

Lexical Analysis

we will use flex

CMPS-104A

ch 3.1

Dragon 3:109

task: read chars
make lexemes
produce tokens

$y = y_{lex}()$
 $y \text{ text} \rightarrow$

Separate from Parser

- simplify design: ex: parser can ignore comments
- buffering
- portability: input separation
ex: ASCII, ISO Latin 1, Unicode
UTF8
UTF16

Tokens

~~Token = pair { token name, lexeme }~~

Token = pair { token name, attribute }

Pattern = description^{V_T} of form (regex)

Lexeme = seq of chars

ex:

Token	Lexeme
ID	foo bar
STRING	"foo" "\n"
IF	if
LE	<=
'<'	<

attributes

each token has a lexeme
plus source coordinates

file
line
offset

$G = \langle V_T, V_N, P, S \rangle$

36 "foo" ~

Lex Errors

- skip char
- ins/repl/swap

we
will
do.

CMPS-104A
ch 3.2

Dragon 3:114

note -3 is two tokens not 1

$x = 3 - 4;$ $x = 3 + -4;$
 ↑ ↑
 binop unop

lexer can't tell.

Buffering

- need to look ahead (usu. 1 char, more?)

ex Fortran DO 5 I = 1.25
 DO 5 I = 1, 25

C: < <= <<
 - -> --
 foo <=

x+++y
x+ ++y

buffer | if (den < num) foo = f(bar); |

 ↑ ↑
 lex lex
 begin end.
 (yytext)

yy lex
scan () {
 replace null plug.
 yytext lexbegin = old lexend
 advance lexend
 ~~*lexend = '\0'~~
 t = *lexend
 *lexend = 0
 return lexbegin
}

two ptrs

lex begin \rightarrow 1st char of lexeme

lex end \rightarrow 1st char after lexeme

end of buffer: shift back.

CMPS-104A
ch 3.3

Dragon 3: 115

Specifications

Strings & Languages.

alphabet: finite set of symbols.

EBCDIC

ASCII: 0..127

flex \rightarrow ISO Latin 1: 0..255*

Unicode: 0..0x10FFFF

never use 7-bit charset. $\uparrow = 1, 114, 111$

UTF8 } warning
UTF16 }

string: over an $\alpha\beta$ is a seq of sym \in that $\alpha\beta$

empty string: ϵ (~~ϵ vs ϵ~~) (ϵ vs ϵ)

language: set of strings over an $\alpha\beta$
may be \exists

$t = \text{prefix}(s) \iff \exists u: tu = s$

$t = \text{suffix}(s) \iff \exists u: ut = s$

$t = \text{substring}(s) \iff \exists u, v: utv = s$

proper means \neq

Operation on Lang

CMP5-104A
ch 3.4

Dragon 3:119

~~union, concat, closure~~

Union $L \cup M = \{s \mid s \in L \vee s \in M\}$

concat $LM = \{st \mid s \in L \wedge t \in M\}$

~~closure~~

Kleene closure $L^* = \bigcup_{i=0}^{\infty} L^i$

Positive closure $L^+ = \bigcup_{i=1}^{\infty} L^i$

Regular Exprs

~~is~~ a c ident $= \lambda (\lambda \mid \delta)^*$

if $\lambda \in a \dots z A \dots Z$

$\delta \in 0 \dots 9$

ex

in flex: $[a-zA-Z_][a-zA-Z_][0-9]^*$

or: $[a-zA-Z_][a-zA-Z_0-9]^*$

~~is a~~

each regex r denotes a lang $L(r)$

~~given an alphabet Σ~~

~~Basis 1. ϵ is a regex, $L(\epsilon) = \{\epsilon\}$~~

~~2. if $a \in \Sigma$ then a is a regex~~

↑ empty string

LETTER

$[A-Za-z_]$

DIGIT

$[0-9]$

IDENT

$\{LETTER\} \left(\{LETTER\} / \{DIGIT\} \right)^*$

each regex r denotes a
lang $L(r)$
defined recursively

CMPS-104A

ch 3.5

Dragon 3:121

given an $\alpha \beta \Sigma$

Basis 1. ϵ is a r.e.; $L(\epsilon) = \{\epsilon\}$

2. if $a \in \Sigma$ then a is a ~~RE~~ RE
and $L(a) = \{a\}$

ϵ
 \in

Induction

assume r, s are RE denoting $L(r), L(s)$

1. $(r)|(s)$ is a RE denoting $L(r) \cup L(s)$
2. $(r)(s)$ $L(r)L(s)$
3. $(r)^*$ $(L(r))^*$
4. (r) $L(r)$

Precedence

- (a) $*$ highest left assoc
- (b) concat medium left assoc
- (c) $|$ lowest left assoc.

Laws of RE

commutative
 $r|s = s|r$
x

associative
 $r|(s|t) = (r|s)|t$
 $r(st) = (rs)t$

identity
 $\epsilon r = r \epsilon = r$

distributive
 $r(s|t) = rs|rt$
 $(s|t)r = sr|tr$

closure
 $r^* = (r|\epsilon)^*$
 $r^{**} = r^*$

Definitions (flex)

CMP5-104A
ch 3.6

Dragon 3:123

LETTER $[A-Za-z_]$

DIGIT $[0-9]$

IDENT $\{\text{LETTER}\}(\{\text{LETTER}\}|\{\text{DIGIT}\})^*$

INTEGER $\{\text{DIGIT}\}^+$ *~~~~~ fails with 00.123*

FRACTION $\{\text{DIGIT}\}^+ \backslash . ? | \{\text{DIGIT}\}^* \backslash . \{\text{DIGIT}\}^+$

EXPONENT $[Ee][+-]? \{\text{DIGIT}\}^+$

FNUMBER $\{\text{FRACTION}\}(\{\text{EXPONENT}\})?$

Extensions to Regex

$(r)^+$ ~~means~~ denotes $(L(r))^+$
?, +, * same prec, assoc.

