



Compiler - GATE

* Symbol table is updated during lexical analysis, syntax analysis and semantic analysis

* Top down parser - Leftmost Derivation
Classification

- Predictive Parser
1. LL(1)
 2. Recursive Descent Parsing
 3. Non Recursive Predictive Parsing
 4. Backtracking

* Bottom up parser - Reverse of Rightmost Derivation
also called Shift Reduce (SR) Parser

Classification

1. Operator Precedence Parser - no 2 consecutive variables
2. LR(K) parsing \rightarrow LR(0), SLR(1), LALR(1), CLR(1)

(Shift Reduce Parser's)

* $LR(0) \subseteq SLR(1) \subseteq LALR(1) \subseteq CLR(1)$
 \downarrow
 also called LR(1)

High level language

LL(1)

Lexical Analyzer

Syntax Analyzer

Semantic Analyzer

Intermediate code generation

Code optimizer

Target Code Generation

Assembly Code / Machine code

- INPUT - Character Stream
- read the program character by character
- INPUT - Token stream
- INPUT - Syntax tree
- INPUT - Intermediate Representation

DFA minimization
Key word recognition
Regular Expression
Finite Automata

PDA

Type checking

Dataflow Analysis

DAG

Register Allocation

Syntax Tree

Graph Back coloring and

Expression evaluation

Abelian Group

Type checking

Production Tree

Post order traversal
Push Down Automata

Syntax Directed Translation
Parsing

PTD

Activation record
Leftmost derivation

Runtime environment

Top Down Parsing

Symbol Table - Info about variables and their attributes



* LL(1) Parser / Predictive Parser

- most widely used Top down parser
- ~~also~~ parses the given input string with the help of LL(1) parsing table
- two functions First() and Follow() are used to construct LL(1) parsing table.

~~also~~ FIRST(X) for a grammar symbol X , is the set of terminals that begin the strings derivable from X .

Rules to compute FIRST set -

1. If x is a terminal, then $FIRST(x) = \{ 'x' \}$
2. If $x \rightarrow \epsilon$ is a production rule, then add ϵ to $FIRST(x)$.
3. If $X \rightarrow Y_1 Y_2 Y_3 \dots Y_n$ is a production
 1. $FIRST(x) = FIRST(Y_1)$
 2. If $FIRST(Y_1)$ contains ϵ , then
$$FIRST(x) = \{ FIRST(Y_1) - \epsilon \} \cup \{ FIRST(Y_2) \}$$
 3. If $FIRST(Y_i)$ contains ϵ for all $i=1$ to n , then add ϵ to $FIRST(x)$

FOLLOW(X) to be the set of terminals that can appear immediately to the right of Non-Terminal X in some sentential form. \rightarrow any string derivable from start symbol

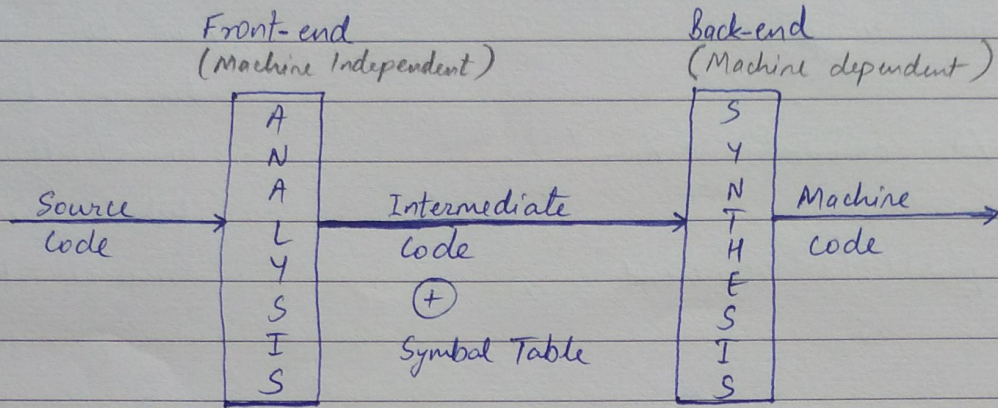
Rules to compute FOLLOW set -

1. If B is start symbol of G \rightarrow Grammar
$$FOLLOW(B) = \{ \$ \}$$
 - $FOLLOW(A) \subseteq FOLLOW(B)$
 - However reverse may not be true
2. If $A \rightarrow pB$ is a production, then everything in $FOLLOW(A)$ is in $FOLLOW(B)$
3. If $A \rightarrow pBq$ is a production, where p, B and q are any grammar symbols, then everything in $FIRST(q)$ except ϵ is in $FOLLOW(B)$
4. If $A \rightarrow pBq$ is a production and $FIRST(q)$ contains ϵ , then $FOLLOW(B)$ contains $\{ ~~FIRST(q)~~ - \epsilon \} \cup FOLLOW(A)$
 ϵ never comes in FOLLOW set



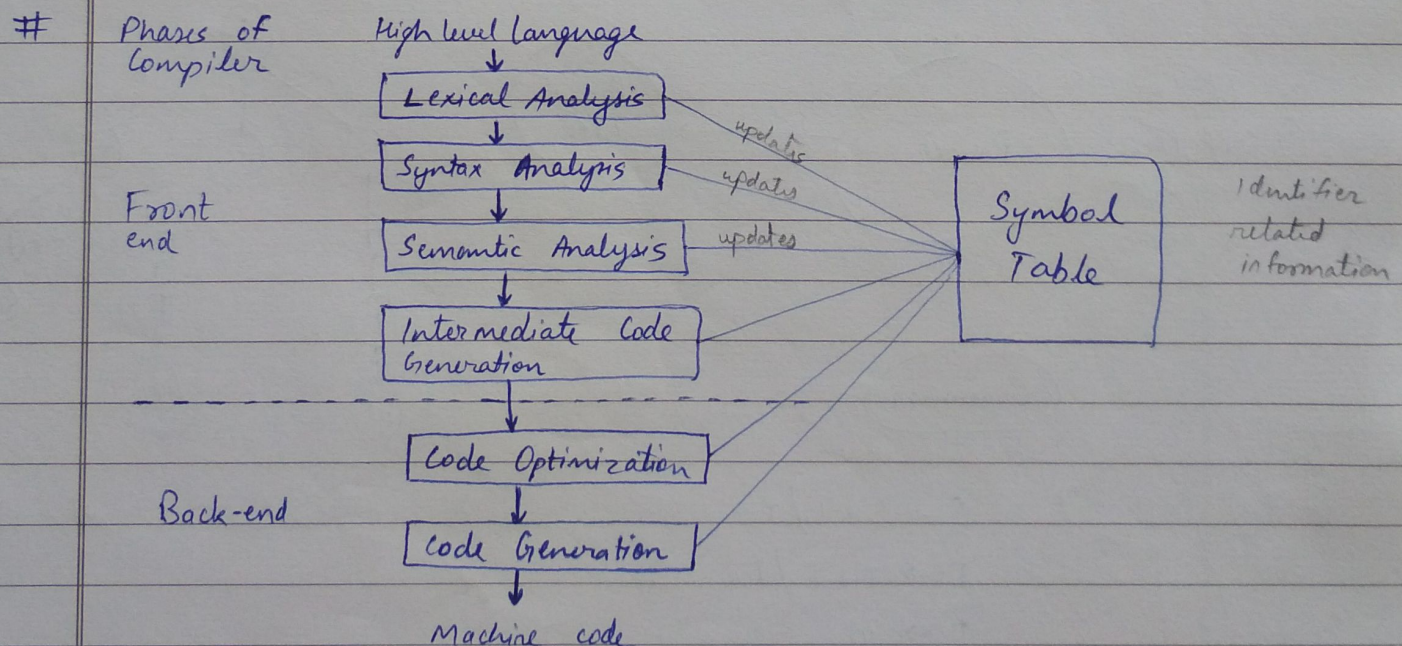
Compiler (Unit 1)

