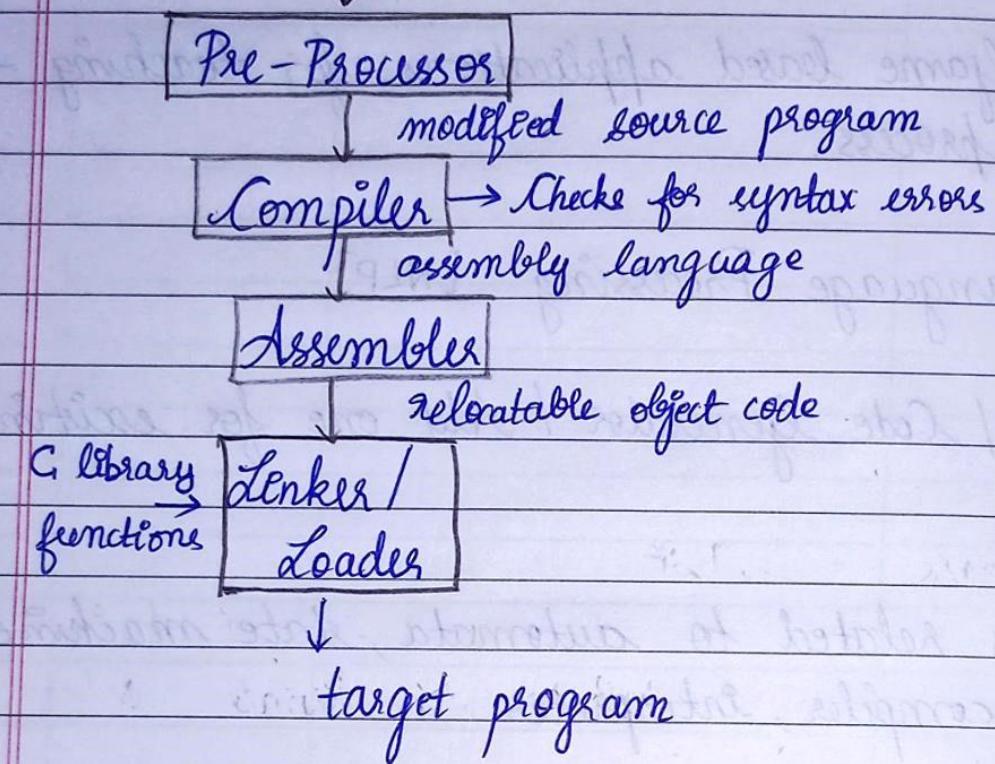
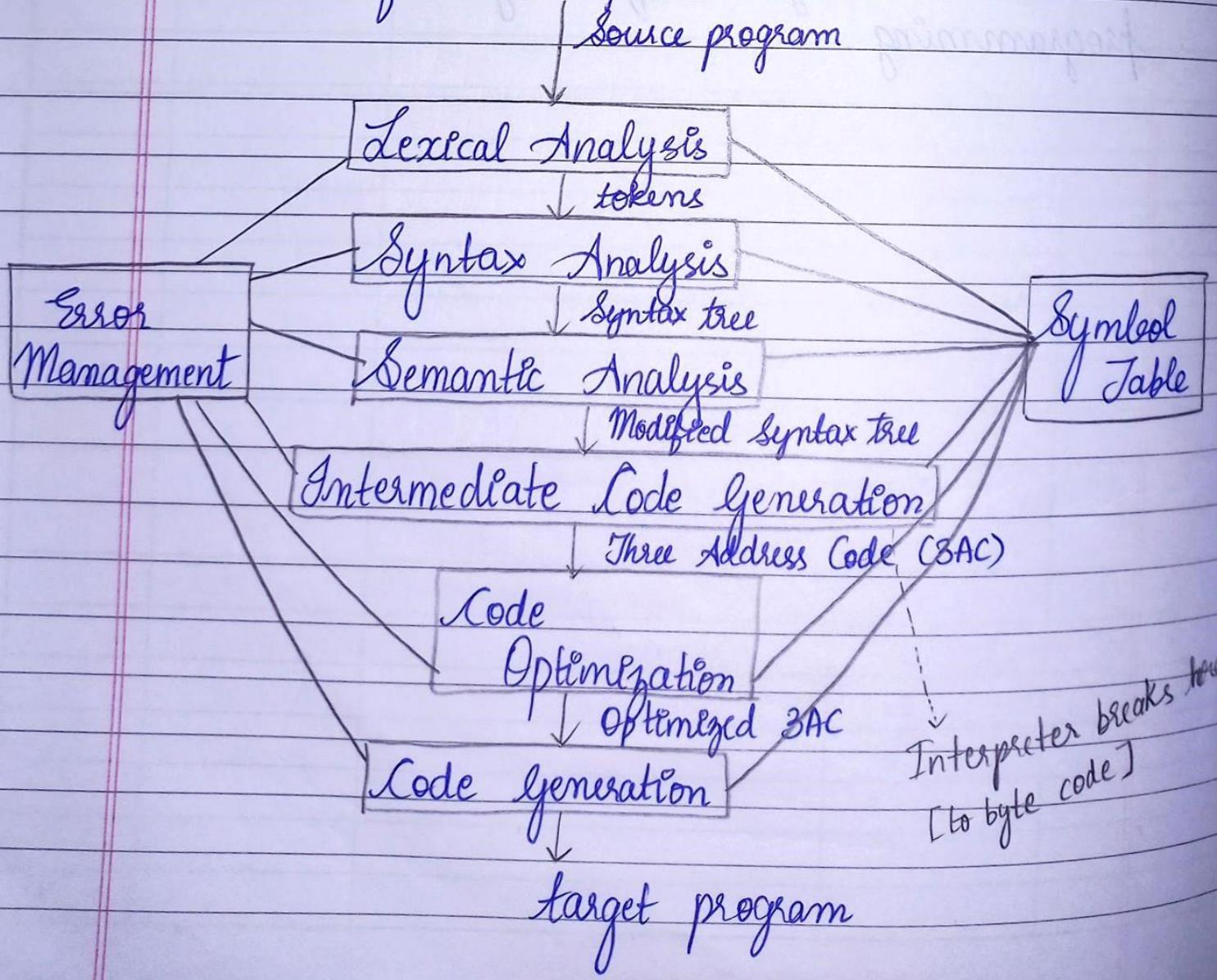