$(r)?$ denotes $L(r) \cup \{\epsilon\}$

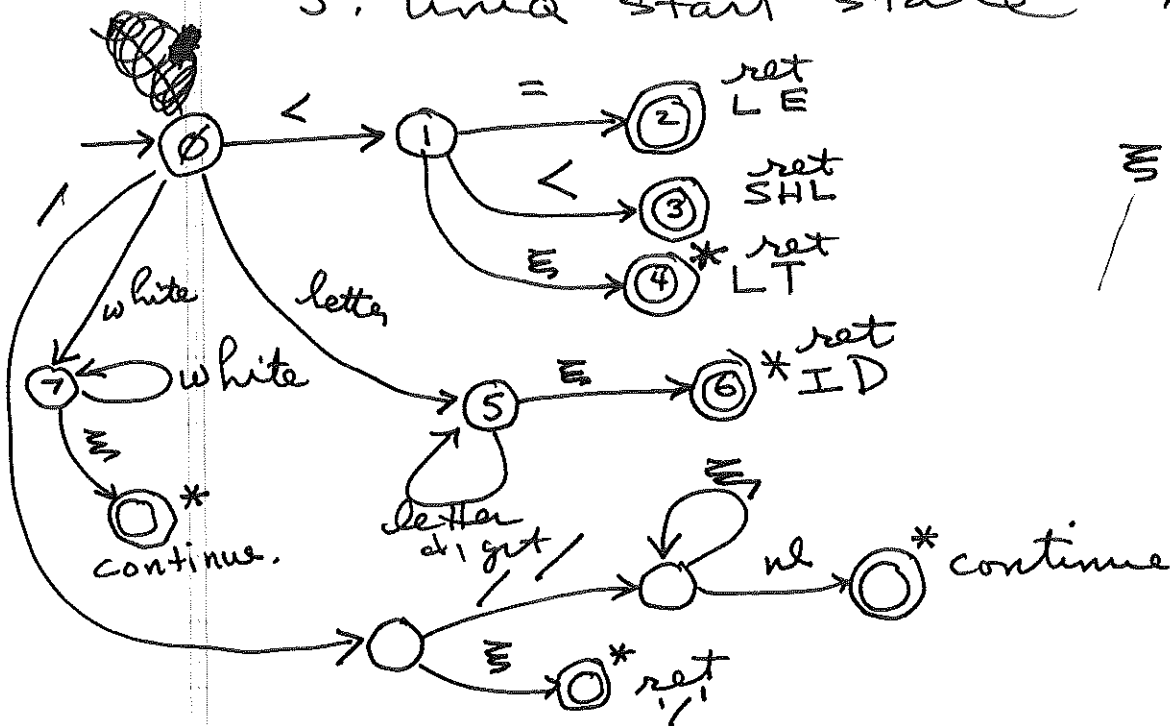
$[abc\dots]$ denotes $a|b|c|\dots$

Dragon 3:128

Transition diagrams

① ~~1. start state~~
~~2. final state~~

1. some states are final \odot
2. some final are retracting \odot^*
3. unia start state $\rightarrow \odot$



Recog Reg. words

CMPS - 104A

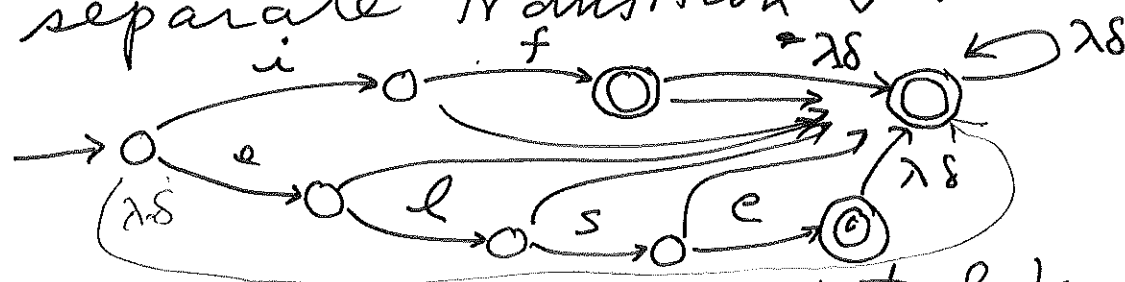
ch 3.9

Dragon 3:132

- scan as ~~IDENT~~
IDENT

1. install res wds in string table.
before scan begins
- extra field in token node

2. separate transition & kw.



3. flex cost: many extra states & trans.
pats.

