

Automata theory & Compiler design

Source program (HLL)

Compiles

I/P → target program (LL) → O/P

assembly lang. (mnemonics)

I/P → link-expreter → O/P
source program

target
Program

Source program

Preprocessor

modified source program

Compiler

assembly program

Assembler

→ LLL / Assembly lang.

→ relocatable machine code

→ lib files | relocatable

→ gcc -lm | obj. files

→ linker

→ for math fn.'s

(cos sin)

[loader] → loads all the relocatable

machine code to memory

for execution

Hybrid compiler

source program

Translator | compiled

→ intermediate
(class) program (byte code) → virtual
by machine → O/P

Java combines compilation & interpretation

byte code compiled on one machine can be interpreted on another machine

Structure of a compiler

compiler is mapping source program to semantically equivalent target program

two phases

- analysis (front end)
- synthesis (back end)

→ lang dependent → lang indep.

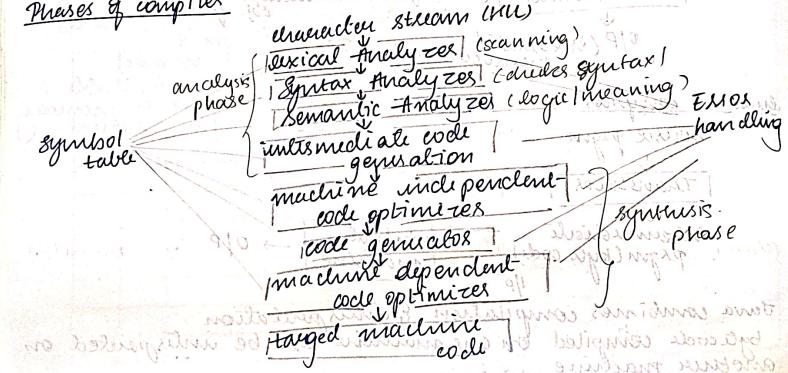
→ machine independent → machine dep.

→ breaks up source program into constituent pieces & imposes grammatical structure on them

→ detects source program is semantically ill formed then corrective step should be taken

→ collects info about source program & stores it in a data struc. called symbol table

Phases of compiler



Lexical analysis / Scanning

→ reads stream of chars making up the source program & groups the chars into meaningful sequences called lexemes (tokens)

→ token is of the form $\langle \text{token name}, \text{attr name} \rangle$ which passes on to syntax phase

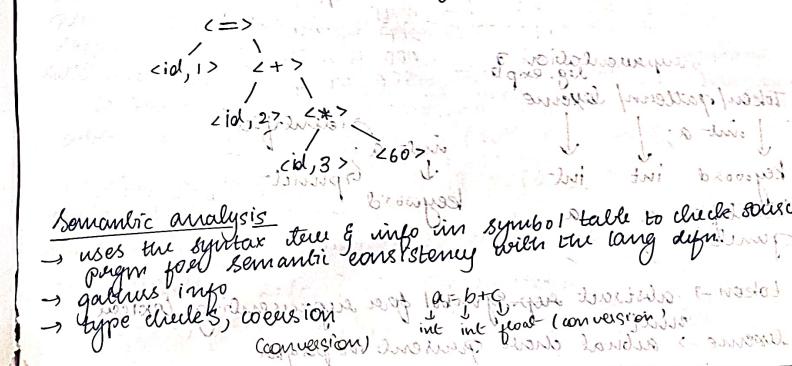
Eg. position = initial + size ≈ 60] this step converts stream into tokens
 $\langle \text{id}, 1 \rangle \Leftrightarrow \langle \text{id}, 2 \rangle \langle \text{id}, 3 \rangle \dots \langle \text{id}, n \rangle$ identifying lexeme
 identifies spaces
 identifier
 $\langle \text{token}, \text{attr} \rangle$ points to an entry in the symbol table for this token
 abstract symbol used during syntax analysis
 finds errors

→ info from symbol table needed for semantic analysis

Syntax analysis / Parsing

→ uses first component of token to create a tree like intermediate representation that depicts grammatical structure of token stream.

→ syntax tree → interior node → operation
 → children → args. of operation



Intermediate code generation

- after syntax & semantic phase, the source program, many compilers generate an intermediate representation → program for an abstract machine
- two properties
 - ↳ easy to produce
 - ↳ easy to translate
- uses three address instn

$t_2 = \text{id}_3 * t_1$

$t_2 = \text{id}_3 + t_2$

$\text{id}_4 = t_3$

$t_3 = \text{id}_3 * t_1$

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

$t_2 = \text{id}_3 * 60.0$

$\text{id}_1 = \text{id}_3 + t_2$

$\text{mul } \text{id}_3 \#60.0 \quad t_2$

$\text{add } \text{id}_2 \quad t_2 \quad \text{id}_1$

Code generation takes ans 1/p an intermediate representation of the source pgm. & maps it to target lang. target lang is machine code, registers for mem loc are selected for each of the variable used by the program. Then they're translated into sequences of machine inst. that perform same task

$t_1 = id_3 * 60.0$
 $id_1 = id_2 + t_1$] \rightarrow MULF R₂, R₃, #60
 \rightarrow LDF R₁, id₂
 \rightarrow ADD R₁, R₂, R₃
 \rightarrow STF id₁, R₁

→ representation of reg. exp b

Token/Pattern Lexeme \rightarrow Identifier
 ↓ int a; ↓ ↓ int a:
 keyword int int keyword → Identifier