Unit - I

* Learning Processing
Input source program



* Structure of a Compiler



* Grammars

↳ checking 'C' expressions for arithmetic operations

* Types :-

1) Context Free Grammar

2) Context Sensitive Grammar

(Statement) $S = E \mid id$

(Expression) $E \rightarrow E + T \mid E - T \mid T$

(Term) $T \rightarrow T * F \mid T / F \mid T \% F \mid F$

(Factor) $F \rightarrow (E) \mid id \mid num$

Non-terminals :- E, T, F

Terminals / Tokens :- (,), id, num, +, -, *, /, %

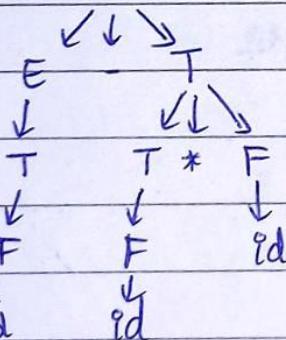
* Tokens - has its own definition in symbol table
(user can't define)

Eg:

$a - b * c$

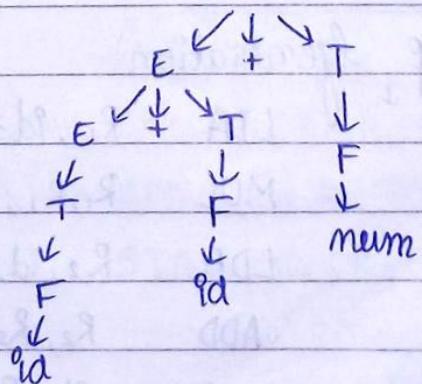


E



$a = b + c + 10 \rightarrow$

$S = E$



* $ab = b + c * 10$



Lexical Analysis



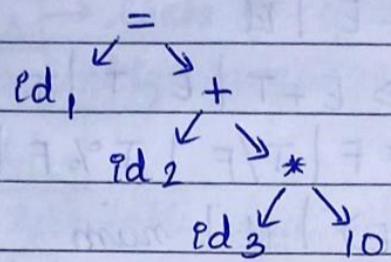
$\langle id, 1 \rangle \quad \langle = \rangle \quad \langle id, 2 \rangle \quad \langle + \rangle \quad \langle id, 3 \rangle \quad \langle * \rangle \quad \langle 10 \rangle$



Syntax Analysis



Syntax tree:



Semantic Analysis



Same syntax tree



Intermediate Code Generation



Three

$$t_1 = id_3 * 10$$

Address

$$t_2 = id_2 + t_1$$

Code (3AC)

$$id_1 = t_2$$



Code Optimization



$$t_1 = id_3 * 10$$

$$id_1 = id_2 + t_1$$



Code Generation

Assembly language code

LDA R₁, id₃

MUL R₁, R₁, #10

LDA R₂, id₂

ADD R₂, R₂, R₁

STA id₁, R₂

* Translate HLL $a = b + c + 10$ to assembly language [$\frac{a,b,c}{\text{float}}$]

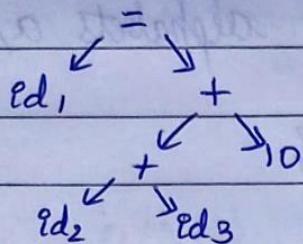
$$a = b + c + 10$$

Lexical analysis

$\langle \text{id}, 1 \rangle \quad \langle = \rangle \quad \langle \text{id}, 2 \rangle \quad \langle + \rangle \quad \langle \text{id}, 3 \rangle \quad \langle + \rangle \quad \langle 10 \rangle$

Syntax analysis

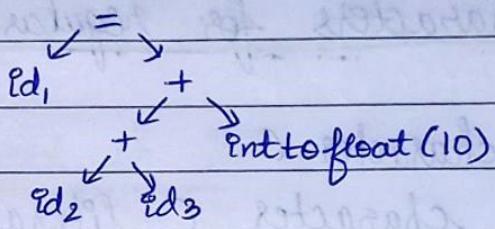
Syntax tree :



Semantic Analysis

Modified

Syntax tree :



Intermediate Code Generation

3AC :

$$t_1 = id_2 + id_3$$

$$t_2 = \text{inttofloat}(10)$$

$$t_3 = t_1 + t_2$$

$$id_1 = t_3$$

Code Optimization

$$t_1 = id_2 + id_3$$

$$id_1 = t_1 + 10.0$$

Code generation

Assembly
language code

LDA R₁, id₂

LDA R₂, id₃

ADD R₂, R₂, R₁

ADD R₃, R₂, #10.0

STA id₁, R₃

+ Precedence: Closure > Concat. > Union
* Operator: |+, Char.; +

* Regular Expression

Union, Concatenation, Closure \rightarrow Positive - atleast 1 occurrence of d
(L|D) (LD) \rightarrow Kleene - occurrence can be null
(+)

Language - A set of finite symbols consisting of alphabets

Eg: L is a set of alphabets a,b,c

$$\Sigma = \{a, b, c\}$$

B is binary language B = {0, 1}

* Lex characters for regular expressions

c - character c

|c - character c literally

"s" - string s

. - any character other than newline

^ - beginning of line

\$ - end of line

[s] - any one of the characters in s

[^s] - any character not s

r* - 0 or more occurrences

r+ - 1 or more occurrences

r? - 0/1 occurrence

r{m,n} - b/n m and n occurrences

r₁r₂ - r₁ followed by r₂

r₁|r₂ - r₁ or r₂

(1) - same as r

r₁/r₂ - r₁ when followed by r₂

[a-zA-Z] - character class

* Let $L = \{a, b\}$

$$\text{i) } \underbrace{(a \mid b)}_c \underbrace{(a \mid b)}_c$$

$$= \{aa, ab, ba, bb\}$$

$$\text{ii) } a^*$$

$$= \{\epsilon, a, aa, aaa, \dots\}$$

$$\text{iii) } (a \mid b)^*$$

$$= \{\epsilon, a, aa, aaa, \dots, b, bb, bbb, \dots\}$$

$$\text{iv) } a \mid a^*b$$

$$= \{a, b, ab, aab, aaab, \dots\}$$

* Represent the statement in RE

i) All strings of lowercase letters that contain five vowels in order.

