

GATE CSE NOTES

by
Joyoshish Saha



Downloaded from <https://gatetcsebyjs.github.io/>

With best wishes from Joyoshish Saha

- **Lexical errors:** Lexical error is a sequence of characters that does not match the pattern of any token.
- For unambiguous grammars, Leftmost derivation and Rightmost derivation represent the same parse tree.
- For ambiguous grammars, Leftmost derivation and Rightmost derivation represent different parse trees.
- Concatenating the leaves of a parse tree from the left produces a string of terminals. This string of terminals is called the **yield of a parse tree**.

- **Eliminating Left Recursion:** <https://www.gatevidyalay.com/left-recursion-left-recursion-elimination/>

> Removing left recursion (as top down parser can't handle it)

a) Immediate/Direct (by Moore's proposal)

$$A \rightarrow A\alpha / \beta \Rightarrow \begin{cases} A \rightarrow \beta A' \\ A' \rightarrow \alpha A'/\epsilon \end{cases}$$

(starts in β , then any number of α)

In general, $A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m | \beta_1 | \dots | \beta_n$ and β_i not starting with A

$$\downarrow$$

$$A \rightarrow \beta_1 A' | \dots | \beta_n A'$$

$$A' \rightarrow \alpha_1 A' | \dots | \alpha_m A' | \epsilon$$

b) Indirect: Arrange non-terminals in some order A_1, \dots, A_n

Alg for i from 1 to n
 for j from 1 to $i-1$
 replace each prodⁿ $A_i \rightarrow A_j \alpha$ by
 $A_i \rightarrow \alpha_1 \alpha_2 \dots \alpha_{i-1}$ where $A_j \rightarrow \alpha_1 \dots \alpha_k$
 eliminate left recursions among A_i prods.

eg $E \rightarrow E + T | T$
 $T \rightarrow T * F | F$
 $F \rightarrow a | (E)$

removing LRs \downarrow

2 immediate LRs

$$\begin{array}{l|l} E \rightarrow E + T & T \rightarrow T * F \\ E \rightarrow TE' & T \rightarrow FT' \\ E' \rightarrow +TE'/\epsilon & T' \rightarrow *FT'/\epsilon \end{array}$$

$(E \rightarrow TE' \quad T \rightarrow FT' \quad F \rightarrow a | (E))$
 $E' \rightarrow +TE'/\epsilon \quad T' \rightarrow *FT'/\epsilon$

eg $S \rightarrow Aa/b$
 $A \rightarrow Ac / Sd/f$

Has indirect LR.
 Rename $A_1 \rightarrow A_2 a/b$ $A_2 \rightarrow A_2 c / A_2 d / f$

$i=1 j=1$ no prodⁿ as $A_1 \rightarrow A_1 \omega$
 $i=2 j=1$ $A_2 \rightarrow A_1 d \Rightarrow A_2 \rightarrow A_2 a/d / bd$.
 Prodⁿ for A_2 becomes $A_2 \rightarrow A_2 c / A_2 ad / bd / f$
 Removing direct LR of A_2 prods $\rightarrow \begin{cases} A_2 \rightarrow bd A_2' / f A_2' \\ A_2' \rightarrow c A_2' / ad A_2' \end{cases}$
 Now, change name to original.

right recursion in its production.

<https://www.gatevidyalay.com/grammar-ambiguity-ambiguous-grammar/>

<https://www.gatevidyalay.com/removing-ambiguity-grammar-ambiguity/>

- For ambiguous grammars,
 - More than one leftmost derivation and more than one rightmost derivation exist for at least one string.
 - Leftmost derivation and rightmost derivation represents different parse trees.
- For unambiguous grammars,
 - A unique leftmost derivation and a unique rightmost derivation exist for all the strings.
 - Leftmost derivation and rightmost derivation represents the same parse tree.
- There may exist derivations for a string which are neither leftmost nor rightmost.
- Leftmost derivation and rightmost derivation of a string may be exactly the same. In fact, there may exist a grammar in which leftmost derivation and rightmost derivation is exactly the same for all the strings.
- For a given parse tree, we may have its leftmost derivation exactly the same as the rightmost derivation.
- If for all the strings of a grammar, leftmost derivation is exactly the same as rightmost derivation, then that grammar may be ambiguous or unambiguous.
- **Left factoring(eliminating non-determinism):** Left factoring is a process by which the grammar with common prefixes is transformed to make it useful for Top down parsers. >>We make one production for each common

prefixes. The common prefix may be a terminal or a non-terminal or a combination of both. Rest of the derivation is added by new productions.

