

Introduction

Compiler is a program which converts High-level program to Low-level program

There are multiple phases of Compilation:

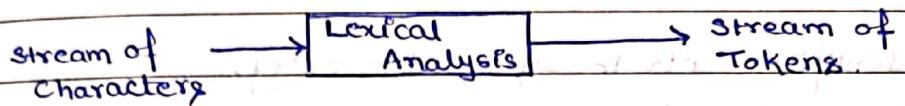
- Lexical Analysis
- Syntax Analysis
- Semantic Analysis
- Intermediate Code Generation
- I_C Optimization
- Code Generation
- Code Optimization

Analysis Phase

Synthesis
Phase

Lexical Analysis

- Token: A logical unit of a program.
- The source program is a sequence of characters.
- The job of Lexical Analysis phase is to convert the stream of characters to stream of tokens:



A token mainly consists of Token Name, but it can additionally contain attributes. Ex: Name of identifier, value of numeral, etc.

Ex. while ($x < 0$) $x++;$

while	(x	<	0	x	$++$;
Token Name	kW.White	L-paren	Ident	L-than	Int	Ident	SColon
Attr:			x		0	x	

Lexeme: The piece of original program from which the token is made.

Ex. In above program x is Lexeme, Identifier is Token.

Tokens are considered as the terminal of Grammar (Useful in Syntax Analysis).

* Symbol Table

A data structure maintained by the Compiler which is used to store & query information about various identifiers (variables, functions, etc), through different phases of compilation. In the lexical analysis phase a lexeme can be identified as identifier, and its

entry can be created in the symbol table, but its type can only be determined in later phases.

Ex. $x += y + 3;$

< identifier, ptr to sym. table entry >

< incr-assign-op, >

< identifier, ptr to symbol table entry >

< add-op >

< integer-const, 3 >

* Recognizing Tokens

- Tokens in a program are recognized based on RE rules specified for different types of tokens.

Ex. KU-FOR

Ident

[A-zA-Z] [A-zA-Z0-9_]*

So, based on these rules, matching is done.

- Longest Match (Maximal Munch): The Lexical Analyzer matches the longest matching lexeme with a rule from start.

Ex. [In C]

a++;

Starting from pos 0, the maximal matching lexeme is 'a' which is a identifier. Move to pos 1.

Starting from pos 1, one match is '+' w/ add-op & other match is '++' which is incr-op. So, since "++" is

- longer, that's matched. Move to pos 3
- Only ';' matches.

at + ;

Reading forward, no match found for 'double'.

- These can't be other problem. A lexeme can match w/ more than one rule.

Ex. KW-DOUBLE double

Identifier: [A-Za-z_][A-Za-z0-9_]*

double x; int i; double d;

matches w/ both rules.

In such situations, match w/ the rule defined first (priority by position).

- In case Lexical Analyzer starting from some position, can't match the lexemes with any rule, then there is a lexical error.

Ex. x="apdb;

x is matched.

= is matched.

"apdb;" isn't matched w/ any token type. This throws error.

Ex. int 3abc;

int matches KW-INT

3abc; doesn't matches any token type.

- Lexical Analyzer ignores comments.

* Token Counting For C programs.

- LA analyzer / scanner removes / ignores comments before converting to tokens.

Before LA runs, preprocessing happens, hence any substitution by #define, any inclusion (copy-paste) by #include, etc must be done before converting to tokens.

- String literals: Anything b/w starting " " & ending " " is counted as one token.

Ex. `x = "abc";` \Rightarrow [x] = ["abc"] ;

Ex. `#define OP +`

[int] main[OP]

`#define M 100`

[]

`int main()`

\Rightarrow

`x = OP / M + OP;`

$\therefore 11 \text{ tokens.}$

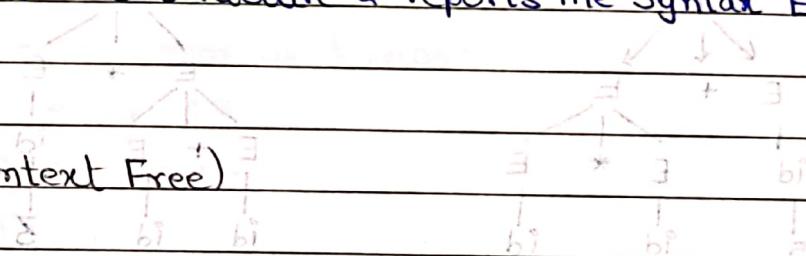
Possible Lexical Errors:

- Unterminated Comment
 - Unterminated String
 - Unidentified token
- etc.

Syntax Analysis

- After Lexical Analysis, we have a stream of tokens, so this is the input for Syntax Analyzer.
- Syntax Analysis phase tries to structure the stream of tokens into "parse Tree", by following the Grammar rules of the language. While converting to Parse Tree, it also has to check if there's any token(s) which don't follow the Grammatical rules & don't properly convert to the structure & reports the Syntax Errors.

* Grammar (Context Free)



- A ~~ex~~Grammar is used to define a language (set of strings) finitely. It is a four tuple: (V, T, P, S)

V : Set of Variables (Non-terminals) {Non empty, finite}

T : Language Alphabet (Terminals) {Nonempty, finite}

P : Set of production rules {finite}

S : Start symbol (variable). $S \in V$

- CFG is used to describe CFLs. Each production of a CFG is of the form:

$$V \rightarrow (V+T)^*$$

Left Most Derivation: At any step of derivation, resolve the leftmost NTerm.

Right Most Derivation: At any step of derivation, resolve the rightmost NTerm.

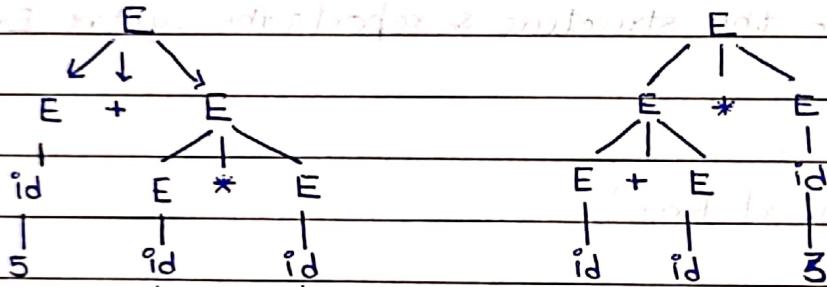
Lang(G)

- For a string $w \in \text{CFG}_i G_i$, the no. of LMDs = no. of RMDs = no. of Parse Trees.

* Ambiguous Grammars

- A CFG G is ambiguous iff there exists some string $w \in L(G)$ such that there ~~are~~ is more than one parse tree for it w.r.t. G .

Ex. $E \rightarrow E+E \mid E * E \mid id$, with $w = 5+2*3$



So, PT_1 results in the op^n: $5+(2*3) = 5+6 = 11 \neq 7$

PT_2 results in the op^n: $(5+2)*3 = 7*3 = 21 \neq 7$

So, Ambiguous grammars can create serious problems, as the meaning of the program becomes unpredictable & may differ from what it was intended to.

- The major problems w/ Ambiguous grammars:

(1) There is no algorithm to detect if a given arbitrary grammar is ambiguous or not.

(2) There is no algorithm to convert a given ambiguous grammar to unambiguous grammar.

Ex. Try to analyze if the following grammar can be made unambiguous.

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

① For a given string $id + id * id$, we'd want to give higher priority to $*$.

② For $id + id + id$, we'd want to have left association for $+$. Similarly for $*$.

We can modify the grammar as follows:

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

- So, since we want left associativity, so we have left recursion, for the $+$ & $*$ ops.
- Since we want higher priority for $*$, we keep it below the level of $+$.

Recursion: For a non-terminal X in CFG G , $\bullet X$ has a recursion iff $X \xrightarrow{*} \dots X \dots$, where " $\dots X \dots$ " is some sentential form w/ X .

Leftmost Recursion: X has leftmost rec. iff $X \xrightarrow{*} X \dots$ (or $X \xrightarrow{*} X \alpha$)

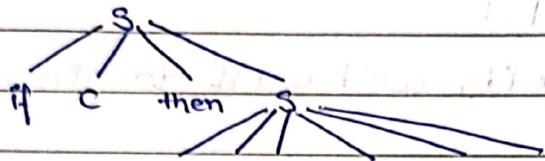
Rightmost Recursion: X has rightmost rec. iff $X \xrightarrow{*} \dots X$ (or $X \xrightarrow{*} \alpha X$)

Ex. Analyze: $S \xrightarrow{①} \text{if } C \text{ then } S \mid \text{if } C \text{ then } S \text{ else } S \mid \text{other}$
 Say $w = \text{if } C \text{ then if } C \text{ then } S \text{ else } S$

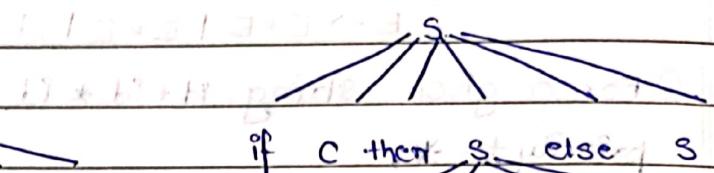
- ① (a) $S \xrightarrow{①} \text{if } C \text{ then } S \xrightarrow{②} \text{if } C \text{ then if } C \text{ then } S \text{ else } S$
- ① (b) $S \xrightarrow{①} \text{if } C \text{ then } S \text{ else } S \xrightarrow{①} \text{if } C \text{ then if } C \text{ then } S \text{ else } S$.

So, this grammar is ambiguous.

Now the question is which of the two is desired.



PT1



PT2

In C language the rule is: an else of this form would associate with closest if preceding it.

So, PT1 is desired.

In Prod. ② of grammar, the problem is caused by having an option to derive Prod ①. If we eliminate it, we can possibly eliminate the ambiguity.

$$S \rightarrow M \mid U$$

$$M \rightarrow \text{if } C \text{ then } M \text{ else } M \mid \text{other}$$

$$U \rightarrow \text{if } C \text{ then } S \mid \text{if } C \text{ then } M \text{ else } U$$

This is unambiguous grammar.

* Left Recursion & Left Factoring

Apart from Ambiguity, two major problems faced by CFGs are:

(1) Left Recursion

(2) Non-left-factored production rules.

Non-left-factored productions: A grammar may have productions of form $A \rightarrow \alpha B_1 | \alpha B_2$

$$A \rightarrow \alpha B_1 | \alpha B_2$$

These are non-left-factored, since they have a prefix α common.

This can be left factored as:

$$A \rightarrow \alpha B$$

$$B \rightarrow B_1 | B_2$$

$$\text{Ex. } S \rightarrow aS | aaA | c$$

$$A \rightarrow b$$

$$\Rightarrow S \rightarrow ab | ac$$

$$B \rightarrow S | aA$$

$$A \rightarrow b$$

$$\text{Ex. } S \rightarrow aS | abA | abc$$

$$A \rightarrow a$$

$$\Rightarrow S \rightarrow aB$$

$$B \rightarrow bS | bc$$

$$A \rightarrow a$$

$$B \rightarrow aA$$

- The problem with having ~~non-left-factored~~ non-left-factored productions in Grammar would become apparent while discussing Parsers, but basically, while trying to substitute a variable A , with production rules of above form, confuses a parser.

- Left Recursion: A ~~production~~ variable S has left recursion iff $S^* \rightarrow S\alpha \{ \alpha \in (V+T)^* \}$.

- The idea to resolve left recursion is to convert it to right recursion, since it doesn't cause any problem for parsers.

- For a variable A in CFG, with left recursion of form:

$$A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m | B_1 | B_2 | \dots | B_n$$

It can be converted to:

$$A \rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_n A' | B_1 | B_2 | \dots | B_n$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_m A' | \epsilon$$

don't forget to add this
 $\alpha_{m+1} A'$

$$A \rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_n A' | B_1 | B_2 | \dots | B_n$$

$$\text{Ex. } S \rightarrow Sa | Sb | b | c$$

$$\Rightarrow S \rightarrow bS' | cS' | \dots | \epsilon$$

$$S' \rightarrow aS' | bS' | \epsilon$$

$$\text{Ex. } S \rightarrow Sa | SA | Ab$$

$$A \rightarrow aA | b$$

$$S' \rightarrow aS' | AS' | \epsilon$$

$$A \rightarrow aA | b$$

$$\text{Ex. } E \rightarrow E + T | T \Rightarrow E \rightarrow TE'$$

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

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

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

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' | \epsilon$$

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

This method works for direct left recursion.

- The grammar may also have indirect left recursion.

$$\text{Ex. } S \rightarrow Aa | b$$

$$A \rightarrow Sa | a$$

$$S \rightarrow Aa \rightarrow Saa$$

$$A \rightarrow Sa \rightarrow Aaa$$

$$\text{Ex. } S \rightarrow Aa | b$$

$$A \rightarrow Ac | Sd | b$$

To remove left recursion in case of (direct + indirect):

$$\{ \text{Say, } G: \quad S \rightarrow Aa \mid b \\ \} \quad \quad \quad A \rightarrow Ac \mid Sd \mid c$$

(1) Decide an order of non-terminals.

$$\{ \text{Say, } S, A. \}$$

for each x_i , if

(2) In order, x_1, x_2, \dots (above order),[↑] in the productions of x_i , if any of x_1, x_2, \dots, x_{i-1} appear, substitute them. Then eliminate any direct left recursion in x_i , if there.

So, in order:

β : S has no possibility for substitution.

S has no direct left recursion.

A:

$$A \rightarrow Acl | Sd | \epsilon$$

↓
substitution

$$\Rightarrow A \rightarrow Acl | Aad | bd' | \epsilon$$

↓
direct left recursion

↓

$$\Rightarrow A \rightarrow bdA' | A'$$

$$A' \rightarrow cA' | adA' | \epsilon$$

$$\therefore G': \quad S \rightarrow Aa \mid b$$

$$A \rightarrow Aad \mid bdA' \mid A'$$

$$A' \rightarrow cA' \mid adA' \mid \epsilon$$

Ex. $A \rightarrow Cd$ $B \rightarrow Ce$
 $C \rightarrow A|B|f$

S1: Order A, B, C. down from left to right

S2:

A: first $A \rightarrow Cd$

No subs. No direct left rec.

B: $B \rightarrow Ce$. A recursive left branch structure
 No subs. No direct left rec.

C: $C \rightarrow A|B|f$ \uparrow \wedge \downarrow
 subs. non-branching and D

$C \rightarrow Cd|Ce|f$ \uparrow \wedge \downarrow
 direct.left.rec.

$C \rightarrow fC'$ \uparrow \wedge \downarrow
 $C' \rightarrow dC' | eC' | e$ \uparrow \wedge \downarrow

$\therefore G': A \rightarrow cd$
 $B \rightarrow Ce$
 $C \rightarrow fC'$
 $C' \rightarrow dC' | eC' | e$

So, before using a grammar G_1 for a parser check ~~for~~:

- (1) Is the grammar ~~ambiguous?~~ to tell at p. 10
- (2) Does any variable in G_1 have left recursion? at p. 10
- (3) ^{Are} ~~Is~~ there any non-left-factored production?

Ex. $S \rightarrow aaA \mid aB \mid S\alpha \mid \beta$ Difficult to tell if it's ambiguous

$A \rightarrow a$

$B \rightarrow b$.

- ① Doesn't seem to be ambiguous.

② Left Factoring:

$S \rightarrow aX \mid S\alpha \mid \beta$ possible ways to do it

$X \rightarrow aA \mid B$ left factoring in X

$A \rightarrow a$

$B \rightarrow b$ base and now expand p. 72 for details

- ③ Removing left recursions: from left to right p. 72

$S \rightarrow aXs' \mid \beta s'$

$s' \rightarrow \epsilon$

$X \rightarrow aA \mid B$

$A \rightarrow a$

$B \rightarrow b$.

- * • Topdown Parsing: Beginning from the start symbol, try to estimate the productions to apply to end up at the string {source program} in this case { }.
- Bottom Up Parsing: Beginning w/ the string {source program} in this case { }, try to apply productions in reverse to convert the program back to start symbol.
- Parsing: The process of constructing the parse tree for a sentence generated by a given grammar.
- For any CFG_i, G_i , any string $w \in L(G_i)^*$ can be parsed using CYK algorithm in $O(n^3)$.
- Subsets of CFG_i however can be used & in conjunction w/ specific algo. to make parsing more efficient.

* Top Down Parsing

- There are two types of TopDown parsers :
- Backtracking Used: For every variable substitution if multiple options, then choose an option & see if it leads to the string, else backtrack & take different string.
- Non Backtracking (Predictive Parsing): For a variable substitution, choose a production, if it works, great, else you can't backtrack.

Top Down Parsing

with Backtracking

Recursive Way (Recursive Descent Parser)

Non recursive way

Without backtracking

Recursive Way (RDP)

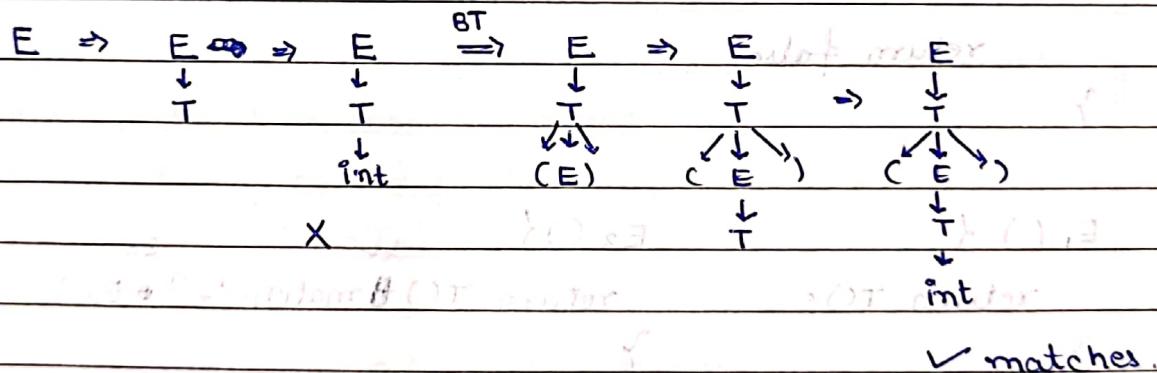
Non recursive way.

Recursive Descent Parser : A parser which uses recursion & descent, i.e., goes top to bottom (TopDown).

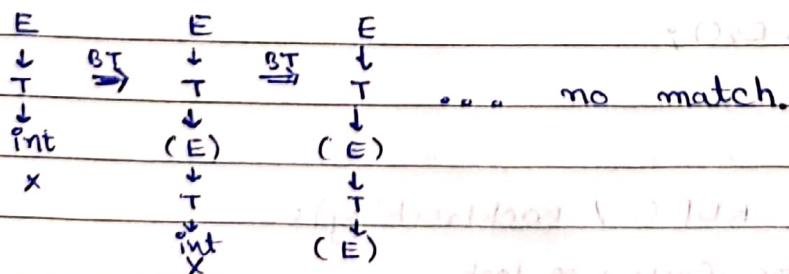
TD Parsing w/ Backtracking (RDP) :

Ex. $E \rightarrow T \mid T+E$

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



$w = \text{int} +$



• RDP (w/ Backtracking) ::

To implement it, for each non-terminal, write a procedure which explores all options.

Ex. $E \rightarrow T \mid T+E$
 $T \rightarrow id \mid (E)$

~~test E~~ ~~check if current symbol~~

(current char) omitted at first step -

match(t) {
 if (currSymbol == t) {
 move to next symbol;
 return true;

}

return false;

}

$E_1()$ {

return T();

}

$E_2()$ {

return T() ~~&~~ match('+') ~~&~~ E();

}

$E()$ {

return $E_1() \parallel E_2()$;

}

• Some pitfalls for RDP (w/ Backtracking):

• Left Recursion. Causes no loop.

• Expensive due to backtracking.

• Sometimes may not even be able to parse a valid string w/ L(G). {Ex shown later}

Predictive Parsing: The parser at every step tries to predict the production based on some criteria & never backtracks from it.

(The criteria generally is next "few" symbols, Lookahead.)

. It can also be done both in Recursive & Non-Recursive way.