$$\text{conso} \rightarrow [b-d-f-h-j-n-p-t-v-z]$$

$$\text{Ans: st} \rightarrow \{\text{conso}\}^* a \{\text{conso}\}^* e \{\text{conso}\}^* i \{\text{conso}\}^* o \{\text{conso}\}^* u \{\text{conso}\}^*$$

2) Any string of $\Sigma = \{a, b\}$ where it starts with 'a' and ends with 'a'

$$= a(a^* \mid b^*)^* a \quad \text{or} \quad a(a^* b^* \mid b^* a^*)^* a$$

~~$a a^* b^*$~~

\$ a (a^* b) \$ a

(S|s)(

12.10.2023

classmate
Date _____
Page _____

* Write regular expression for SQL keywords.

SELECT (Case insensitive) Eg: SeLEct, Select etc

+ 2) Comments consisting of string surrounded by /* and */.

* 3) Describe the language description for the following:
R.E.s

a) $a^* b a^* b a^* b a^*$ = $\Sigma_0, a, aa, aaa \dots b \Sigma_0, a, aa, \dots b \Sigma_0, a, aa, b \Sigma_0, a, aa$

b) $(a|b)^* a (a|b)(a|b)$

c) $^* [^aeiou]^* \$$

Answers

1) SELECT

(S|s)(E|e)(L|l)(E|e)(C|c)(T|t)
[Ss][Ee][Ll][Ee][Cc][Tt]

2)

"/* "[^(*/*)]*/*") (. | "\m")]* /*/"

3)

a) $a^* b a^* b a^* b a^*$

A language $L = \{a, b\}$ with string of any length containing 3 b's in order

b) $(a|b)^* a (a|b)(a|b)$

A string with characters a & b such that last 3rd character is always a

Write RE to identify arithmetic, log, relational
op. in C.

classmate

Date _____

Page _____

c) ${}^{\wedge}[^{\wedge}aeiou]^* \$$

A line consisting of all characters except a,e,i,o,u.

* R.E to identify arithmetic, log. & relational op. in C.

$[- + * / \% = < > ! \& | ^] | = = | ! = | < = | > =$

Tutorial - 1

I Write character class for the following set of characters.

- 1) The first 10 letters (upto j) in either uppercase or lowercase.
- 2) Lowercase consonants
- 3) Characters that appears at the end of English sentence (Eg: Exclamation point)

II Write a regular expression that matches string "\\"

III Translate the foll. assignment statement:-

$\text{Eq}^n = a*a + 2*a*b + b*b$
(Assume all variables are int)

Answers

I

- 1) $[a-jA-J]$
- 2) $[b-df-hj-np-tv-z]$
- 3) $[\.,\?,\!,\;]$

II

" " " "

" " "

III

$$\text{eqn} = \underline{\underline{a * a}} + 2a * b + b * b$$

Lexical Analysis

$\langle \text{id}, 1 \rangle \Rightarrow \langle \text{id}, 2 \rangle \langle * \rangle \langle \text{id}, 2 \rangle \langle + \rangle$

$\langle 2 \rangle \langle * \rangle \langle \text{id}, 2 \rangle \langle * \rangle$

$\langle \text{id}, 3 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle * \rangle$

$\langle \text{id}, 3 \rangle$

$$E \xrightarrow{E + T}$$

$$E \xrightarrow{E + T} T \xrightarrow{T * F}$$

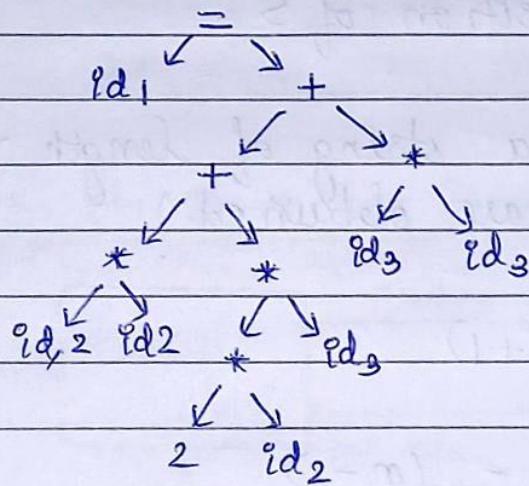
$$E \xrightarrow{E + T} T \xrightarrow{T * F} F$$

$$E \xrightarrow{E + T} T \xrightarrow{T * F} F$$

$$E \xrightarrow{E + T} T \xrightarrow{T * F} F$$

context free grammar tree

ALL:



Intermediate Code Generation

$$t_1 = \text{id}_2 * \text{id}_2$$

$$t_2 = 2 * \text{id}_2$$

$$t_3 = t_2 * \text{id}_3$$

$$t_4 = t_1 + t_3$$

$$t_5 = \text{id}_3 * \text{id}_3$$

$$t_6 = t_4 + t_5$$

$$\text{id}_1 = t_6$$

Code Optimization

$$t_1 = \text{id}_2 * \text{id}_2$$

$$t_2 = 2 * \text{id}_2$$

$$t_3 = t_2 * \text{id}_3$$

$$t_4 = t_1 + t_3$$

$$t_5 = \text{id}_3 * \text{id}_3$$

$$\text{id}_1 = t_4 + t_5$$

- * Prefix of string - string obtained by removing 0 or more characters (symbols) from the end of S.
Eg: Σ , hello, h, he, hel, hell
- * Suffix - string obtained by removing 0 or more characters from beginning of S.
- * Substring - obtained by deleting any prefix and any suffix from S
- * Proper prefix, suffix substring - not Σ or S
Eg: h, he, hel, hell
- * Subsequence - any S formed by deleting 0 or more not necessarily consecutive position of S.