- The presence or absence of left recursion does not impact left factoring and ambiguity anyhow.
- The presence or absence of left factoring does not impact left recursion and ambiguity anyhow.
- The presence or absence of ambiguity does not impact left recursion and left factoring anyhow.
- **Parsers:**
 - Top down parsers
 - w/ full backtracking (brute force)
 - w/o backtracking
 - Recursive descent parser(Goes to infinite loop if grammar is not free of left recursion)
 - Non-recursive descent parser or LL(1) parser - Predictive Parser(No ambiguity, No Left Recursion, No nondeterminism)
 - Bottom up parsers
 - Operator precedence parser
 - LR Parsers (LR grammar can never be ambiguous)
 - LR(0)
 - SLR(1)
 - LALR(1)
 - CLR(1)

- **First and Follow calculation**(eliminate left recursion before calculating, *follow* never contains *eps*)

- **LL(1) Parsing**

- Using LMD, left → right scanning, #lookahead = 1
- Constructing LL(1) parsing table [Table size = #V * (#T + 1), 1 for \$]
 - All null productions put under *follow* set of the symbol, rest under the *first* set.
- In case of >1 entries in a cell of table, we can't parse using LL(1) parsing.
- LL(1) parsing examples (using stack, put the production in reverse into the stack, match top of the stack with input symbol, if matched pop, then look for the appropriate production by seeing the top and the input string symbol)
- Checking whether grammar is LL(1) or not(>2 entries => not LL(1)).
- If a grammar contains Left recursion/Common prefixes/Ambiguity, it can't be LL(1).
- LL grammar is a subset of CFG with some restrictions on it.

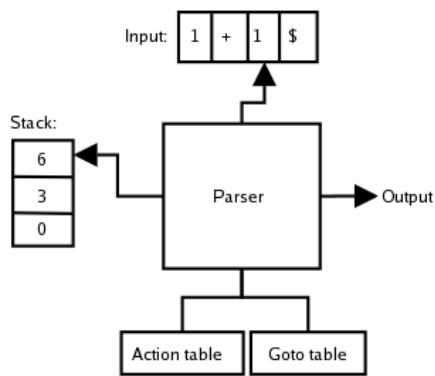
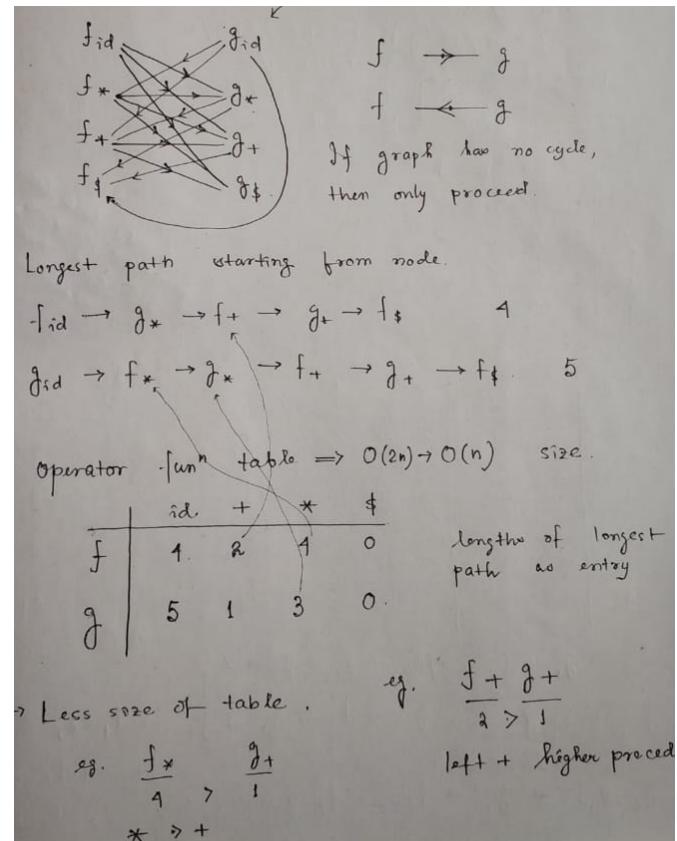
- **Recursive descent parsing**