. It is easy to check & match next few, say one, symbols when the productions start with ~~non-terminal~~ terminal.

$$\text{Ex. } E \rightarrow +E \mid id$$

① ②

$$w = ++ id \quad \uparrow$$

$(+) (+) id = (++) id$

at cur pos., lookahead symbol = '+'

trying to match with prods. we get a match with prod ①. Use it.

$$(A) T2917 \cup (B) = (B) T2917 \quad (B, id) = (A) T2917$$

. There can be two problems here:

① Non left factored productions.

$$\text{Ex: } E \rightarrow +E \mid +T \mid id$$

$$T \rightarrow = P$$

This can be solved by left factoring.

② Major problem is when the productions don't have non-terminals at start.

$$\text{Ex. } E \rightarrow TE'$$

Here E can't match ~~non-terminal~~ terminal

$E' \rightarrow +TE' \mid \epsilon$ with lookahead directly.

. So, for Predict Parsing, we need two concepts:

(1) First Set

(2) Follow Set.

• First Set:

- For a non-terminal A, the First Set, $\text{FIRST}(A)$, is the set of terminals that can be at the start of strings produced by A.

$$\text{FIRST}(A) = \{ t \mid A \xrightarrow{*} w, w \in T^*, t = w[0] \}$$

- Main Idea: For any production $A \rightarrow B_1 B_2 \dots B_K$, keep processing B_1, B_2, \dots, B_K until you encounter a NonTerm. that doesn't produce ϵ . Then go to next rule.

- $a \in \text{FIRST}(S)$, iff $S \xrightarrow{*} a$

Ex. $S \rightarrow ab \mid cd \mid de$

$$\text{FIRST}(S) = \{a, c, d\}$$

Ex. $S \rightarrow A b \mid a$

$$A \rightarrow a A \mid \epsilon$$

$$\text{FIRST}(A) = \{a, \epsilon\} \quad \text{FIRST}(S) = \{a\} \cup \text{FIRST}(Ab)$$

so $\text{FIRST}(S) = \{a\}$ but $S \not\xrightarrow{*} \epsilon$ so ϵ is not included but $S \xrightarrow{*} \epsilon$ is the rule (1)

$$\text{So, } \text{FIRST}(S) = \{a\}$$

Since $S \xrightarrow{*} \epsilon b \equiv S \xrightarrow{*} b$

Ex. $S \rightarrow a \mid b A$

$$A \rightarrow AB \mid a \mid \epsilon$$

$$B \rightarrow c \mid d \mid e \mid f$$

$$\text{FIRST}(B) = \{c, d, e, f\} \quad \text{FIRST}(A) = \{a, \epsilon, c, d, e, f\}$$

$$\text{FIRST}(S) = \{a, b\}$$

$$a \xrightarrow{*} c$$

• Formal defⁿ: $G(V, T, P, S)$

- For $t \in T$:

$$\text{FIRST}(t) = \{t\}$$

- For a var. X , with prod., $X \rightarrow Y_1 Y_2 Y_3 \dots Y_K \mid K \geq 1$:

- $a \in \text{FIRST}(X)$, if $t \in \text{FIRST}(Y_i)$ AND $Y_1 Y_2 \dots Y_{i-1} \xrightarrow{*} a$

- $a \in \text{FIRST}(X)$, if $Y_1 Y_2 \dots Y_K \xrightarrow{*} a$.

- If $X \rightarrow a$, then $a \in \text{FIRST}(X)$.

Ex. $S \rightarrow ab \mid ABSd \mid BC$ so $\text{FIRST}(S) = \{a, b, f\}$

$A \rightarrow eS \mid Se \mid a$ so $\text{FIRST}(A) = \{a, e, d, f\}$

$B \rightarrow dS \mid Sd \mid a$ so $\text{FIRST}(B) = \{a, d, e, f\}$

$C \rightarrow AB \mid Sf \mid a$ so $\text{FIRST}(C) = \{a, e, a, d, f\}$

Ex. $E \rightarrow TX$ so $\text{FIRST}(E) = \{c, int\}$

$T \rightarrow (E) \mid \text{int} \mid Y$ so $\text{FIRST}(T) = \{c, int\}$

$X \rightarrow +E \mid a$ so $\text{FIRST}(X) = \{+, a\}$

$Y \rightarrow *T \mid a$ so $\text{FIRST}(Y) = \{*, a\}$

Follow Set: The follow set, for some non-terminal X , $\text{FOLLOW}(X)$ is the set of all terminals that might come after X .

$$\text{FOLLOW}(X) = \{t \mid S^* \xrightarrow{*} X t \beta, t \in T; X, \beta \in (V+T)^*\}$$

{ S is start symbol of grammar}

Ex. $S \rightarrow Sa \mid Sb \mid e$

$$\text{FOLLOW}(S) = \{a, b, \$\}$$

In Parsing, we always assume a string w is ~~not~~ always followed by a special symbol on input tape called End Marker symbol ($\$$). So, the start symbol S , $\text{FOLLOW}(S)$ will always

contain \$, since we assume $S \xrightarrow{*} w$, so S is followed by \$ after S derives w .

- Formal defⁿ: Apply the following rules for all NonTerm. A, until no new symbol can be added to Δ their Follow sets:

- ~~Rules~~ $\$ \in \text{FOLLOW}(S)$, where S is start symbol.
- If there is a production $X \rightarrow \alpha A \beta$, $\{\alpha, A \in V; \alpha, \beta \in (V \cup T)^*\}$, then everything in $\text{FIRST}(\beta)$ except A is in $\text{FOLLOW}(A)$.
- If there is a prod. $(X \rightarrow \alpha A) \text{ or } (X \rightarrow \alpha A \beta \text{ AND } A \in \text{FIRST}(\beta))$ then, everything in $\text{FOLLOW}(X)$ is in $\text{FOLLOW}(A)$.

		FIRST	FOLLOW
$E \rightarrow 9 \text{int} b + a \cdot B$	E	int, (\$, +, *,)
$O \rightarrow + 1 *$	O	+,*	int, (

	FIRST	FOLLOW
$S \rightarrow a S S a b B C$	S	a, b, e, A
$A \rightarrow A S a E$	A	a, a, b, e
$B \rightarrow B A A B C$	B	A, a, b, e
$C \rightarrow c G$	C	a, e
		\$, a, b, e

	FIRST	FOLLOW
$E \rightarrow T X$	E	C, int
$T \rightarrow (E) \text{int} y$	T	C, int
$X \rightarrow + E A$	X	+, \$, A
$y \rightarrow * T A$	y	*, A
		\$, +, *,)

		FIRST	FOLLOW
Ex. $E \rightarrow TE'$	E	C, id	$\$,)$
$E' \rightarrow +TE' \mid \epsilon$	E'	$+, E$	$\$,)$
$T \rightarrow FT'$	T	C, id	$+, \$,)$
$T' \rightarrow *FT' \mid \epsilon$	T'	$*, E$	$+, \$,)$
$F \rightarrow (E) \mid id$	F	$(, id)$	$*, +, \$,)$

Ex. $S \rightarrow aBAS \mid \epsilon$	FIRST	FOLLOW
$A \rightarrow ba \mid SB$	a, b, S	$\$, a, c, A$
$B \rightarrow CA \mid S$	C, A	$a, \$, c, A$
	B	c, a, A

Predictive Parsing (contd.)

- TDParser which looks at next few symbols in string & predicts which production to use.

- It has linear parsing time, since no backtracking.

- Predictive Parsers can parse language of LL(K) grammars:

- L: left-to-right scan of inputs.

- L: leftmost derivation

- K: means predict based on next K tokens (lookahead)

- If we are looking at only one next token, then the parser can parse LL(1) grammars.

LL(1) Parser: Parsing Table

- Construct a Parsing Table T for CFG G. For each production $A \rightarrow \alpha$ in G:

mark cells that will be filled later

- $\forall t \in FIRST(\alpha)$, do $T[A, t] = A \rightarrow \alpha$

- If $A \in FIRST(\alpha)$, then $\forall t \in FOLLOW(A)$, do $T[A, t] = A \rightarrow \alpha$.

Ex. $S \rightarrow aS$ $\text{FIRST}(S) = \{a, b, c\}$

$S \rightarrow bS$ $\text{FOLLOW}(S) = \{\$\}$

$S \rightarrow c$

T	a	b	c	\$
S	$s \rightarrow aS$	$s \rightarrow bS$	$s \rightarrow c$	

LL(1) Parsing table.

With help of this table we can perform parsing. Obtain the order of LMD. Check if string $\in L(G_1)$ or not.

say, $w = abac$.

We start with Start Symbol S. Two pointers, one on current input tape symbol & other on current derivation symbol.

Input Tape	Derivation	Action
abac \$	S	Init state. Take T [S, a].
abac \$	as	Match. Move pointer.
abac \$	a\$	Take T [S, b]
abac \$	ab\$	Match. move pointer.
abac \$	ab\$	Take T [S, a]
abac \$	abas	Match. move.
abac \$	aba\$	Take T [S, c]
abac \$	abac	Match. move.
abac \$	abac	

So, the string is successfully matched.

(This is not the full algorithm, btw.)

	FIRST	FOLLOW
Ex. $S \rightarrow SAa \mid E$	$S \rightarrow E, C, b, a$	$\$, a, b, c$
$A \rightarrow Sb \mid Ac \mid C$	$A \rightarrow E, C, b, a$	$a, c, \$$

T	t, a, x, T	$b, y, int, C, \$$	$\$$
S	$S \rightarrow Sa,$ $S \rightarrow A$	$S \rightarrow Sa,$ $S \rightarrow A$	$S \rightarrow A$
A	$A \rightarrow Sb,$ $A \rightarrow Ac, A \rightarrow A$	$A \rightarrow Sb,$ $A \rightarrow Ac, A \rightarrow A$	$A \rightarrow Ac, A \rightarrow A$

- If the Parsing Table T for a CFG G has a cell w/ more than one entry then G is not LL(1) grammar.

	FIRST	FOLLOW
Ex. $E \rightarrow \overset{①}{T} X (A) T \overset{②}{\rightarrow} Y \mid \overset{③}{int}$	$E \rightarrow C, int$	$\$,)$
$T \rightarrow (E) \mid \overset{④}{int} Y \mid \overset{⑤}{T}$	C, int	$+, \$,)$
$X \rightarrow +E \mid \overset{⑥}{A} \mid \overset{⑦}{int}$	$+, A$	$\$,)$
$Y \rightarrow *T \mid A \mid \overset{⑧}{A} \mid \overset{⑨}{int}$	\ast, A	$+, \$,)$

T	$s, (,), +, *, \ast, int, \$, A, T, Y$	int	$+, \ast$	$*$	$\$$
E	1	1	1	1	1
T	2	3	2	2	2
X	5	4	4	5	5
Y	7	7	6	7	7

- In case a CFG G has only one production per variable, it will always be LL(1). {since for each row, any column will have atmost one entry, so conflict is impossible.}

• So, a CFG G is LL(1) iff ; for all prod's; $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$

• $\forall i, j \quad \text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \emptyset$.

• AND If for some α_i ; $\alpha_i \xrightarrow{*} A$,

then, $\forall j, \underset{j \neq i}{\text{FOLLOW}(A)} \cap \text{FIRST}(\alpha_j) = \emptyset$.

• So, a CFG G with a left recursive variable can never be LL(1). Since, the production will be of the form:

$$A \rightarrow A\alpha_1 | B_1 | B_2 \dots$$

\downarrow
 $\text{FIRST}(A\alpha_1) \cap \text{FIRST}(B_1) \neq \emptyset$

So, $\text{FIRST}(B_1) \subset \text{FIRST}(A\alpha_1)$

(since $\text{FIRST}(B_1) \subset \text{FIRST}(A)$).

• Also, a non-left factored CFG can't be LL(1) either.

Since there would be a variable A with productions of the form:

$$A \rightarrow \alpha_1 | \alpha_2 | \dots$$

So, ~~a T[A, x]~~ would contain two or more entries.

$$T[A, \boxed{x}]$$

where $x \in \text{FIRST}(\alpha)$

{ in case $\text{FIRST}(\alpha) = \{\alpha\}$, then conflict in $T[A, \alpha]$ where }

{ $\alpha \in \text{FOLLOW}(A)$ }

Ex. $E \rightarrow T + E \mid T$ Is this LL(1) grammar?

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

No, not left factored for E .

$$\text{FIRST}(T+E) \cap \text{FIRST}(T) \neq \emptyset$$

If a CFG G_1 is ambiguous, then its not LL(1).

If G_1 is ambiguous, then for some w there are two LMDs.

$$L_1: S \rightarrow \dots \rightarrow | \dots \rightarrow \dots \rightarrow \dots \rightarrow w$$

$$L_2: S \rightarrow \dots \rightarrow | \dots \rightarrow \dots \rightarrow \dots \rightarrow | \dots \rightarrow \dots \rightarrow w$$

Suppose till here both
LMDs are same.

At this point there is a confusion as to which path to predict.

The Predictive parser can be implemented in both way: Recursive Descent & Non-recursive.

For LL(1) Parser : Non Recursive Algorithm

Input: String w , Parsing Table M .

Initially on input buffer, we have $w\$$.

On stack we mark bottom $w\$\$$ & have the start symbol S on top.

Output: If $w \in L(G)$, a leftmost derivation of w .

Else, appropriate error.

Algorithm:

Let a be first symbol on input buffer;

let X be Top of Stack Value;

while ($X \neq \$$) {

if ($X = a$) pop stack top & let a be next symbol on buffer;

else if (X is terminal) error();

else if ($M[X, a]$ is error entry) error();

else if ($M[X, a]$ is $X \rightarrow Y_1 Y_2 \dots Y_k$) {

output the prod. $X \rightarrow Y_1 Y_2 \dots Y_k$;

pop stack;

push $Y_k, Y_{k-1}, \dots, Y_2, Y_1$ on stack, Y_1 on top;

}

let X be current TOS value;

}

Ex. $E \rightarrow TX^*$	FIRST	FOLLOW
$T \rightarrow (E) \mid \text{int}$	$(, \text{int})$	$(, \$)$
$X \rightarrow +E \mid a$	(C, int)	$(+, \$)$
$y \rightarrow *T \mid \lambda$	$(X, +, \lambda)$	$(\$, \lambda)$
$w = (\text{int} * \text{int}) + \text{int}$	y	$(+, \$)$

String formed of $\text{int} * \text{int} + \text{int}$ (2nd part) followed by $\$$ with FA

E 1 1

T 3 2 1st part ended with first part.

X 4 5 5 ends & from part 1

Y 7 6 7 7

$w = (\text{int} * \text{int}) + \text{int}$ didn't predict new predict. of input

3rd part found goes as final token no update

Buffer	Top Stack	Action
$(\text{int} * \text{int}) + \text{int}$	$E \$$	Init config.
$(\text{int} * \text{int}) + \text{int} \$$	$T X \$$	Take 1 (T + trigger).
$(\text{int} * \text{int}) + \text{int} \$$	$(E) - X \$$	Take 2.
$(\text{int} * \text{int}) + \text{int} \$$	$E) X \$$	Match. Pop. Incr on buffer
$(\text{int} * \text{int}) + \text{int} \$$	$T X) X \$$	Take 1. st. in fa
$(\text{int} * \text{int}) + \text{int} \$$	$\text{int} Y X) X \$$	Take 3 rd. X. fa
$(\text{int} * \text{int}) + \text{int} \$$	$Y X) X \$$	Match. Pop. Incr
$(\text{int} * \text{int}) + \text{int} \$$	$* T X) X \$$	Take 6.
$(\text{int} * \text{int}) + \text{int} \$$	$T(X) X \$$	Match. Pop. Incr.
$(\text{int} * \text{int}) + \text{int} \$$	$\text{int} Y X) X \$$	Take 3.
$(\text{int} * \text{int}) + \text{int} \$$	$Y X) X \$$	Match. Pop. Incr.
$(\text{int} * \text{int}) + \text{int} \$$	$X) X \$$	Take 7.
$(\text{int} * \text{int}) + \text{int} \$$	$X \$$	Match. Pop. Incr.
$(\text{int} * \text{int}) + \text{int} \$$	$+ E \$$	Take 5.
$(\text{int} * \text{int}) + \text{int} \$$	$E \$$	Match. Pop. Incr

Ex. The grammar $S \rightarrow aSa \cup a$ generates all even length strings of 'a's. A recursing-descent-parser (with backtracking) can be devised for the grammar.

Show that the RDP (w/ BT) recognizes a^2, a^4 & a^8 but not a^6 . (Suppose that order followed by RDP is: $s \rightarrow asa$, then $s \rightarrow aa$).

- Suppose the order followed is: ① $S \rightarrow aa$ ② $S \rightarrow aSa$. Let's try writing code for it.

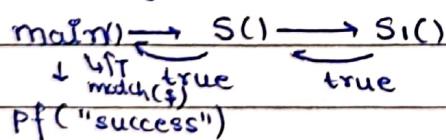
```

main() {
    sc();
    if (sc == match('$')) {
        printf("Success");
    }
}

sc() {
    if (cs1() && match('a')) {
        if (cs2() && match('a')) {
            return true;
        }
    }
    return false;
}

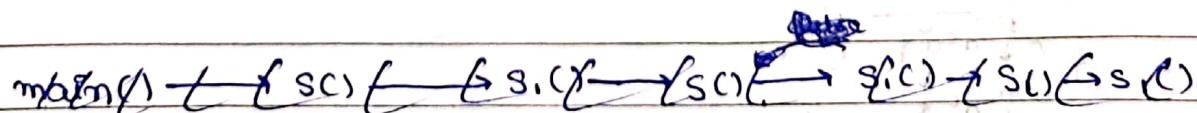
```

Now, try $w = aa$.

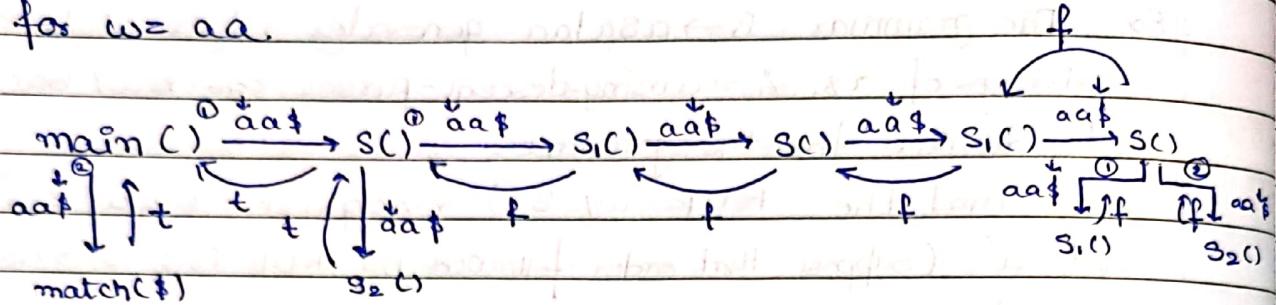


So, for any a^{2n} , $n > 1$, S would return true just after matching first two a 's. The next 'a' won't match '\$', so parsing unsuccessful even though $w \in L(G)$.

- Now suppose we have orders: ① $S \rightarrow aSa$ ② $S \rightarrow aa$
 $w = aa$.