Q: * If there is a string of length "n", How many of the foll. are obtained?

a) Prefix - $(n+1)$

b) Proper Prefix - $(n-1)$

c) Substring - $n \frac{(n+1)}{2} + 1$

$$\downarrow \quad 3 \left(\frac{4}{2} \right) + 1 = 7$$

$3 \rightarrow a, b, c, ab, bc, ac, abc$

* Digits $\rightarrow [0-9]^+$

Num \rightarrow digits (floating point value) ? (E [+ -] ? digits) ? } Regular Definition
Regular Expression

letter $\rightarrow [A-Za-z-]$

Id \rightarrow letter (letter | digits)* }

R.E for identifier

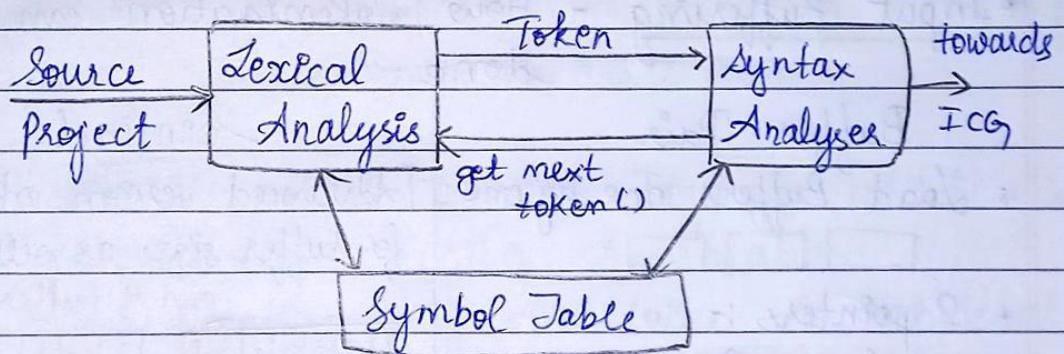
relational :- relop $\rightarrow < | \leq | <= | > | \geq = | =$
operators

arithmetic :- arithop $\rightarrow [+|-|^*|/| \%|=]$ } "-" character literally
operators or
 $+|-|*|/| \%|$

* Finite Automata (Deterministic)

* Transition Diagram

* Interaction b/w Lexical and Syntax Analyser



Tasks :

- 1) Tokenization
- 2) Stripping comments, tab spaces and white spaces.
- 3) Keep tracking of line nos.
- 4) Correlating error messages.

* Terms:

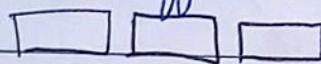
- 1) Pattern - Description of the form that lexemes of a token may take.
- 2) Lexeme - Sequence of character in a source program that matches pattern for a token.
- 3) Token - pair consisting of token names and attribute values.

* Lexical Analyzer Error Recovery Strategy can do the foll : (Lexical Errors Handler)

- 1) Delete any extra character
- 2) Insert any one character
- 3) Swap adjacent characters

* Input Buffering - How Tokenization can be done

Buffer Pair
* Load Buffers one by one



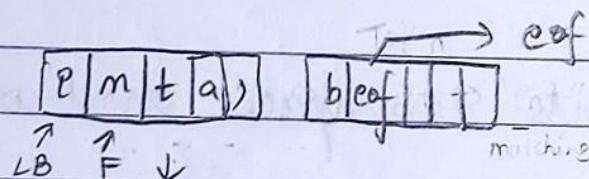
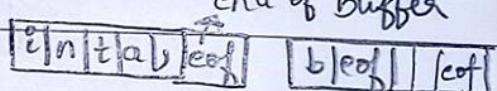
* 2 pointers :-

1) LB - Lexeme Begin

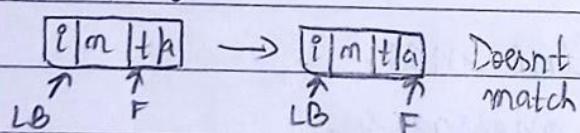
2) F - Forward

Sentinel

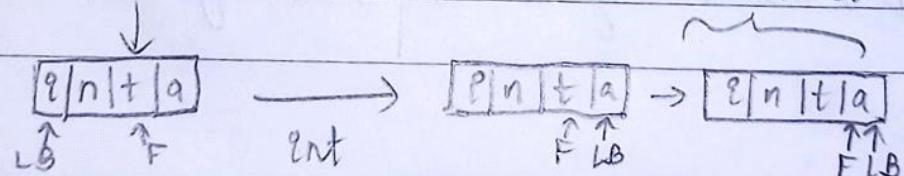
[Advanced version which checks for buffer size as well]