Compiler design Architecture



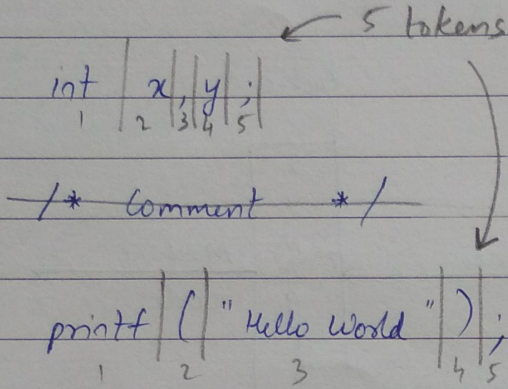
- Lexical Analysis
- Syntax Analysis
- Semantic Analysis
- Intermediate Code Generation
- Code Optimization
- Code Generator

Compiler - without changing the meaning of the high level language, it converts the source code into machine code





Lexical Analysis



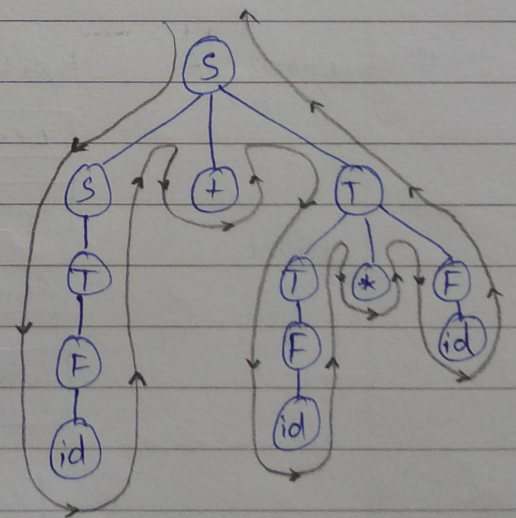
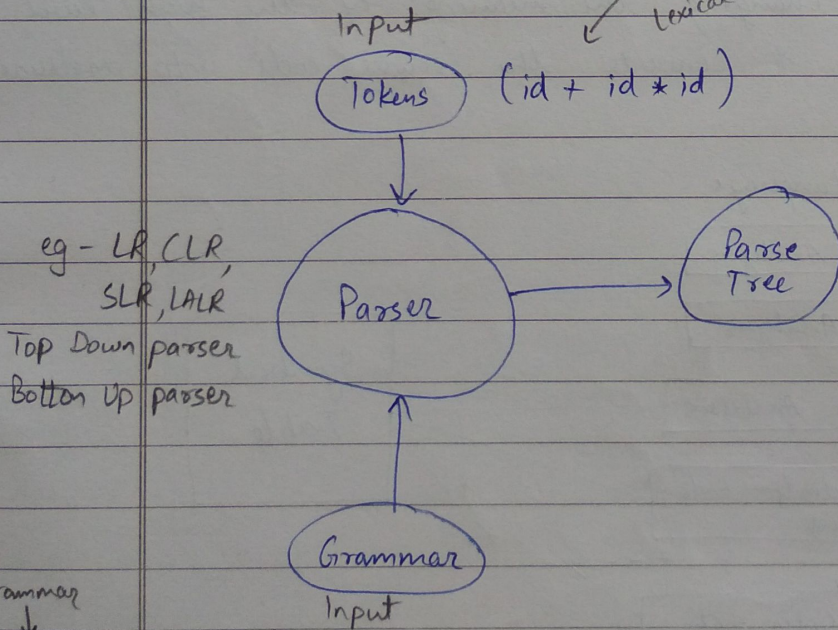
Lexical Analyser -

- Read the High Level language (HLL) line by line
- Discards all comments
- White space elimination
- Generate tokens

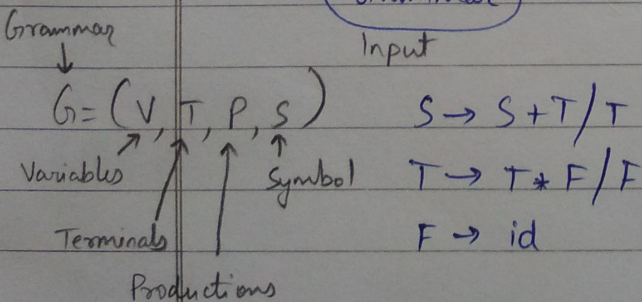
Syntax Analysis

Syntax Analyser / Parser - evaluation of parse tree must map/match to the input tokens. Evaluation is top-down, left-to right