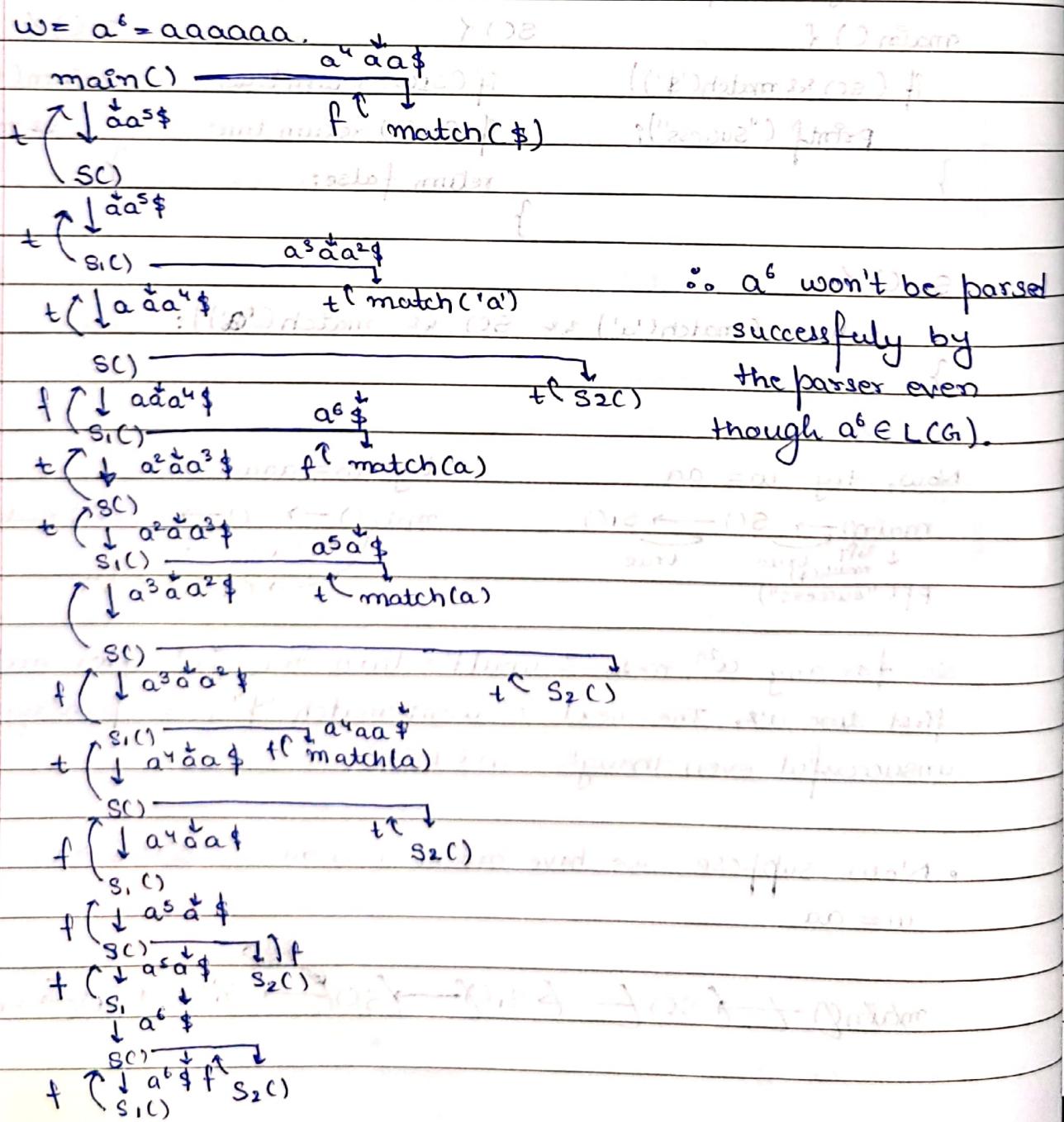


for $w = aa$.



so, aa is parsed successfully.

$w = a^6 = aaaaaa$.



- So, RDP (w/ Backtracking) are very sensitive to the order of productions' implementation. It can happen that a string $w \in L(G)$ doesn't get successfully parsed.

- A CFG₁ even if unambiguous & left factored & non-left-recursive can still be not LL(1).

Ex. $S \rightarrow AaB$ $A \rightarrow aB$ $B \rightarrow a\mid b$

$$\text{FIRST}(Aa) \cap \text{FIRST}(B) \neq \emptyset$$

$$= \{a\}$$

Left factoring A is not left-linear

Input same for right-linear transformation

- A grammar is non-left-linear. A right-linearization of it is not necessarily unique for most cases. Acceptance of a right-linearization of a grammar may not be unique.

Non-left-linear grammar is not right-linear.

$S \rightarrow a \mid AaB$ is not left-linear

- Non-left-linear grammar is not right-linear. Ex. $a^k b^l c^m$ and $a^k b^l c^m d^n$ are not right-linear.

Non-left-linear grammar is not right-linear.

$S \rightarrow a \mid AaB$ is not left-linear

Non-left-linear grammar is not right-linear.

$S \rightarrow a \mid AaB$ is not left-linear

Non-left-linear grammar is not right-linear.

$S \rightarrow a \mid AaB$ is not left-linear

Non-left-linear grammar is not right-linear.

$S \rightarrow a \mid AaB$ is not left-linear

Non-left-linear grammar is not right-linear.

* Bottom Up Parsing

- Sentential Form: Any α , such that $S \xrightarrow{*} \alpha$, where S is start symbol & $\alpha \in (V+T)^*$ is called sentential Form.
- Sentence: A sentential form α , such that $\alpha \in T^*$
- Left Sentential Form: A sentential form α that occurs in the leftmost derivation of some string.
- Right Sentential Form: A sentential form α that occurs in the rightmost derivation of some string.
- Handle: For a sentential form $\alpha\beta\gamma$, such that $S \xrightarrow{*} \alpha A\gamma \rightarrow \alpha\beta\gamma$, handle is two tuple $\langle A \rightarrow \beta$, at pos. following $\alpha \rangle$, i.e. handle consists of the production & the pos. at which the production rule applies in the sentential form.

Ex. $S \rightarrow a A c B e$
 $A \rightarrow A b b b$
 $B \rightarrow d$

List the handles in the rightmost derivation of $w = abcd e$

$S \rightarrow a A c B e \rightarrow a A b c B e \rightarrow a b b c B e \rightarrow a b b c d e$.

So, handle for:

$ab c d e : \langle B \rightarrow d, \text{ pos } 4 \rangle$ { we can also denote $B \rightarrow d$ by just the }
 $RHS : "d", \text{ pos } 4$
 $a A b c B e : \langle A \rightarrow b, \text{ pos } 1 \rangle$
 $a A c B e : \langle S \rightarrow a A c B e, \text{ pos } 0 \rangle$

Ex. $E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$ & if $w = id * id$, RMD.

$F \rightarrow id$.

Find handle of: (a) $id * id$ (b) $T * id$.

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

Handle for: $id * id : \langle F \rightarrow id, pos 0 \rangle$

$T * id : \langle F \rightarrow id, pos 2 \rangle$

Informally, handle is substring of some sentential form.

Ex. $S \rightarrow OS_1 | O_1$. Indicate ill-formed handle in each of these right sentential forms.

(a) 000111 (b) $00S11$

Since right sentential form, ~~RMD~~.

$S \rightarrow OS_1 \rightarrow 00S_11 \rightarrow 000111$

Handle for: $000111 : \langle S \rightarrow OS_1, 2 \rangle$

$00S_11 : \langle S \rightarrow OS_1, 1 \rangle$ or OS_1

"Informally, handle is substring of some sentential form such that it is on RHS of some production".

X This is misconception.

Ex. $S \rightarrow AB; A \rightarrow aa; B \rightarrow a$ & $w = aaa$.

$S \rightarrow AB \rightarrow Aa \rightarrow aaa$

So, in "aaa", 'a' at pos 1 & 2 are not handles even though they are substring of sentential form & appear on RHS of ~~aaa~~ $B \rightarrow a$.

"aa" at pos 0 is handle.

Ex. Which of the following statements are true:

- (A) Handle is any substring appearing on RHS of the prod.
- (B) If grammar is unambiguous, then we have only one handle in any right sentential form.
- (C) For any string $w \in L(G)$, if we keep on finding handles then we will reach to start symbol.
- (D) If $\alpha\beta\gamma$ is right sentential form & β is handle then γ is string of terminals.

A is false. Handle must be substring of some sentential form & also represent a step in some derivation, i.e.

$$S \xrightarrow{*} \alpha A \gamma \rightarrow \alpha \beta \gamma$$

\uparrow 1st step \downarrow substr. of sentential form.

B is true.

If G is unambiguous, for any string we have only one RMD, so every step in derivation is unique for RMD.

{ We can show by contradiction }

C is true. By def? $S \xrightarrow{*} \alpha A \gamma \rightarrow \alpha \beta \gamma$

D is true. If $\alpha\beta\gamma$ is right sentential form, then

$$S \xrightarrow{*} \alpha A \gamma \rightarrow \alpha \beta \gamma$$

If γ had some non-terminal, then we'd have first taken that production, otherwise it violates RMD.

Viable Prefix: For a sentential form $\alpha\beta^*$, where β is handle of the sentential form, Viable prefix is any prefix of $\alpha\beta$ substring.

So, to a viable prefix, it is always possible to add some appropriate sequence of terminals & non-terminals to obtain sentential form.

$$Ex. E \rightarrow E + T \mid T$$

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

Find viable prefix of $E + (id)$

$$E \rightarrow E + T \rightarrow E + (E) \rightarrow E + (\underline{E}) \rightarrow E + (id).$$

So, handle of $E + (id)$ is $\langle T \rightarrow id, 3 \rangle$. So, any ~~prefix~~ prefix of ' $E + (id)$ ' is viable prefix.

$$\{E, E+, E+(, E+(id), \}$$

$$Ex. E \rightarrow E + E \mid E * E \mid id. \text{ Find all handles in derivation of}$$

$w = id + id * id$. Also find viable prefixes of $id + id * id$.

(RMD)

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

$$\text{Handles: } \{E+E\}, \{E * id\}, \{E+E\}, \{id\}, \{id\}.$$

So, viable prefix of $id + id * id$ = pref. of $id = \{id, \epsilon\}$.

$$Ex. \text{ Find all viable prefixes of the grammar: } S \rightarrow 0S1101.$$

Earlier we had a string for the RMD of which we could find viable prefixes, but to find all the viable prefixes of a grammar, we first need the concept of LR(K) grammars.

Ex. Suppose we have a string $w = ab$, which belongs to grammar: $S \rightarrow A B$

$$A \rightarrow a$$

$$B \rightarrow b$$

If we try to go from string w to start symbol while reading it from left to right we get the following steps:

$$S_1: \underline{a}b \rightarrow Ab \quad \{ \text{Since } a \text{ can be reduced to } A \}$$

$$S_2: Ab \rightarrow A B \quad \{ A \text{ can't, } Ab \text{ can't. } b \text{ can be reduced} \}$$

$$S_3: \underline{AB} \rightarrow S$$

If we just reverse the steps, we essentially get RMD for w .

$$S \rightarrow \underline{AB} \rightarrow \underline{Ab} \rightarrow ab$$

So, this is what BottomUp parsers do. They perform reverse of rightmost derivation to go from string to start symbol. They don't do it blindly & apply some strategies. If they are able to parse the string back to start symbol, $w \in L(G)$, else not.

- If you know the handle at every step of "Reverse of RMD" we can simply replace w/ ~~LHS~~ LHS of handle at pos. of handle, & get the previous sentential form in RMD of w . Eventually we'd reach S .

- Idea of Bottom Up parsing: $S \rightarrow \gamma_0 \rightarrow \gamma_1 \rightarrow \dots \rightarrow \gamma_m \quad \{ w = \gamma_m \}$
for $i = m \text{ down to } 0$:

find the handle $\langle A_i \rightarrow \beta_i, pos \rangle$ in γ_i

replace β_i in γ_i at 'pos' to get γ_{i-1}

So, the crucial step is to find the handle at each step.

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

$w = id + id * id$. Parse in LtoR BottomUp way

Suppose $\alpha \mid \beta$ means, have read α , & haven't read β .

S1: $| id + id * id$.

S2: $id | + id * id$.

[Shift]

We have two options:

(1) Try to reduce what we've read.

(2) Shift $|$ to right.

{ say we reduce whenever we can. }

S3: $E | + id * id$.

[Reduce]

We can't reduce.

S4: $E + | id * id$.

[Shift]

S5: $E + id | * id$.

[Shift]

We can reduce 'id'.

S6: $E + E | * id$.

[Reduce]

S7: $E | * id$.

[Reduce]

S8: $E * | id$.

[Shift]

S9: $E * id |$.

[Shift]

S10: $E * E |$.

[Reduce]

S11: $E |$.

[Reduce]

If we have a choice to reduce, then we have handle in α portion of $\alpha \mid \beta$.

Also, if I have handle in α portion of $\alpha \mid \beta$, I'll always reduce.

So, the question, "Can I reduce?" is ~~same as~~ "Is there a handle in α portion of $\alpha \mid \beta$?".

But this doesn't mean, if we reduce some β in $\alpha \mid \beta$, then β is handle. We may reduce wrongly.

∴ If handle in α of $\alpha \mid \beta$ then reduce.

handle in α of $\alpha \mid \beta \rightarrow$ reduce.

- Data Structure for Bottom Up Parsing: We'll be putting the α portion of $\alpha|\beta$ on stack.

Say, we have $w = \alpha_1 \alpha_2 \alpha_3 \alpha_4 \alpha_5$ & we are at:

$\alpha_1 \alpha_2 \alpha_3 | \alpha_4 \alpha_5$

So, what is the question at this state?

The possible questions are: $[\text{Id}]$ or $b_1 + b_2 + b_3$

(a) α_3 OR (b) $\alpha_2 \alpha_3$ or (c) $\alpha_1 \alpha_2 \alpha_3$

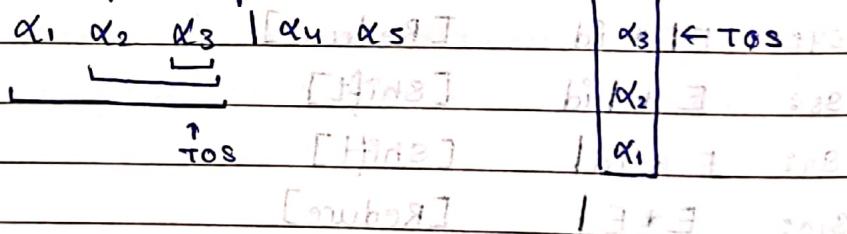
handle?

We can't ask "is α_2 a handle?" at this state, because we'd already have asked that question on state:

$\alpha_1 \alpha_2 | \alpha_3 \alpha_4 \alpha_5$

Since we found out that the answer was no, that's why we shifted the boundary.

So at any ~~state~~ state $\alpha|\beta$ only the substrings ending at TOS can be option for handles.



Ex. $E \rightarrow E+E | E*E | Id$ $w = \text{id} * \text{id}$

Stack	Action	Inp. Buffer	Remark/Question
\$	\bar{S}	\downarrow $\text{id} * \text{id} \$$	Init Config.
\$ Id	S	\downarrow $\text{id} * \text{id} \$$	Q: Is 'id' handle?
\$ E	R	\downarrow $\text{id} * \text{id} \$$	'id' is handle. Q: Is E handle?
\$ E*	S	\downarrow $\text{id} * \text{id} \$$	Q: Is '*' handle? OR Is 'E*' handle?
\$ E*Id	S	\downarrow $\text{id} * \text{id} \$$	Q: Is 'id' OR '*' handle? OR Is 'E*Id' handle?
\$ E+E	R	\downarrow $\text{id} * \text{id} \$$	'Id' is handle. Q: Is 'E+E' handle?
\$ E*	R	\downarrow $\text{id} * \text{id} \$$	Q: Is 'E' OR '*' handle? OR Is 'E+E' handle?

So, we reached start symbol. Found RMD. Verified that $w \in L(G)$.

- Note that this strategy is dependent on correct detection of handle. This blind strategy might go wrong.

Ex. $S \rightarrow ABA$

$A \rightarrow a$

$w = aba$

$B \rightarrow b$

s1: $|aba$ \xrightarrow{RMD} $a|ba$ \xleftarrow{RMD} $a|b|a$

s2: $a|ba \xrightarrow{R} A|ba$.

s3: $A|ba \xrightarrow{R} A|B|a$.

s4: ABA

at this stage if we first ask, is 'a' handle {before 'Ba' & 'Aba'} we get yes. If we blindly reduce:

$\xrightarrow{R} ABA$

Now we are stuck.

So, blindly reducing doesn't guarantee we are reducing handle.

Being able to reduce doesn't mean we are reducing a handle.

So, the most important question is : "How to find handle", since if we find handle on top of stack {in a part of $\alpha | B$ }, then we can reduce.

{ Since in sentential form, if we keep reducing, we'll eventually get to start symbol. }

Strategies to find handles:

- Strategy 1: Whenever possible to reduce, just reduce.

Ex. $S \rightarrow aA ; A \rightarrow a$, $w = aa$.

$aa\$ \xrightarrow{S} a|a\$ \xrightarrow{R} A|a\$ \xrightarrow{S} A|a\$ \xrightarrow{S} Aa\$|$

So this strategy fails easily.

This assumes that if reduction is possible (by say $A \rightarrow \alpha$) then it is handle.

Strategy 1 is called LR(0). Since we are Left to Right reading string, while performing Reverse RMD & looking at $_0$ symbols after $|$ in $\alpha | \beta$. \therefore LR(0)

\downarrow
0 lookahead.

- Strategy 2: Reduce $\alpha\beta\gamma$ to $\alpha A\gamma$ only if $a \in FOLLOW(A)$

Ex. $S \rightarrow aA ; A \rightarrow a$ where $w = aa$ and $\$$ is absent
 $FOLLOW(S) = \{\$\}$
 $FOLLOW(A) = \{\$\}$

$1aa\$ \xrightarrow{S} a1a\$ \xrightarrow{a} aa1\$ \xrightarrow{a} a11\$ \Rightarrow S1\$ \Rightarrow 1\$ \#$
 \uparrow
 don't red. with
 $A \rightarrow a, \because$
 $a \notin FOLLOW(A)$

Successfully parsed.

This makes sense. If the next symbol ~~is~~ after A after reduction is not in $FOLLOW(A)$, i.e. $\alpha A a \$ \notin FOLLOW(A)$ that means this is not a sentential form. So why take this derivation. And since $\alpha A a \$$ is not sentential form, so by defⁿ of handle B is not the handle of $\alpha B a \$$.

This strategy is called: SLR(1) parsing.

S stands for Simple.
 Simple LR(1).

- Strategy 3: Reduce $\alpha B a \$$ to ~~$\alpha B \gamma$~~ only if $a \in FOLLOW(\alpha A)$.

This strategy is called CLR(1) or just LRC(1).
 C stands for canonical (i.e. standard).

	FIRST	FOLLOW
$S \rightarrow Ba$	$w = aaba$	$S \quad a, b \quad \$$
$B \rightarrow aAa \mid aab \mid bAb$		$B \quad a, b \quad a$
$A \rightarrow Ba$		$A \quad a, b \quad a, b$

Let's apply SLR(1) strategy.

~~Alaba\$ $\xrightarrow{S} A\overline{lab}\overline{a\$}$ $\xrightarrow{R} \overline{lab}\overline{a\$}$ $\xrightarrow{S} \overline{lab}\overline{ba\$}$ $\xrightarrow{R} \overline{lab}\overline{ba\$}$ $\xrightarrow{S} \overline{lab}\overline{ba\$}$ $\xrightarrow{R} \overline{lab}\overline{ba\$}$~~

~~Alaba\$ $\xrightarrow{S} A\overline{lab}\overline{a\$}$ $\xrightarrow{R} A\overline{lab}\overline{a\$}$ $\xrightarrow{S} A\overline{lab}\overline{a\$}$ $\xrightarrow{R} AA\overline{ba\$}$~~

~~AAba\$ $\xleftarrow{S} A\overline{Ab}\overline{a\$}$ $\xleftarrow{S} A\overline{Ab}\overline{a\$}$ $\xleftarrow{S} A\overline{Ab}\overline{a\$}$~~

So, parsing fails.

① Since $a \in \text{FOLLOW}(A)$ ② Since $b \in \text{FOLLOW}(A)$

Applying LR(1) A forward lookahead strategy

Suppose state is $aa|ba\$ \xrightarrow{S} aab|a\$ \xrightarrow{R} B|a\$$
and lookahead symbol is b forward

$s\$ \xleftarrow{S} s\$ \xleftarrow{R} B\$$

So, parsing successful.

There is another strategy called LALR(1), which we'll look at later.

LR(0) & SLR(1) use something called "LR(0) items".

LR(1) & LALR(1) can also use LR(0) items, but they use LR(1) items to make their job easier.

- LR(0) Item : For any production $X \rightarrow \alpha\beta$, we have following

LR(0) Items :

- $X \rightarrow \cdot\alpha\beta$ { same as $\mid\alpha\beta$ }
- $X \rightarrow \alpha\cdot\beta$ { same as $\alpha\mid\beta$ }
- $X \rightarrow \alpha\beta.$ { same as $\alpha\beta\mid$ }

- An item $X \rightarrow \alpha\cdot\beta$ means that:

- Have α on TOS { i.e. have already seen α }
- Expecting to find string derived from β on the input buffer.
- The parser is looking to reduce to X , if $\alpha\beta$ is found.

- CLOSURE (I) { I is LR(0) item } : For a LR(0) item I,

say $X \rightarrow \alpha\cdot A\beta$, where A is non-terminal, its closure contains itself, all prods. of A with the dot/border at start of LHS (i.e. $A \rightarrow \cdot Y$) and recursively if any prod has non-terminal B after dot/border then all prod of B with dot/border at start of RHS.