Coding

```
begin = forward;
```

```
state = 0;
```

```
for (j) {
```

```
state = trans[state][*forward++];  
nstate = trans(state, *forward++);
```

```
if (nstate == ERROR) fail();
```

```
else if (final(nstate)) {
```

```
if (iskeyword)  
return token(nstate)
```

```
} else {
```

```
state = nstate
```

```
}
```

```
}
```


actions

- trans to non final \rightarrow continue
- trans to final space/comment
begin = forward
- trans to final token
retract if needed
return token, lexeme
- end of buffer - refill
- graph repr = 2D array
but SPARSE
- use char classes.
instead of n states & 256 cols
(Unicode ~~1000000~~?)
- init state very branchy
- others simpler
- or code directly.

scan tokens example

jflex
cup

CMPS-104A
ch3.10

Dragon 3:13)

Java idents

kitty
mačka
кошка
پیشی
猫

\$*£€

'A' ≠ 'A' ≠ 'A'

0x10FFFF
1114111

10/18/07
19:17:16

Regular expressions with a HUGE number of states
/afs/cats.ucsc.edu/courses/cms104a-wm/News/01642

1

```
1:
2: Article: 1642 of ucsc.class.cms104a
3: Newsgroups: ucsc.class.cms104a
4: Subject: Regular expressions with a HUGE number of states.
5: Distribution: ucsc
6: Organization: UC Santa Cruz CIS/CE
7: From: mackey@cse.ucsc.edu (W. Mackey)
8: NNTP-Posting-Host: asphalt3.cse.ucsc.edu
9: Message-ID: <47181434$1@darkstar>
10: Date: 18 Oct 2007 19:19:32 -0800
11: X-Trace: darkstar 1192760372 128.114.58.48 (18 Oct 2007 19:19:32 -0800)
12: Lines: 45
13: Path: darkstar!mackey
14: Xref: darkstar ucsc.class.cms104a:1642
15:
16: CS154 - Introduction to Automata and Complexity Theory
17: Stanford University
18: Jeffrey Ullman (Dragonbook author)
19: http://InfoLab.Stanford.EDU/~ullman/ialc/win00/win00.html
20:
21:
22:
23: Here are some test files for your regular-expression processor
24: (Project Part 1).
25:
26: 1. test1.txt: a RE whose DFA has about 1000 states.
27:
28: (0|1)*1(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)
29:
30:
31: 2. test2.txt: a RE whose DFA has about 32,000 states.
32:
33: (0|1)*1(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)
34: (0|1)(0|1)(0|1)(0|1)(0|1)(0|1)
35:
36:
37: 3. test3.txt: a RE whose DFA has about 1,000,000 states.
38:
39: (0|1)*1(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)
40: (0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)
41:
42:
43: 4. test4.txt: a RE whose DFA has about 32,000,000 states.
44:
45: (0|1)*1(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)
46: (0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)(0|1)
47: (0|1)(0|1)(0|1)(0|1)(0|1)(0|1)
48:
49:
50: 5. test5.txt: a RE with a big NFA but a very small DFA.
51:
52: [a-z]?[a-z]?[a-z]?[a-z]?[a-z]?[a-z]?[a-z]?[a-z]?[a-z]?[a-z]?
53:
54:
55: 6. test6.txt: an even more extreme version of test5.txt.
56:
57: [0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*
58: [0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*[0-9]*
59:
60:
61:
62:
```

$[0|1]\{2^5, \}$

3.5 Lex

```
int yylex() {  
    for (;;) {
```

$n = \text{scan pattern}$

```
        switch(n) {
```

```
            case: {act} break
```

```
            case: {act} brk
```

```
        }
```

```
    }
```

return tokencode.

FILE *yyin — fin

yytext → flex buffer.

yyval = semantic info.

yylen = #ch. scanned.

token codes $\emptyset = \text{EOF}$

1..255 = 'x'

256.. named tokens.

ambiguity

1. always longest char match

2. else use earlier part.

3 evil thing

goto

global

pointer.

$a[i] \equiv *(a+i)$
 $\equiv *(i+a)$
 $\equiv i[a]$

CMPS-104A

~~3.12~~ 3.12

Dragon 3:140

3.6 Finite Automata

CMP5-104A

3.13

Dragon 3:147

NFA: nodes & directed edges

DFA: at most one edge leaving a state for given symbol.

$$NFA = \{S, \Sigma, t, s_0, F\}$$

S = finite set of states

Σ = finite input alphabet, $\epsilon \notin \Sigma$

t : trans function.

$$t: S \times (\Sigma \cup \{\epsilon\}) \rightarrow S$$

s_0 : start state

F : set of final states, $F \subseteq S$

ϵ

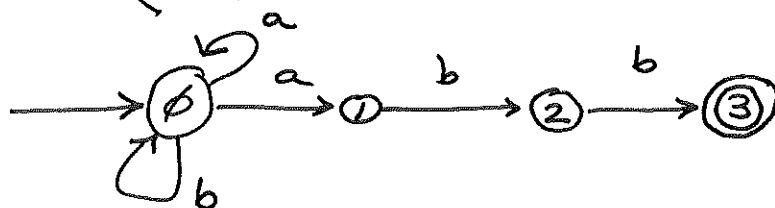
ϵ

transition diagram

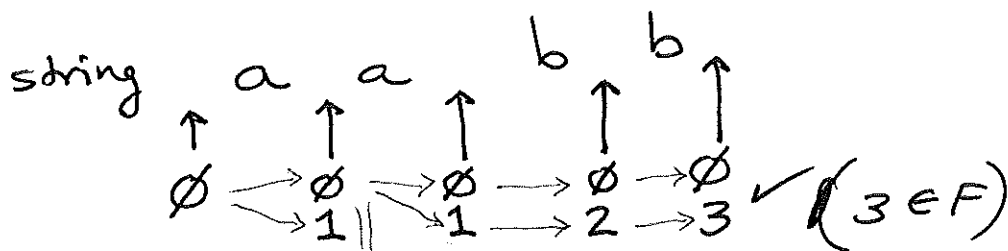
$\textcircled{s} \xrightarrow{a \in \Sigma} \textcircled{t}$ eats a when move