- Uses stack
- If contains left recursion will go to infinite loop
- If contains common prefixes may generate parsing error
- Recursive descent is a top-down parsing technique that constructs the parse tree from the top and the input is read from left to right. It uses procedures for every terminal and non-terminal entity. This parsing technique recursively parses the input to make a parse tree, which may or may not require back-tracking. But the grammar associated with it (if not left factored) cannot avoid back-tracking. A form of recursive-descent parsing that does not require any back-tracking is known as predictive parsing.
- This parsing technique is regarded recursive as it uses context-free grammar which is recursive in nature.
- For Recursive Descent Parser, we are going to write one recursive function/procedure for every variable.
<https://www.geeksforgeeks.org/recursive-descent-parser/>

- **Operator precedence parsing**

- Operator grammar should not have 2 variables adjacent to each other and null productions(~~X~~).
- Can also work on ambiguous grammar. Only parser that works on ambiguous grammar.

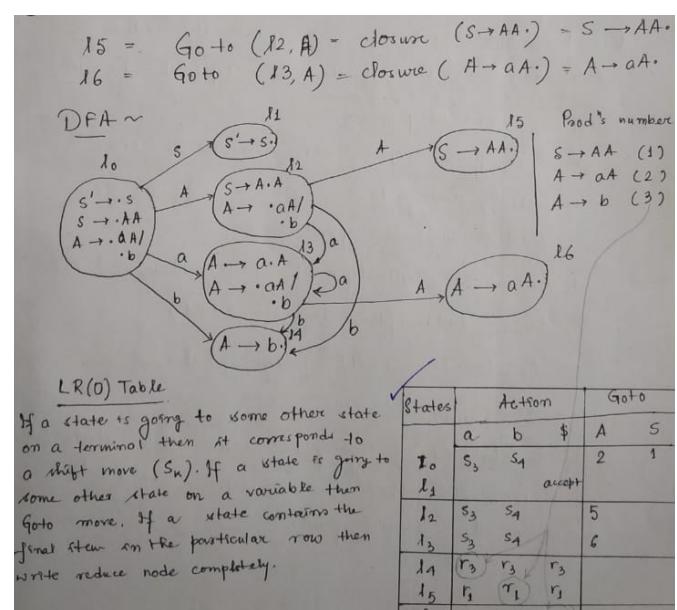
- Uses operator relation table(having associativity and precedence description)
- Size of operator relation table = n^2 [n = #operators]
 - So, we have operator function table. Size = $2n$ [$O(n)$]
 - Building Operator function table(from operator relation graph): >Proceed only if graph is acyclic. >Length of longest path from the node.
 - Disadvantage of the function table is there's no blank entries here(as in operator relation table). So, error detecting capability of the function table is less.



• LR(0) Parser

- Uses Canonical collection of LR(0) items: Add augmented production $S' \rightarrow S$ and put “.” at first position of every production's RHS. (“.” signifies what it has seen so far in the string)
- CFG → Check ambiguity → Remove ambiguity if possible → Add augment production and make canonical collection of LR(0) items → Draw DFA → Construct LR(0) parsing table
- Making the DFA: Put all productions in the initial state. Then, compute Goto and the closure of that production, thereby, giving birth of new states.

- Making the LR(0) table:
 - If a state is going to some other state on seeing a terminal, it's a action *shift move*(specified with the state number where it's going in DFA).
 - If a state is going to some other state on seeing some variable, then it's a *goto move*(specified with the state number where it's going in DFA).
 - If a state contains some final item(like $A \rightarrow Aa.$) then in the action part of the table put *reduce move*(reduce using specified production) throughout the row.