Ex. $S \rightarrow AB$

$A \rightarrow a\mid Aa\mid Ba$.

$B \rightarrow a$.

$$\text{Closure}(S \rightarrow \cdot AB) = \{ S \rightarrow \cdot AB, A \rightarrow a, A \rightarrow \cdot Aa, A \rightarrow \cdot Ba, B \rightarrow \cdot a \}$$

$$\text{Closure}(S \rightarrow A \cdot B) = \{ S \rightarrow A \cdot B, B \rightarrow \cdot a \}$$

Including $B \rightarrow \cdot a$ in closure of $S \rightarrow A \cdot B$, means that we've already seen A, & we expect to see string derived from B on input buffer & the options to get that include $B \rightarrow \cdot a$.

Similarly, for LR(0) item $S \rightarrow \cdot AB$, next symbol is non-terminal A, so we can get A through $A \rightarrow \cdot a$ or $A \rightarrow \cdot Aa$ or ~~$A \rightarrow \cdot B$~~ $A \rightarrow \cdot Ba$.

Building LR(0) automaton:

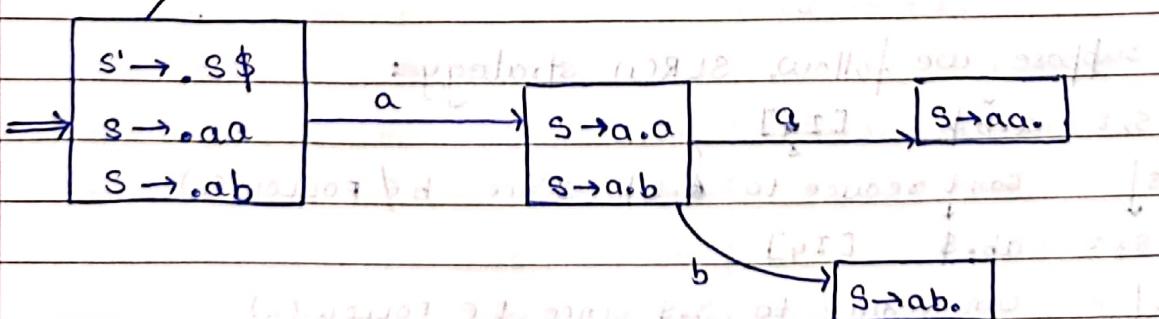
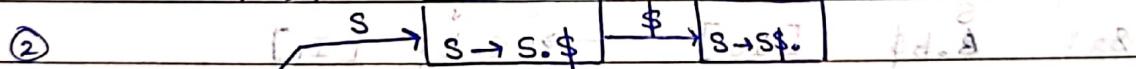
Before building LR(0) automaton for a Grammar G w/ start symbol S , we have to augment it, with the following production: $S' \rightarrow S\$$, meaning that we've seen nothing & expect to see string derived from ' $S\$$ '. If we successfully see a string derived from ' $S\$$ ' we can reduce to S' & accept the string.

The initial state for automaton is the closure of ' $S' \rightarrow S\$$ ' & every state is closure of a LR(0) item.

Ex. L $S \rightarrow aa|ab$

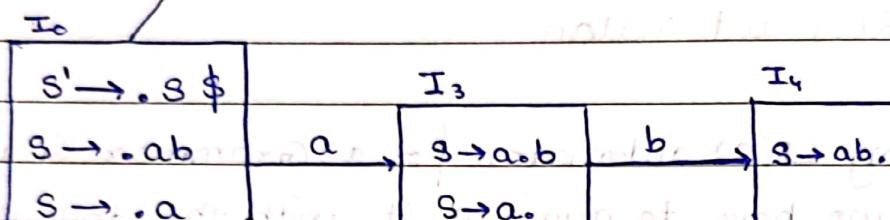
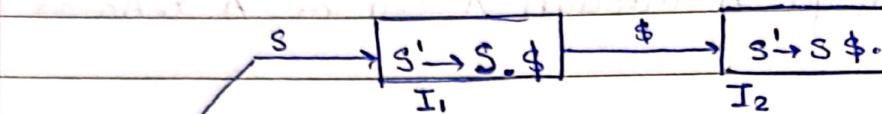
① $S' \rightarrow S\$$

$S \rightarrow aa|ab$



This automaton captures all the possibilities. Based on Parsing strategy either reduction is taken or not. Each state represents a closure: what we've seen, what we expect to see, if we see what we'll reduce to, & what are the possibilities to see it. $\{ \text{closure}(S \rightarrow \alpha \beta) \}$

Ex. $S \rightarrow ab\$$



Say $w = ab$.

$$\begin{array}{l} s_1: \overset{+}{a} b \$ [I_0] \\ s_2: a \overset{+}{b} \$ [I_3] \end{array}$$

Now, as per DFA, we've two choices.

Shift $\Rightarrow a.b.\$$

Reduce $\Rightarrow S.b\$$

Say we are following LR(0) strategy, i.e. always reduce when possible.

$$s_3: \overset{+}{S} b \$ [I_0] \longrightarrow \overset{+}{S} . b \$ [I_1]$$

Now we have no outgoing \xrightarrow{b} transition, nor can we reduce.

Suppose, we follow SLR(1) strategy:

$$s_2: a \overset{+}{b} \$ [I_3]$$

$\downarrow S$ can't reduce to $\overset{+}{a} b \$$ since $b \notin \text{FOLLOW}(S)$

$$s_3: a b . \$ [I_4]$$

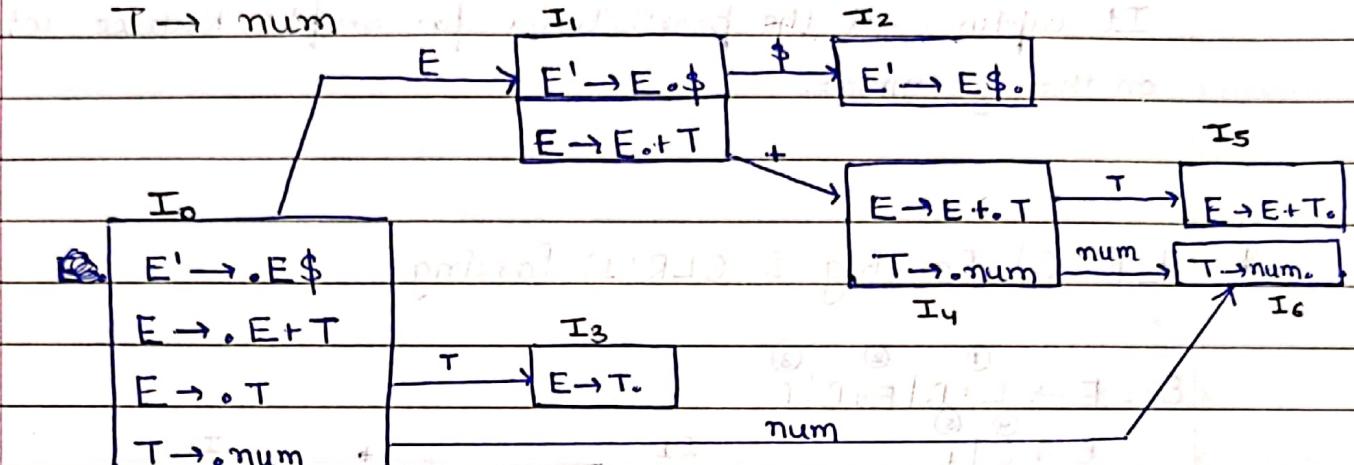
$\downarrow R$ can reduce to $S.\$$ since $\$ \in \text{FOLLOW}(S)$.

$$s_4: S . \$ [I_0] \longrightarrow S . \$ [I_1] \longrightarrow S \$. ^+ [I_2]$$

So, we can accept the string.

Ex. $E \rightarrow E + T \mid T$

$T \rightarrow \cdot \text{num}$ set I_1 to $E \rightarrow E \cdot \text{num}$ and I_2 to $E \rightarrow E + \cdot T$



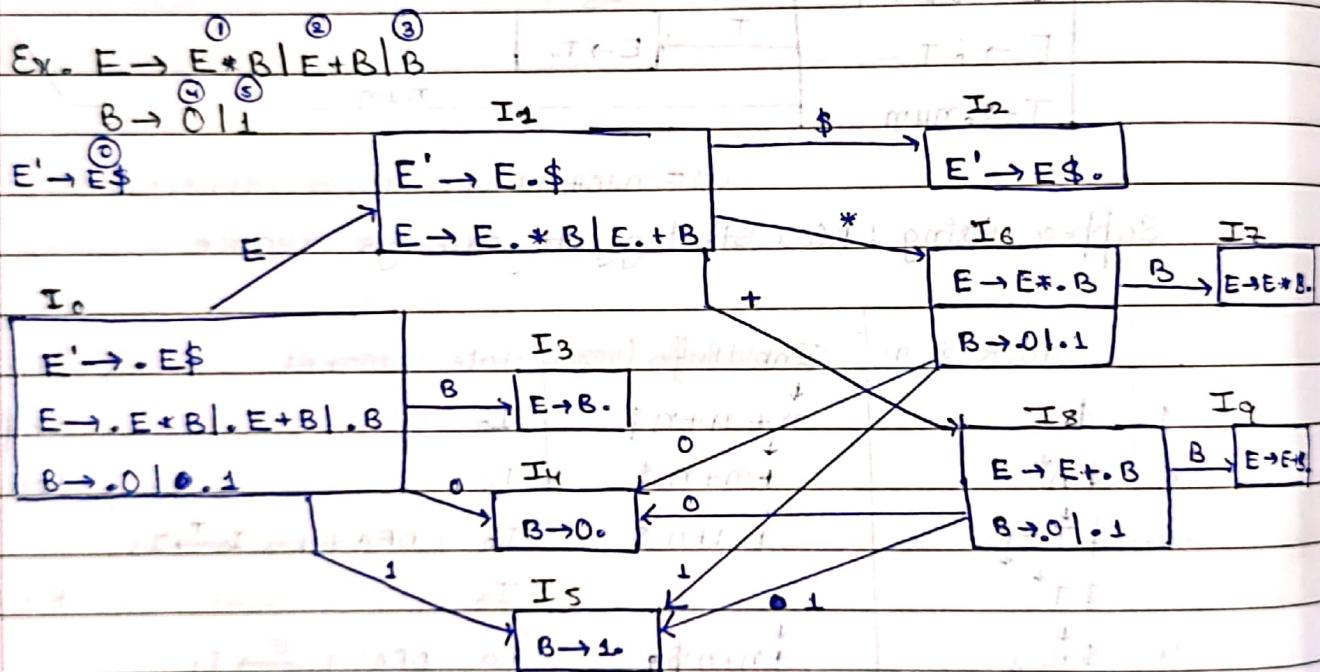
$$w = \text{num} \cdot \text{num} + \text{num} \text{ or } n \cdot n + n$$

Suppose using LR(0) strategy to always reduce.

Stack (seen)	Input Buffer (unseen)	State	Comment	Action
1. \$	↓ n + n + n \$	I ₀		shift
2. \$ n ⁺	↓ + n + n \$	I ₆		Reduce.
3. \$ T ⁺	↓ + n + n \$	I ₀	DFA takes I ₀ $\xrightarrow{T} I_3$	
\$ T ⁺		I ₃		Reduce.
4. \$ E ⁺	↓ + n + n \$	I ₀	DFA: I ₀ $\xrightarrow{E} I_1$	
\$ E ⁺		I ₁		shift
5. \$ E + ⁺	n + n \$	I ₁ $\xrightarrow{+} I_4$		
\$ E + ⁺		I ₄		shift
6. \$ E + n ⁺		I ₄ $\xrightarrow{n} I_6$		
\$ E + n ⁺	+ n \$	I ₆		Reduce.
7. \$ E + T ⁺	+ n \$	I ₀ $\xrightarrow{E} I_1 \xrightarrow{T} I_4$		
\$ E + T ⁺		I ₁ $\xrightarrow{T} I_5$		
\$ E + T ⁺		I ₄		
\$ E + T ⁺		I ₅		Reduce.
8. \$ E ⁺	+ n \$	I ₀		
\$ E ⁺		I ₁		shift
9. \$ E + ⁺	n \$	I ₁		
\$ E + ⁺		I ₄		shift
10. \$ E + n ⁺	\$	I ₄		

So the Automaton is independent of the Parsing Method.
It captures all the possibilities for shift & Reduce actions on the grammar.

* LR(0) Parsing & SLR(1) Parsing



We can convert this automaton to Tabular form:

States	Action					Goto
	+	*	0	1	\$	
0			S4	S5	I	E S1
1	S8	S6			(Accept) S2	B S3
2	R0	R0	R0 + R0	R0		
3	R3	R3	R3	R8	R8	
4	R4	R4	R4	R4	R4	
5	R5	R5	R5	R5	R5	
6		S4	S5			
7	R1	R1	R1	R1	R1	
8		S4	S5			
9	R2	R2	R2	R2	R2	

• Si entry: For action part, Si means Shift & Go to state i.

• Ri entry: For action part,

$T[I_i, \$]$: The entry $T[I_i, \$]$, where I_i is the current state & $\$$ is the current symbol on top of stack, means what action to take.

- The only difference b/w LR(0) & SLR(1) parsers is that LR(0), whenever it can reduce will reduce, while SLR(1) whenever it can reduce, first checks the follow of the reduced non-terminal.

So, in the Parsing Table, the only actions that would differentiate b/w LR(0) & SLR(1) are the R_i {reduction entries}.

$T[I_i, t] = \text{Reduce } A \rightarrow \alpha B$ means, when you are at state I_i & the next symbol is t , then reduce αB to A in $\text{Follow}(A)$.

Since SLR(1) always ensures that ' t ' is in $\text{FOLLOW}(A)$ before reducing, only those entries in Table T, for which

(1) I_i contains LR(0) item " $A \rightarrow \alpha B$ ".

(2) AND t is in $\text{FOLLOW}(A)$

can have $R_{A \rightarrow \alpha B}$ entry.

Ex. Previous Example

$$T[I_2, \$] = R_0 \quad \{ \because \text{FOLLOW}(E') = \{\$\} \}$$

$$T[I_3, +] = T[I_3, *] = T[I_3, *\$] = R_{E \rightarrow B}$$

$$T[I_4, +] = T[I_4, *] = T[I_4, \$] = R_{B \rightarrow 0}$$

$$T[I_5, +] = T[I_5, *] = T[I_5, \$] = R_{B \rightarrow L}$$

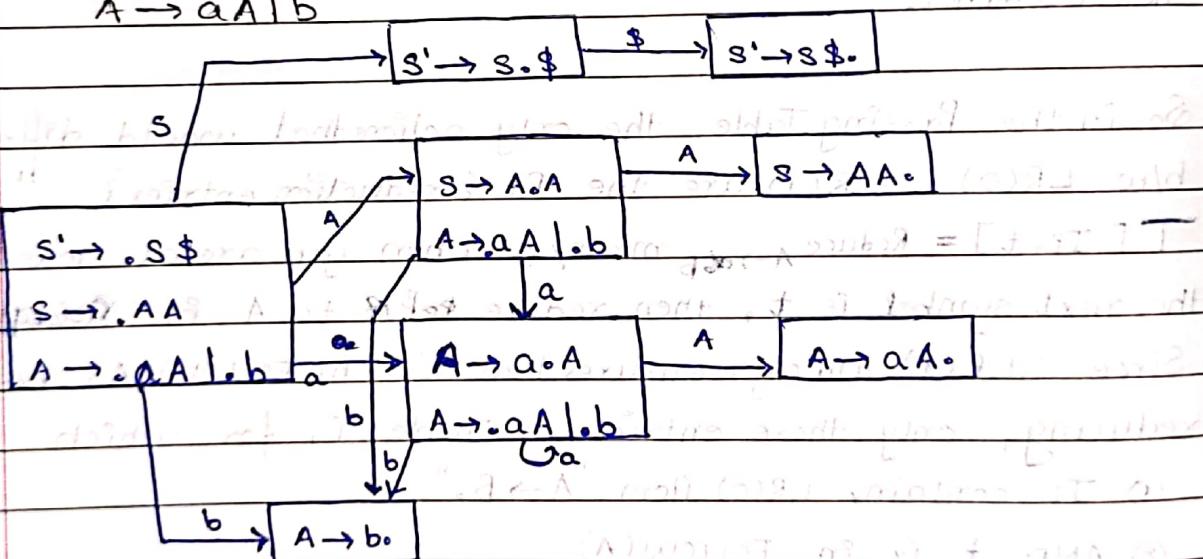
$$T[I_7, +] = T[I_7, *] = T[I_7, \$] = R_{E \rightarrow EKB}$$

$$T[I_9, +] = T[I_9, *] = T[I_9, \$] = R_{E \rightarrow F+B}$$

- So, every CFG suitable for LR(0) parsing is also suitable for SLR(1) parsing, i.e. every LR(0) grammar is SLR(1) grammar.
- Since we are only reducing entries from LR(0) table to obtain SLR(1) table.

Ex. Check if the following grammar is (1) LR(0) (2) SLR(1):

$$S \rightarrow AA \quad A \rightarrow aA \mid b$$



We know that in no cell of LR(0) & SLR(1) parsing table we can have two shift entries since on one symbol we get only one state.

For non-terminals however it is possible to have:

- ① One shift entry
- ② One reduce entry
- ③ One shift, one reduce entry } conflicts
- ④ Two reduce entries.

The last two {③ & ④} are problem for parsing.

since the parser gets confused as to which action to take.

So, for any parser its parsing table shouldn't have multiple entries.

For this example, we don't have any state with any entry $T[I_i, t]$ with multiple entries.

So, this is $LR(0)$ grammar.

{ Also, since $LR(0)$, hence $SLR(1)$ }

• Checking for $SLR(1)$ conflicts when a state has :

1. Two Reduce entries :

$$I_i : A \rightarrow \alpha.$$

$$B \rightarrow \beta.$$

In $SLR(1)$ table, entry $T[I_i, t]$ can have both

$R_{A \rightarrow \alpha}$ & $R_{B \rightarrow \beta}$ entry iff $t \in FOLLOW(A) \cap t \in FOLLOW(B)$.

So, iff $FOLLOW(A) \cap FOLLOW(B) = \emptyset$ then no $T[I_i, t]$ for any t can have multiple entries, i.e., no RR conflict.

2. One Reduce entry, One Shift entry :

$$I_i : A \rightarrow \alpha.$$

$$B \rightarrow \beta. tY$$

In $SLR(1)$ table, entry $T[I_i, t]$ can have both $R_{A \rightarrow \alpha}$ &

S_j entry entry iff $t \in FOLLOW(A)$ & t is terminal. So iff

$FOLLOW(A) \cap \{t\} = \emptyset$ then no SR conflict.
(or $t \notin FOLLOW(A)$)

• For $LR(0)$, any state having "two reduce entries" or "one shift, one reduce" entry has conflict.

Ex. Check if LL(1), LR(0), SLR(1): $A \rightarrow Cb_1 aCa$

$C \rightarrow aB$

$B \rightarrow Ca_1 G$

①	FIRST	FOLLOW
A	a	\$
C	b, a	
B	a, a	ba

Not LL(1), since, $\text{FIRST}(Cb) \cap \text{FIRST}(aCa) \neq \emptyset$

$I_1: A \rightarrow A\$ \xrightarrow{\$} A \rightarrow A\$.$

$I_2: A' \rightarrow A\$ \xrightarrow{a} A' \rightarrow A\$.$

$I_3: A \rightarrow C.b \xrightarrow{b} A \rightarrow Cb.$

I_0	C	I_5	I_6	I_7
$A' \rightarrow A\$$	$a \rightarrow A \rightarrow A\$.$	$A \rightarrow a.Ca$	$C \rightarrow a.a \xrightarrow{a} A \rightarrow aCa.$	
$A \rightarrow Cb_1 aCa$		$C \rightarrow a.B$	$B \rightarrow a.B \xrightarrow{B} C \rightarrow a.B$	
$C \rightarrow a.B$		$B \rightarrow a.Ca \xrightarrow{a} B \rightarrow a.Ca.$	$B \rightarrow a.Ca \xrightarrow{a} B \rightarrow a.Ca.$	

Not LR(0). Since I_5 has both shift & reduce entry for 'a'.

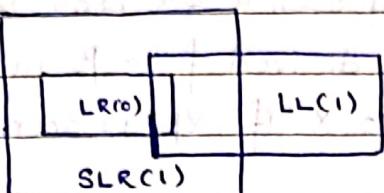
$T[I_5, a] = S_q \& R_{B \rightarrow e}$

similarly $T[I_9, a] = S_q \& R_{B \rightarrow C}$.