← Tokens from Lexical Analyser



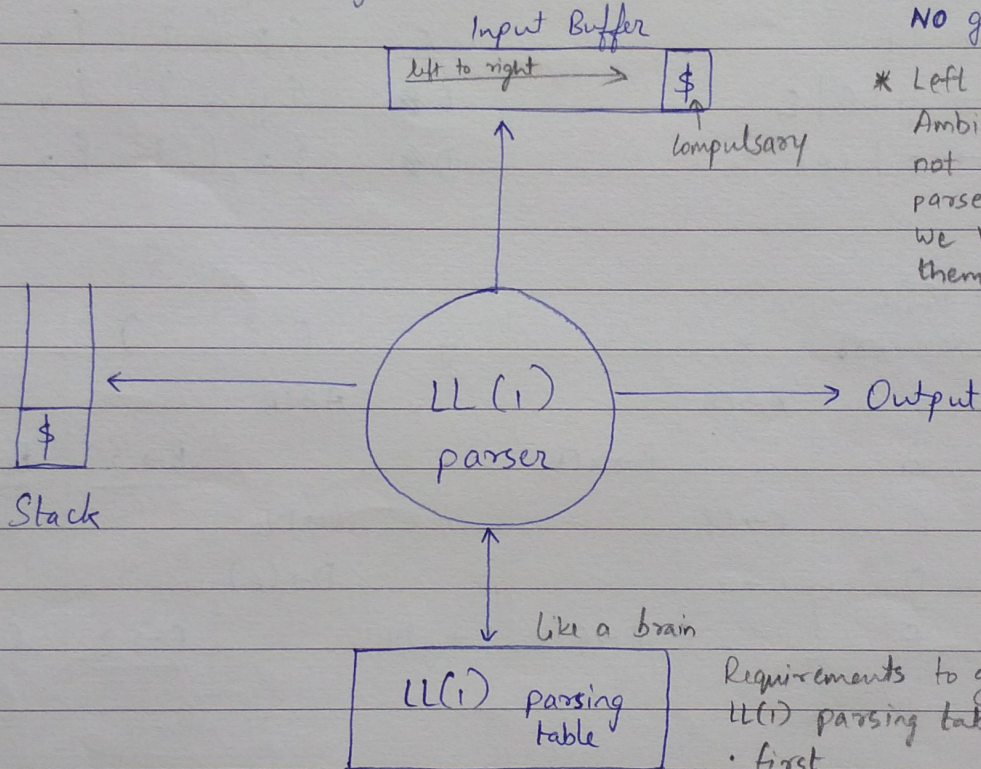
evaluation = id + id * id



only 1 symbol lookahead

LL(1) Parser / Predictive Parser

- Top down parser
- Left to right parsing
- Left most derivation - left most variable is expanded first
- follows LL grammar - subset of CFG, under LL grammar NO grammar is ambiguous



* Left recursion and Ambiguous grammar is not allowed for LL(1) parser. We have to first remove them and then solve

Requirements to generate LL(1) parsing table

- first
- follow

Always Variable

#

First(α) and Follow(α) $\xrightarrow{\text{output}}$ set of terminals only

It is a set of terminal symbols that begin in strings derived by " α "

$$A \rightarrow ab|de|gh$$

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

It is a set of terminal symbols that appear immediately to the right of " α "

TIP - always try to locate α on the RHS of other productions

$$S \rightarrow aAb$$

$$A \rightarrow c$$

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

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

$$\text{For } A \rightarrow \alpha B, \text{ Follow}(A) \subseteq \text{Follow}(B)$$

Follow of start symbol



Construction of LL(1) parsing ~~table~~ table

Production rule	Variables	First	Follow
$A \rightarrow CB$	A	{ z, (}	{ \$,) }
$B \rightarrow xCB \epsilon$	B	{ x, \epsilon }	{ \$,) }
$C \rightarrow DE$	C	{ z, (}	{ x,), \$ }
$E \rightarrow yDE \epsilon$	E	{ y, \epsilon }	{ x,), \$ }
$D \rightarrow z (A)$	D	{ z, (}	{ x, y,), \$ }

Variables ↓ Terminals →	z	x	y	()	\$
A	A → CB			A → CB		
B		B → xCB			<u>B → ε</u>	<u>B → ε</u>
C	C → DE			C → DE		
D	D → z			D → (A)		
E		E → ε	E → yDE		<u>E → ε</u>	<u>E → ε</u>

Synchronizing entries - those entries which can be obtained with the help of FOLLOW set. All synchronizing entries are underlined in the above table.

- All productions of "A" will be in the row of "A". Same applies for other variables.
- Each production will go into the column of terminals which are present in the set generated using $first(x)$ where x is the first variable/terminal of that respective production. Same applies for all production of that particular variable's productions, and for all variables as well, eg- for $A \rightarrow CB$, row is A & column = $first(C)$.
- For production like $B \rightarrow \epsilon$, first of ϵ is ϵ .
 ∴ We use $follow(B)$ to get the column for ~~B~~ $B \rightarrow \epsilon$
 We do the same for $E \rightarrow \epsilon$

$A \rightarrow X_1 X_2$
 ↑ ↑
 row FIRST(x)
 columns

#

How to check if a Grammar is LL(1) or not

- For a Grammar to be LL(1), each cell or LL(1) parsing table should have only one production.
- If a cell has two productions in it, then it is ambiguous and NOT a LL(1) grammar.
- Example - the table on the previous page has only one production in each cell. Hence the respective grammar is a LL(1) grammar.
- Example - the following grammar is not LL(1) grammar.

Productions	Variable	first	follow
$X \rightarrow Y \mid a$	X	{a}	{ $\$$ }
$Y \rightarrow a$	Y	{a}	{ $\$$ }

Terminals \rightarrow Variables \downarrow	a	$\$$
X	$X \rightarrow Y$ $X \rightarrow a$	
Y	$Y \rightarrow a$	

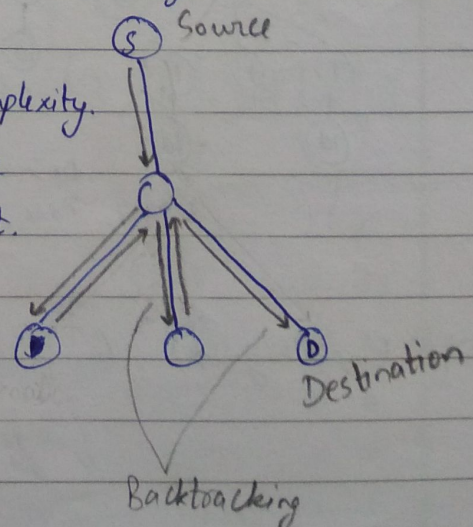
As two productions in one cell, this is not a valid LL(1) parse table

#

Recursive Descent Parser

- Top-down
- Left to right
- Disadvantage - exponential time complexity.
- Uses `match(...)` function to see if the path chosen is correct or not.