• SLR(1) Parser

- Also uses canonical collection of LR(0) items.

- Only difference with LR(0) is here we put *reduce move* only in the *follow* of LHS.
- Makings DFA(same as LR(0)), SLR(1) table: Same as LR(0) except the rule that when a state contains a final item(as $A \rightarrow ab.$), fill with *reduce entry* only in the *follow* of A.
- Conflicts:**
 - R/R conflict: State containing 2 *reduce items*.
 - S/R conflict: State containing both *shift* and *reduce moves*.
 - Conflicts can happen for LR(0), SLR(1).
 - If conflict is there, we say it's not LR(0), or SLR(1), in which case it applies.
 - If grammar is LR(0), it's SLR(1).
 - No relation of a grammar of being LL(1) and [LR(0) or SLR(1)].
 - Number of **inadequate states** is nothing but the number of states which have S/R or R/R conflicts.

• CLR(1) Parser (Canonical LR) or LR(1)

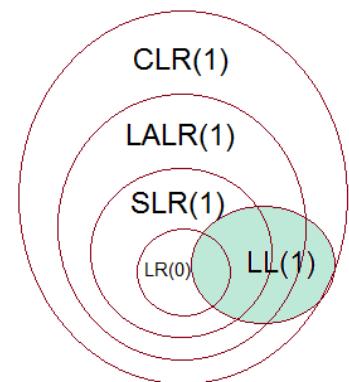
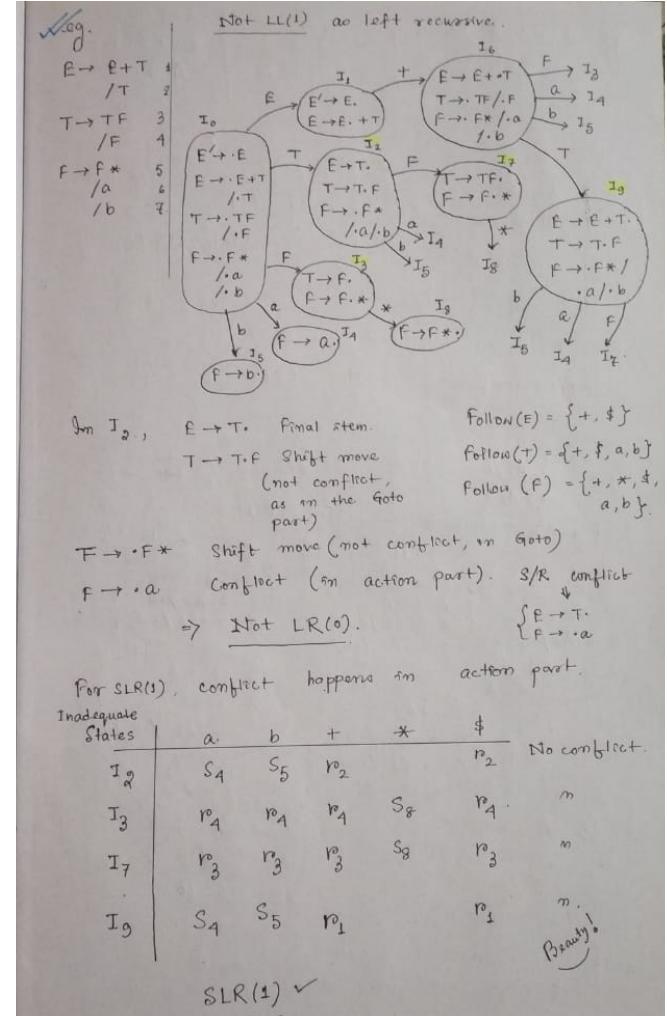
- LR(1) item:** Collection of LR(0) item along with a lookahead symbol.
- If $A \rightarrow \alpha.B\beta$, a is in closure of some state and $B \rightarrow \gamma$ is a production rule of the grammar, then $B \rightarrow \gamma$, b will be in the closure of that state also for each terminal b in *first*(βa)
- Lookahead is used to determine where we place the final item. For the augment production, lookahead is always the \$.
- Placement of *shift* nodes is the same as LR(0) or SLR(1), only difference is for the placement of *reduce* nodes. While making the CLR(1) table, we place the *reduce* only in the lookahead symbols.
- Example: Next page