Only $I_5 \& I_9$ have conflicting entries: $B \rightarrow .$ & $C \rightarrow a.B$. since $a \in \text{FOLLOW}(B)$, hence conflict. Not SLR(1).

CFG, if

- Since every $^T LR(0)$ is $SLR(1)$, we can easily say:



$LLC(1)$ has no relation w/ $LR(0)$ & $SLR(1)$. They are independent.

i.e. CFG G is: ~~SLR(1) + LR(0) + LLC(1)~~

$$LLC(1) \leftrightarrow SLR(1)$$

$$LLC(1) \leftrightarrow LR(0).$$

but $LR(0) \rightarrow SLR(1)$.

• Limitation of $SLR(1)$:

- While reducing $\alpha\beta\gamma$ to $\alpha A\gamma$, $SLR(1)$ only checks if $a \in FOLLOW(A)$. It may be possible, $a \in FOLLOW(A)$ but ~~a~~ $a \notin FOLLOW(\alpha A)$. This may lead to the parser not reaching the start state of grammar.
- $LRC(1)$ parser takes into account the whole context, i.e., αA while checking if the next symbol a can follow αA or not.
- So basically ~~SLR(1)~~ parser can't remember ~~the~~ the left context α in $\alpha\beta$.

* $LR(1)$ Parsing & $LALR(1)$ Parsing

- For $LR(1)$ parsing we'll be required to remember the entire left context. To achieve this we'll use the concept of lookahead to keep track of ~~the~~ terminals which can follow a state. So, basically lookahead ~~is~~ the ~~a~~ for a $LRC(1)$ state is the set of all terminals that can follow the state, i.e. FOLLOW set of the state.

- LR(1) Item: $A \# LR(1)$ item $X \rightarrow \alpha.\beta, \alpha$ means:
 - The parser has already seen α (α is on T₀S)
 - Expects to find string derived from β in the next input.
 - Looking to reduce to X if $\alpha\beta$ is seen and $\alpha\beta$ is followed by 'a'.

- LR(1) parsing also, doesn't always work.

Ex. $S \rightarrow Xab | aaa$

$X \rightarrow a$

$|aaa \xrightarrow{S} a|aa \xrightarrow{R} X|aa \xrightarrow{S} x|a|a \xrightarrow{S} x|aa$

Show even LR(1) parsing strategy can fail.

- Since SLR(1) allows reduction based on just FOLLOW(A) & LR(1) on FOLLOW(αA), & FOLLOW(αA) \subseteq FOLLOW(A), every CFG which is SLR(1) is also LR(1) parseable.

- Closure of LR(1) Item: For a LR(1) item $S \rightarrow \alpha.AB, L$, where A is non-term & L is set of follow terminals, its closure contains:
 - Itself
 - All prod of A, $A \rightarrow \alpha$, include $A \rightarrow \alpha.L'$, where $L' = \bigcup_{\alpha \in L} FIRST(\alpha.B)$.

Ex. $A \rightarrow BC$, $B \rightarrow d|b$, $C \rightarrow a$. Find closure of $A \rightarrow .BC$, $\{a\}$.

closure ($A \rightarrow^* BC, \{a\}$) =

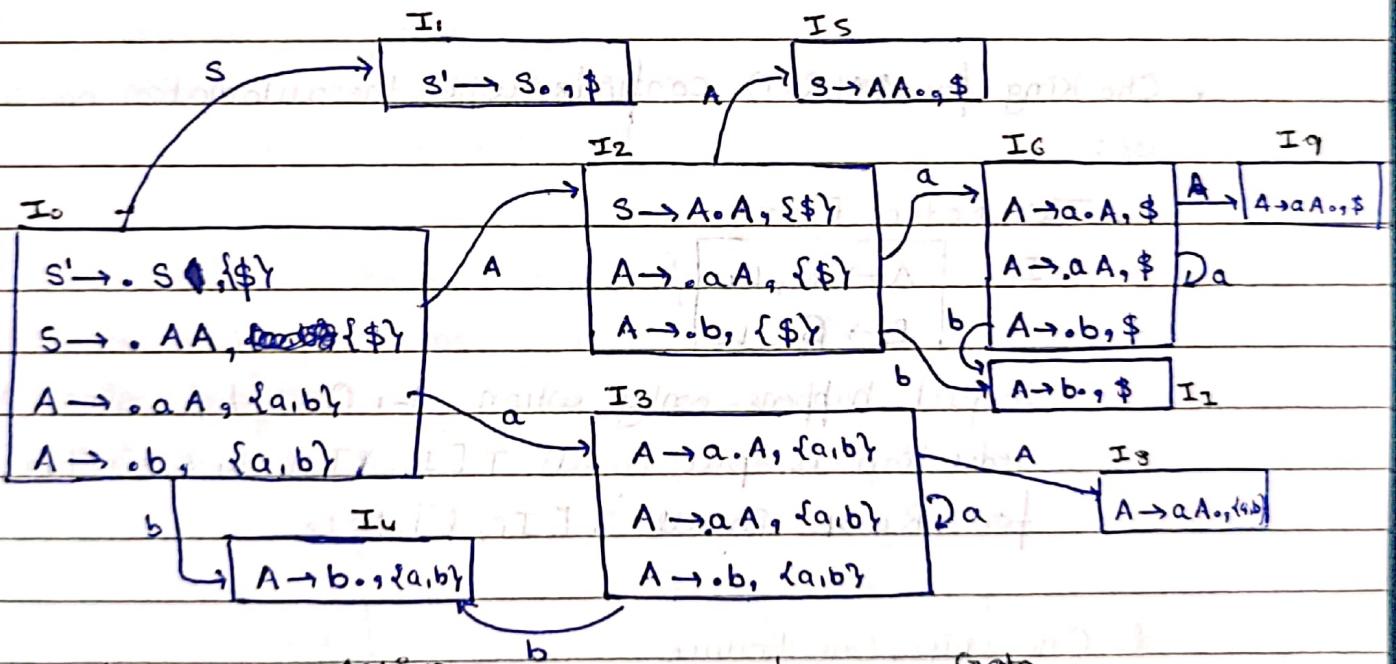
$A \rightarrow^* BC, \{a\}$ followed by a for $\text{FIRST}(B)$ to $\text{FIRST}(C)$

$B \rightarrow^* a, \{b, a\}$ followed by b for $\text{FIRST}(B)$ to $\text{FIRST}(C)$

$C \rightarrow^* b, \{a, b\}$ followed by a for $\text{FIRST}(B)$ to $\text{FIRST}(C)$

LR(0) Automaton

Ex. $S \xrightarrow{①} AA$ $\text{FIRST}(S) = \{a, b\}$ $\text{FIRST}(A) = \{a, b\}$
 $A \xrightarrow{②} aA1b$ removes a from $\text{FIRST}(A)$ because a is part of $\text{FIRST}(B)$



	action			state
	a	b	\$	
I0	S_3	S_4		1
I1			(R_0) Accept	2
I2	S_6	S_7		5
I3	S_3	S_4		8
I4	R_3	R_3		
I5			R_1	
I6	S_6	S_7		9
I7			R_3	
I8	R_2	R_2		
I9			R_2	

- So all entries of LR(1) parsing table are same as that of LR(0) & SLR(1) parsing table, except for reduction entry.
- For LR(1) parsing table, for a state I_i having LR(1) item of the form: $X \rightarrow \alpha \cdot, L$ where L is the set of Lookahead, $R_{X \rightarrow \alpha}$ is added only to $T[I_i, l]$ where $l \in L$.
- Again for a CFG, if the LALR(1) PT contains more than one entry in any cell, then the grammar is not LALR(1).

- Checking for LALR(1) conflicts when the automaton has a state w_1 :

1. Two Reduce Entry

$I_i:$	$A \rightarrow \alpha_0, L_1 \cdot \alpha_1 A$
	$B \rightarrow \beta_0, L_2 \cdot \beta_1 B$

Conflict happens only when $L_1 \cap L_2 \neq \emptyset$ since $R_{A \rightarrow \alpha} = R_{B \rightarrow \beta}$.
 reduction is put in all $T[I_i, l], l \in L_1$ & similarly for $R_{B \rightarrow \beta}$, in all $T[I_i, l], l \in L_2$.

2. One Shift, One Reduce

$I_i:$	$A \rightarrow \alpha_0, L_1$
	$B \rightarrow \beta \cdot a \beta, L_2$

Conflict happens when, $a \in L_1$.

Since $T[I_i, a] = I_j$ for some j

& if $a \in L_1$ then $T[I_i, a] = R_{A \rightarrow \alpha}$.

- LALR(1) Parsing means: Look Ahead Left to Right Reverse rightmost derivation.

LALR(1) Automaton can be created from LRC(1) automaton by combining all the states which only differ by lookahead.

Ex.

I _i :	$A \rightarrow \alpha \cdot \beta, l_{11}$
	$B \rightarrow \gamma \cdot S, l_{12}$

I _j :	$A \rightarrow \alpha \cdot \beta, l_{21}$
	$B \rightarrow \gamma \cdot S, l_{22}$

can be merged as:

I _{ij} :	$A \rightarrow \alpha \cdot \beta, l_{11} \cup l_{21}$
	$B \rightarrow \gamma \cdot S, l_{12} \cup l_{22}$

- Same rules to check conflict apply as LRC(1)
- Notice that for some grammar G_1 , if its ~~not~~ LRC(1) DFA has no ~~conflict~~ shift-reduce conflict, then its LALR(1) DFA would also not have any SR conflict.

I _i :	$A \rightarrow \alpha \cdot l_{11}$
	$B \rightarrow \beta \cdot a \gamma, l_{12}$

I _j :	$A \rightarrow \alpha \cdot l_{21}$
	$B \rightarrow \beta \cdot a \gamma, l_{22}$

if $\alpha \notin l_{11} \wedge \alpha \notin l_{12}$

then $\alpha \notin l_{11} \cup l_{12}$.

- Notice that if a CFG is not LRC(1) then it's not LALR(1) either, since we merge states to get LALR(1) DFA from LRC(1) DFA, so this can never reduce conflict {rather it may introduce conflicts}.
- $\therefore G_1$ is not LRC(1) $\rightarrow G_1$ is not LALR(1).

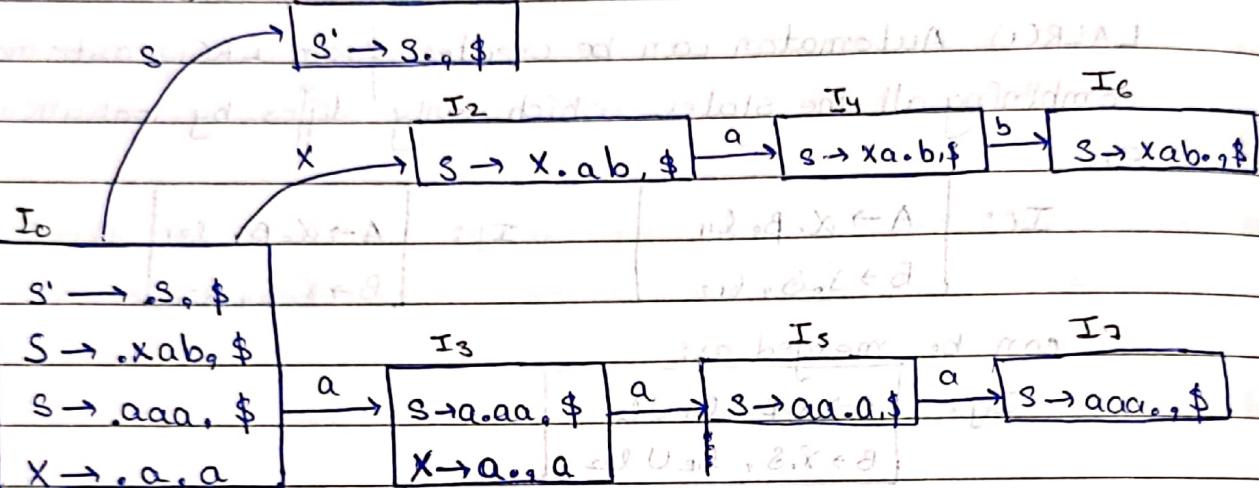
So,

If G_1 is LALR(1) then G_1 is LRC(1).

Ex. Check if LRC(1)

$$S \rightarrow Xab\alpha a$$

$$X \rightarrow a.$$

 I_1 

In I_3 : $\text{# } a \in \text{Lookahead}(X \rightarrow a.b, a)$ at value a .

Not LRC(1). Parsing stage will stuck with a .

So, the confusion at I_3 is: ① If 'a' is already seen & next symbol is 'a', we can reduce. Also ② If 'a' is already seen & 'a' is next symbol. If 'a' we can shift to I_5 .

Ex. $w = aab$

$\overset{\dagger}{aab} \Rightarrow \overset{\dagger}{a}ab$ At this state the parser is confused.
(I_0) (I_3)

If the lookahead was of length 2, this confusion could've been avoided. We could say, if lookahead is 'aa', shift. If lookahead is 'ab' reduce.

Ex. A CFG G , if G is LRC(1) but not LALR(1), then there's a SR conflict in Raw LALR(1) parsing table of G .

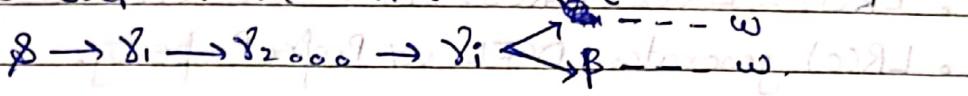
No. If G is LRC(1), then there's no SR conflict in LRC(1) PT.

If no SR conflict in LRC(1) PT, then no SR conflict in LALR(1) PT possible. So, if G is a LRC(1) but not LALR(1), there must be a RR conflict in LALR(1) table.

for account of backtracking for non-deterministic parser

- Ambiguous grammar can't be LR grammar, since for some string there would be two rightmost derivations (or more).

So, at some step $\gamma_i \in P_m RMD:$



there would be two paths in the RMD of w .

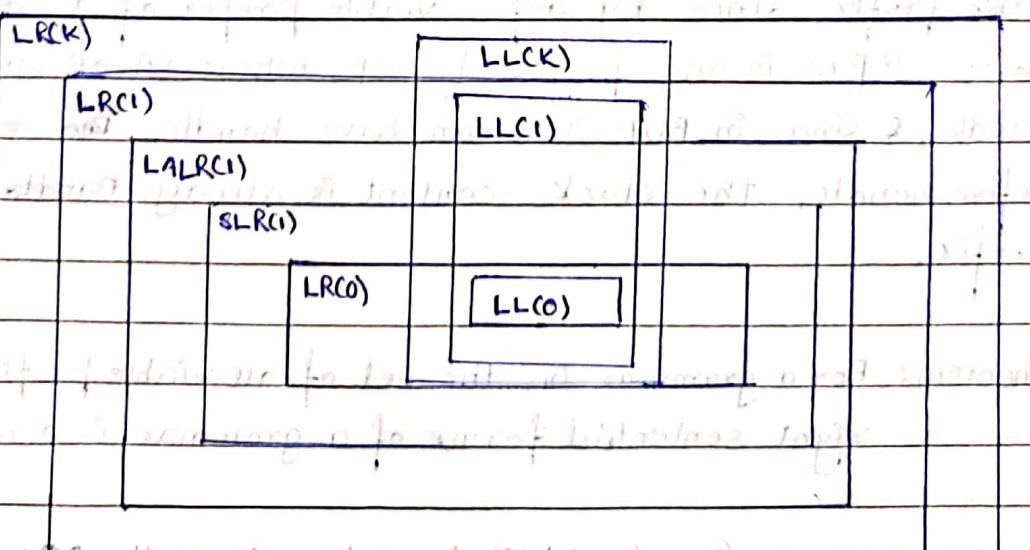
So in reverse RMD of w , the parser can't know either to reduce to α or reduce to β at a certain step.

- Since every SLR(1) is also LR(1) & LALR(1) is performing opⁿ similar to LR(1) by considering the whole left context while deciding to reduce, every SLR(1) grammar is also LALR(1).

- LL(1) is less powerful than LR(1) parsing strategy, because LL(1) doesn't really care about context of derivation.

Also, left recursion & non-left factoring are not problems for LR grammars.

* Types of Grammars



Classification of CFGs (Not CFLs)

Note: Power of a parser is measured in terms of no. of grammars they allow

- LL(1) can generate some subset of DCFL languages.
- $\text{LL}(K) \subset \text{LL}(1) \subset \text{LL}(2) \subset \text{LL}(3) \dots \subset \text{LL}(K)$
- $\text{LR}(0) \subset \text{LR}(1) \subset \text{LR}(2) \subset \dots \subset \text{LR}(K)$ {In terms of #grammars they allow}
- $\text{LR}(1) = \text{LR}(2) = \text{LR}(3) = \dots = \text{LR}(K)$ {In terms of lang. they produce}
- LR(0) generates DCFL w/ Prefix Property.
- LR(K) generates DCFL languages
- Although $\text{LR}(1)$ is more powerful than $\text{SLR}(1)$ since it can parse all the grammar that $\text{SLR}(1)$ can & more. But both $\text{LR}(1)$ & $\text{SLR}(1)$ can generate DCFL languages, i.e., for a $\text{LR}(1)$ grammar G w/ lang. $L(G)$ there exists some $\text{SLR}(1)$ grammar G' w/ lang $L(G')$ such that $L(G) = L(G')$.
- If a CFG G is $\text{LL}(1)$ & $w \notin L(G)$, then G is $\text{SLR}(1)$ also.
- For ~~exists~~ a CFG G to be $\text{LL}(K)$, it means that for any ~~symbol~~ non-terminal A , deciding which of A 's production to use for the next derivation, takes atmost K symbols look-ahead on the input buffer.

In bottom up parsing of a string $w \in L(G)$, whatever that appears on the stack {or α part of $\alpha\beta$ } must be a viable prefix, since by def, viable prefix of a sentential form $\alpha\beta\omega$ is any prefix of $\alpha\beta$, where $\langle A \rightarrow \beta, \text{after } \alpha \rangle$ is a handle, & since in BUP TOS can have handle, the ~~can~~ or string before handle, the stack content is always ~~handle~~ viable prefix.

Theorem: For a grammar G , the set of all viable prefixes of all the right sentential forms of a grammar is a regular language

For any grammar G , the LR(0) Automata acts as the ~~DO~~ NFA to recognize all viable prefixes of all sentential forms of G . Just make every state as final.

Semantic Analysis

- Syntax Analysis checks for structural or syntactical problems with the program.

Semantic Analysis checks if the meaning of program is clear or not & correct or not.

Ex. $a = b + c;$

In C, this is perfectly OK syntax. $id = id + id;$

But suppose that b is int, c is char array & b is double.

Does it make sense? $double = float + char[];$

It's the job of Symantic Analyzer to check the meaning of program.

- * Syntax Directed Definition: SDD[†] specifies the value of attr by associating semantic rules w/ grammar productions.

Ex. $E \rightarrow E_1 + T$ $E.val = E_1.val + T.val$

- * Syntax Directed Translation: SDT scheme[†] embeds program fragments called semantic actions within production body.

Ex. $E \rightarrow E_1 + \{ point '+' \} T$ Ex. $E \rightarrow E_1 + T \{ E.val = E_1.val + T.val \}$

- These are two famously used Semantic Analysis schemes.
- SDD has the rule separate from production body & the order of the execution of rule is not specified.
- SDT ~~can't~~^{actions} have the rule embedded within the production, could be in start, middle or end. Their position specifies the order of execution.

- Attribute Grammar: A Syntax Directed Defⁿ (SDD) in which the func. in the semantic rules can't have side effect i.e. the rules only are allowed to calculate the attribute values?

* Attributes

For a node, if the

- Synthesized Attr: Value of an ~~synthesized~~ attr. is computed from the values of only ~~its~~ the child nodes of that node & itself then the attr. is called Synthesized.
- Inherited Attr: For a node, if the value of an attribute is computed only in terms of its parent, sibling & itself, then the attr. is called Inherited.
- Although Inherited attrs. of a node N can't be defined in terms of N's children, but Synthesized attrs. of N can be defined in terms of Inherited attrs. of N.
- A leaf node (terminal) can only have synthesized attrs, the value of which is supplied by the lexical analyzer.
- An Attribute Grammar with only synthesized attrs, is called S-attributed grammar.