Backtracking





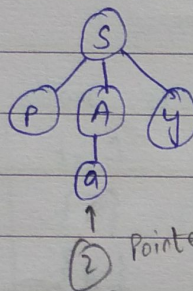
• example -

Grammar:

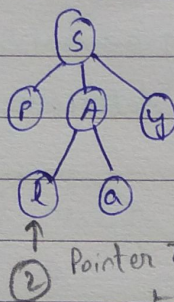
$S \rightarrow pAy$

$A \rightarrow a \mid la$

Input: play
↑
① pointer 1



match(...) returns false. Hence backtrack



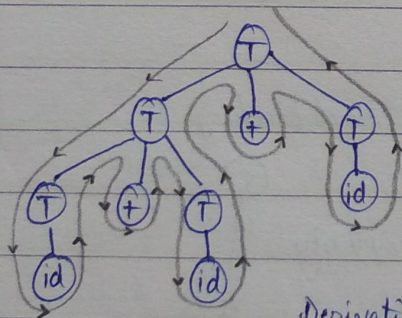
match returns true. Hence proceed

#

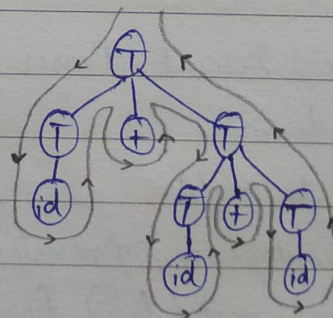
Ambiguous Grammar - when there exists more than 1 derivation tree for a given valid input string.

Production: $T \rightarrow T+T \mid id$

Input: id + id + id



Derivation Tree 1



Derivation Tree 2

#

Unambiguous Grammar - there exists a unique derivation tree for every given valid input string.

- * Write \$ at the end of input string
- * Default top of stack is \$



Operator Precedence Parser

- Bottom-up parser
- Reduction operation is performed on the input string to get the starting symbol

• Restriction on Operator Precedence Grammar -

(1) ϵ should **not** be on RHS of the productions

(2) 2 non-terminals (variables) can **NOT** be adjacent to each other

$$\left. \begin{array}{l} A \rightarrow AB \\ A \rightarrow BB \end{array} \right\} \text{Not valid} \quad \left\{ \begin{array}{l} S \rightarrow SAS/a \\ A \rightarrow bSb/b \end{array} \right.$$

↓ fix the problem

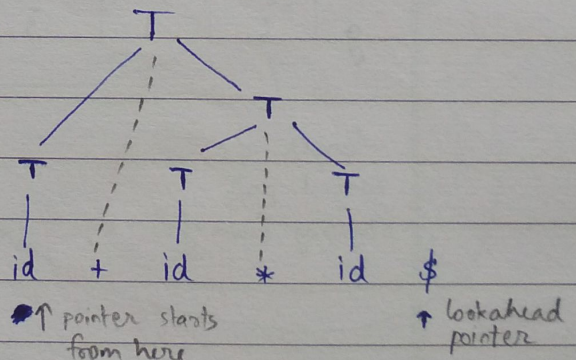
$$\left. \begin{array}{l} T \rightarrow T+T/T * T/id \end{array} \right\} \text{valid} \quad \left\{ \begin{array}{l} S \rightarrow S b S b S/a/S b S \\ A \rightarrow b S b/b \end{array} \right.$$

Operator Precedence Relation

$$T \rightarrow T+T \mid T * T \mid id$$

Precedence: $id \triangleright * \triangleright + \triangleright \$$

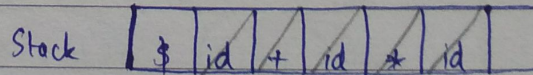
Terminals = $\{ id, +, *, \$ \}$



Disadvantage of this relation - its size is n^2

	id	+	*	\$
id	-	\triangleright	\triangleright	\triangleright
left associative +	\triangleleft	\triangleright	\triangleleft	\triangleright
left associative *	\triangleleft	\triangleright	\triangleright	\triangleright
\$	\triangleleft	\triangleleft	\triangleleft	-

parse tree successfully built



if (precedence of Stack symbol \triangleleft precedence of operator pointed by lookahead pointer) { push } increment pointer

if (precedence of Stack symbol \triangleright precedence of operator pointed by lookahead pointer) { pop } reduce the popped element

Successfully built parse tree ← stack \$ ← pointer

New Stack Comparisons

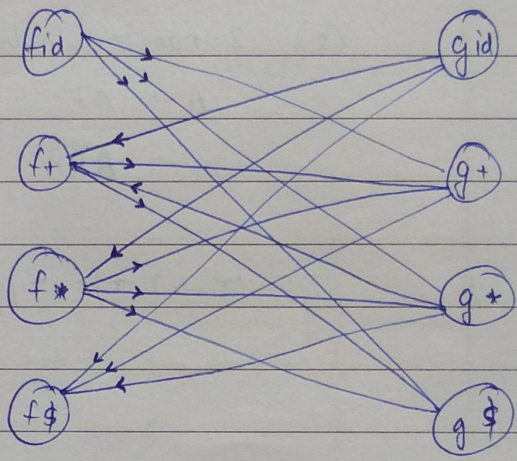
\$ id	\$ \triangleleft id	push & increment
\$	id \triangleright +	pop & reduce
\$ +	\$ \triangleleft +	push & increment
\$ + id	+ \triangleleft id	push & increment
\$ +	id \triangleright *	pop & reduce
\$ + +	+ \triangleleft *	push & increment
\$ + + id	* \triangleleft id	push & increment
\$	id \triangleright \$	pop
\$	+ \triangleright \$	pop
\$	* \triangleright \$	pop

Operator Function Table

- There should not be any cycle in the graph generated. If there is any cycle, then the function table can not be generated for that ~~graph~~ graph.
- New size = $2n$

space = n^2

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



space = $2n$

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

$fid \rightarrow g* \rightarrow f+ \rightarrow g+ \rightarrow f\$$
 $gid \rightarrow f* \rightarrow g* \rightarrow f+ \rightarrow g+ \rightarrow f\$$

length of longest path starting from gid

Compare $fid, g*$ $\rightarrow 4 > 3$
 $f+, gid$ $\rightarrow 2 < 5$



Shift Reduce Parser

- Bottom up parser
- Actions
 - Shift - push current input symbol to stack
 - Reduce - replace symbols by non-terminals (variables) i.e. replace RHS of a production by LHS