• LALR(1) Parser (Lookahead LR)

- The LR(1) items(states in the DFA) having the same productions but different lookahead are combined to form a single set of items.
- If there are no states with same productions and different lookahead in CLR(1), then LALR(1) parsing table is identical to CLR(1).
- If a grammar is CLR(1) it may or may not be LALR(1). After combining items it may give rise to conflicts. We just reduce the size of the parsing table in LALR(1) compared to CLR(1), but the conflicts can't go away.
- If a grammar is not CLR(1) it cannot be LALR(1).
- Example: next page

• Some pointers:

- $LR(0) \subset SLR(1) \subset LALR(1) \subset CLR(1)$
- $LL(1) \subset LR(1)$
- Every ϵ free LL(1) grammar is SLR(1) grammar, thereby LALR(1) and CLR(1) too.
- (#states(SLR(1)) = #states(LALR(1))) \leq #states(CLR(1))
- CLR(1) parser most powerful among the bottom-up parsers.
- Based on power, OPP < LL(1) < LR(0) < SLR(1) \leq LALR(1) \leq CLR(1)
- If a LL(1) grammar's each variable can derive a non null string then it is LALR(1) i.e. an LL(1) grammar with symbols that have both empty and non-empty derivations is also an LALR(1) grammar.



- If a LL(1) grammar has all the variables producing only null strings, then it may or may not be LALR(1).
- Every LL(k) grammar is also an LR(k) grammar.
- <https://gateoverflow.in/blog/6215/parsing-notes>

eg. $S \rightarrow AA \quad (1)$
 $A \rightarrow aA/b$
 $\quad \quad (2)$
 $\quad \quad (3)$

$$\left\{ \begin{array}{l} S' \rightarrow S, \$ \\ S \rightarrow \cdot AA, \$ \\ A \rightarrow \cdot aA, a/b \\ A \rightarrow \cdot b, a/b \end{array} \right.$$

I0

$$\begin{aligned} I0 &= \text{closure } (S' \rightarrow S) \\ &= S' \rightarrow \cdot S, \$ \\ &= S \rightarrow \cdot AA, \$ \\ &= A \rightarrow \cdot aA, a/b \\ &= A \rightarrow \cdot b, a/b. \end{aligned}$$

I1

$$\begin{aligned} I1 &= \text{Goto } (I0, S) = \text{closure } (S' \rightarrow S \cdot, \$) \\ &= S' \rightarrow S \cdot, \$ \end{aligned}$$

I2

$$\begin{aligned} I2 &= \text{Goto } (I0, A) = \text{closure } (S \rightarrow A \cdot A, \$) \\ &= S \rightarrow A \cdot A, \$ \\ &= A \rightarrow \cdot aA, \$ \\ &= A \rightarrow \cdot b, \$ \end{aligned}$$

I3

$$\begin{aligned} I3 &= \text{Goto } (I0, a) = \text{closure } (A \rightarrow a \cdot A, a/b) \\ &= A \rightarrow a \cdot A, a/b \\ &= A \rightarrow \cdot aA, a/b \\ &= A \rightarrow \cdot b, a/b \end{aligned}$$

Goto (I3, a) = closure (A → a · A, a/b) = (same as I3)

Goto (I3, b) = closure (A → b ·, a/b) = (same as I1)

I4

$$\begin{aligned} I4 &= \text{Goto } (I0, b) = \text{closure } (A \rightarrow b \cdot, a/b) \\ &= A \rightarrow b \cdot, a/b \end{aligned}$$

I5

$$\begin{aligned} I5 &= \text{Goto } (I2, A) = \text{closure } (S \rightarrow AA \cdot, \$) \\ &= S \rightarrow AA \cdot, \$ \end{aligned}$$