Lexical analysis process
Syntax analysis



int a



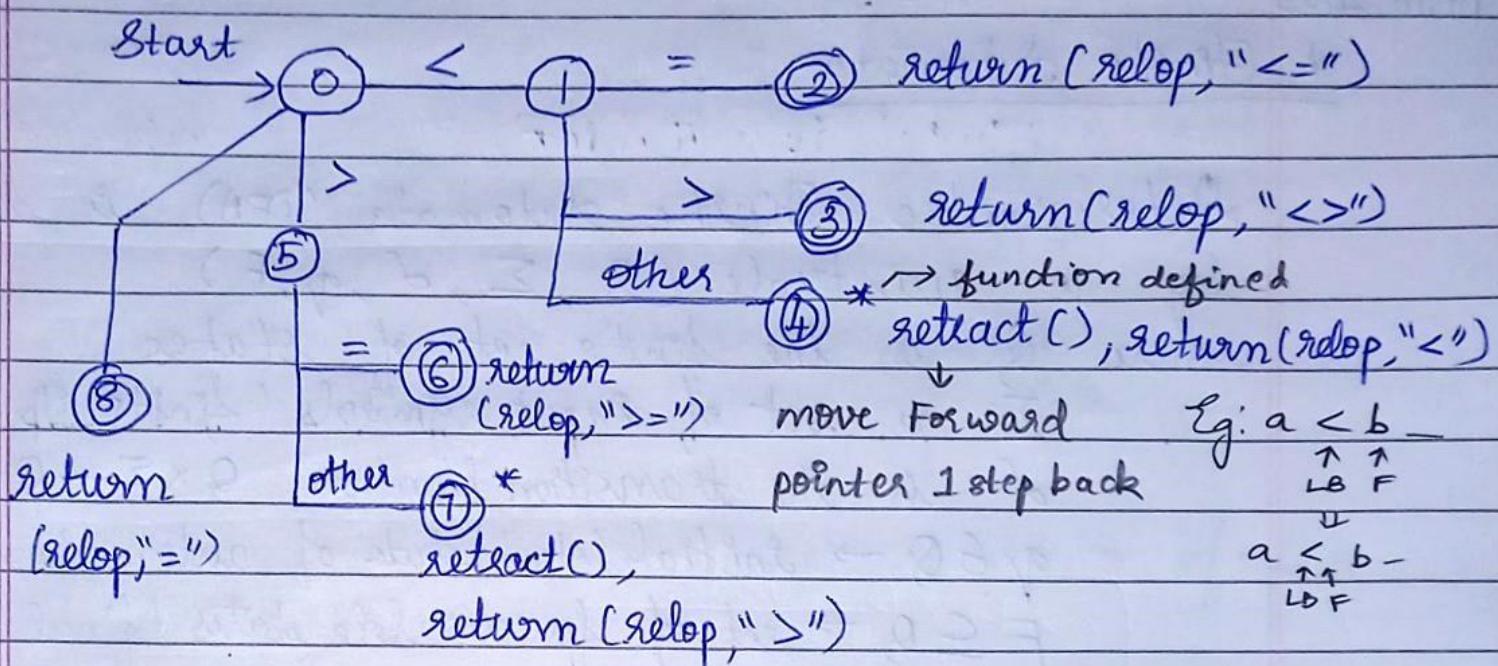
token name att. value

- * Defined on R.E
 - { float → "float"
 - int → "int"
 - id → $l(lid)^*$
 - d → [0-9] $^+$

- * Some scripting languages → panic-mode recovery
(errors in all lines after an error in particular line)

- * relop → <= | <> | < | > | >= | =

Start state, Accepting state



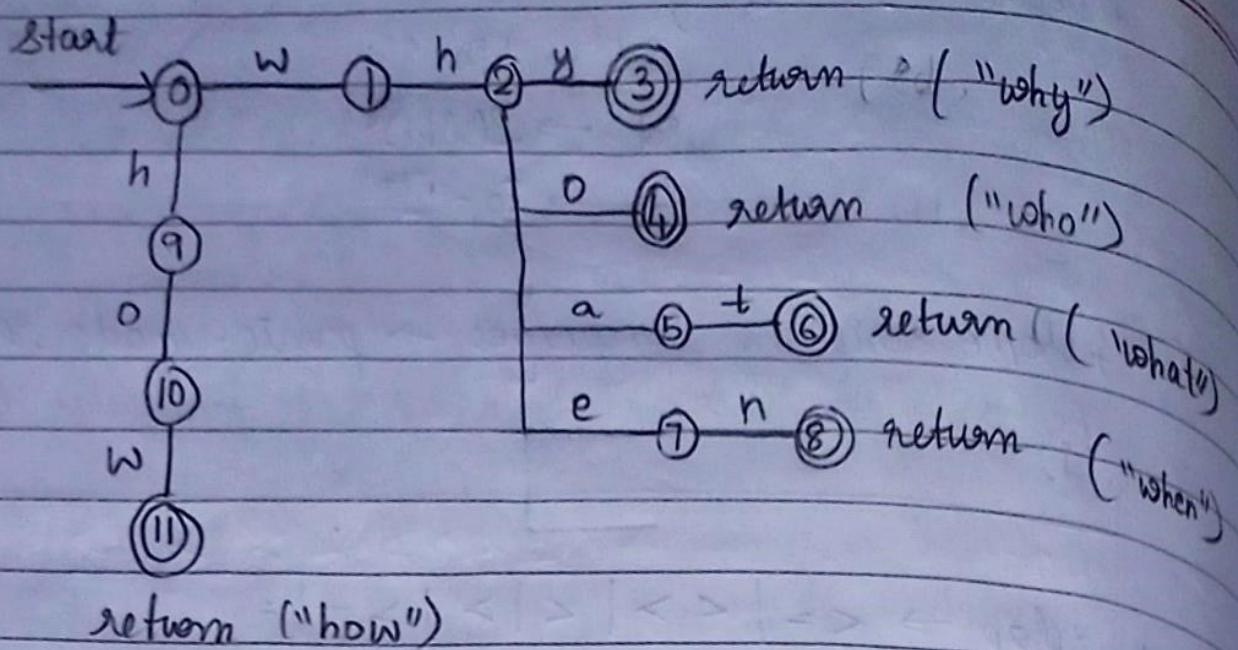
This is called Transition Diagram which is Deterministic Finite Automata.

[Can identify only one token at each state]

- * Draw Transition diagram to accept why, who, what, when, how.

Σ_i
 $+ \int_{\text{state} = 0, n}^{\text{input}} \rightarrow \downarrow \text{new state}$

CLASSMATE
 Date _____
 Page _____



19.10.2023

* Finite Automata

Deterministic Finite Automata (DFA) is denoted by a quintuple $(Q, \Sigma, \delta, q_0, F)$ where Q is the finite set of states Σ is set of input symbols (finite I/P alphabet) δ is the transition function $Q \times \Sigma \rightarrow Q$ $q_0 \in Q \rightarrow$ initial (start) state of automation $F \subseteq Q \rightarrow$ set of favourable states (final state)

* Deterministic means every move of a Finite Automata is completely determined by the input and current state