↓
keyword int int
↓
identifier a a
↓
Punct / ,
↓
keyword

tokens are specified using regular exp. or words (set of chars.)
 specifications of token: alphabet, string, lang.
 finite set of symbols
 sequence of symbol drawn from
 countable set of strings over
 fixed alphabet

Operations on string

prefix: banana → ban (elimination at end)
 suffix: banaa → ana
 Substring: x → can't be an empty string / string → ban, ana
 ↳ empty string = ϵ (string with length = 0)
 subsequence: → bna → bnana (not continuous)

Operations on languages. $L = \{0, 1\}$ $S = \{a, b, c\}$

union: $L \cup S = \{0, Y, a, b, c\}$

concatenation: $L_5 = \overline{L_1} \cup L_2 \cup L_3 \cup L_4$... $\cup L_n$

$$f^* = \{ f_0, f_1, f_2, f_3, f_4, f_5, f_6, f_7 \}$$

$$\begin{aligned} S^* &= L(S)^*, \text{ e.g., } a^* \\ S^+ &= L(S) + a \text{ (one more)} \\ S/S &= L(S) \cup L(S) \text{ (two)} \\ SS &= L(S) L(S) \text{ (three)} \end{aligned}$$

Regular defn. vs

giving name to the regular exp.

$$\left. \begin{array}{l} d_1 \rightarrow s_1 \\ d_2 \rightarrow s_2 \\ \vdots \\ d_n \rightarrow s_n \end{array} \right\} d_1, d_2 \rightarrow \text{distinct s-name}$$

3.9 E + 2

Shortend

One or more occurrence $\rightarrow +$
zero or one instance $\rightarrow ?$
char class $\rightarrow []$
 \rightarrow unsigned
digit $\rightarrow 0/1 \dots 9$
digits \rightarrow digit(digit)* \rightarrow 1 or more digits
optional fractions \rightarrow digits/.
optional exponents $\rightarrow (E (+/-)(digits)/e)$
unsigned no. \rightarrow digits, optional fraction optional exponent

Expression	matches	Example
'c'	the one non-operator character 'c'	a2, a2, a2, a2
'*' or '^'	character c literally	*
"S"	string S literally	abc
."	any char but newline	a.b
^	beginning of a line	abc
\$	end of a line	abc\$
[s]	any line char starting	a[b]c

Exercise "abc" matches:
normal match: abc
Word up: abc
"": abc
"abc": abc
"abc": abc

classifies:
G(letter) (letter/digit)

$\cap A - z \ a - z \ - \ \cup$ union
letter: A1B1 ... z/a/b/.../z
(digit: 0/1 ... 9)

letter = [A...z, a...z]

digit = [0...9]

expression

["S"]

a, b

? , *

, ^

, \$

, []

, { }

, ()

matches

Example as do ()
[a/b/c]?

and more occurrance than 3 is after each

more of elements than 3

as well as more than 3

as well as more than 3

Write a complete line that doesn't contain a lowercase vowel

^ [aeiou]*\$

relational operators $\rightarrow 2 > / 2 = / 2 \leq / 2 \geq$

lexeme token attr value

\leq (less than or equal to) LT (less than or equal to)

\geq (greater than or equal to) GT (greater than or equal to)

$<$ (less than) LT (less than)

$>$ (greater than) GT (greater than)

$=$ (equal to) EQ (equal to)

\neq (not equal to) NE (not equal to)

single char at least

start ① a → ② at least

starting node ending node

return (relop LT)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

return (relop GE)

return (relop GT)

return (relop EQ)

return (relop NE)

return (relop LE)