Ex. $L \rightarrow E_n$

$L.val = E.val$

$E \rightarrow E_1 + T$

$E.val = E_1.val + T.val$

$E \rightarrow T$

$E.val = T.val$

$T \rightarrow T_1 * F$

$T.val = T_1.val * F.val$

$T \rightarrow F$

$T.val = F.val$

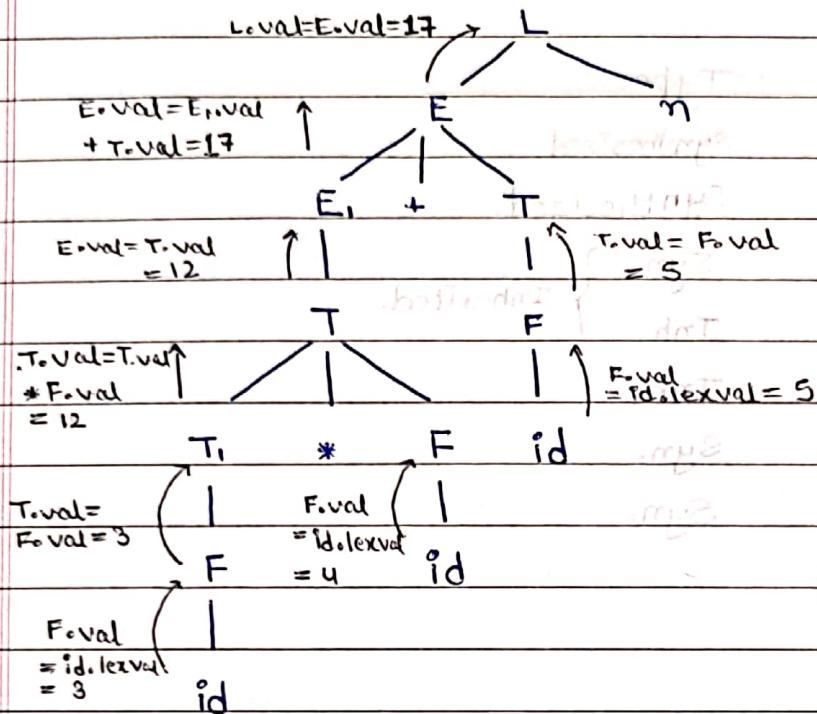
$F \rightarrow (E)$

$F.val = E.val$

$F \rightarrow \text{digit}$

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

$$w = 3 * 4 + 5n$$



Since the SDD rules have no side effect & only compute attr. values, the Grammar is Attribute Grammar.

Since all attributes are synthesized, the AG is S-Attributed Grammar

- The above approach is Bottom Up Left to Right. For evaluation of Parse Tree of S-Attributed Grammar, we can have:

(1) Bottom Up ~~Attributed~~ Left to Right

(2) Topdown

Ex. Consider the grammar & define attr. grammar for its evaluation: ~~convert binary no. to decimal no.~~

Num \rightarrow Sign List

Sign \rightarrow + | -

List \rightarrow List Bit | Bit

Bit \rightarrow 0 | 1

Num \rightarrow S L

S \rightarrow +

S \rightarrow -

L \rightarrow L, B

L \rightarrow B

B \rightarrow 0

B \rightarrow 1

L.pos = 0;

Num.val = $(-1)^{S.\text{neg}} \times L.\text{val}$

S.neg = 0

S.neg = 1

B.pos = L.pos; L.pos = L.pos + 1; L.v = L.v + B.v

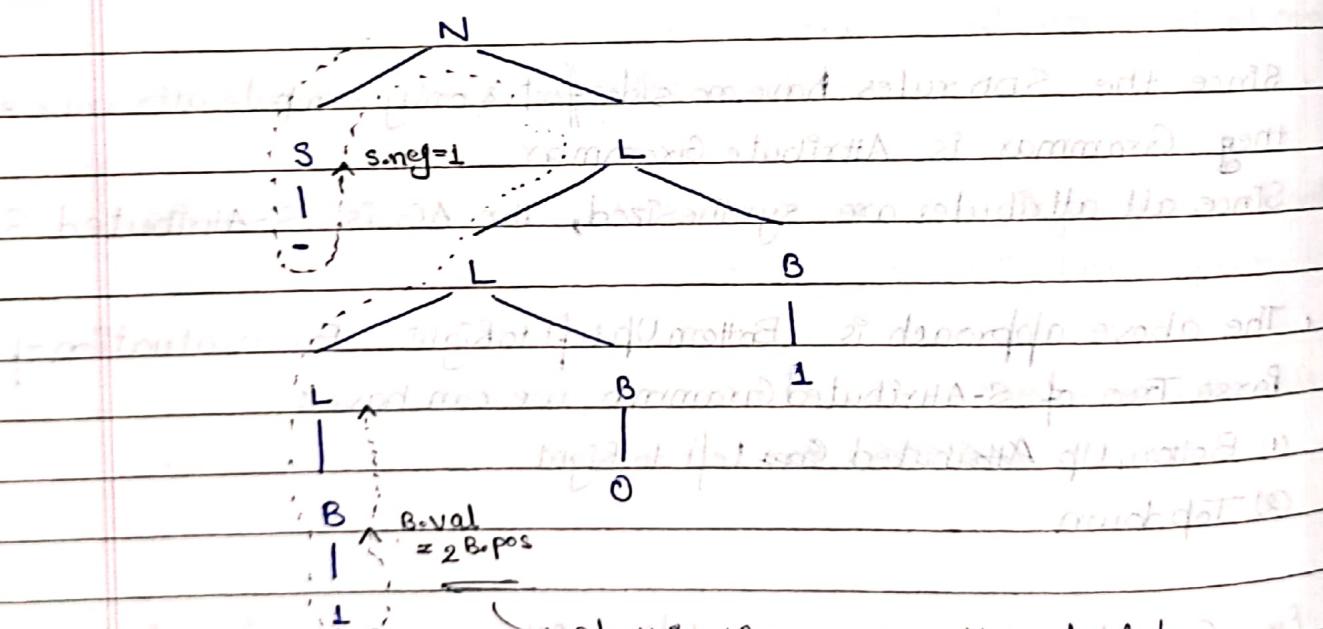
B.pos = L.pos; L.val = B.val

B.val = 0

B.val = $2^{B.\text{pos}}$

Attribute	Type
Num. val	Synthesized.
s.neg	Synthesized.
List. pos	Syn. } Inherited.
L. pos	Inh.
B. pos	Inh.
L. val	Syn.
B. val	Syn.

Say, $w = -101$



at this time, we can't calculate $B.val$ since we don't have $B.pos$.

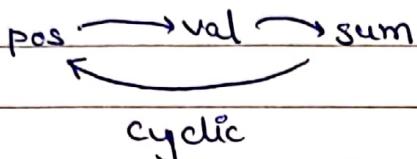
So, bottomup left to right approach doesn't work too well.

The attr. 'pos' flows downwards. Attr. 'val' flows upwards.

So, we can do one topdown pass & one bottom up pass in that order since, val depends on pos.

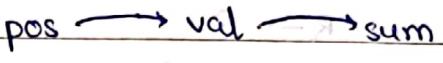
- In general we can evaluate SDDs with attribute dependency as long as the dependency isn't cyclic.

Ex.



$\{ \text{sum} = \text{val} \} \cap T \neq \emptyset$

$\{ \text{sum} = \text{val} \} \cap T = \emptyset$



$\{ \text{sum} = \text{val} \} \cap T = \emptyset$

- Detecting cyclic dependency is exponential complexity problem. So, we restrict our discussion to the following SDDs:

- S-attributed
- L-attributed

Since they always guarantee acyclic attribute dependency.

- In general the evaluation of attributes involve:

- Construct Parse Tree.
- Construct attr. dependency graph & perform topological sort.
- Evaluate attrs in the topological sort order.

- * S-Attributed SDD: A SDD is S-attributed iff it only has synthesized attributes.

- L-Attributed SDD: A SDD is L-attributed iff each inherited attr. of X_i , where $A \rightarrow X_1 X_2 \dots X_n$, $1 \leq i \leq n$, the attr. depends only on:

- The attributes of X_1, X_2, \dots, X_{i-1} {i.e left siblings of X_i }
- The inherited attributes of A {i.e. parent}.

- So, by def", every "S-attributed def" is also L-attributed.
- Note that "L-attr. SDD" only imposes restriction on inherited attributes.

Ex. [SDT scheme] give structures and run bottom up

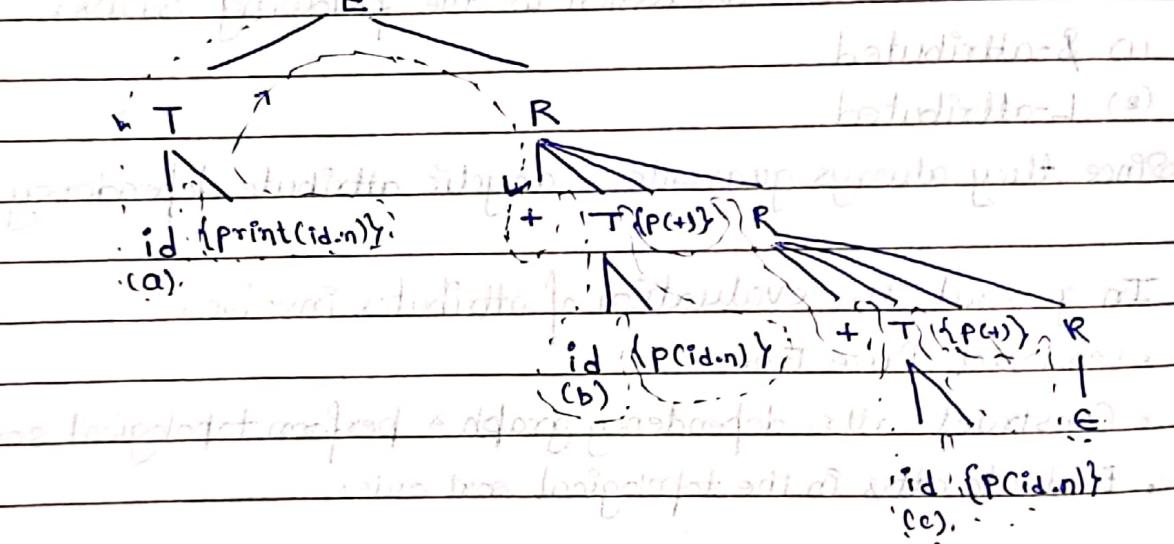
$E \rightarrow TR$

$R \rightarrow + T \{ \text{print}('+) \} R_1$

$R \rightarrow E \{ \text{print}(E) \} R_2$

$T \rightarrow id \{ \text{print}(id.name) \}$

Suppose $w = a+b+c$

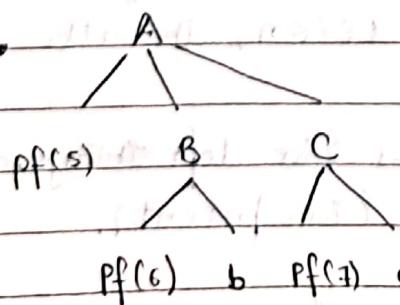


since no attr dependency we can go a single pass

Bottom Up or Top Down.

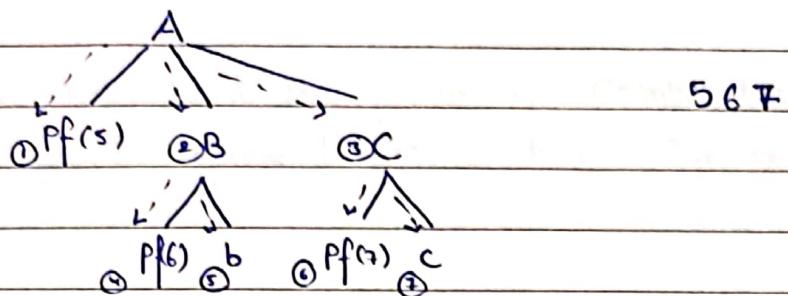
Output: ab+c

Ex.



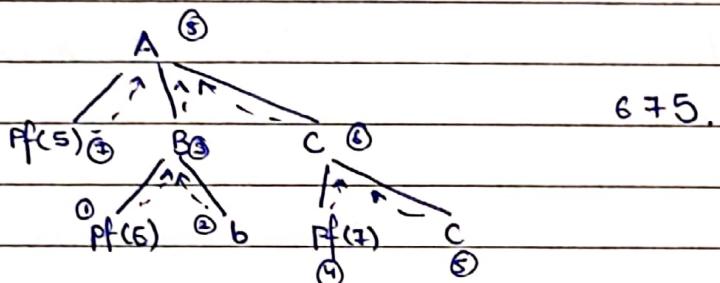
Perform both Top-down & Bottom-up parsing & record the output of semantic actions. In order to

TD:



567

BUp:



675.

- In case the semantic actions are embedded before RHS of production (at start) the output of Topdown & bottom can differ.
{or it says that SR parser performs actions, before reducing}
just

* SDD to SDT.

3-attributed.

Just embed the rules as action in the end of corresponding production. Works both for Topdown & bottom up evaluation.

L-attributed

For a production $A \rightarrow X_1 X_2 X_3 \dots X_n$, for the ~~inherited~~^{inherited} attr. 'in' of X_i , embed the rule right before X_i . For ~~synthesized~~ synthesized attr., put the rule as action in the end. Works for TopDown parsing.

Ex. $A \rightarrow BCD$ $B.i = A.i ; C.i = B.i + 2, A.s = B.s + 2$

$\Rightarrow A \rightarrow \{B.i = A.i\} B \{C.i = B.i + 2\} C D \{A.s = B.s + 2\}$

Note: L attrs can be used for BottomUp parsing as well, but it is not to be discussed.

Intermediate Code

- Throughout the phases of compilation, there are various intermediate representation, such as Syntax Analyzer's output, Parse Tree.
- After the Analysis phase of compiler is over, the synthesis phase turns the representation of the Analysis phase to some "intermediate code" so as to perform possible optimizations on it.

* Three Address Code

- In three address code atmost one operator is there on the right side of an instruction.
- Ex. $[x + y * z] \Rightarrow t_1 = y * z$ $t_2 = x + t_1$
where t_1 & t_2 are compiler generated temp names.
- In symbol table we have entries for each variable, so in any expression, the pointer to symbol table entry is used. For convinience, we use names instead of address.
- In Three Address Code representation, we treat all three as addresses : Variable name, Constants & Compiler generated variables.
- Assignment is also called Definition.

Standard TAC forms

- $x = y \text{ op } z$
- $x = \text{op } y$ (A value is done unconditionally after it)
- $x = y$
- goto L
- $\text{if } x \text{ goto L}$ (if true, go to L; if false, continue)
- $\text{if False } x \text{ goto L}$ (if false, go to L; if true, continue)
- $\text{if } x \text{ op } y \text{ goto L}$ (if condition on x, then go to L; if not, continue)
- Function call $p(x_1, x_2, x_3, \dots, x_n)$:

param x_1

param x_2

:

param x_n (n terms, then see below point)

call p, n

Indexed addressing mode:

$x = y[i]$ { meaning $x \leftarrow [y + i]$ in Assembly}

- $x = \&y$
- $x = *y$ (bottoming address of s at memory)
- $*x = y$

Ex. Convert to TAC: $i = 1; \text{do } a[i] < v \text{ while } i < 10;$

Assume size of elements of a is 8B.

$L_1: t_1 = p+1$

$i = t_1$ (initial value, sum of three 8B values within 10)

$t_2 = 9 * 8$

$t_3 = a[t_1]$

$\text{if } t_3 < v \text{ goto } L_1$

- There are multiple ways to represent TAC in system:
- Quadruple: $\langle \text{result}, \text{op1}, \text{op2}, \text{operator} \rangle$
- Triple: $\langle \text{op1}, \text{op2}, \text{operator} \rangle$
 $\text{op}_1 \leftarrow \text{op}_2, \text{operator } \text{op}_2$
- Indirect Triple

$$L + S = M$$

* Static Single Assignment

- The 'static' part means that SSA works on "text" part of the code.
- Each variable in SSA has exactly one definition. This becomes very helpful in performing optimizations such as Constant Propagation.

Ex. Convert to TAC & SSA:

```
int dist(int x, int y) {
    x = x * x;
    y = y * y;
    return sqrt(x+y);
}

dist(x0, y0):
    x1 ← x0 * x0
    y1 ← y0 * y0
    t0 ← x1 + y1
    t1 ← sqrt(t0)
    return t1
```

(TAC)

(SSA)

- SSA makes use of "def-use information".
 - SSA doesn't restrict loops. Remember that it restricts a variable to have one "def" in program text (static) not while running.
- Ex. $\text{while } (x < 5) \{$

$$y = x + 1$$

}

Here, in the text of the program, y appears to be assigned exactly once, although when the code executes, y may get assigned multiple times. So statically y has single def.

$$"y = x + 1"$$

- SSA allows the use of a special function: ϕ -func., which allows for choosing one of the options as the value to be used.

Ex. $x = 0$

$\text{while } (x < 10) \{$

$$y = 2 * x$$

$$x = y + 1$$

}

$$x_0 = 0$$

if $(x \geq 10)$ goto L₁, then it should always be x_1 . So, we

$$y_0 = 2 * (x_0)$$

$$x_1 = y_0 + 1$$

L₀:

If we try to convert this to SSA, for $y = 2 * x$, what should we use, x_0 or x_1 ?

Similarly at $x < 10$.

(phi) func.

For the first time it should be x_0 ,

use ϕ func. here:

$$x_0 = 0$$

if $\phi(x_0, x_1) \geq 10$ goto L₁

$$y_0 = 2 * \phi(x_0, x_1)$$

$$x_1 = y_0 + 1$$

L₁:

(N/A)

* IC Optimization

Code transformations to improve program, mainly in terms of the execution time, (and sometimes in program size), without affecting the desired result of the source program.

- We are mainly concerned with static optimizations, i.e. optimizations on program text, i.e. analysis & transformation is done on what the inst's are, not how they function.

Ex. Suppose the following code:

- ① $x = y + 1$ which inst's / statements can be removed
- ② $y = 2 * z$ without affecting program execution & result?
- ③ $x = y + z$
- ④ $z = 1$ Inst. ① & ④ are dead code, for sure.
- ⑤ $z = x$.

Ex. Same as above:

- ① $x = y + 1$; ① Not deadcode. If cond. at ③ might result to False
- ② $y = 2 * z$; & the value of x at ① may get used at ④, ⑤
- ③ $\text{if } (d) \ x = y + z$; ②, ③ we can't guarantee dead code by looking
- ④ $z = 1$; at just this block. ④ is definitely dead code
- ⑤ $z = x$; since ⑤ overwrites it immediately.

Ex. Same as above:

- ① $\text{while } (c) \{$
 - ② $x = y + 1;$
 - ③ $y = 2 * z;$
 - ④ $\text{if } (d) \ x = y + z;$
 - ⑤ $z = 1;$
 - ⑥ $z = x;$
- No dead code. $z = 1$ in ⑤ ^{can be} used at ③.

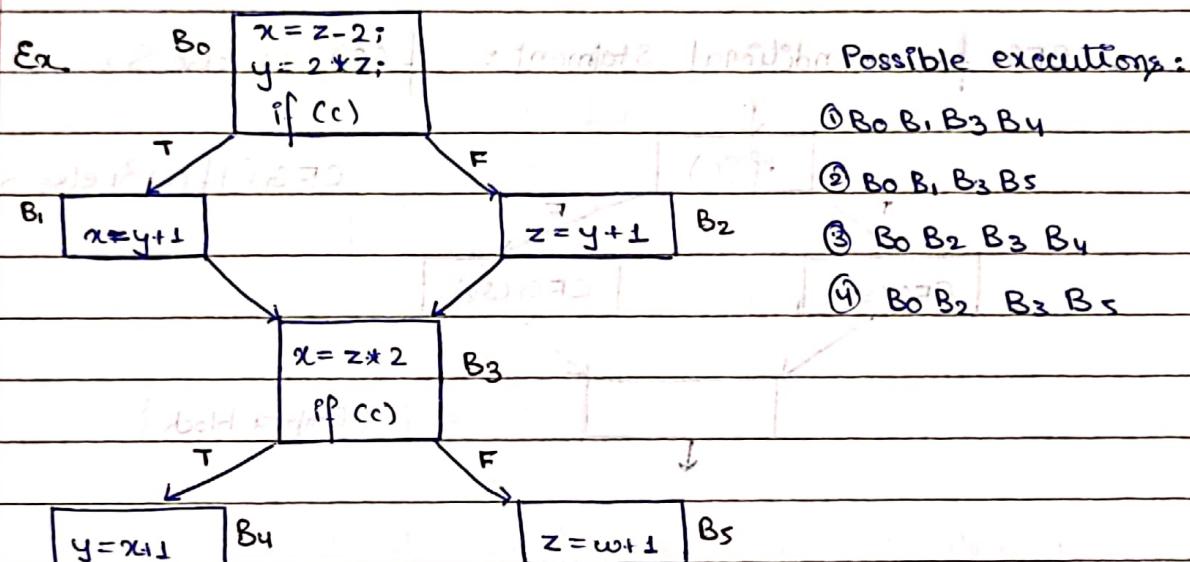
- Applying certain optimizations on low-level assembly code is tedious & can differ from vendor to vendor, so, optimizations are performed on IC.
- The static optimization should ensure that all possible executions can occur & the optimization doesn't affect any dynamic execution.

