

CMPT 379 Compilers

Anoop Sarkar

<http://www.cs.sfu.ca/~anoop>

9/8/05

1

Token Attributes

- Some tokens have attributes
 - T_IDENT “sqrt”
 - T_INTCONSTANT 1
- Other tokens do not
 - T_WHILE
- *Token*=T_IDENT, *Lexeme*=“sqrt”, *Pattern*
- Source code location for error reports

9/8/05

3

Lexical Analysis

- Also called *scanning*, take input program *string* and convert into tokens
- Example:

```
double f = sqrt(-1);
```

T_DOUBLE	("double")
T_IDENT	("f")
T_OP	("=")
T_IDENT	("sqrt")
T_LPAREN	("(")
T_OP	("-")
T_INTCONSTANT	("1")
T_RPAREN	(")")
T_SEP	(";")

9/8/05

2

Lexical errors

- What if user omits the space in “doublef”?
 - No lexical error, single token
T_IDENT(“doublef”) is produced instead of
sequence T_DOUBLE, T_IDENT(“f”)!
- Typically few lexical error types
 - E.g., illegal chars, opened string constants or
comments that are not closed

9/8/05

4

Implementing Lexers: Loop and switch scanners

- Ad hoc scanners
- Big nested switch/case statements
- Lots of `getc()/ungetc()` calls
 - Buffering
- Can be error-prone, use only if
 - Your language's lexical structure is simple
 - Tools don't do what you want
- Changing or adding a keyword is problematic
- Key idea: separate the defn from the implementation
- Problem: we need to reason about patterns and how they can be used to define tokens (recognize strings).

9/8/05

5

Regular Languages

- The set of regular languages: each element is a regular language
- Each regular language is an example of a (formal) language, i.e. a set of strings
e.g. $\{ a^m b^n : m, n \text{ are +ve integers} \}$

9/8/05

7

Formal Languages: Recap

- Symbols: a, b, c
- Alphabet : finite set of symbols $\Sigma = \{a, b\}$
- String: sequence of symbols bab
- Empty string: ϵ Define: $\Sigma^\epsilon = \Sigma \cup \{\epsilon\}$
- Set of all strings: Σ^* cf. *The Library of Babel*, Jorge Luis Borges
- (Formal) Language: a set of strings
 $\{ a^n b^n : n > 0 \}$

9/8/05

6

Regular Languages

- Defining the set of all regular languages:
 - The empty set and $\{a\}$ for all a in Σ^ϵ are regular languages
 - If L_1 and L_2 and L are regular languages, then:
 - $L_1 \cdot L_2 = \{xy \mid x \in L_1 \text{ and } y \in L_2\}$ (concatenation)
 - $L_1 \cup L_2$ (union)
 - $L^* = \bigcup_{i=0}^{\infty} L^i$ (Kleene closure)are also regular languages
 - There are no other regular languages

9/8/05

8

Formal Grammars

- A formal grammar is a concise description of a formal language
- A formal grammar uses a specialized syntax
- For example, a **regular expression** is a concise description of a regular language
 $(a|b)^*abb$: is the set of all strings over the alphabet $\{a, b\}$ which end in abb

9/8/05

9

Regular Expressions: Examples

- Alphabet $\{0, 1\}$
- All strings that represent binary numbers divisible by 4 (but accept 0) $((0|1)^*00)|0$
- All strings that do not contain “01” as a substring 1^*0^*

9/8/05

11

Regular Expressions: Definition

- Every symbol of $\Sigma \cup \{\epsilon\}$ is a regular expression
- If r_1 and r_2 are regular expressions, so are
 - Concatenation: $r_1 r_2$
 - Alternation: $r_1 | r_2$
 - Repetition: r_1^*
- Nothing else is.
 - Grouping re's: e.g. $aalbc$ vs. $((aa)lb)c$

9/8/05

10

Regular Expressions

- To describe all lexemes that form a token as a *pattern*
 - $(0|1|2|3|4|5|6|7|8|9)^+$
- Need decision procedure: to which token does a given sequence of characters belong (if any)?
 - Finite State Automata

9/8/05

12

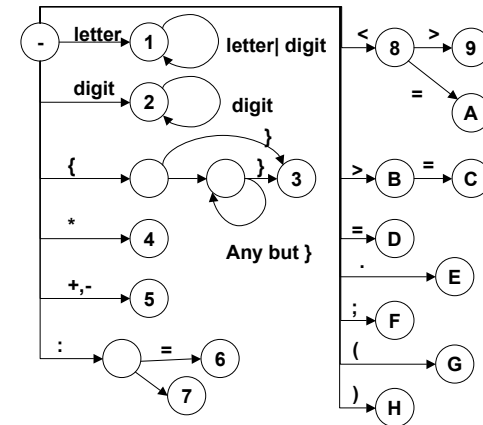
Finite Automata: Recap

- A set of states S
 - One start state q_0 , zero or more final states F
- An alphabet Σ of input symbols
- A transition function:
 - $\delta: S \times \Sigma \Rightarrow S$
- Example: $\delta(1, a) = 2$

9/8/05

13

FA: Pascal Example

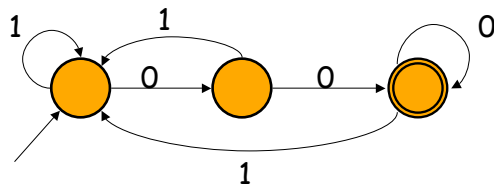


9/8/05

15

Finite Automata: Example

- What regular expression does this automaton accept?