Eg: $R.E = (a|b)^*$

Apply DFA on RE,

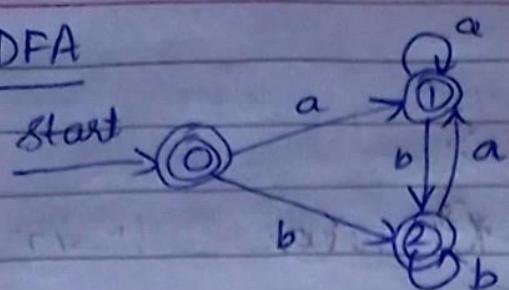
$$\Sigma = \{a, b\}$$

$$(a|b)^* = \{\epsilon, a, b, aa, \dots, bb, \dots, ab, \dots, ba, \dots\}$$

$$Q = \{q_0, q_1, q_2\}$$

$$F = \{q_0, q_1, q_2\}$$

$$q_1 = q_0$$

DFA

$$\delta(q_0, a) = q_1$$

$$\delta(q_0, b) = q_2$$

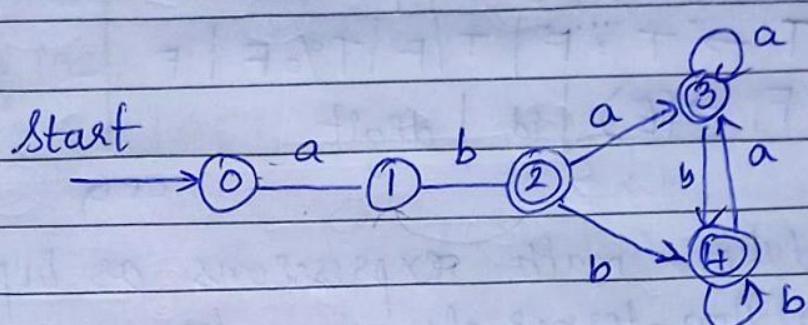
$$\delta(q_1, a) = q_1$$

$$\delta(q_1, b) = q_2$$

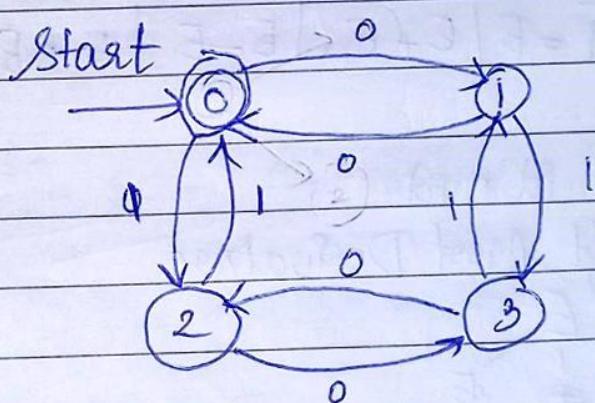
$$\delta(q_2, a) = q_1$$

$$\delta(q_2, b) = q_2$$

* RE = $ab(a \mid b)^*$ = ab, aba, abb

DFA

* RE Construct DFA to accept even no. of 0s and 1s.

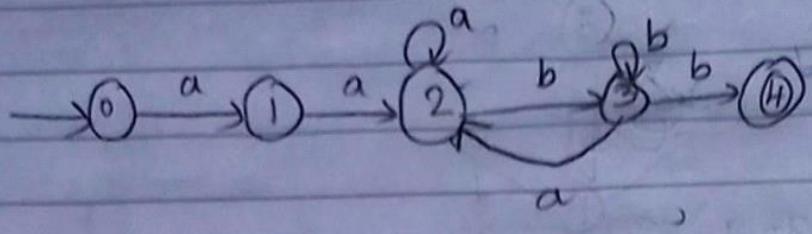
DFA

0011

0110

1111

* $RE = aa \underset{12}{(a|b)}^* bb \underset{34}{(a|b)}^*$ $1 \rightarrow M^n$



26.10.2023

* Context Free Grammars

$$G = (V, T, P, S) \rightarrow \begin{array}{l} \text{terminals - tokens : +, -, *, id} \\ \text{start symbol} \end{array}$$

$$G = (V, T, P, S) \rightarrow \begin{array}{l} \text{set of grammar symbols - production} \\ \text{non-terminals} \end{array}$$

\hookrightarrow can be expanded further: E, T, F

	$S \rightarrow id = E E$	
Exp.	$E \rightarrow E + T E - T T$	
Term	$T \rightarrow T * F T / F T \% F F$	$\rightarrow 1 \text{ production set}$
factor	$F \rightarrow (E) id \text{digit}$	CFG

* It takes math expressions as input.

* V → Non-terminals - can be expanded further:

E, T, F

* T → terminals - tokens : +, -, *, id

* Grammars:

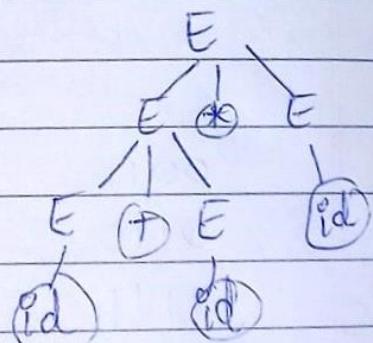
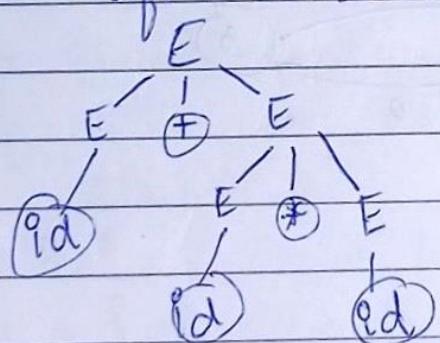
$$E \rightarrow id = E | E + E | E - E | E * E | E / E | E \% E | (E) | id$$

digit

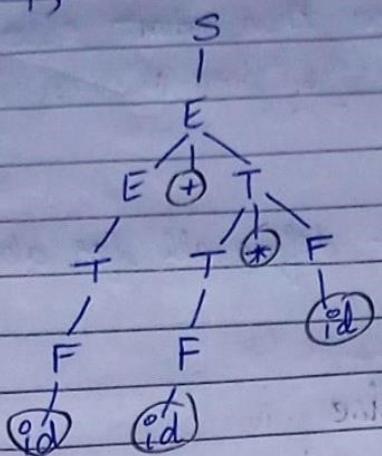
Q: id + id * pd

LMD - Left Most Derivative

LMD



* string :- $id + id * id$
using CFG,



* Design a & CFG to accept strings
0011, 000111, 00001111...