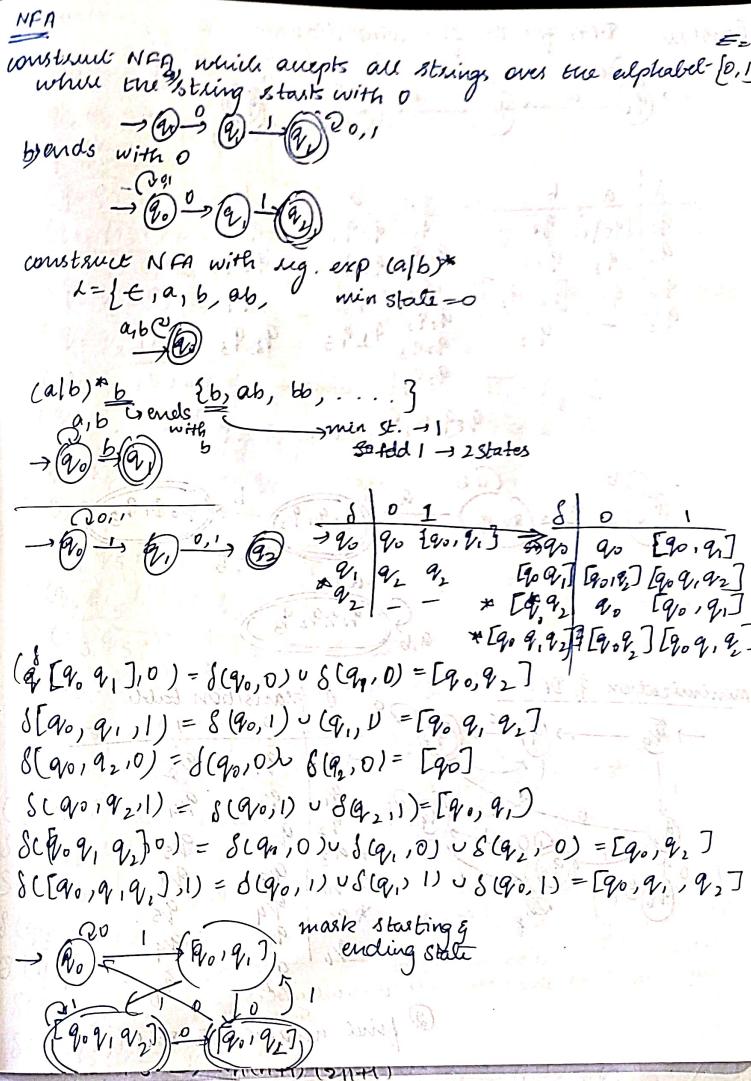
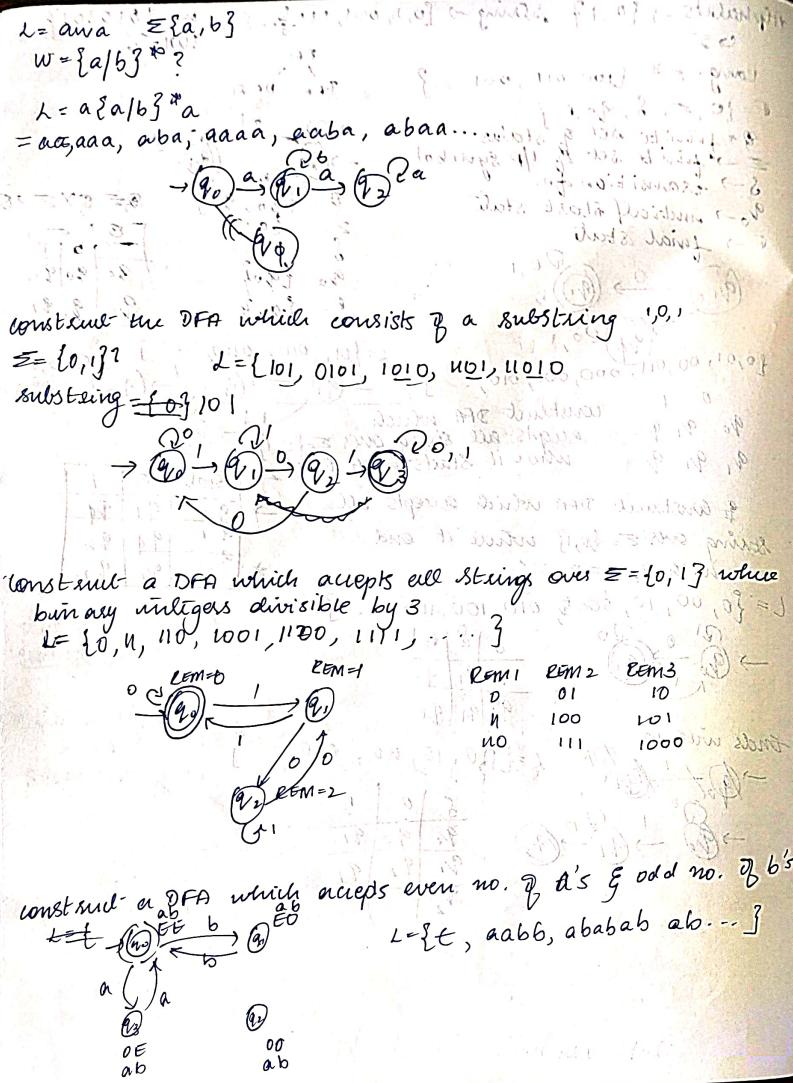
return (relop GE)

return (relop GT)

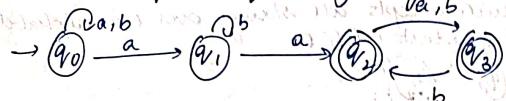
return (relop EQ)

return (relop NE)

return (relop LE)