• Example -

* Productions -

$$E \rightarrow E + E$$

$$E \rightarrow a$$

$$E \rightarrow b$$

* Input string -

a + b

Stack	Input String	Action
\$	a + b \$	Shift a
\$ a	+ b \$	Reduce $E \rightarrow a$
\$ E	+ b \$	Shift +
\$ E +	b \$	Shift b
\$ E + b	\$	Reduce $E \rightarrow b$
\$ E + E	\$	Reduce $E \rightarrow E + E$
\$ E	\$	Accept

LR parser

LR(k)

Left to right scanning

Rightmost derivation in reverse

number of input symbols. When k is omitted, k is assumed to be 1, i.e. CLR(1)/LR(1)

Bottom up parser

LR parser

LR(0)

Simple LR(1)

SLR(1)

LALR(1)

lookahead LR(1)

CLR(1)

Conical LR(1)

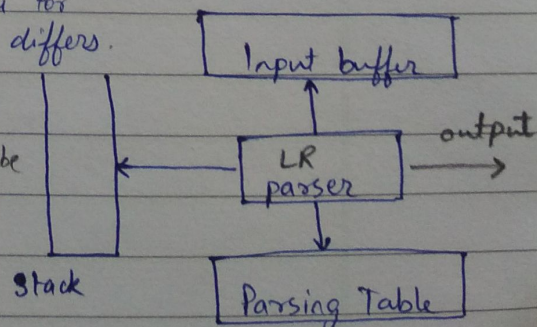
Note - Parser, stack & input buffer are same for all parsers
 - Only the parsing table construction differs.

Dot means that the input has been read till that point.

→ LR(0) items (LR(0), SLR(1))
 $E \rightarrow \cdot aA$
 Requirement - 'dot' should be present on RHS

→ LR(1) items (LALR(1), CLR(1))

$E \rightarrow \cdot aA$, a/b
 Same as LR(0) Lookahead



LR(0) item

$$E \rightarrow TT$$

$$T \rightarrow aT/b$$

Convert normal grammar to ~~augmented~~ augmented grammar. Steps -

1. Add an augmented production

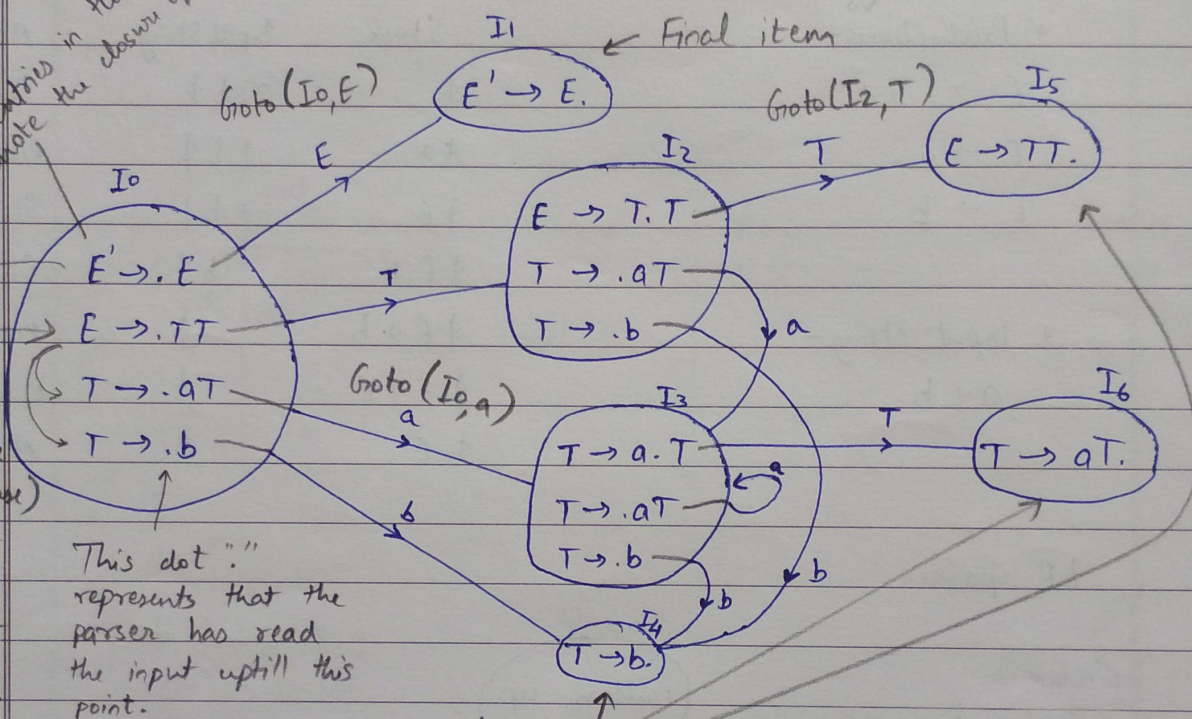
$$E' \rightarrow E$$
 ← Augmented Production

$$E \rightarrow TT$$

$$T \rightarrow aT/b$$

The entries in this set denote the closure of $E' \rightarrow E$

For each variable to the right of ".", add the productions of that variable (i.e. E in this case) and put the dot "." in the first place of the RHS. Similarly for T



This dot "." represents that the parser has read the input up till this point.

Final Item (The RHS of the first production is completely read)

This is the most important and prerequisite structure to construct the LR(0) parsing table



Construction of LR(0) parsing table
 Refer the diagram and grammar on pg 10

$E' \rightarrow E$
 $E \rightarrow TT$ ①
 $T \rightarrow aT / b$ ② ③

Necessary - We number the productions for use in the table when denoting the reduce operation

Perform shift on seeing "a" when in I_0 , and move to I_3

The state in which the augmented production becomes the final item, that is when we use "accept"

On seeing E when in I_0 , move to I_1

	Action (terminals)	Goto (variables)				
		a	b		\$	E
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			

Represents states I_0, I_1, \dots

This is augmented relation

Final Items

Note - we fill all columns under Action with r_n where n is the production number in the augmented grammar

only this will change in SLR(1) table

Construction of SLR(1) parsing table

Production:

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

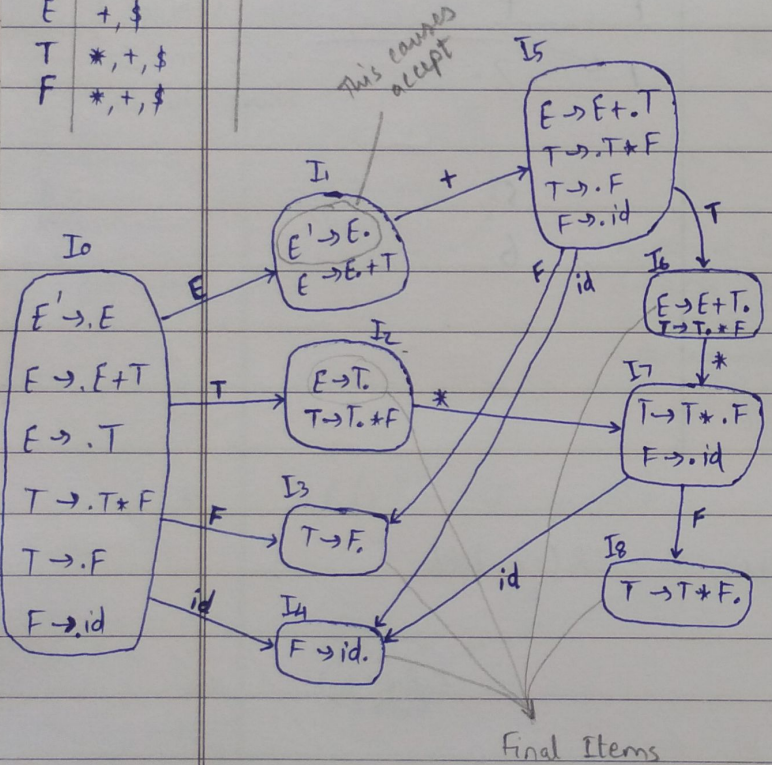
- $E' \rightarrow \cdot E$
- $E \rightarrow \cdot E + T$ ①
- $E \rightarrow \cdot T$ ②
- $T \rightarrow \cdot T * F$ ③
- $T \rightarrow \cdot F$ ④
- $F \rightarrow \cdot id$ ⑤

Convert to augmented grammar and make them LR(0) items by adding "." to the left most side of RHS

Necessary - we number the productions / LR(0) items for use in the table while writing the reduce operation.

Follow(x)

E	+, \$
T	*, +, \$
F	*, +, \$



	Action	Go to						
		id	+	*	\$	E	T	F
0	S_4					1	2	3
1	S_5		accept					
2	r_2	S_7	r_2					
3	r_4	r_4	r_4					
4	r_5	r_5	r_5					
5	S_4					6	3	
6	r_1	S_7	r_1					
7	S_4						8	
8	r_3	r_3	r_3					

Steps

- Find all the states where final item has come. Except the state with augmented production $E' \rightarrow \cdot E$.
- For example, consider I_4 item/state in the transition diagram

this is 5th production in the converted grammar.
 Hence write r_5
 $F \rightarrow id$
 \downarrow
 $Follow(F) = \{ *, +, \$ \}$

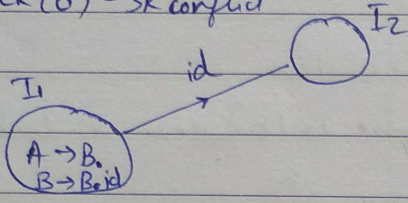
Write r_5 under these terminals for the 4th row



SR conflict and RR conflict

* LR(0) - SR conflict

Assume $A \rightarrow B$ is 3rd production in the augmented grammar.



SR conflict

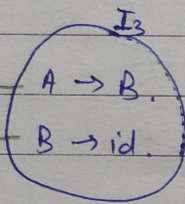
Action	id	+	*	\$...
1	s_2/s_3	r_3	r_3	r_3	...

(s_2 Because I_1 has a transition to I_2 for "id")

(r_3 Because I_1 has a final item)

* LR(0) - RR conflict

Assume - $A \rightarrow B$ as 2nd production & $B \rightarrow id$ as 3rd production in the augmented grammar



RR conflict

Action	t_1	t_2	t_3	Goto
3	r_2/r_3	r_2/r_3	r_2/r_3	...

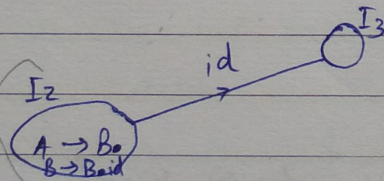
r_2 because I_3 has

$A \rightarrow B$ as final state which is 2nd production in the augmented grammar

r_3 because I_3 has $B \rightarrow id$ as final state which is 3rd production in the augmented grammar

* SLR(1) - SR conflict

Assume $A \rightarrow B$ is 3rd production in the augmented grammar & $\text{Follow}(A) = \{id\}$



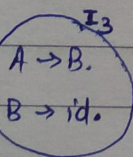
SR conflict

Action	id	+	*	\$...
2	s_3/r_3				...

because $A \rightarrow B$ is final item

* SLR(1) - RR conflict

Assume $A \rightarrow B$ as 2nd production & $B \rightarrow id$ as 3rd production in the augmented grammar



RR conflict

Action	t_1	t_2	t_3	Goto
3	r_2	r_2/r_3	r_3	...