I6

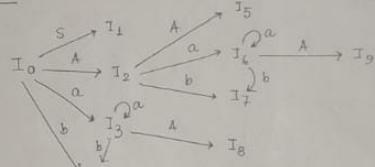
$$\begin{aligned} I6 &= \text{Goto } (I2, a) = \text{closure } (A \rightarrow a \cdot A, \$) \\ &= A \rightarrow a \cdot A, \$ \\ &= A \rightarrow \cdot aA, \$ \\ &= A \rightarrow \cdot b, \$ \end{aligned}$$

Goto (I6, a) = closure (A → a · A, \\$) = (same as I4)

Goto (I6, b) = closure (A → b ·, \\$) = (same as I7)

$$\begin{aligned} I7 &= \text{Goto } (I2, b) = \text{closure } (A \rightarrow b \cdot, \$) \\ &= A \rightarrow b \cdot, \$ \\ I8 &= \text{Goto } (I3, A) = \text{closure } (A \rightarrow aA \cdot, a/b) \\ &= A \rightarrow aA \cdot, a/b \\ I9 &= \text{Goto } (I4, A) = \text{closure } (A \rightarrow aA \cdot, \$) \\ &= A \rightarrow aA \cdot, \$ \end{aligned}$$

DFA



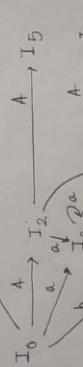
Data flow diagram

CLR(1) parsing table.

States	Action			Goto
	a	b	\$	
I0	S3	S1		1 2
I1			Accept	
I2	S6	S7		5
I3	S3	S1		8
I4	R3	R3		
I5			R1	
I6	S6	S7		9 action(I6, a) = S6
I7			R3	action(I7, a) = R3
I8	R2	R2		action(I8, a) = R2 action(I8, b) = R3
I9			R2	action(I9, a) = R2

{ Placement of shift nodes in CLR(1) parsing table is same as the SLR(1) parsing table. Only difference is in the placement of reduce node.

Reduced no. of states than CLR(1).



LALR(1) parsing table.

States	a	b	\$	A
I0	S36	S17		1 2
I1	S26	S17		
I2	S36	S17		5
I36	R3	R3		89
I47	R3	R3		R1
I5				
I89	R2	R2		R2

- Any string derivable from the start symbol is a **sentential form** — it becomes a sentence if it contains only terminals.
- A sentential form that occurs in the leftmost derivation of some sentence is called **left-sentential form**.
- A sentential form that occurs in the rightmost derivation of some sentence is called **right-sentential form**.
- A **handle** of a right sentential form ' γ ' ($\gamma = \alpha\delta\beta$) is a production $E \rightarrow \delta$ and a position in γ where δ can be found and substituted by E to get the previous step in the right most derivation of γ — previous and not "next" because we are doing rightmost derivation in REVERSE.
- **Viable prefixes** are the prefixes of right sentential forms that do not extend beyond the end of its handle, i.e., a viable prefix either has no handle or just one possible handle on the extreme RIGHT which can be reduced.
- **Viable prefixes** can be recognized by using a FINITE AUTOMATA. Using this FINITE AUTOMATA and a stack we get the power of a Push Down Automata and that is how we can parse context-free languages.
- In LR(1) we see the input from left till we get a handle. After this, we see one more lookahead symbol and determine the parser action. i.e., the parser has MORE information to decide its action than in LL(1) and this makes it more powerful than LL(1). More powerful means, any grammar that can be parsed by LL(1) can also be parsed by LR(1). This holds in general case too – for any k , LL(k) set of grammars is a PROPER subset of LR(k) set of grammars.

Just revising on the previous part we have the following relation. Here Lang denote the set of languages defined by the given set of grammars.