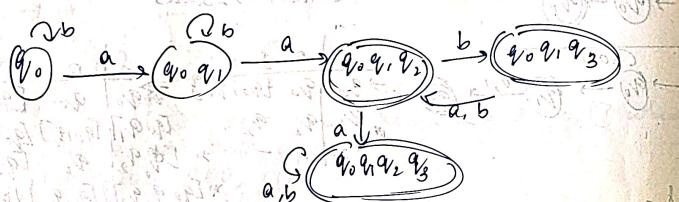


Construct DFA for the following diagram

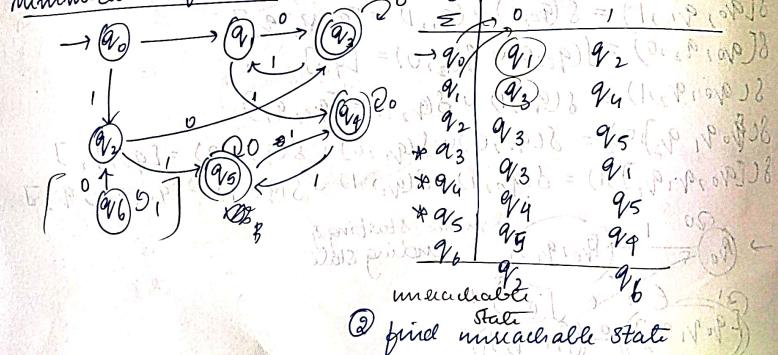


δ	a	b
q_0	$\{q_0, q_1\}$	q_0
q_1	q_2	q_0, q_1
q_2	q_3	q_0, q_1
q_3	$-$	q_2

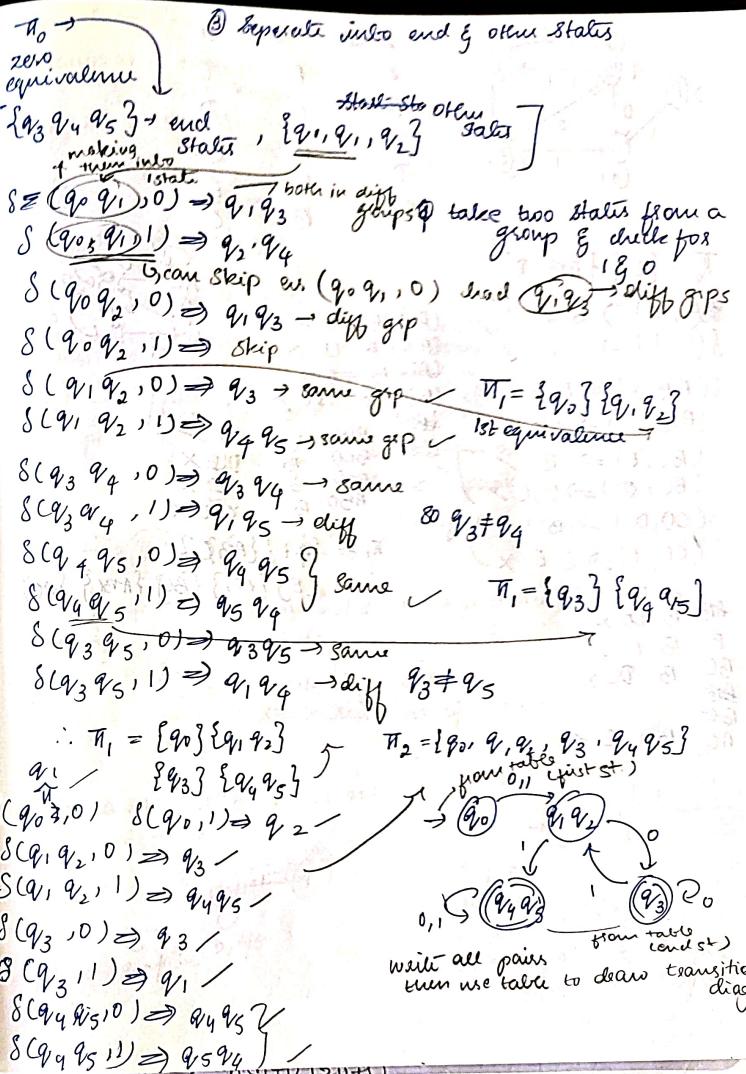
δ	a	b
q_0	q_0, q_1	q_0
q_1	q_0, q_1	q_0, q_1
q_2	q_1	q_2
q_3	q_2	q_2, q_3