$$S \rightarrow 00EI1 \\ E \rightarrow \epsilon \mid 0EI$$

* Consider CFG $S \rightarrow SS+ \mid SS^* \mid a$ for
string $aa + aa^* +$

a) Give LMD

b) Give RMD

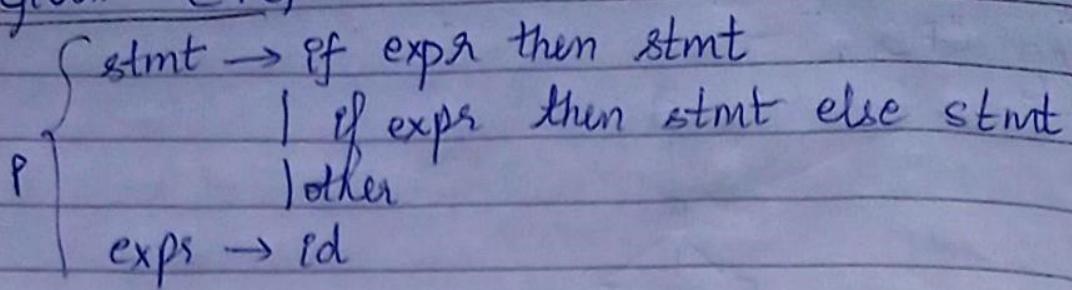
c) Draw Parse tree

d) Is Grammar ambiguous? Justify.

a) LMD

$$S \xrightarrow{lm} SS+ \xrightarrow{lm} SS+S+ \xrightarrow{lm} aS+S+ \xrightarrow{lm} aa+S+ \xrightarrow{lm} aa+SS^*+ \xrightarrow{lm} aa+aa^*+$$

* Given CFG :



a) $G = V, T, P, S$ - define all for above CFG

$V = \text{stmt, expr}$

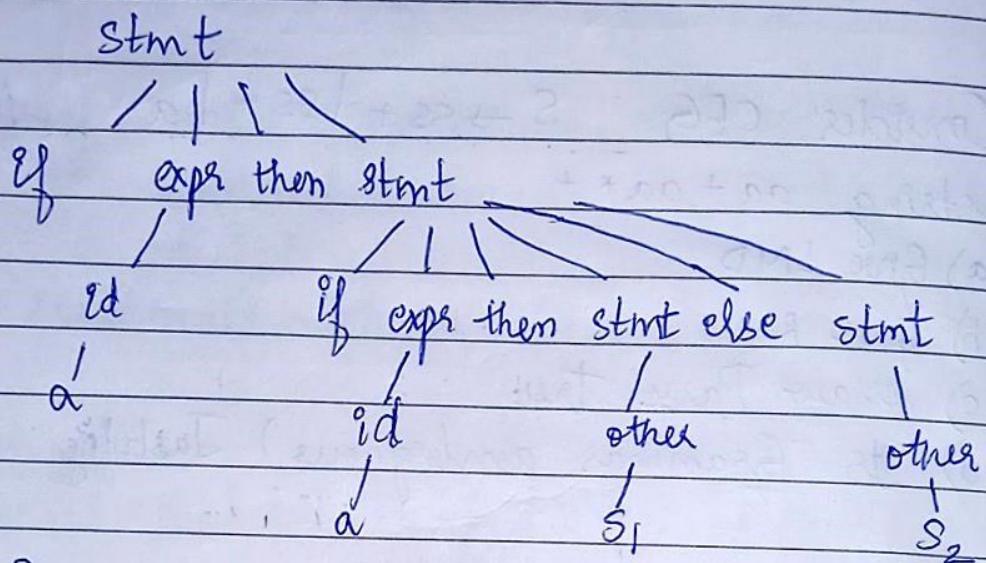
$T = \text{id, if, then, else}$

$P = \text{whole CFG}$

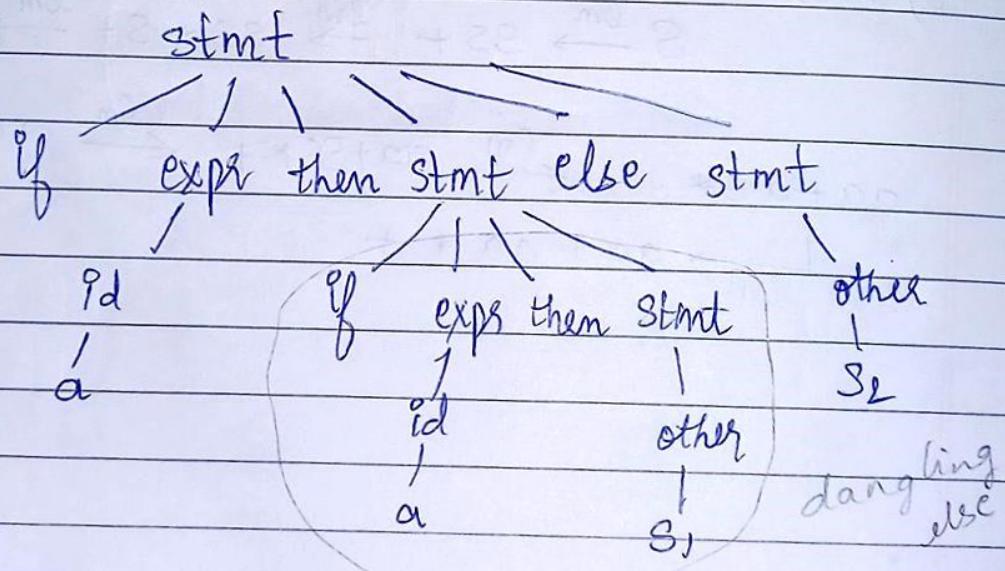
$S = \text{stmt}$

b) p/p :- if a then if a then s_1 else s_2
 ↳ dangling else

LMD - 1



LMD - 2



* Removing ambiguity from dangling else grammar:

stmt → if expr then stmt
| if expr then stmt else stmt
| other

Can be converted to stmt → matched - stmt
 →
 | open stmt
matched - stmt → if expr then matched - stmt,
 stmt else matched - stmt
 | other
open - stmt → if expr then stmt
 | if expr then matched - stmt else
 | open - stmt
expr → id

* Eliminating ambiguity for TDP - Top Down Parser:

$$* \quad S \rightarrow SS * \mid ss + |a$$

* Not suitable for NR-PP : $A \rightarrow A\alpha_1 | \beta$
 $(TDP) \qquad \qquad A \rightarrow \alpha P_1 | \alpha P_2$

as backtracking is not allowed

* For TDP,

If left recursion is there $\Rightarrow \underline{s} \rightarrow \underline{s} s^* | ss + la$, then
eliminate LR