$\textcircled{s} \xrightarrow{\epsilon} \textcircled{t}$ move w/o ~~consume~~ eat a char.

ex: $(a|b)^*abb$



not
DFA



simulate NFA

CMPS-104A

3.14

Dragon 3:148

~~code~~
 $C \leftarrow \{S.\}$

loop {

$a \leftarrow \text{getc}();$

 if ($a == \text{EOF}$) break;

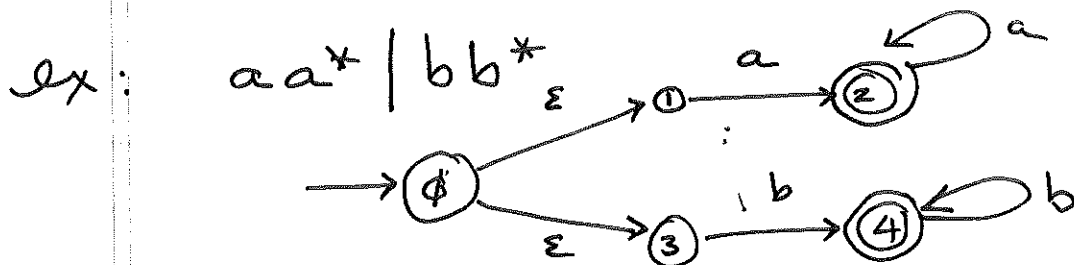
$C \leftarrow \{t' \mid t \in C \wedge \exists t' = \text{move}(t, a)\}$

}

~~return $\exists t \mid t \in C \wedge t \in F$~~

return $\exists t \mid t \in C \wedge t \in F$

$\Phi \xrightarrow{a} \Phi'$

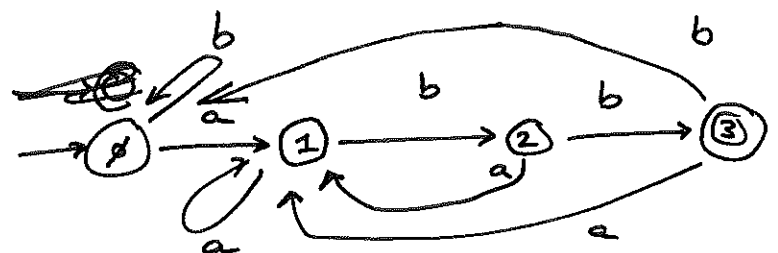


DFA

a DFA is an NFA except:

- no ϵ moves
- each outgoing edge ~~is unique~~ has unique label

ex: $(a/b)^* abb$



\uparrow \uparrow \uparrow \uparrow \uparrow
 \emptyset 1 1 2 3 ✓

simulate DFA

CMPS-104A
3.15

Dragon 3 : 151

```
s = S0
loop {
    a =getc()
    if (a == EOF) break
    s = move(s, a)
}
return s ∈ F
```

NOTE: we don't use dead states in actual scanner.

```
s = S0
loop {
    a =getc()
    t = move(s, a)
    if (!t) break.
    s = t
}
if ungetc(a);
if s ∈ F OK else error -
```

note: we scan until stuck
then exit
don't scan to EOF.

theory DFA stops @ EOF
scanner stops @ end of token

~~ex: strict DFA~~

~~so~~

CMPS-104A
3.16

Dragon 3:152

simulate NFA:

$S = \epsilon\text{closure}(s_0)$

loop {

$a = \text{getc}();$

 if ($c == \text{EOF}$) break;

$S = \epsilon\text{closure}(\text{move}(S, a))$

}

if ($S \cap F == \emptyset$) error else OK

3.7 Regex \rightarrow Automata

Thompson's construction

RE \rightarrow NFA

input: RE r over Σ

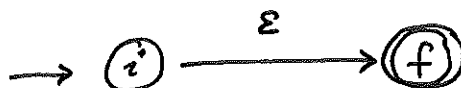
output: NFA accepting $L(r)$

~~Induction~~

Basis

(1)

ϵ



(2)

$a \in \Sigma$



for simple symbols.

Induction

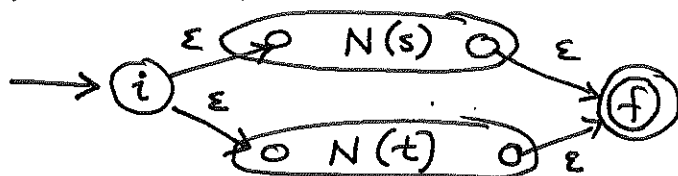
let $N(s)$, $N(t)$ be the NFA for reg s , and t