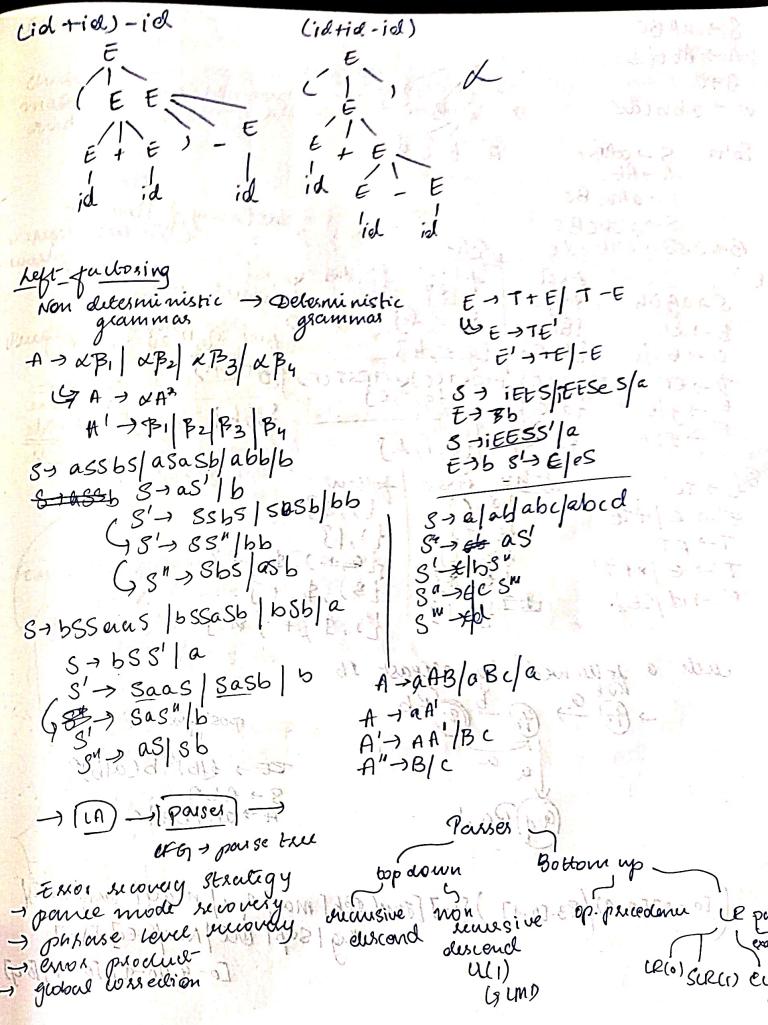
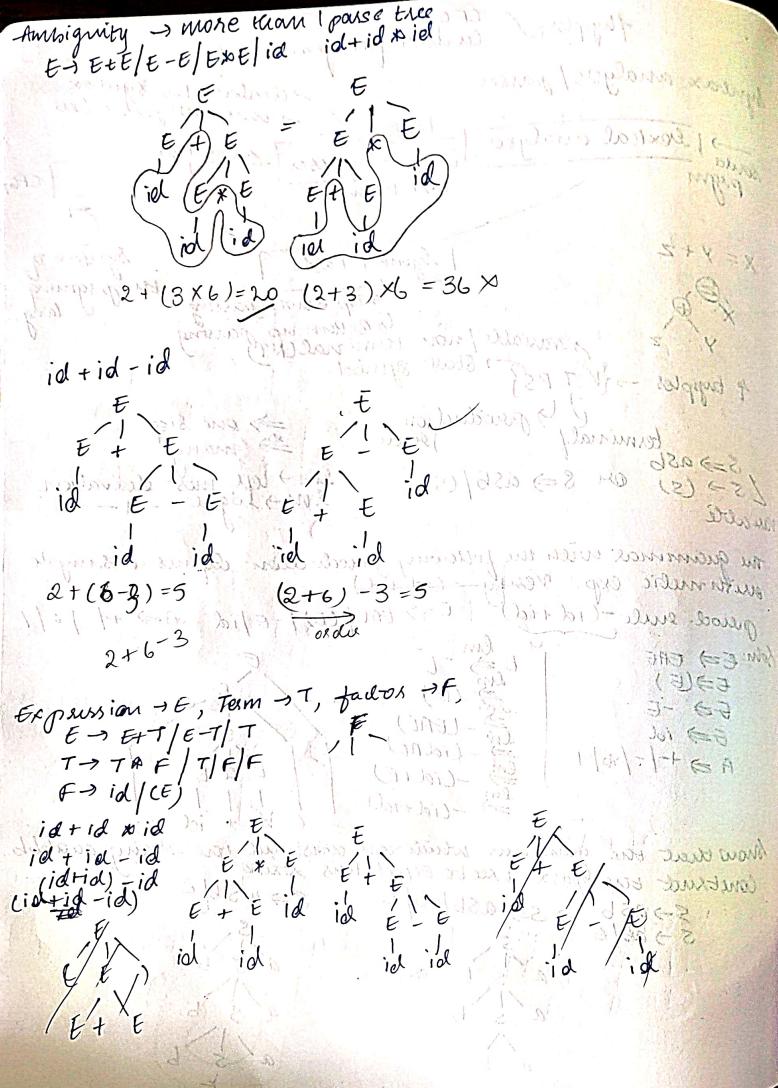


Minimization of DFA



② find unreachable state





$S \rightarrow aABe$
 $A \rightarrow Abc/b$
 $B \rightarrow d$
 $w \rightarrow abbcde$

Soln. $S \rightarrow aABe$

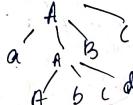
$A \rightarrow Abc$

$S \rightarrow AbcBe$

$S \rightarrow abbCBBe$

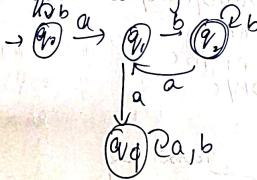
$B \rightarrow a$

$\Rightarrow S \rightarrow abbcde$



Non-terminal	Production	First	Follow
S	$S \rightarrow aABe$	$\{a\}$	$\{\$, b, e\}$
A	$A \rightarrow Abc$	$\{a, b\}$	$\{b, \$\}$
B	$B \rightarrow d$	$\{d\}$	$\{\$, e\}$
E	$E \rightarrow EF$	$\{F\}$	$\{E, F\}$
F	$F \rightarrow f E$	$\{f, E\}$	$\{f, E\}$
T	$T \rightarrow FT'$	$\{T'\}$	$\{T, \$\}$
T'	$T' \rightarrow ET' FT' \epsilon$	$\{E, T'\}$	$\{E, T', \$\}$
E'	$E' \rightarrow ET' \epsilon$	$\{E, \$\}$	$\{\$\}$
E	$E \rightarrow id CE TE ET'$	$\{id, C, T, E\}$	$\{\$\}$
C	$C \rightarrow id (E)$	$\{id, (\}\}$	$\{id, \$\}$
T	$T \rightarrow id FT ET FT ET \epsilon$	$\{id, F, E, T\}$	$\{\$\}$
F	$F \rightarrow id CE TE ET'$	$\{id, C, T, E\}$	$\{\$\}$

each a followed by atleast 1 b



postfix

$\rightarrow EE \rightarrow (ab)^*$

$\rightarrow S \rightarrow ABA$

$\rightarrow A \rightarrow aA / bA / \epsilon$

$([0-2][0-9]/[0-3][0-1])[1-][Jan/Feb/Mar/Apr/May/Jun/Jul]$
 $[Aug/Sep/Oct/Nov/Dec][1-][0-9][0-9]/[0-9][0-9]$

consider the CFG given

$S \rightarrow 0\$ / 01$
 check whether the g is suitable for predictive parsing. If not make it suitable and write modified to CFG. Compute FIRST and FOLLOW set of all non-terminals in the CFG