* Control Flow Graph

- It is a graphical representation of computation & control flow over the static program text.
- Nodes of the graph are "basic blocks" & edges represent the possible flow of control from end of a block to beginning of others.
- Basic Block: A sequence of Instⁿs, which guarantee that if the execution happens for them, then they will always execute in the sequence from first to last. If Instⁿ & every Instⁿ of the block will execute. So, execution of Instⁿs in a basic block always starts at the first Instⁿ & all the Instⁿ in the block execute till the last Instⁿ.
So, flow of prog. enters at ^{beginning} of block and can only exit from end of the block.
- CFG makes implementation of various optimization techniques convenient.

- CFG models all possible executions of the program. All paths in CFG trace a possible execution.

It might happen that when you run the program some of the execution path never gets taken, but we can't take any risk while performing static analysis, so CFG models all possible executions.



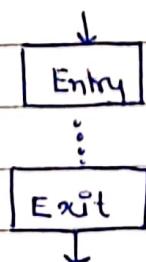
But if c is True, only B₀ B₁ B₃ B₄ is possible.

If c is False, only B₀ B₂ B₃ B₅ is possible.

So, CFG covers superset of "possible execution during runtime".

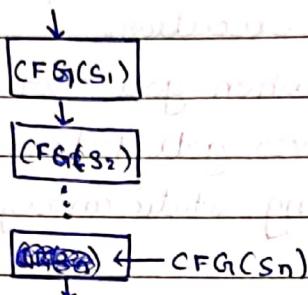
CFG for high-level IR

- CFG(S): Represents flow graph of high-level stmt S. It is a single entry, single exit graph, i.e. one entry node, one exit node.

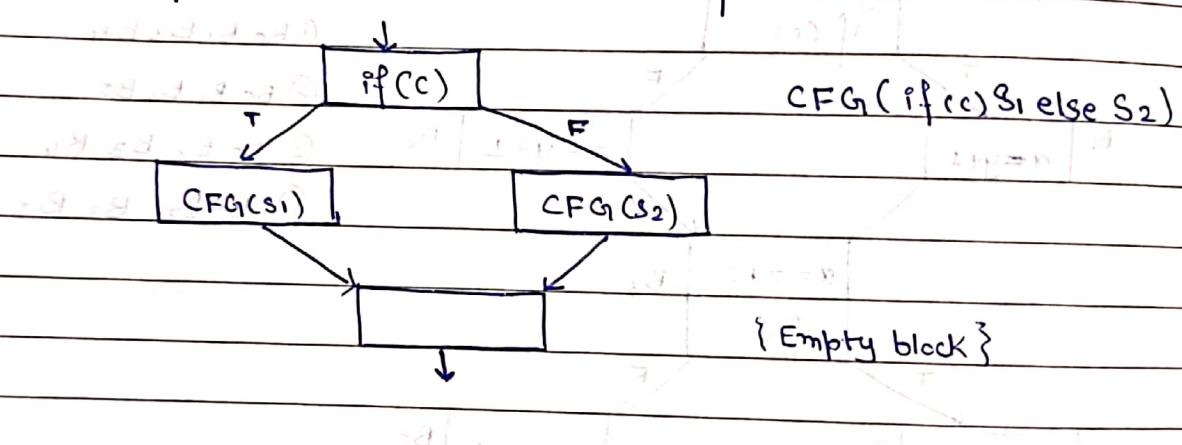


- We recursively define the func. CFG(S)

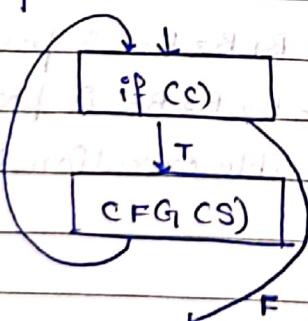
- CFG for block statements: $S_1, S_2, S_3, \dots, S_n$



- CFG for conditional statement: $\text{if } (c) S_1 \text{ else } S_2$



- CFG for while(c) S :



Ex. Make CFG for:

$\text{while } (c) \{$

$$x = y + 1;$$

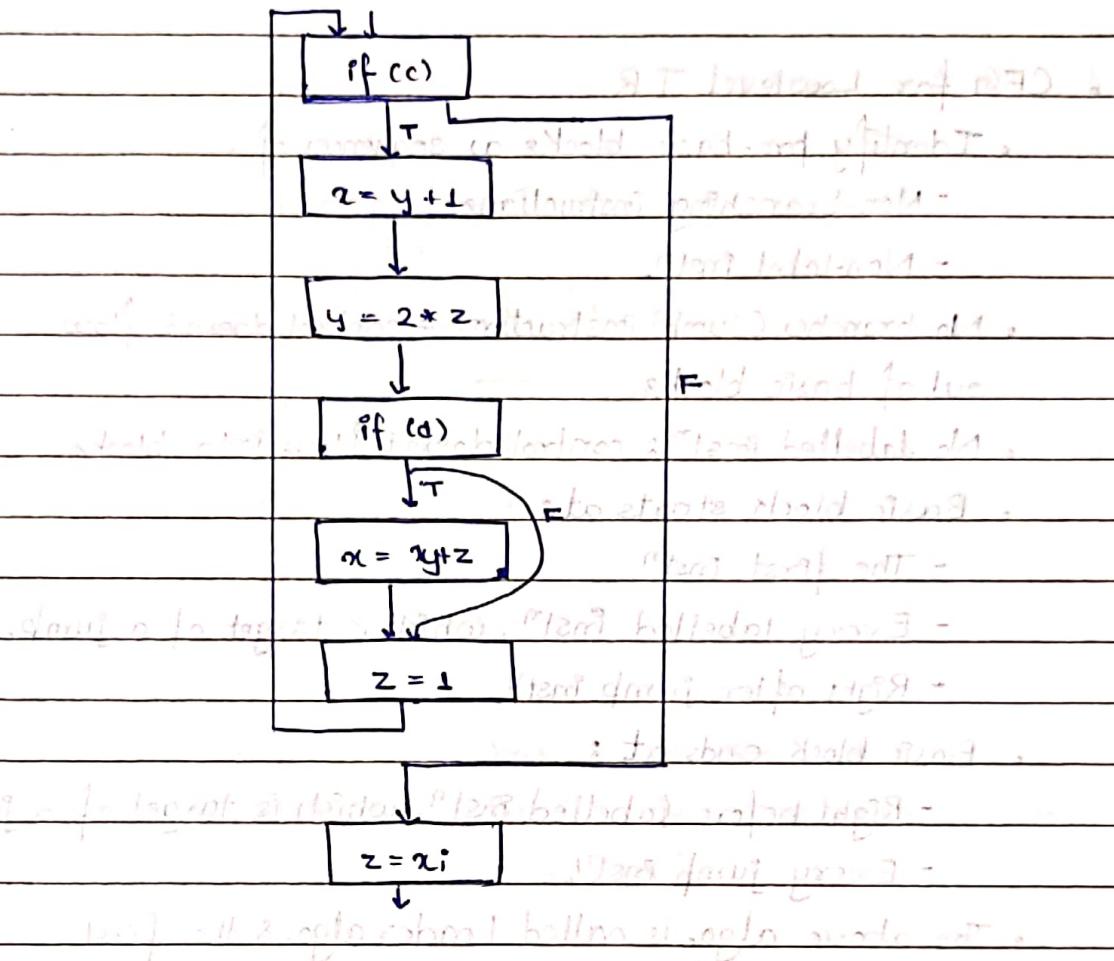
$$y = 2 * z$$

$$\text{if } (d) x = y + z;$$

$$z = 1;$$

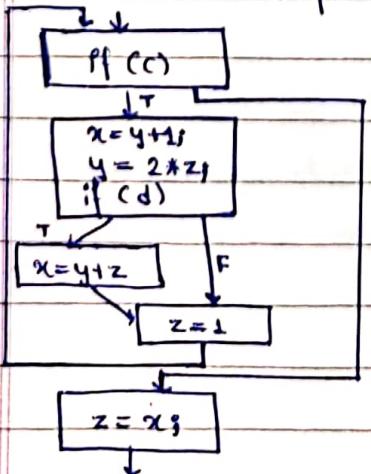
γ

$$z = x^y \quad \text{(Note: This part is handwritten and may be incorrect.)}$$



- This is simple algo, but it is inefficient. We want algo. which minimizes the number of basic blocks.
- There should be no pair of basic blocks (B_1, B_2) such that:
 - B_2 is successor of B_1 .
 - B_1 has one outgoing edge.
 - B_2 has one incoming edge.
- There should be no empty basic blocks.

Ex. So, we can optimize above CFG:



- CFG for Lowlevel IR

- Identify pre-basic blocks as sequences of :

- Non-branching instructions.
- Non-label instn.

- No brancher (jump) instructions = control doesn't flow out of basic blocks.

- No labelled instn: control doesn't flow into blocks.

- Basic block starts at :

- The first instn.

- Every labelled instn, which is target of a jump.

- Right after jump instn.

- Basic block ends at :

- Right before labelled instn, which is target of a jump.

- Every jump instn.

- The above algo. is called Leader algo. & the first instn of every basic block is called leader.

- In case of Unconditional Jump there would be single successor, and for conditional jump two successor based on the condition.

Note: Some author also put special Entry & Exit nodes.
So, check the options of questions accordingly.

Ex. (1) $P = 0$

(2) $I = 1$

(3) $P = P + I$

(4) IF $P \leq 60$ GOTO C7

(5) $P = 0$

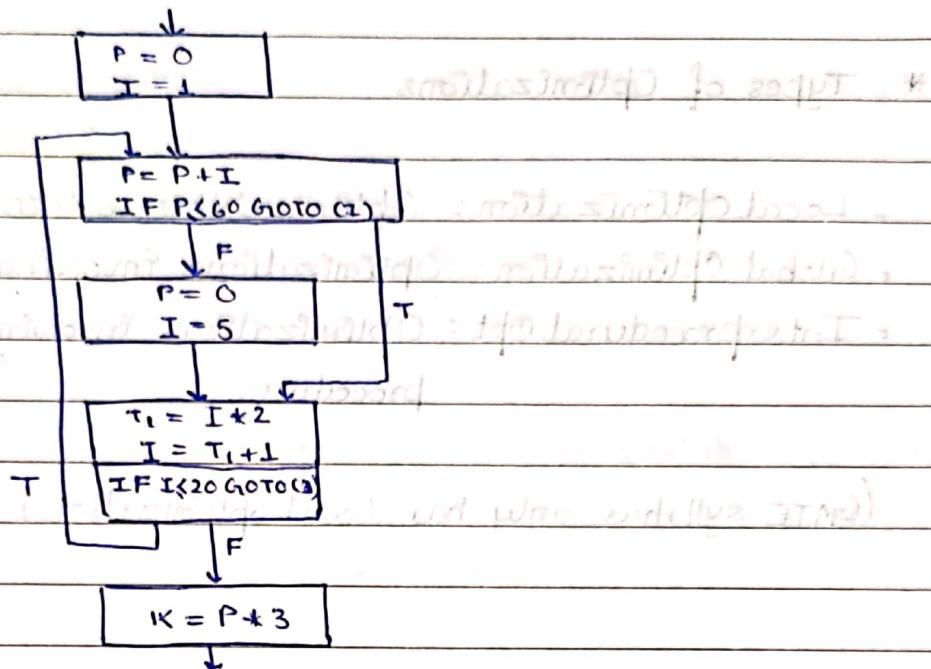
(6) $I = 5$

(7) $T_1 = I * 2$

(8) $I = T_1 + 1$

(9) IF $I \leq 20$ GOTO B

(10) $K = P * 3$



So, we have 5 leaders. One basic block per leader.

Total no. of basic blocks = 5 (OR 7 if one ENTRY & one EXIT block is also counted.)

- So, the leaders are:

- The first instⁿ
- Target of branch (jump) instⁿ
- Instⁿ that follows a branch (jump) instⁿ.

- CFG is created for every procedure in the program.

- For a block B_i , any block B_j , such that $B_j \rightarrow B_i$, is called predecessor of the block B_i .

- For a block B_i , any block B_j , such that $B_i \rightarrow B_j$, is called successor of the block B_i .

Note: Most authors include Entry & Exit block along with others. So by default count two extra blocks.

- HALT instⁿ equals exit block.

* Types of Optimizations

- Local Optimization: Optimizations within a single basic block.
- Global Optimization: Optimizations involving a single CFG.
- Interprocedural Opt: Optimizations involving CFGs of multiple procedures.

(GATE syllabus only has local optimization.)

* Local Optimizations

- Optimization performed within a single basic block.
- Common Subexpression Elimination: If an expression E is calculated & stored in a variable T_1 & then later in T_2 , & the value of T_1 hasn't changed while T_2 is assigned value of E we can replace E with T_1 in " $T_2 = E$ ", i.e.,

$$T_1 = a_1 \text{ op } a_2 \quad T_2 = a_1 \text{ op } a_2$$

$$\vdots \qquad \Rightarrow \qquad \vdots$$

$$T_2 = a_1 \text{ op } a_2$$

- Copy Propagation: If we have a variable assignment, such as " $v_1 = v_2$ ", then as long as v_1 & v_2 are not reassigned, any statement of the form " $a = v_1$ " can be rewritten as " $a = v_2$ ", provided that such a rewrite is legal.

$$v_1 = v_2$$

 \vdots

$$a = v_1$$

 \Rightarrow

$$a = v_2$$

This helps greatly w/ dead code elimination.

- A variable can be:

(1) Used = read from memory or write to memory.

(2) Defined = Assigned = write to.

Ex. $a = b + c;$ both b, c are used. a is defined.

a is temporary variable.

- Dead Code Elimination: An assignment to variable v is called dead if the value of that assignment is never read anywhere. Such dead assignments can be removed.

→ Dead code elimination for variable a .

Note: There is no fixed order in which optimizations can be applied.

You may have to apply each optimization multiple times.

→ Multiple optimization steps may be required.

$$\text{Ex. } b = a * a$$

$$b = a * a$$

$$b = a * a.$$

$$b = a * a$$

$$c = a * a \xrightarrow{\text{CSElim.}}$$

$$c = b * b \xrightarrow{\text{Copy}} c = b$$

$$c = b * b \xrightarrow{\text{CSElim.}} c = b$$

$$d = b * c$$

$$d = b * c$$

$$d = b * b$$

$$e = b * b.$$

$$e = b * b.$$

$$e = d.$$

- Arithmetic Simplification: Replace complex ops w/ easier ones.

$$\text{Ex. } x = 4 * a \Rightarrow x = a \ll 2$$

- Constant Folding: Evaluate expressions at compile time if they have constant value.

$$\text{Ex. } x = 4 * 5 \Rightarrow x = 20$$

- * Optimization & Analyses: Most optimizations are only possible by analyzing the program's behaviour. In order to implement optimizations, we must also talk about the corresponding program analysis.

- Dead Code Elimination requires Liveness Analysis, which helps us know which variables are alive at which point in program.
- Available Expression Analysis helps w/ Common Subexpression Elimination, & Copy Propagation.

* Available Expression Analysis

- An expression is called available if some variable in the program holds the value of that expression.

- In Common Subexpr. Elimination, we replace available expression by the variable holding its value.

- In Copy Propagation, we replace the use of a variable by the available expression it holds.

Finding available expressions:

- Initially no expressions are available.
- Whenever we execute a statement, $a = b + c$:
 - Any expression holding 'a' becomes invalidated. (a on RHS)
 - The expression $a = b + c$ becomes available.
- If $a = a + c$, then $a = a + c$ doesn't become available.
- . Iterate across the basic block, begin beginning w/ the empty set of expressions & updating available expressions at each

Ex. ① $a = b;$

$$\textcircled{2} \quad c = b;$$

$$\textcircled{3} \quad d = a + b;$$

$$\textcircled{4} \quad e = a + b;$$

$$\textcircled{5} \quad f = a + b;$$

$$\{a = b, c = b, d = a + b, e = a + b\}$$

$$\{a = b, c = b, d = b, e = a + b\}$$

$$\{a = b, c = b, d = b, e = a + b, f = a + b\}$$

Now, we can apply Common Subexpression Elimination using the set of Available Expressions.

$$a = b;$$

$$c = a; \quad \text{with } a \in \emptyset$$

$$d = a + b;$$

$$e = d;$$

$$f = e;$$

$$g = f + b;$$

available at write point

* Liveness Analysis

- It gives us the set of live & dead variables at a point in program.

- A variable is called Live at a point in program if later in the program, its value is read before being overwritten.

- DCE elimination works by computing liveness for each variable & then eliminating assignment to dead variables.

$$\text{Ex. } ① a = b \quad \{ b \}$$

$$③ d = a + b \quad \{ a, b \}$$

$\{ \text{we go in reverse order} \}$

$$\text{② } c = b \quad \{ a, b \}$$

* Computing Live Variables

- To know if a variable will be used at some point, we iterate across the statements, in a basic block in reverse order.

- Initially some set of values are known to be alive. (Given)

- When we see a stmt, $a = b + c$; just before the stmt:

- The value of a is dead, since its value is about to be overwritten.
- The values of b, c are alive, since they are about to be read.

- If the stmt. is of the form: $a = a + b;$, then both a & b are alive right before the stmt, since both will be read before getting rewritten.

~~dead~~
Ex. $a = b;$ $\{b\} - S_6$

a is live after,

$\{a, b\} - S_5$

$c = a;$

c is dead after assignment.
eliminate $c.$

$\{a, b\} - S_4$

$d = a + b;$

d is live after,

$\{b, d, a\} - S_3$

$e = d;$

Don't elim. e is alive after

$\{b, e, a\} - S_2$

$d = a;$

d is live after assignment. Don't eliminate

$\{b, d, e\} - S_1$ Since b is not live after assignment, it can be removed.

$f = e;$

Since f is not live after assignment,
it can be removed.

$\{b, d\} - \text{Initial}$

$\{b\}$ point to b here is value of a , so

i. $a = b$ we can again apply liveness analysis.

$d = a + b$ $\{b, a\}$

so now d is dead after assignment.

$e = d$ $\{b, g\}$

$d = a$ $\{b, d\}$

So, d is dead in stmt 2.

$d = a + b$ $\{b, a\}$

$d = a$ $\{b, d\}$

i. $a = b$ $\{b\}$
 $d = a + b$ $\{b, a\}$
 $d = a$ $\{b, d\}$

So all vars are live right after
assignment.

- We can combine (interleave) LA & DCE together. (It will require just single pass)

* Other Intermediate Code

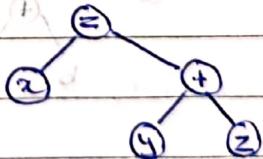
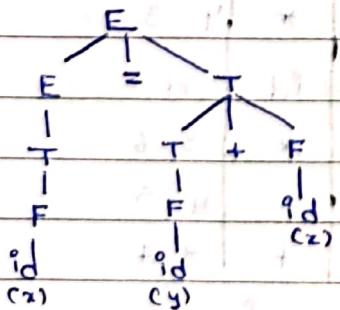
- There are various other Intermediate Representation along with Three Address Code & Static Single Assignment, such as, Parse Tree, Abstract Syntax Tree, Directed Acyclic Graph, Control Flow Graph, Two Address Code, etc.

- Interestingly, not only the above data structures, but other programming languages may also be used as IRs, such as C language.

Abstract Syntax Tree (or Syntax Tree):

- It is an abstract representation of Parse Tree.
- The operation base root & intermediate nodes & operands are on leaf.