* If left factoring is there $\Rightarrow S \rightarrow 0S1 | 01 \Rightarrow$ do LF
 $A \rightarrow \alpha B_1 | \alpha B_2$
 $S \rightarrow 0 S1 | 01$

$$\begin{array}{l|l} A \rightarrow x.A' & \} \\ A' \rightarrow B_1 \quad | \quad B_2 & \} \end{array}$$

$$\begin{array}{l} S \rightarrow OS' \\ S' \rightarrow SI \end{array} / /$$

* Eliminate LR (left recursion)

$$\begin{array}{l}
 \text{eliminate LR} \\
 \begin{array}{c}
 \begin{array}{c}
 E \rightarrow E + T \mid E - T \mid T \quad - \textcircled{1} \\
 T \rightarrow T * F \mid T / F \mid T \% F \mid F \quad - \textcircled{2} \\
 F \rightarrow (E) \mid \text{id} \mid \text{num} \quad - \textcircled{3}
 \end{array}
 \quad \left\{ \begin{array}{l} E \rightarrow TE' \\ E' \rightarrow +TE' \mid -TE' \mid \epsilon \\ T \rightarrow FT' \\ T' \rightarrow *FT' \mid /FT' \mid \%FT' \\ F \rightarrow (E) \mid \text{id} \mid \text{num} \end{array} \right.
 \end{array}
 \end{array}$$

To check for LR, see if the below format exists ,

$$A \rightarrow A\alpha \mid B$$

$$\begin{array}{c}
 \text{eliminate} \\
 \begin{array}{c}
 A \alpha_1 \quad A \alpha_2 \quad B \\
 E \rightarrow E + T \mid E - T \mid T \quad - \textcircled{1}
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 A \rightarrow A\alpha \mid B \Rightarrow \left\{ \begin{array}{l} A \rightarrow BA' \\ A' \rightarrow \alpha A' \mid \epsilon \end{array} \right. \\
 \text{eliminate} \\
 \text{LR}
 \end{array}$$

Similarly for $\textcircled{1}$,

$$\begin{array}{l}
 E \rightarrow TE' \\
 E' \rightarrow +TE' \mid -TE' \mid \epsilon.
 \end{array}$$

Now for $\textcircled{2}$,

$$\begin{array}{l}
 T \rightarrow FT' \\
 \textcircled{2} \Rightarrow T' \rightarrow *FT' \mid /FT' \mid \%FT' \mid \epsilon. \\
 F \rightarrow (E) \mid \text{id} \mid \text{num}
 \end{array}$$

* Left Recursion

$$\begin{array}{c}
 \text{Direct LR} \\
 A \rightarrow A\alpha \mid B
 \end{array}$$

$$\begin{array}{c}
 \text{Indirect LR} \\
 S \rightarrow Aa \mid b \\
 A \rightarrow Sa \mid a
 \end{array}$$

- * Input :- Grammar G
O/P :- G , without LR

Method :- Apply algorithm to G

Arrange the non-terminals in same order A_1, A_2
for each i from 1 to $n\}$

for each j from 1 to $i-1\}$

replace each production of the form

$A_i \rightarrow A_j \alpha$ by the production

bring it to this format $\rightarrow A_i \rightarrow f_1 \alpha_1 | f_2 \alpha_2 | \dots | f_n \alpha_n$ where

$A_j \rightarrow f_1 | f_2 | \dots | f_k$ A_j production

eliminate immediate LR from A_i^0 production

$$G: S \rightarrow Aa|b$$

$$A \rightarrow Ac|Sd|\varepsilon_0$$

Iteration

$$\begin{cases} i=2 \\ i=1 \end{cases}$$

$$A_1 = S$$

$$i=1 \rightarrow 2^n$$

$$A_i \rightarrow A_j \alpha$$

$$A_2 = A$$

$$j=1 \rightarrow 0$$

$$A_2 \rightarrow A_1 \alpha$$

$$A \rightarrow S \alpha$$

$$① S \rightarrow Aa|b$$

$$A \rightarrow Ac|Aad|bd|\varepsilon_0$$

$$② A \rightarrow bdA' | A'$$

$$A \prec_1 A \prec_2 P_1 B_2$$

$$③ A \rightarrow cA' | adA' | \varepsilon_0$$

$$* G: S \rightarrow Aa|b$$

$$A \rightarrow Ac|Sd|\varepsilon_0$$

$$1^{st} \text{ iteration: } i=1$$

$$S \rightarrow Aa|b$$

$$j=1 \rightarrow 0$$

$$A_j$$

2nd

$$: i=2$$

$$S \rightarrow Aa|b$$

$$j=1$$

$$A_p$$

$$A \rightarrow Ac|Sd|\varepsilon_0$$

$$\bar{A}_j$$

$$A \rightarrow Ac | Aad | bd | \varepsilon_0$$

↓

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

$$A' \rightarrow cA' | adA' | \varepsilon_0$$

* Check whether the G are suitable for top-down parsers & if not, make it suitable.

$$1) S \rightarrow S+S \mid SS \mid (S) \mid S* \mid a$$

$$2) S \rightarrow (L) \mid a$$
$$L \rightarrow L, S \mid S$$

$$3) \begin{array}{c} A \\ \hline S \end{array} \rightarrow \begin{array}{c} A \alpha_1 \\ S+S \end{array} \mid \begin{array}{c} A \alpha_2 \\ SS \end{array} \mid \begin{array}{c} P_1 \\ (S) \end{array} \mid \begin{array}{c} A \alpha_3 \\ S* \end{array} \mid \begin{array}{c} P_2 \\ a \end{array}$$

$$\begin{array}{c} S \rightarrow aS' \mid (S)S' \\ S' \rightarrow _+ SS' \mid SS' \mid *S' \mid \epsilon. \end{array}$$

$$2) S \rightarrow (L) \mid a$$
$$L \rightarrow L, S \mid S$$

$$\begin{array}{c} S \rightarrow (L) \mid a \\ L \rightarrow \begin{array}{c} L, S \mid (L) \mid a \\ \Downarrow \\ L \rightarrow (L)L' \mid aL' \\ L' \rightarrow , SL' \mid \epsilon. \end{array} \end{array}$$