$S \rightarrow 0\$$

$S' \rightarrow 0\$$

Non-terminal

S

S'

First

$\{0\}$

$\{0, 1\}$

Follow

$\{\$, 1\}$

$\{1, \$\}$

$\{0\}$

$\{0, 1\}$

$\{1, \$\}$

Predictive parse table

NT / P	P		
	0	1	\$
Non-terminal	$S \rightarrow 0\$$		
S'	$S' \rightarrow 0\$$	$S' \rightarrow 1\$$	

CFG represented as
 $\frac{A \rightarrow x}{A \rightarrow x}$
 no context

for every production $\frac{A}{B}$ form

$A \rightarrow x$ (check if its of $A \rightarrow \alpha$ form)

$\textcircled{1} (S \rightarrow 0\$)$ (find first of RHS)

find first of RHS

$\textcircled{2} 0$

$M[S, 0]$

$\textcircled{3} (S' \rightarrow 1\$)$ first(S') $\textcircled{2} 0$

$M[S', 0]$

$\textcircled{3} (S' \rightarrow 1)$

$M[S', 1]$

- For each terminal 'a' in FIRST(α), add $A \rightarrow a$ to M.
 - If ϵ is in FIRST(α), then for each terminal 'b', add $A \rightarrow \epsilon$ to M.
 - If ϵ is in FIRST(α) and $\$$ is in FOLLOW(α), add $A \rightarrow \$$ to M.
 - All blank entries are removed.

$S \rightarrow AaAb / BbBa$ will be in the form $A \xrightarrow{B} X$

NT	λP	a	b	\$
S	$S \rightarrow AaAb$		$S \rightarrow BbB\alpha$	
A		$A \rightarrow \epsilon$		$A \rightarrow \epsilon$
B		$B \rightarrow \epsilon$	$B \rightarrow \epsilon$	

$$S \rightarrow A a A b \{$$

first($\{$) = a

$$m[S, a]$$

$$S \rightarrow B b B a$$

first($\{$) = b

$$m[S, b]$$

$A \rightarrow E$ $\text{first } E(A) = E$
 so check follow $E(A) = a, b$

$m[A, a], m[A, b] \quad S \rightarrow A \# B$

$B \rightarrow E$
 $\text{first } E(B) = E$
 so check follow $(B) = b, a$

$m[B, a] \quad m[B, b]$

$$\begin{aligned}
 E &\rightarrow T E' \\
 E' &\rightarrow + T E' | E \\
 T &\rightarrow F T' \\
 T' &\rightarrow * F T' | E \\
 F &\rightarrow (E) | \text{id}
 \end{aligned}$$

compute first & follow of all

NT \ P	+	*	()	id	\$
E	E → E + B	B → a	E → TE'		E → TE	
T		T → FT'			T → F T'	
E'	E' → E + TE'			E' → E	B	E' → E
T'	T' → E	T' → FT'		T' → E		T' → E
F			F → (E)		F → id	

$$\begin{array}{l}
 E \rightarrow TE' \\
 \text{first}(TE') \rightarrow (, \text{id}) \\
 M[E, ()], M[E, \text{id}] \\
 T \rightarrow FT' \\
 M[T, ()], M[T, \text{id}]
 \end{array}$$

$T' \rightarrow \#FT'B$
 $\text{follow}(T') \rightarrow \$, +,)$
 $M[T', \$] \leftarrow M[T', +] \quad M[T', +] \leftarrow M[T', \#])$
 $E^I \rightarrow E^I$
 $\text{for } \text{follow}(E^I) \rightarrow \text{follow}(E^I)$
 $F \rightarrow \#(E)$
 $F \rightarrow \text{id}$
 $m[F, \{ \}] \quad m[F, \text{id}]$
 $\text{follow}(T') \rightarrow \$$
 $\text{for } \text{follow}(T')$
 $\text{for } \text{follow}(E^I)$
 $\text{for } \text{follow}(E)$

svd to create a simple desk/calculator for eval of exp / annotated false
true for $3k^5 + qn$

Production	Semantic rules
$1) L \rightarrow E_N$	$L.val = E.val$
$2) E \rightarrow E_1 + T$	$E.val =$

- Production

 - 1) $L \rightarrow EN$
 - 2) $E \rightarrow E1 + T$
 - 3) $E \rightarrow \text{digit} T$
 - 4) $T \rightarrow TI * F$
 - 5) $T \rightarrow \text{digit} F$
 - 6) $F \rightarrow (E)$
 - 7) $F \rightarrow \text{digit}$

$L.\text{val} = E.\text{val}$

$E.\text{val} =$

$F.\text{val} = E.\text{val}$

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

80 $3 * 5 + 4n$
~~Even = 19~~
~~Even is T. Val = 6~~
~~Even is F. Val = 4~~
~~Even is F. Val = 5~~
~~T. Val = 3~~
~~F. Val = 1~~
 Fixing F. Val digit digit. Val = 5
 $\text{digit} = \text{lexval} = 103$