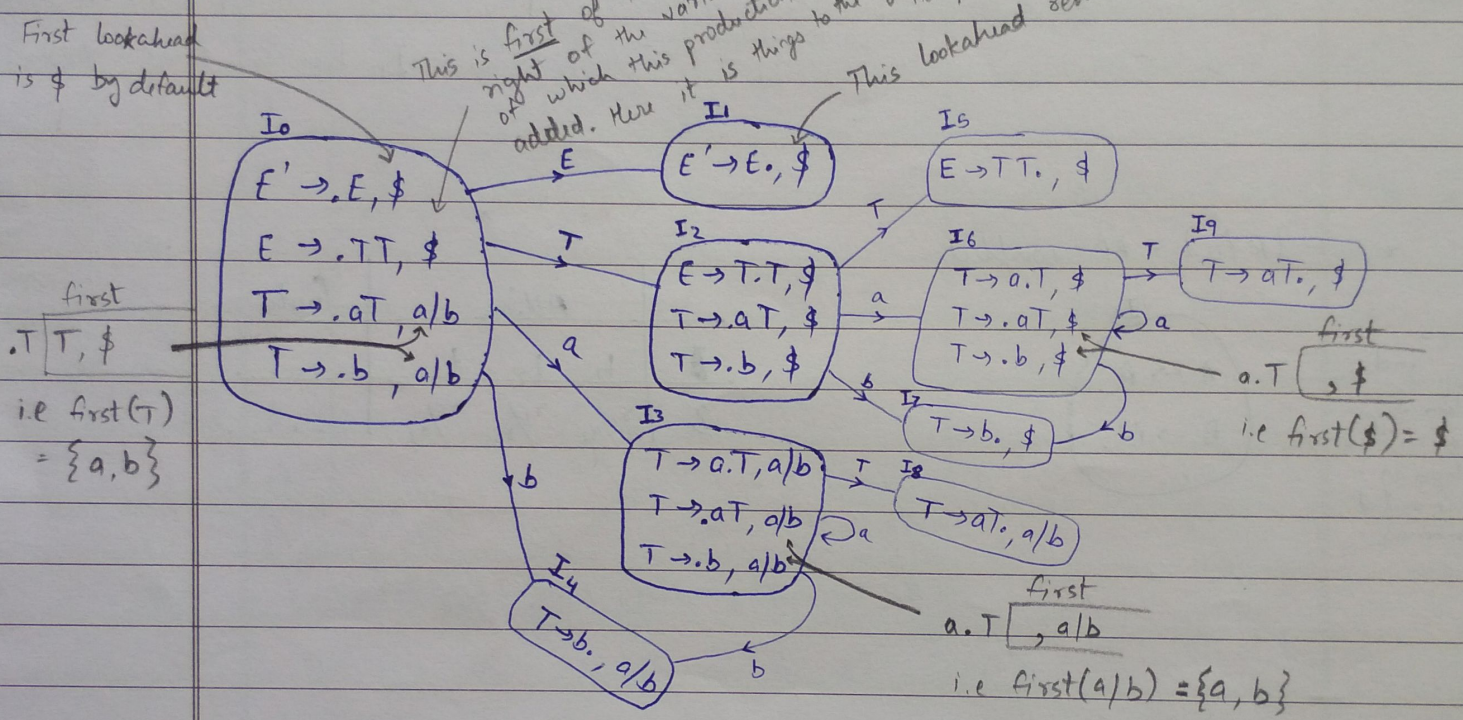
RR conflict

and $\text{Follow}(A) = \{t_1, t_2\}$
 $\text{Follow}(B) = \{t_2, t_3\}$

CLR(1) / LR(1) parsing table - more states as compared to LR(0) and SLR(1)

$E' \rightarrow E$
 $E \rightarrow TT$ ①
 $T \rightarrow aT/b$ ② ③

LR(1) item \hookrightarrow LR(0) item + lookahead



Steps

- 1 Write the first production $E \rightarrow \cdot E, \$$, and add other production if there is a variable to the right of ".", which is same as LR(0) + SLR(1)
- 2 Lookahead remains the same after transition from one state to another, for the first production only. i.e only the dot shifts towards the right. The lookahead ~~look~~ lookahead may change for production generated due to first production of the same state.



The format of the parsing table remains the same.

Same as LR(0), SLR(0)

It is 3 because on 'd', I₀ has transition to I₃

I₄

T → b, a/b

under a and b we write r₃

because T → b is 3rd production in the augmented grammar

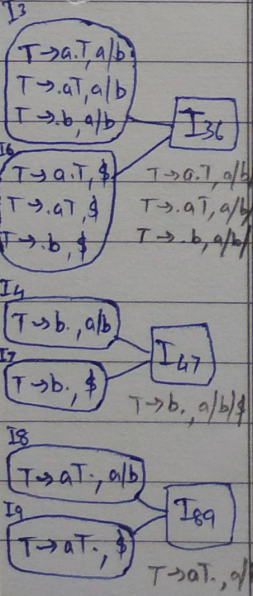
	Action			Goto	
	a	b	\$	E	T
0	S ₃	S ₄		1	2
1			accept		
2	S ₆	S ₇			5
3	S ₃	S ₄			8
4	r ₃	r ₃			
5			r ₁		
6	S ₆	S ₇			9
7			r ₃		
8	r ₂	r ₂			
9			r ₂		

Note - filling of shift operation and values in Goto column is same as LR(0) and SLR(1) parsing table like SLR(1).
Instead of using follow(X)_a, we use the lookahead to decide the columns.

#

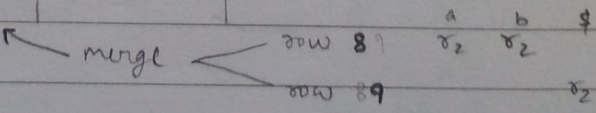
LA LR(1) parsing table (or LR(1))

- We use the same example as CLR(1)
- Take those states whose base content (i.e LR(0) items) is same but lookahead is different, and merge them to generate new states. Replace the old state numbers under Goto and ~~reduce~~ shift operation with the new state numbers.

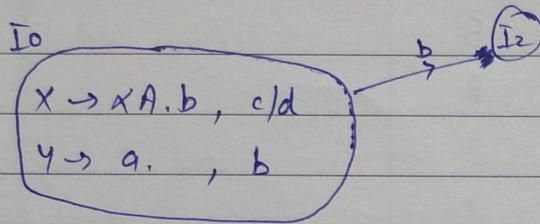
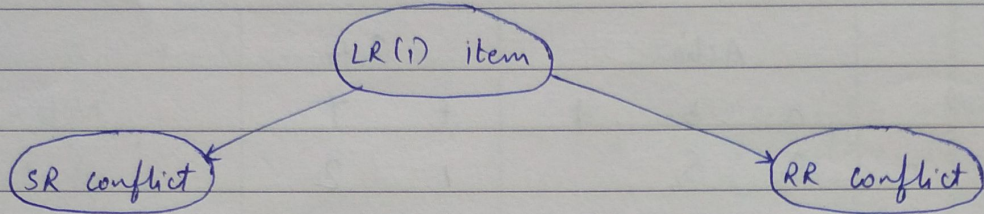


	Action			Goto	
	a	b	\$	E	T
0	S ₃₆	S ₄₇		1	2
1			accept		
2	S ₃₆	S ₄₇			5
36	S ₃₆	S ₄₇			89
47	r ₃	r ₃	r ₃		
5			r ₁		
89	r ₂	r ₂	r ₂		

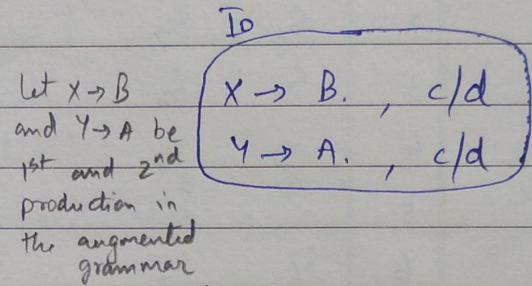
- Do NOT change numbers with reduce as they are production number from the augmented grammar.
- Merge the rows having same row number and take a union of their entries.