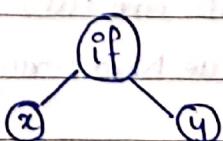
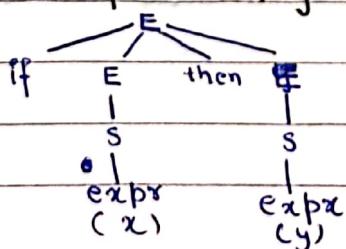
Ex. $x = y + z$



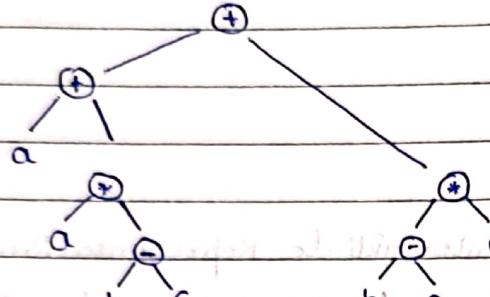
AST

Parse Tree

Ex. If x then y



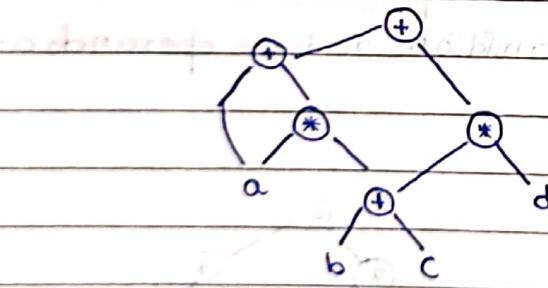
Ex. $a + a * (b - c) + (b - c) * d$
 [By default AST is constructed with Precedence & Assoc. in mind]



- The 'Abstract' in AST means that the syntax is abstracted, it only bothers about semantics & it does so by the structure of tree.

- Directed Acyclic Graph: A variant of AST which doesn't recompute expressions more than once.

Ex. $a + a * (b + c) + (b + c) * d$



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

- The above implementation is called Value Numbering.
- DAG can also be used to express basic blocks, allowing Value Numbering.

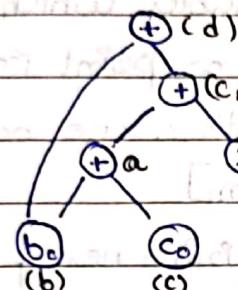
Ex. $a = b + c.$

$a_1 = b_0 + c_0$ (Replace b with b_0)

$$c = a + x \Rightarrow c_1 = a_1 + x_0$$

$$d = b + c. \text{ So } d_1 = b_0 + c_1$$

$$b = a + x. \text{ So } b_1 = a_1 + x_0$$



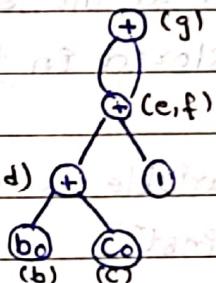
Ex. $a = b + c$ (Now $a_1 = b_0 + c_0$ if b comes with b_0)

$$e = a + 1 \Rightarrow e_1 = a_1 + 1$$

$$d = b + c. \Rightarrow d_1 = b_0 + c_0$$

$$f = d + 1 \Rightarrow f_1 = d_1 + 1$$

$$g = e + f \Rightarrow g_1 = e_1 + f_1$$



* DAG can easily be used for Common Subexpression Elimination.

* There are two types of analysis:

. Control Flow Analysis: Discovering control structure (such as calls, basic blocks, loops)

. Data Flow Analysis: Discovering data flow structure (variable uses, expression evaluation).

Ex. Dead Code elimination eliminates stmts, the result (data) of which isn't used. So, part of Dataflow Analysis.

Unreachable code eliminates the stmts/procedures which are never executed under any circumstance. So, its part of Control Flow analysis.

* Liveness Analysis (contd.)

- Liveness analysis can also be used for register allocation. At every point in prog. we know which variables are live & which are not. So, at a point all the variables which are live must be given different registers. If there's no point in program where variable x & y are live together, we can allocate them same registers. [Graph Coloring Problem]

During IR generation we assume infinite no. of registers, but the system may not have enough registers to hold all the variables at the same time, so some variables may have to be stored in the m/m. This is called Register Spilling.

- A variable is live "at an inst" if it is live "right before the inst".

Ex. $S_1: a = a + c \xrightarrow{\{d, a, b, c\}}$
 $S_2: f = a + b \xrightarrow{\{d, a, b\}}$
 $S_3: a = d + f \xrightarrow{\{d, f\}}$
 $S_4: \text{ret } a, d \xrightarrow{\{a, d\}}$

Find live vars at $S_1, S_2,$

$S_3 \& S_4$. Assume no var is live after S_4 .

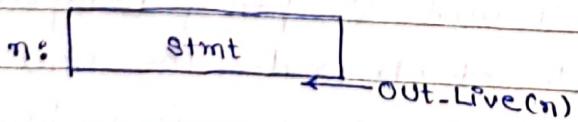
\therefore At $S_1: \{d, a, b, c\}$

$S_2: \{d, a, b\}$

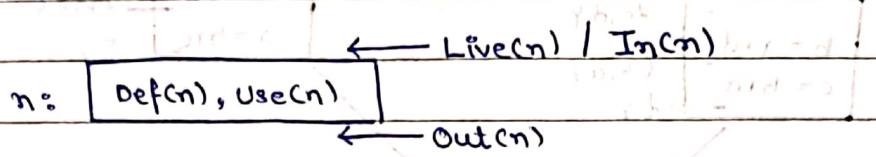
$S_3: \{d, f\}$

$S_4: \{a, d\}$

- For a stmt, the variables live right after it are called OUT or OutLive.



- assigned
- For a stmt, give its Out-Live, its defined variable & its used variables as:



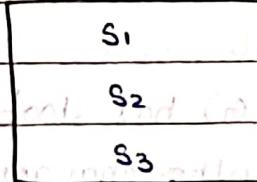
then, $\text{Live}(n) = [\text{Out}(n) - \text{Def}(n)] \cup \text{Use}(n)$.
 $(\text{Live}(n) \equiv \text{In}(n))$

- Live Variable Analysis is backward analysis (for efficiency), but we can also do it in forward way. For each definition, go ~~back~~ to all the instns after it. If the defined var is used before reassignment it is alive at that point in program.

- $\text{Use} \equiv \text{Read} \equiv \text{Reference}$

- $\text{Define} \equiv \text{Assign} \equiv \text{Write} \equiv \text{Kill}$

- In a basic block:



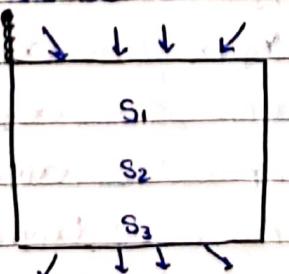
the no. of successor for: $S_1 \Rightarrow 1$, $S_2 \Rightarrow 1$, $S_3 = ?$

For S_3 we can't say there can be 0 or many successors.

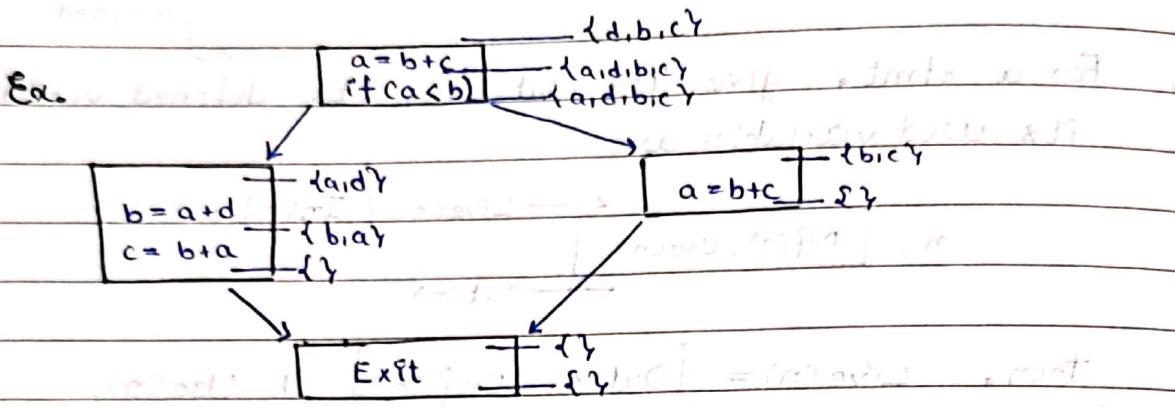
~~Ans.~~ $\text{In}(S_3) = \text{Out}(S_2)$

~~Ans.~~ $\text{In}(S_2) = \text{Out}(S_1)$

~~Ans.~~ But what would be $\text{Out}(S_3)$?

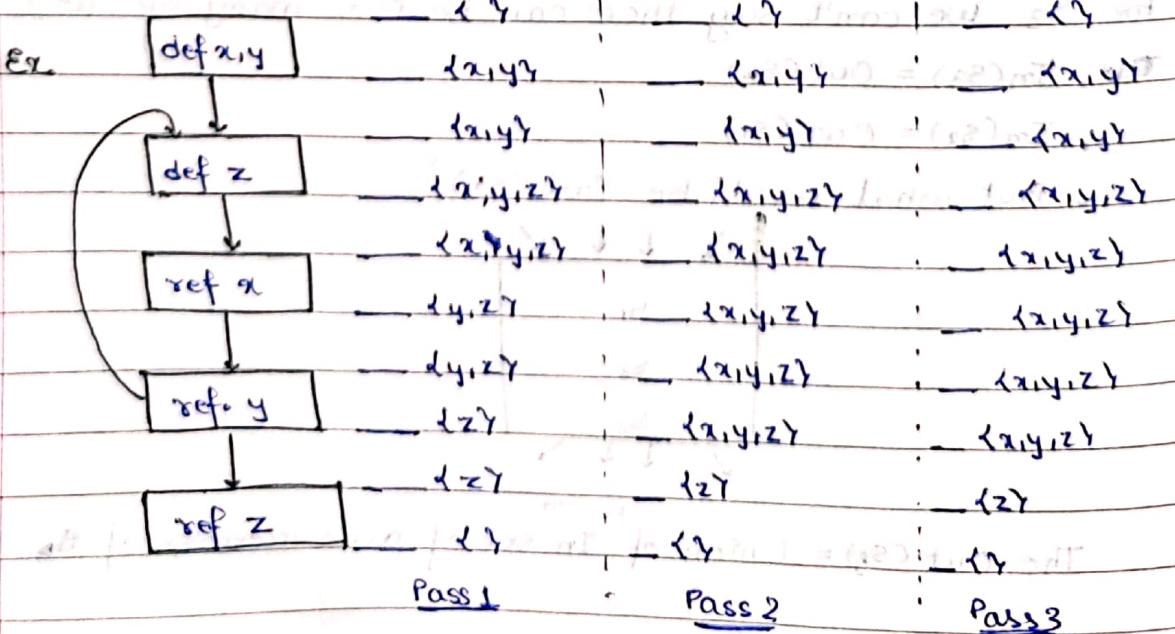


The $\text{Out}(S_3) = \text{Union of } " \text{In set of all successors of } S_3 "$.



- So we can make a general stmt, that for any Stmt S_n , its "Out" set is equal to Union of "In" set of all its successors.
- For a stmt S_n , its "out" set $Out(S_n)$, if :
 - $Def(S_n) \subseteq Out(S_n)$, then the defined variable is alive. Don't remove S_n .
 - Else defined variable is dead. removed stand in ΔT .

- If the program (also CFG) has loops we'd require multiple passes to ensure no further changes to In & Out set of stmts are possible.



Since Pass 2 & 3 are same, no further changes.

Ex.

$$P = q + r$$

$$S = p + q$$

$$U = S * V$$

$$V = r + u$$

$$q = s * u$$

$$q = v + r$$

$$\{v, r, q\}$$

$$\{v, r, p, q\}$$

$$\{v, r, s, u\}$$

$$\{r, u\} \rightarrow \{v, r, s, u\}$$

$$\{v, r\} \rightarrow \{v, r\}$$

$$\{v, r\}$$

Pass 1

$$\{v, r, q\}$$

$$\{v, r, p, q\}$$

$$\{v, r, s\}$$

$$\{v, r, s, u\}$$

$$\{r, u\} \rightarrow \{v, r, s, u\}$$

$$\{v, r\} \rightarrow \{v, r\}$$

$$\{v, r\}$$

$$\{v, r, q\}$$

Pass 2 & 3 would be same, so no further changes.

Ex. For the basic block, find defined & used.

$$a = b + c;$$

$$a = a + b;$$

$$b = b + c;$$

$$x = b + d;$$

a is not used, since it is defined before being used.

b, c, d are used.

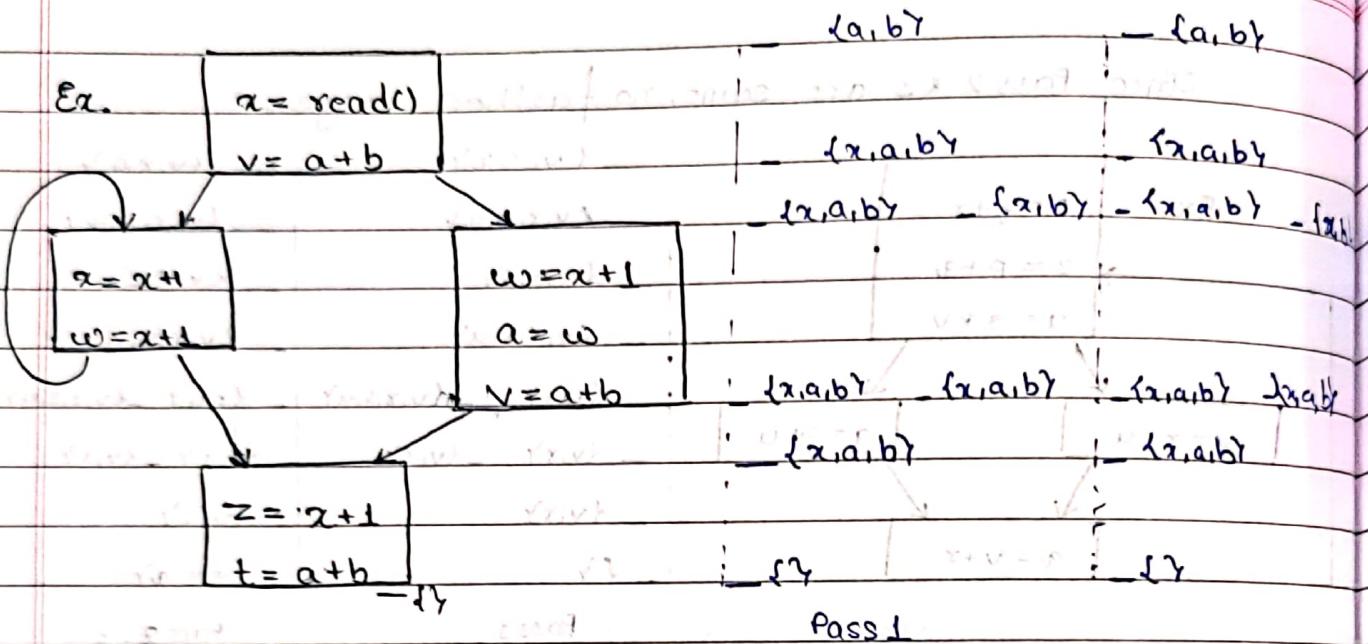
x is defined. b is defined

$$\therefore \text{Def} = \{a, x\}, b\}$$

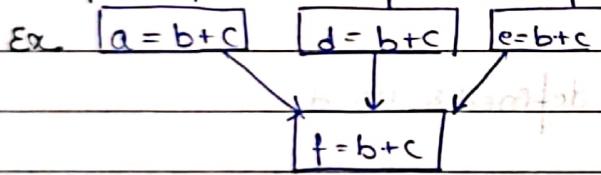
$$\text{Use} = \{b, c, d\}$$

- For a basic block B, the $\text{Use}(B)$ is the set of all the variables 'b' such that b is used in B, before it is defined.

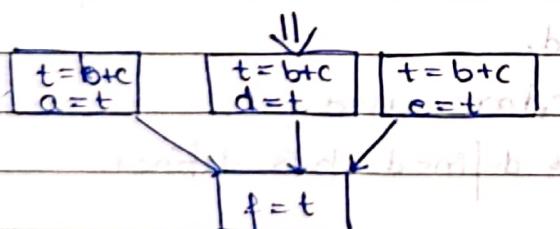
$\text{Def}(B)$ contains all the variables defined within the block.



* Compiler puts computed expressions in temporary variables.



So with which variable
should $f = b + c$ be replaced?



So whichever t reaches
to f on execution path.

- For available expression analysis; for a stmt S_i :

$\text{In}(n)$

S_i

$\text{Out}(n)$

$\text{In}(n)$ is the set of available expr. ~~at~~ at S_i , $\text{Out}(n)$ is the set of available expr. after S_i .

$$\text{Out}(n) = [\text{In}(n) \cup \text{Gen}(n)] - \text{Kill}(n)$$

where $\text{Gen}(n)$ for $x = a+b$ if ~~take b~~^{x=a+b} $\text{Kill}(n)$ contains any expr. in $[\text{In}(n) \cup \text{Gen}(n)]$ with 'x'.

& $\text{In}(n) = \underset{\text{Intersect}}{\text{Intersection}} \text{ of } \text{out}(p) \text{ of all predecessors } p \text{ of } n.$

& For a basic block B, an expression is available if it is in $\text{Out}(p)$ of all predecessors p of B.

Since Common Subexpr. elimination can lead to increased no. of variables, it leads to program size getting increased.
temp

Runtime Environment

Internally, we consider that's a static initialisation.

- * Static: Some issue regarding the program, that can be decided by the compiler (during compilation) is called as "Static".

→ Our point, Let others are question "What's a "Static"?

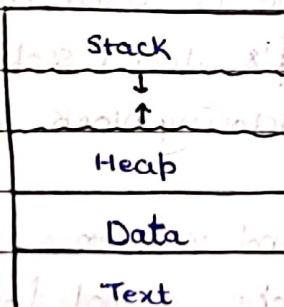
- Dynamic: Some issue regarding the program that can only be decided by running the program (during runtime) is called as "dynamic".

Dynamic

- Compiler only deals w/ Virtual m/m, i.e. it assumes program starts at loc. 0 & grows up. It can compile the program as such.

↑ contiguously after that

- Virtual m/m of a program:



- Static Keyword in C:

- The **static variable** is allocated m/m once in the **'Data'** area of the program & retains its value b/w function calls. The m/m in the **Data** section is allocated at **Compile Time**. The lifetime of such a variable is throughout the program.

- **Global Variable vs Static Variable:** Global variable is visible to every entity in the program, static variable is only visible to the function it is defined in. If static var. is global, then its scope is limited to the file (compilation unit).

Temporary Environment

- So the "Data" section & "Text" sections are allocated during compile time.
- "Stack" & "Heap" sections are allocated during runtime.

* Scoping

Static Scoping: For a symbol reference, look for the definition of the symbol at the physical location of the reference, i.e. look within scope, if not found look in the enclosing block, if not found, repeat.

(Lexical Scoping)

Ex. Suppose that a variable x is used in func. f , then x will be first searched in f 's local scope, if not found, will be searched in f 's enclosing block, i.e. file scope & so on.

Dynamic Scope: For a symbol reference to x in func. f , look for defn. in f 's local scope, if not found, look in f 's caller's scope. if not found, repeat.

Ex. `int x = 5;`

`void g()`

`{`

`f(x);`

`}`

`void f()`

`{`

`g();`

`}`

`void main()`

`{`

`f();`

`}`

global main f g static scoping

global main f g dynamic scoping

- In the activation record of the func., the access link points to the activation record to which the func. must refer for symbol reference resolution.

* Memory Allocation

- Static Allocation:** If the decision of allocation can be made by the compiler looking only at the text of the program.
- Dynamic Allocation:** If the decision of allocation can be made only when the program executes.

* Activation Record

When a func. 'f' calls another func. 'g', an activation record is created on the stack for 'g'. Some of the possible contents of Act. Record (from bottom to TOS):

- Actual Parameters passed by f to g.
- Returned Value: Space for g to place the return value.
- Control Link: Pointing to AR of func. to which g would return control, 'f'.
- Access Link: Points to location for access of non-local data.
- Saved M/C Status: Includes return addr (old PC value) & contents of registers
- Local Data: Belonging locally to g.
- Temporary Vars: Used to store temp vals., for expressions, etc.