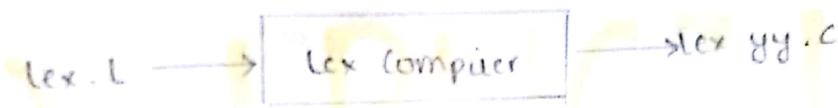
## case study on LEX.

LEX:

LEX is a tool in lexical analysis phase to recognise tokens using a regular expression.

LEX tool itself is a lex compiler.

① A LEX source program is a specification of lexical Analyser consisting of set of regular expressions together with an action for each regular expression.



\* lex.l: is an input file written in a language which describes the generation of lexical analyser.  
lex compiler transforms lex.l to c program known as lex.yy.c

lex.yy.c :- is compiled by the c-compiler to file called a.out.

The attribute value can be numeric code, pointer to symbol table or nothing.

## lex programs:

Program will be in the following form

Declarations

Translation Rules

%%

Auxiliary functions

a) Declarations: This section includes declaration of variables, constants and regular definitions.

b) Translation Rules: It contains regular expression and code segments.

form: Pattern { Action }

c) Auxiliary Functions: holds additional functions which are used in actions. These are compiled separately and loaded with Lexical Analyser.

### \* Conflict Resolution in lex :

Conflict arises when several prefixes of input matches one or more patterns. This can be resolved by the following,

- ① Always prefer a longer prefix than a shorter prefix.
- ② If two or more patterns are matched for the longer prefix, then the first pattern listed in lex program is preferred.

## lexical Analyser

A lexical analyser can either be generated by NFA or by DFA. It is preferable in the implementation of Lex.

Distinguish between top-down and bottom-up  
parsing.

### Top- Down

starts at the root and goes to the leaves.

- ① It performs a left most derivation.
- ② It ends when the stack is empty.
- ③ It expands non-terminals.
- ④ It reads terminals when popped off stack.
- ⑤ Pre Order traversal of Parse tree.
- ⑥ Pops the non-terminals and pushes it's RHS onto stack.
- ⑦ The stack is used for what is expected.
- ⑧ It may requires back tracking.
- ⑨ It requires grammar free of left recursion & left factoring.

### Bottom - Up

- ① Starts at the leaves and goes to the root.
- ② It performs Right Most derivation in reverse.
- ③ It starts when the stack is empty.
- ④ It reduces non-terminals.
- ⑤ It reads terminals when they're pushed on stack.
- ⑥ Post Order traversal of Parse tree.
- ⑦ It tries to recognise RHS. Pops it and pushes corresponding non-terminal.
- ⑧ The stack is used for what is already seen.
- ⑨ It doesn't require back tracking.
- ⑩ It doesn't require such a grammar.

the advantages and disadvantages of top-down  
bottom-up parsing.

### Top-down Parsing:

- Advantages,
- ① Highly flexible
- ② Efficient error recovery.

### Disadvantages,

- ① Semantic actions cannot be performed while making a prediction. The actions must be delayed until the prediction is known to be a successful parse.
- ② Precise error reporting is not possible. A mismatch merely triggers back tracking. A source string is known to be erroneous only after all predictions have failed.

### Bottom-up Parsing:

#### Advantages,

- ① Attribute computation is easy.
- ② Parser can be generated automatically.

#### Disadvantages,

- ① It doesn't support left/right information flows.