CMP5-104A

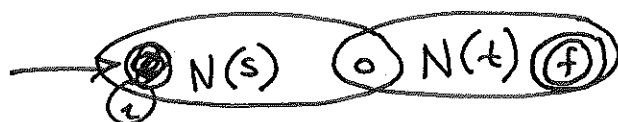
3.17

Dragon 3:160

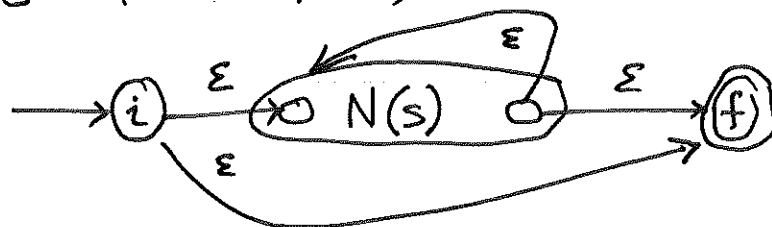
(a) $r = s|t$. then $N(r)$ is



(b) $r = st$. then $N(r)$ is



(c) $r = s^*$ then $N(r)$ is



(d) $r = (s)$ then $L(r) = L(s)$
 $\therefore N(r) = N(s)$

each step consists of at most 2 states

\therefore size ($N(r)$) is $O(n)$ for size (r)
[states] [symbols]

size: $O(|r|)$

construction speed: $O(|r|)$

scanning speed: $O(|r| * |s|)$
string s

Efficiency of NFA simulation

CMPS-104A
3.18

Dragon 3:157

- (1) two sets: old set, new set
- (2) temp bool set already on if in new set.
- (3) array more[s,a] elt is set of states
- (4) perform ε closure on each state
[could precompute]

for each char:
for each state

Conv NFA \rightarrow DFA

alg: subset construction

each state of DFA is a set of states of NFA

worst case: $O(2^n)$

where NFA has n states
usually not so bad.

input: NFA N

output: DFA D accepts same lang

$\text{Eclosure}(s) = \text{set of NFA states reachable from } s \text{ via } \epsilon \text{ trans}$

$\text{Eclosure}(T) = \text{set of NFA states reachable from NFA state } s \in T$
 $= \bigcup_{s \in T} \text{Eclosure}(s)$

move(T, a) set of NFA states
 \exists trans from $a \in \Sigma$
 from $s \in T$

CMPS-104A
 3.19

Dragon 3:153

~~Define~~

subset construction

```

Dstates =  $\epsilon$ closure( $S_0$ )
while ( $\exists$  unmarked  $T \in Dstates$ )
  mark  $T$ 
   $\forall a \in \Sigma$  {
     $U = \epsilon$ closure(move( $T, a$ ))
    if  $U \notin Dstates$  {
      add  $U$  to  $Dstates$ 
      unmarked
    }
     $Dtrans[T, a] = U$ 
  }
}

```

ϵ closure(T)

```

push all states of  $T$  on stk
 $\epsilon$ closure( $T$ ) =  $T$ 
while (stk not empty) {
   $t = \text{pop stk}$ 
  for  $\forall u : \text{edge } t \xrightarrow{\epsilon} u$  {
    if ( $u \notin \epsilon$ closure( $T$ )) :
      add  $u$  to  $\epsilon$ closure( $T$ )
      push  $u$ 
  }
}
}

```

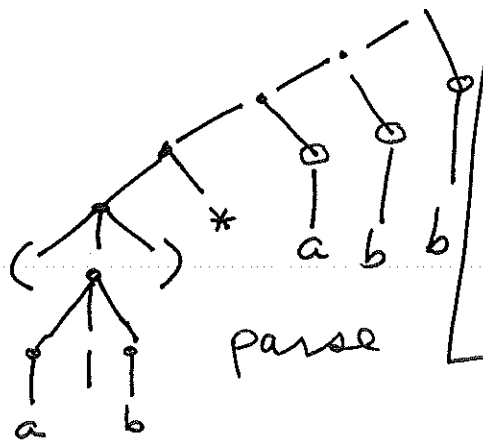
CMPS-104A

3.20

Dragon 3:155

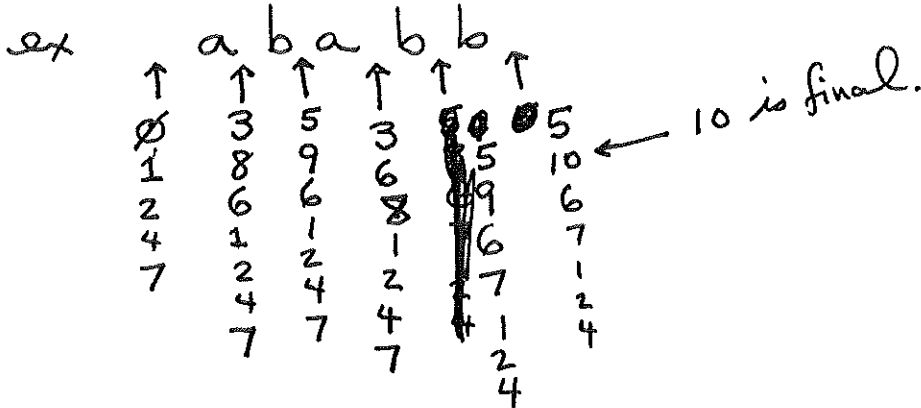
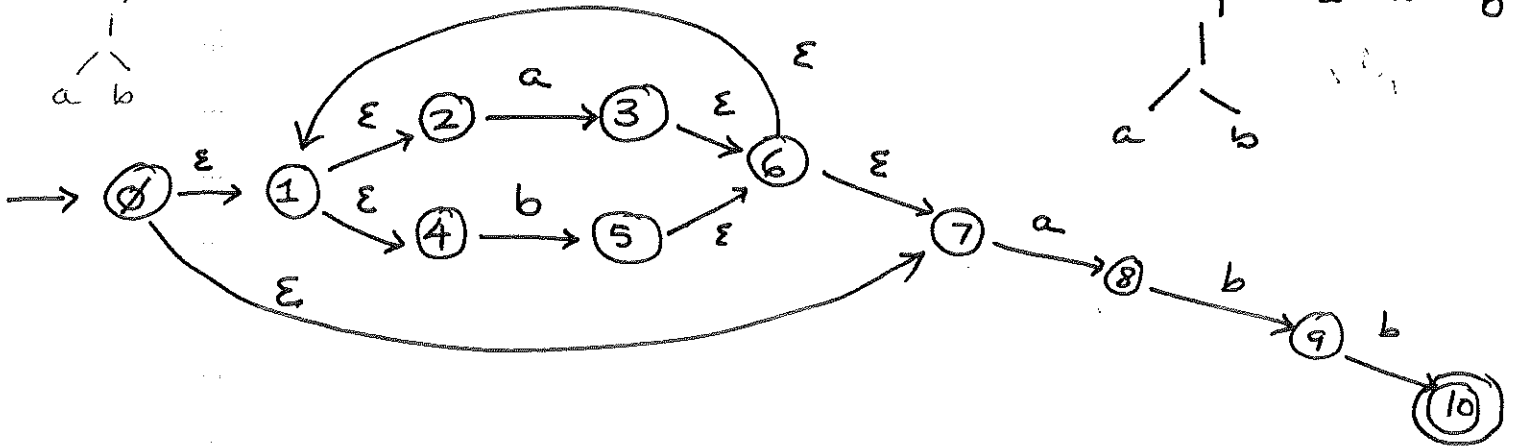
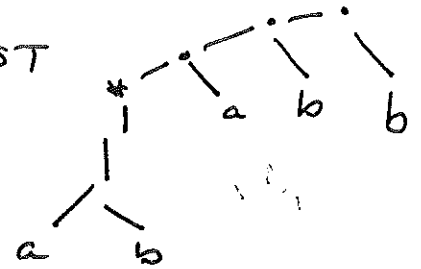
example