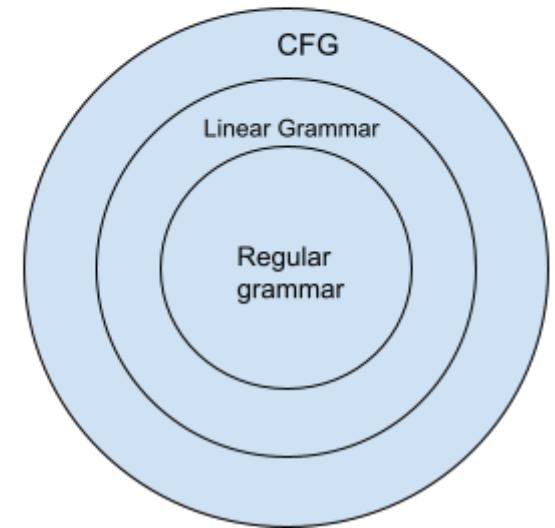
1. $LL(0) \subset LL(1) \subset LL(2) \dots LL(k) \subset LL(k+1) \dots$
2. $\text{Lang}(LL(0)) \subset \text{Lang}(LL(1)) \subset \text{Lang}(LL(2)) \dots \text{Lang}(LL(k)) \subset \text{Lang}(LL(k+1)) \dots$
3. $LR(0) \subset LR(1) \subset LR(2) \dots LR(k) \subset LR(k+1) \dots$
4. $\text{Lang}(LR(0)) \subset \text{Lang}(LR(1)) = \text{Lang}(LR(2)) \dots \text{Lang}(LR(k)) = \text{Lang}(LR(k+1)) \dots$ ($\text{Lang}(LR(0))$ is DCFL with prefix property and $\text{Lang}(LR(k))$; $k \geq 1$ is DCFL (with or without prefix property))
5. $\text{Lang}(LR(0)\$) = \text{Lang}(LR(1)) = \text{Lang}(LR(2)) \dots \text{Lang}(LR(k)) = \text{Lang}(LR(k+1)) \dots$ ($\text{Lang}(LR(0)\$)$ is DCFL as when we add an end delimiter no string becomes a prefix of another)

- Every LL(1) is LR(1)
- If a grammar is LL(k) then it is LL($k+1$) too.
- Every LL(1) is LALR(1)
- Every LL(1) is LR(0)
- LL(k) is subset of LR(k) (LL(2) is subset of LR(2))
- LR(1) table \cong size LL(1) table size
- #states (LR(0) = SLR(1) = LALR(1)) \leq CLR(1)
- Reading terminal on LR parser DFA, if we go out of that state then Shift move is put in the table.
- LR Parsers conflict
 - S/R
 - R/R
 - No S/S conflict as it is DFA, not NFA. If NFA was there then S?S conflicts could be there.
- **Inadequate state:** State having conflicts.
 - Inadequate state must have a minimum of 2 productions.
 - At least 1 reduce move is there.
- In LR(0), SLR(1), LALR(1)
 - Number of states same
 - Shift entries same
 - Goto part same
 - Reduce entries may change
 - Error entries may change
- If CLR(1) doesn't have an R/R conflict then LALR(1) may have an R/R conflict.
- If CLR(1) doesn't have S/R conflict then LALR(1) also does not have S/R conflict.
- If LALR(1) have S/R conflict, then CLR(1) also have S/R conflict.

- **YACC:**

- LALR(1) parser generator: Create LALR(1) parsing table.
 - If having conflict YACC resolves them too.
 - S/R => S Resolving S/R conflict
 - r3/r4 => r3 Resolving R/R conflict

- Every Regular grammar need not be LL(1).
 - Every Regular language is LR(1).
 - Every regular grammar is linear, but every linear grammar need not be regular.
 - Every LR grammar is unambiguous, but every unambiguous grammar is not LR grammar.
 - $\text{Follow}(A) = \text{RFollow}(A)$ [RFollow: Follow in right sentential form]
 - $\text{Follow}(A)$ and $\text{LFollow}(A)$ need not be the same. [LFollow: Follow in left sentential form]
 - For Unambiguous grammar
 - LMDT and RMDT should be same for some string production
 - LMD and RMD need not be same though
 - **Viable prefixes:** The set of prefixes of right sentential forms that appear on the stack of Shift/Reduce parser(or LR parser) are called Viable prefixes.
 - TDP (Uses LMD: Problem in multiple productions)
 - w/ backtracking: RDP
 - w/o backtracking: LL(1) / Predictive parser / NRP
 - BUP (No backtracking for any parser, Uses reverse of RMD: Pr



- Syntax directed translation

- To interleave semantic and syntax analysis
 - SDT = Grammar + Semantic rules

eg $s \rightarrow \text{axw}$ { print (1) }
 y { print (2) }
 $w \rightarrow s2$ { print (3) }