SR conflict and RR conflict with LR(1) items



Let $Y \rightarrow a$ be numbered 2 in augmented grammar



Let $X \rightarrow B$ and $Y \rightarrow A$ be 1st and 2nd production in the augmented grammar

	Action				Goto
	b	c	d	\$...
Conflict I_0	S_2/S_2				...

because of transition from I_0 to I_2

because I_0 has final item which is numbered 2 in augmented production and lookahead is b

	Action			Goto
	c	d	\$...
Conflict I_0	r_1/S_2	r_1/S_2		...



Input- Parse Tree \rightarrow Semantic Analysis \rightarrow Output Semantically Verified Parse Tree

Syntax Directed Translation (SDT)

Grammar + Semantic Rule = SDT

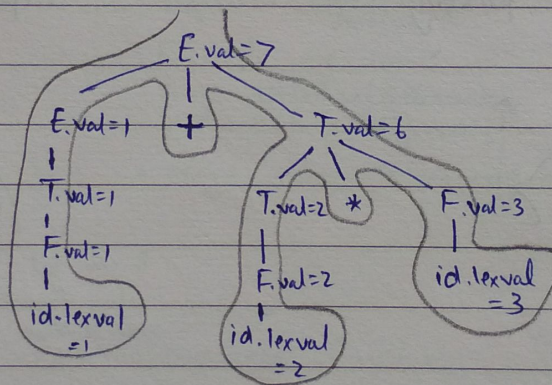
Example-

These are semantic rules associated with the productions. They are executed during the REDUCTION phase.

- $E \rightarrow E + T \quad \{ E.val = E.val + T.val \}$
- $\quad \quad \quad | T \quad \quad \quad \{ E.val = T.val \}$
- $T \rightarrow T * F \quad \{ T.val = T.val * F.val \}$
- $\quad \quad \quad | F \quad \quad \quad \{ T.val = F.val \}$
- $F \rightarrow id \quad \quad \quad \{ F.val = id.lexval \}$

"val" is an attribute associated with Non terminals (variables)

lets solve for $\rightarrow 1 + 2 * 3$



#

Synthesized attribute

- It is an attribute of a non-terminal (i.e. variable) which is at the LHS of the production rule.

$$A \rightarrow BCD \quad \left\{ \begin{array}{l} A.val = B.val, \\ A.val = C.val, \\ A.val = D.val \end{array} \right\}$$

Attribute of A is synthesized attribute as it can get its value from its children B, C, D

Inherited attribute

- It is an attribute of a non-terminal (i.e. variable) which is at the RHS of the production rule.

$$A \rightarrow B(C)D \quad \left\{ \begin{array}{l} c.val = A.val \\ c.val = B.val \\ c.val = D.val \end{array} \right\}$$

Can take its values from its siblings or its parent Same goes for B and D



#

S-attributed SDT

L-attributed SDT

- Synthesized attributes are used (i.e. parent can take its value from its children only)

- Both Synthesized attributes and inherited attribute (left side only) ^{for inherited attributes only}

eg- $A \rightarrow B(C)D$ C can take its value from its parent or left siblings only

✓ $C.val = B.val$
 X $C.val = D.val$

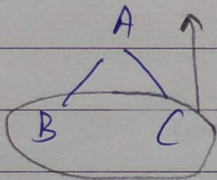
- Semantic actions can only be at the rightmost position

$A \rightarrow BC \{ \}$
 $B \rightarrow C \{ \}$

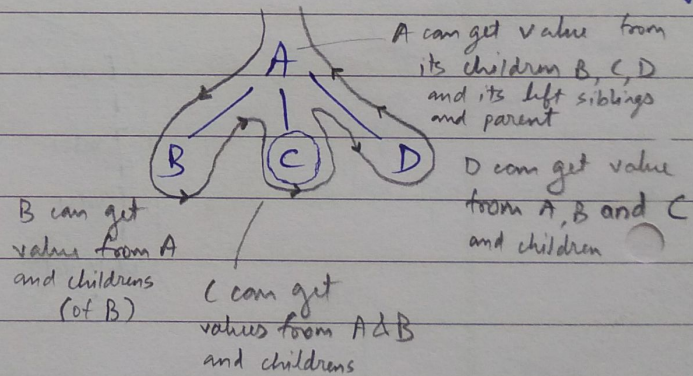
- Semantic actions can be ~~everywhere~~ anywhere on the RHS

$A \rightarrow BCD \{ \}$
 $B \rightarrow C \{ \} D$
 $C \rightarrow \{ \} D$

- Bottom up parsing



- Depth first, left to right parsing



- Every S-attributed definition is also L-attributed definition

- Bottom up evaluation is also possible



Intermediate code generation

3 address code - statement/sequence of variables/instruction which has atmost 3 address locations.

Other representations

(1) Graph & DAG/Syntax tree

(2) Postfix

Most popularly used forms of 3 address code are -

1. Quadruples
2. Triples

**

	Quadruples	Triples
$-(a+b) * (c-d)$	op opr1 opr2 result	Number op opr1 opr2
$t_1 = a+b$	1) + a b t_1	1) + a b
$t_2 = -t_1$	2) - t_1 t_2	2) - (1)
$t_3 = c-d$	3) - c d t_3	3) - c d
$t_4 = t_2 * t_3$	4) * t_2 t_3 t_4	4) * (2) (3)

<ul style="list-style-type: none"> • Efficient optimization possible • Flexible, i.e can swap instruction/statements without any difficulty • More space is required 	<ul style="list-style-type: none"> • NOT Efficient (in terms of instruction swapping) • Less flexible, i.e if we want to swap (1) & (3) above, we will have to change their references as well • Less space required
---	---

Indirect Triples

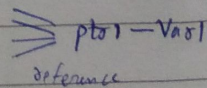
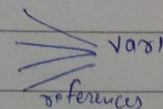
Number	op	opr1	opr2	Statement
1)	+	a	b	1) (11)
2)	-	(11)		2) (12)
3)	-	c	d	3) (13)
4)	*	(12)	(13)	4) (14)

Advantage

- Similar to quadruple
- Easy swapping of instruction as compared to Triples
- Space optimised as compared to quadruples.

Tripl

Indirect Triple



If var1 becomes var2, all references to be updated

if var1 becomes var2, only the ptr1 is updated