$(a|b)^*abb$



parse

AST



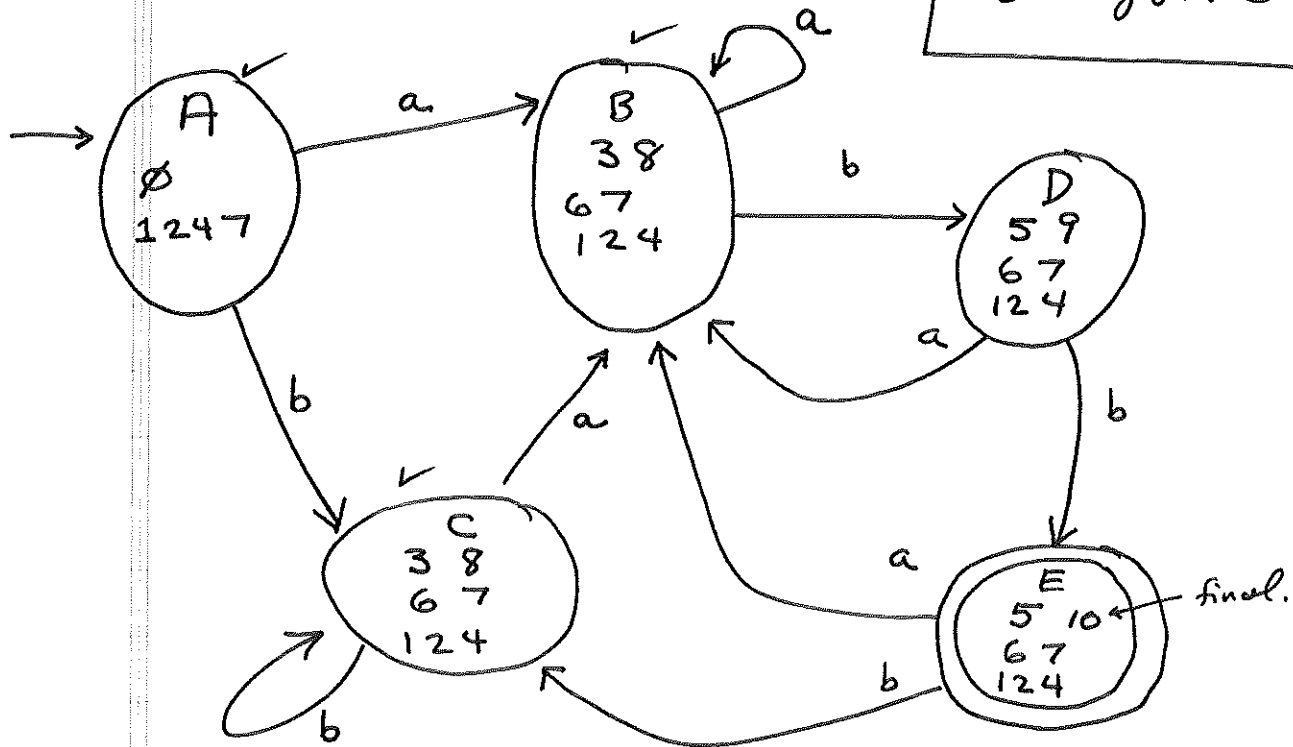
S	ε closure(s)
0	0 1 2 4 7
1	1 2 4
2	2
3	3 1 2 4 6 7
4	4
5	5 1 2 4 6 7
6	6 1 2 4 7
7	7
8	8
9	9
10	10

NFA \rightarrow DFA

CMP5 - 104A

p. 3.21

Dragon 3:155



\uparrow \uparrow \uparrow \uparrow \uparrow \uparrow
 a b a b b
 A B D B D E.