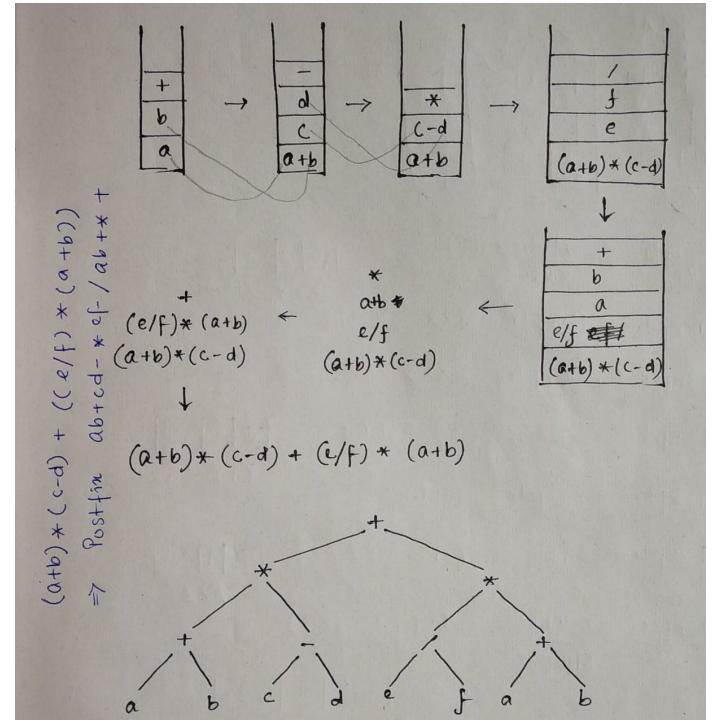
String $axxxy_2z_2$.

Output
 $\# 3 1 3 1$

```

graph TD
    s((s)) -- "a" --> w((w))
    s -- "x" --> w
    w -- "x" --> w5[w]
    w -- "y" --> w2[2]
    w5 -- "y" --> w5
    w2 -- "z" --> w4[4]
    w2 -- "z" --> w1[1]
    style w fill:#ffff00
    style w5 fill:#ffff00
    style w2 fill:#ffff00
    style w4 fill:#ffff00
    style w1 fill:#ffff00
  
```

- Syntax tree: Compact parse tree
 - Construction example: ----->
 - **Attributes**
 - **Synthesized:** Computed from the values of the attributes of the children nodes



- **Inherited:** From both sibling and parent nodes
- **Types of SDT:**
 - **S-attributed definition:** Uses only synthesized attributes
 - Evaluated in bottom-up parsing
 - Semantic actions placed in rightmost place of RHS
 - **L-attributed definition:** Uses both synthesized and inherited attributes with restriction(inherit from parent or left sibling only)
 - Evaluated by depth first order, left to right evaluation
 - Semantic actions placed anywhere in RHS.
 - If a definition is S-attributed then it is L-attributed also. Not vice versa.

- **Static(fixed) storage allocation**

- To make memory in the static area, it should be made in compile time
- 1-time memory created
- Recursion not allowed
- Dynamic data structure not allowed

- **Stack storage allocation**(will grow and shrink in runtime)

- Recursion allowed
- Dynamic DS not allowed

- **Heap memory allocation**(will grow and shrink in runtime)

- Recursion allowed
- Dynamic DS allowed

- **Activation record contains:**

- Local variables: hold the data that is local to the execution of the procedure.
- Temporary values: stores the values that arise in the evaluation of an expression.
- Machine status: holds the information about the status of the machine just before the function call.
- Access link (optional): refers to non-local data held in other activation records.
- Control link (optional): points to activation record of caller.
- Return value: used by the called procedure to return a value to calling procedure
- Actual parameters
- Activation record does not contain global variables.

Lexical phase error can be:

- Spelling error.
- Exceeding length of identifier or numeric constants.
- Appearance of illegal characters.
- To remove the character that should be present.
- To replace a character with an incorrect character.
- Transposition of two characters.