- ~~1) $\frac{2+3}{7} \neq 4$~~

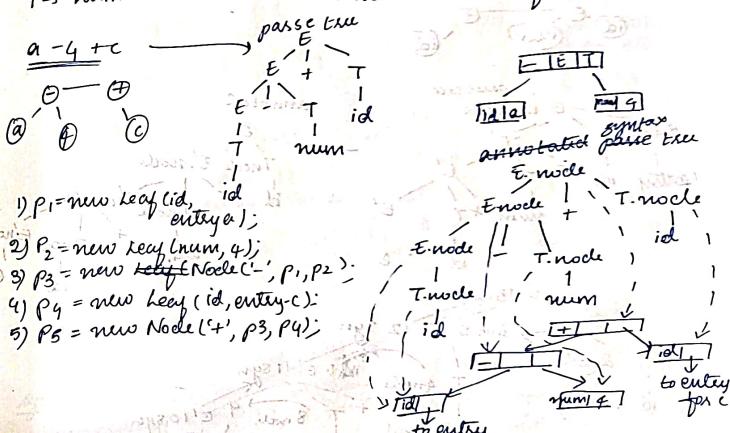
 - 2) $E \rightarrow E_1 + T$
 - 3) $E \rightarrow T$
 - 4) $T \rightarrow T_1 + F$
 - 5) $\textcircled{B} T \rightarrow F$
 - 6) $F \rightarrow (E)$
 - 7) $F \rightarrow \text{digit}$

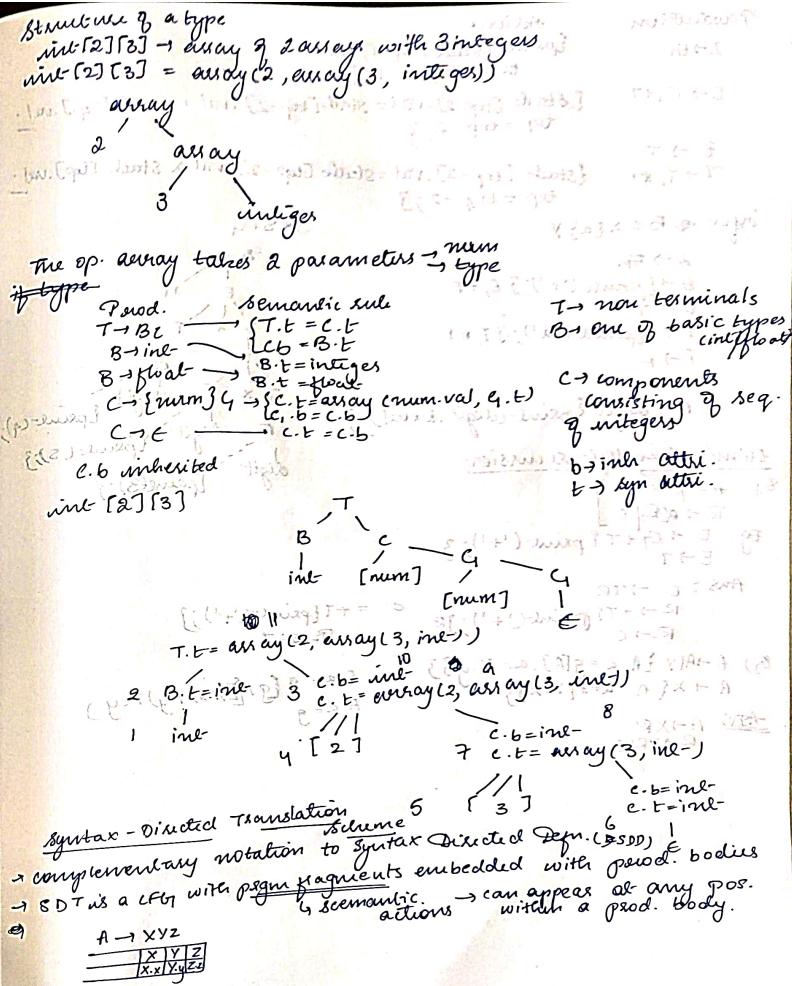
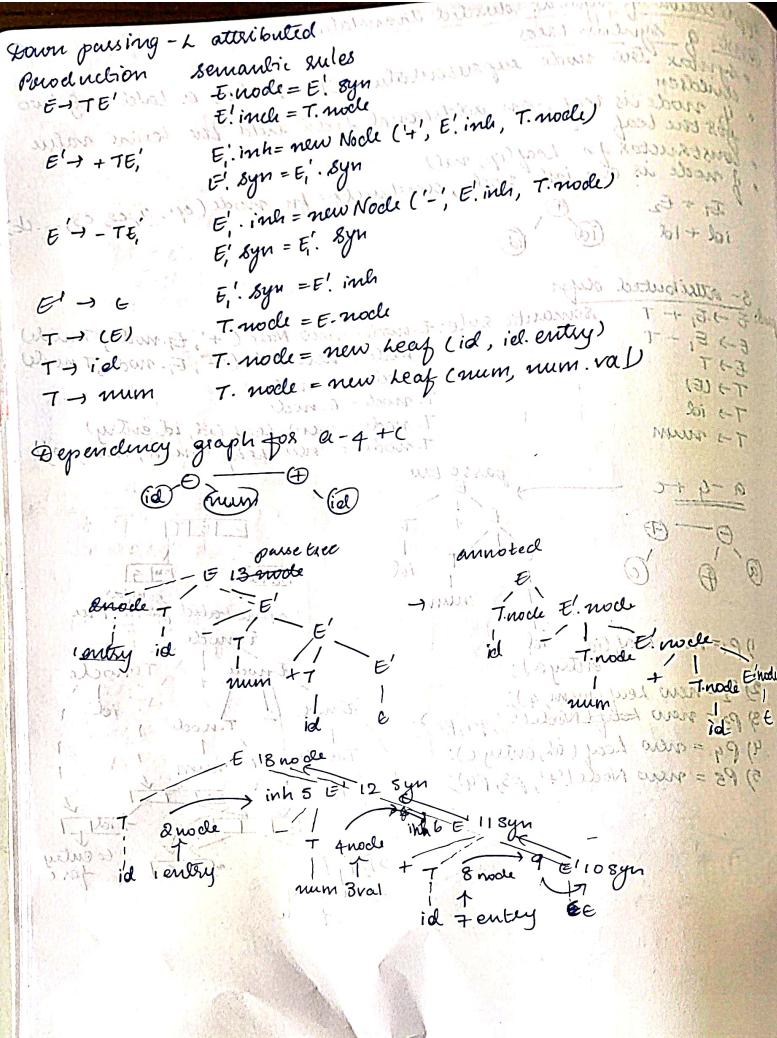
Application of syntax directed translation
const. of syntax trees