~~Minimizing @ DFA~~

Minimizing a DFA ($D \rightarrow D'$)

CMPS-104A

p. 3.22

Dragon 3:180

1. Start with partition π
two groups, F and $S-F$

S = set of all states

F = set of final states

$S-F$ = " " non final states

2. $\pi' = \pi$

$\forall (G \in \pi) \{$

partition G into subgroups where
states $s, t \in$ same subgroup

iff $\forall a \in \Sigma \quad s, t \xrightarrow{a} \text{same subgroup of } \pi$

replace G in π' by new subgroups

$\}$

3. if $(\pi' = \pi)$ let $\pi_{\text{final}} = \pi$ & goto (4)
else $\pi = \pi'$ and repeat (2)

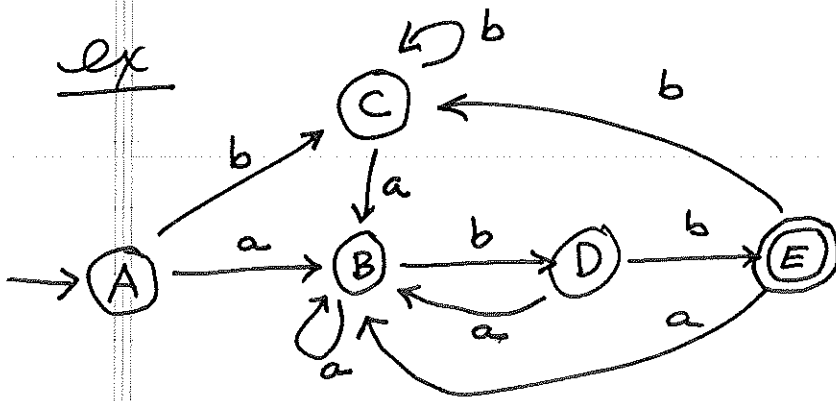
(4) Choose each state of each group $\in \pi_{\text{final}}$
as representative of minimum D'

(a) start s of D' contains start s of D

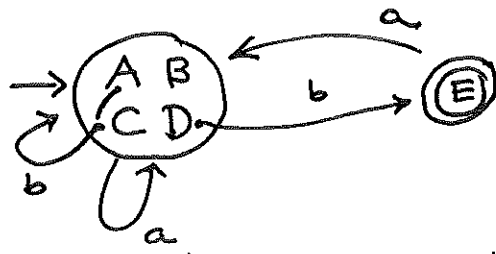
(b) final states of D' are any that have
final of D

(c) Let s be representative of one $G \in \pi_{\text{final}}$
trans $s \xrightarrow{a} t$

ex

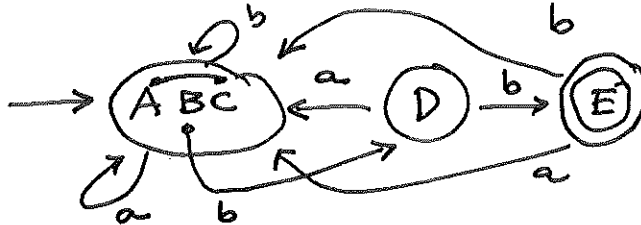


pass 1



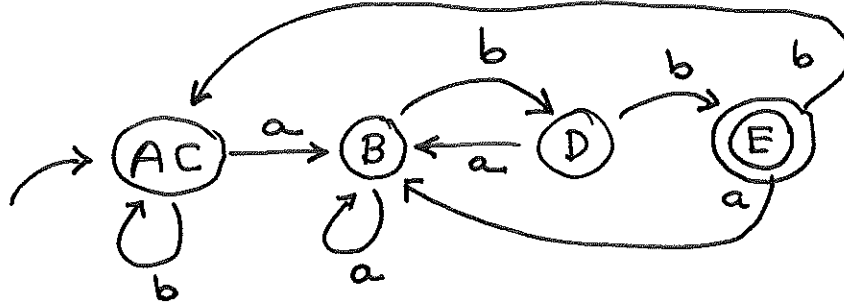
a: agrees.
b: ABC | D

pass 2



a: agree
b: AC | B

pass 3



Efficiency

automaton

space

time

pasting

NFA

$O(|r|)$

$O(|r|)$

$O(|r| \times |x|)$

DFA typical

$O(|r|^3)$

$O(|x|)$

DFA worst

$O(2^{|r|})$

$O(|r|^2 2^{|r|})$

$O(|x|)$

lex - prefer DFA

grep - prefer NFA?

regex r
string x

~~example~~

~~$a(b|c)^*|bb$~~

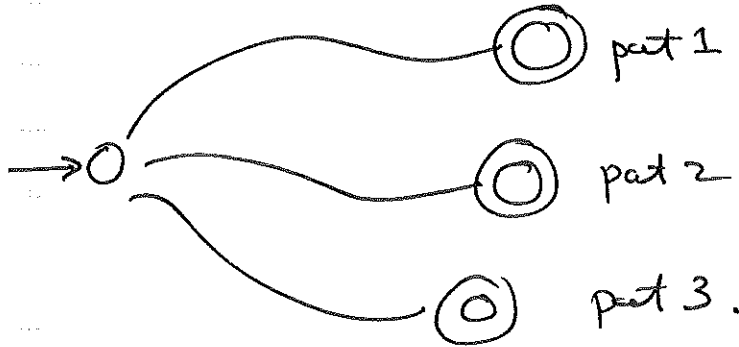
CMP5-104a

p3.24

Dragon 3.