- symbol tree node representation $E1+E2$ has a label + 2 children
 - if node is leaf, an additional field hold the lexical value for the leaf.
 - constructor fn. leaf(op, val)
 - if node is an int. node, constructor fn. node(op, a, c2, c3, ... ct)

prob: $E_3 \rightarrow E_1 + T$ seems

- prob: $E \rightarrow E_1 + T$ semantic rule: $E\text{-node} = \text{new Node} ('+', E_1\text{-node}, T\text{-node})$
 $E \rightarrow E_1 - T$ $E\text{-node} = \text{new Node} ('-', E_1\text{-node}, T\text{-node})$
 $E \rightarrow T$ $E\text{-node} = T\text{-node}$
 $T \rightarrow (E)$ $T\text{-node} = E\text{-node}$
 $T \rightarrow id$ $T\text{-node} = \text{new Leaf} (id, id_entry)$
 $T \rightarrow num$ $T\text{-node} = \text{new Leaf} (num, num, val)$





Production

$L \rightarrow E_n$

$E \rightarrow E_1 + T$

$T \rightarrow T_i \# F$

Actions

{ $\text{push}(\text{stack}[\text{top}-1].\text{val});$
 $\text{top} = \text{top} - 1;$ }

{ $\text{stack}[\text{top}-2].\text{val} = \text{stack}[\text{top}-2].\text{val} + \text{stack}[\text{top}].\text{val};$
 $\text{top} = \text{top} - 2;$ }

{ $\text{stack}[\text{top}-2].\text{val} = \text{stack}[\text{top}-2].\text{val} \times \text{stack}[\text{top}].\text{val};$
 $\text{top} = \text{top} - 2;$ }

infix $\Rightarrow B \rightarrow \times \{a\} Y$

$L \rightarrow E_n$

$E \rightarrow \{\text{push}('+'; 3) E, + T$

$E \rightarrow T$

$T \rightarrow \{\text{push}('*'; 3) T, * F$

$T \rightarrow F$

$F \rightarrow (E)$

$F \rightarrow \text{digit } \{\text{push}(\text{digit}, \text{lexval}); 3$

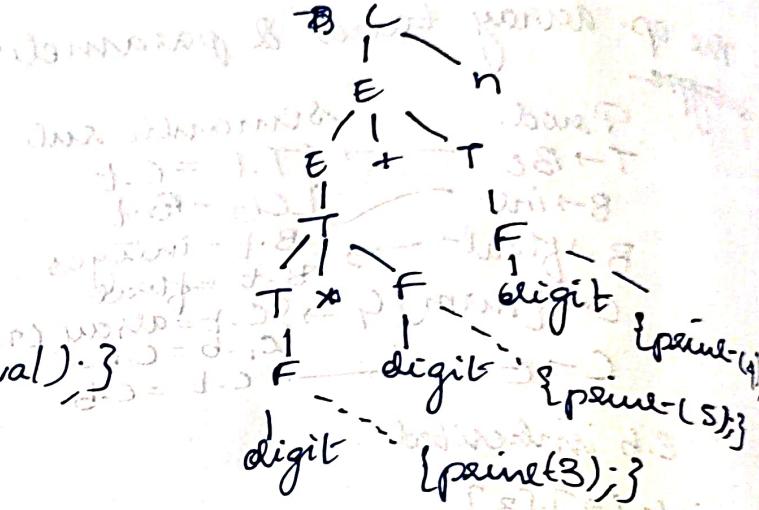
eliminating left-recursion

Q) $A \rightarrow \beta R$
 $R \rightarrow \alpha R | \epsilon$

Eg. $E \rightarrow E_1 + T \{\text{push}('+'; 3)$
 $E \rightarrow T$

Ans: $E \rightarrow TR$

$R \rightarrow + TS \{\text{push}('+'; 3) R$
 $R \rightarrow G$



Q) $A \rightarrow A(y \{ A.a = g[A].a = Y.y \})^3$
 $A \rightarrow X \{ A.a = f[X.x] \}$

$A.a = g(gf(x.x)y.y)^*y.y$

~~AES~~ $A \rightarrow XR$
 $R \rightarrow YR | E$