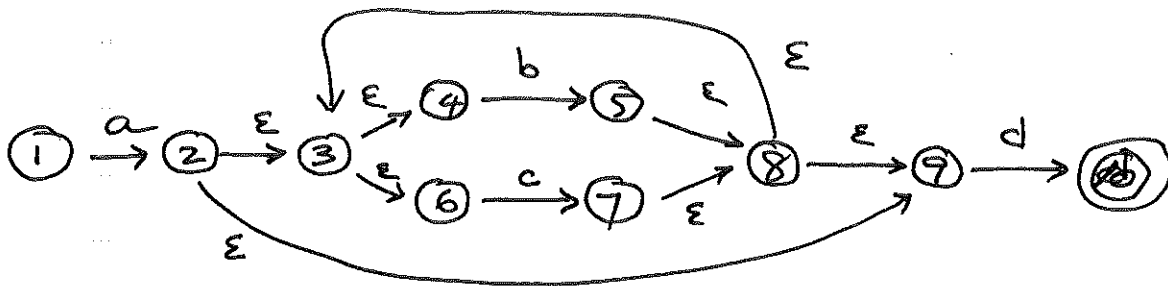
Flex

need sp final state \forall rule



another ex

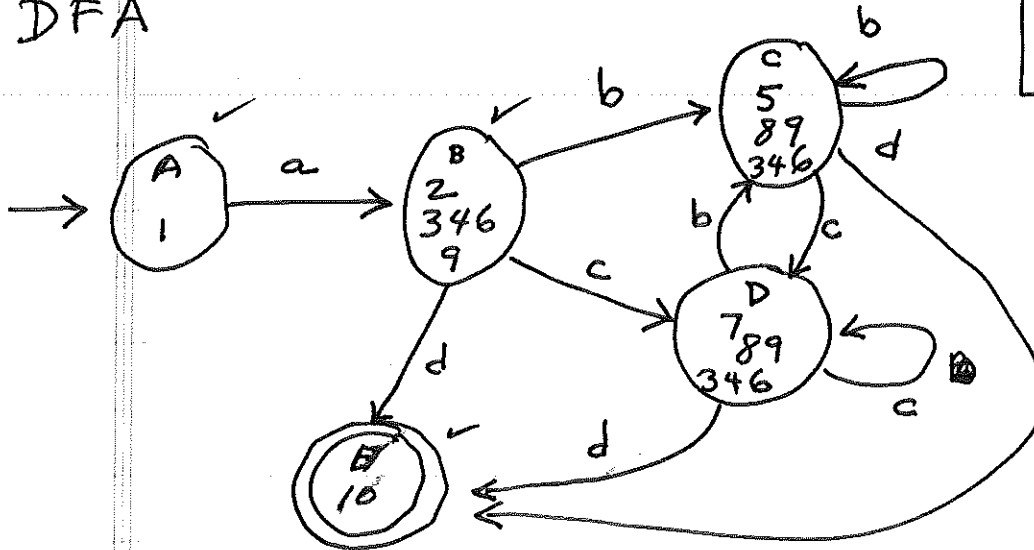
$a(b|c)^*d$



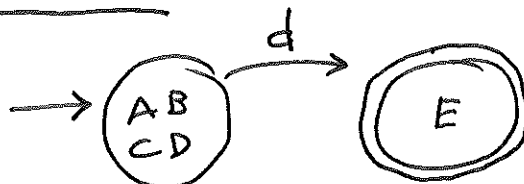
S	ϵ closure(s)
1	1
2	2 3 9 4 6
3	3 4 6
4	4
5	5 8 9 3 4 6
6	6
7	7 8 9 3 4 6
8	8 9 3 4 6
9	9
10	10

DFA

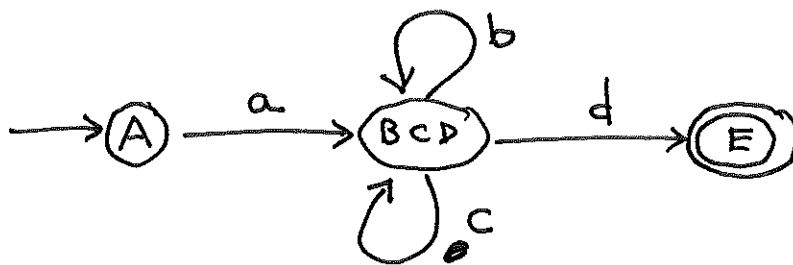
CMPS-104a
p3.25



Minimize

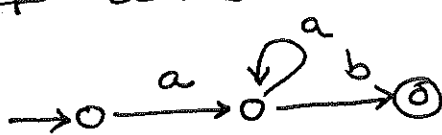


state	A	B	C	D
d:	A	B	C	D
a:	A	B	C	D
b:	A	B	C	D
c:	A	B	C	D

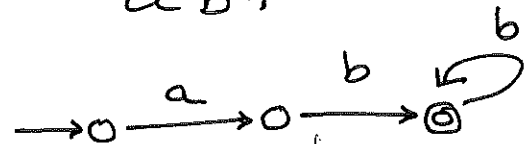


~~a+b~~ a+b

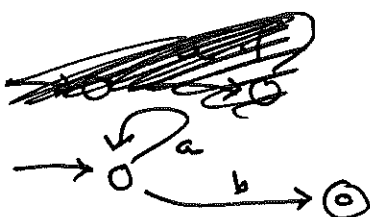
DFA



ab+



a*b



ab*