Answer: $(0|1)^*00$

9/8/05

14

Building a Lexical Analyzer

- Token \Rightarrow Pattern
- Pattern \Rightarrow Regular Expression
- Regular Expression \Rightarrow NFA
- NFA \Rightarrow DFA
- DFA \Rightarrow Lexical Analyzer

9/8/05

16

NFAs

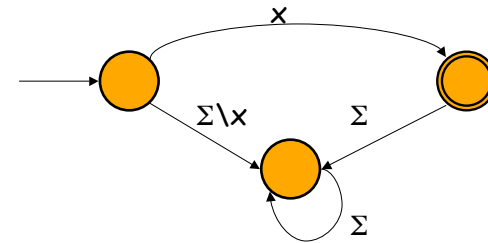
- NFA: like a DFA, except
 - A transition can lead to more than one state, that is, $\delta: S \times \Sigma \Rightarrow 2^S$
 - One state is chosen non-deterministically
 - Transitions can be labeled with ϵ , meaning states can be reached without reading any input, that is, $\delta: S \times \Sigma \cup \{ \epsilon \} \Rightarrow 2^S$

9/8/05

17

Thompson Rule 1

- For each symbol x of the alphabet, there is a NFA that accepts it (include a *sinkhole* state)



9/8/05

19

Thompson's construction

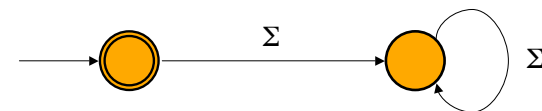
- Converts regexps to NFA
- Five simple rules
 - Symbols
 - Empty String
 - Alternation (r_1 or r_2)
 - Concatenation (r_1 followed by r_2)
 - Repetition (r_1^*)

9/8/05

18

Thompson Rule 2

- There is an NFA that accepts only ϵ

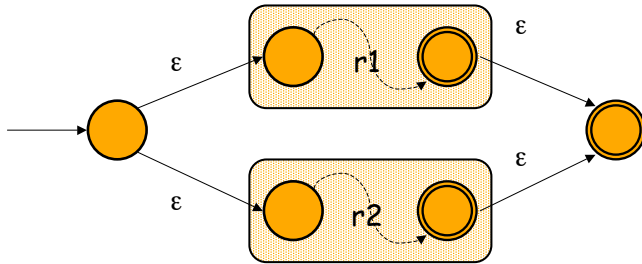


9/8/05

20

Thompson Rule 3

- Given two NFAs for r_1 , r_2 , there is a NFA that accepts $r_1|r_2$

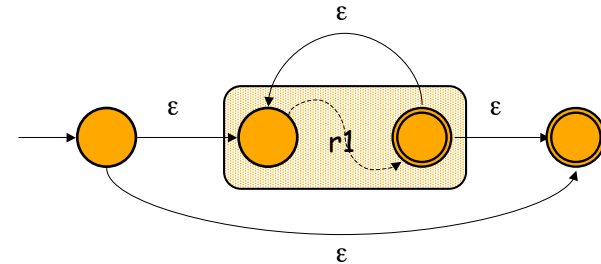


9/8/05

21

Thompson Rule 5

- Given a NFA for r_1 , there is an NFA that accepts r_1^*

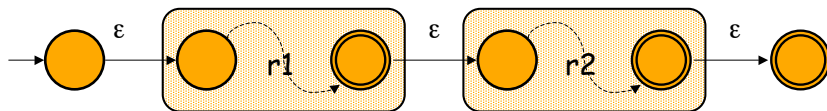


9/8/05

23

Thompson Rule 4

- Given two NFAs for r_1 , r_2 , there is a NFA that accepts r_1r_2



9/8/05

22

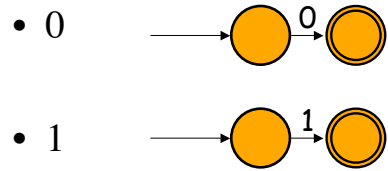
Example

- Set of all binary strings that are divisible by four (include 0 in this set)
- Defined by the regexp: $((0|1)^*00) | 0$
- Apply Thompson's Rules to create an NFA

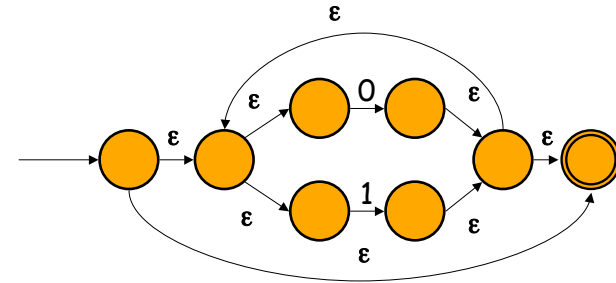
9/8/05

24

Basic Blocks 0 and 1



(this version does not report errors: no *sinkholes*)



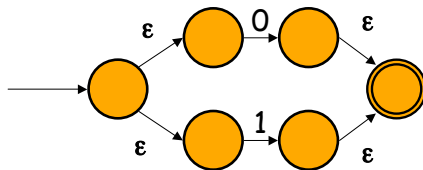
$(0|1)^*$

9/8/05

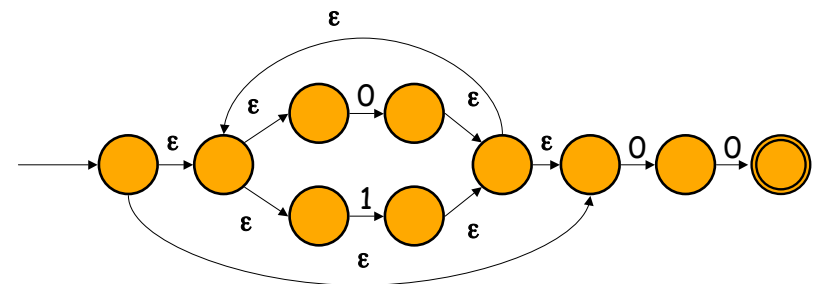
25

9/8/05

27



$0|1$



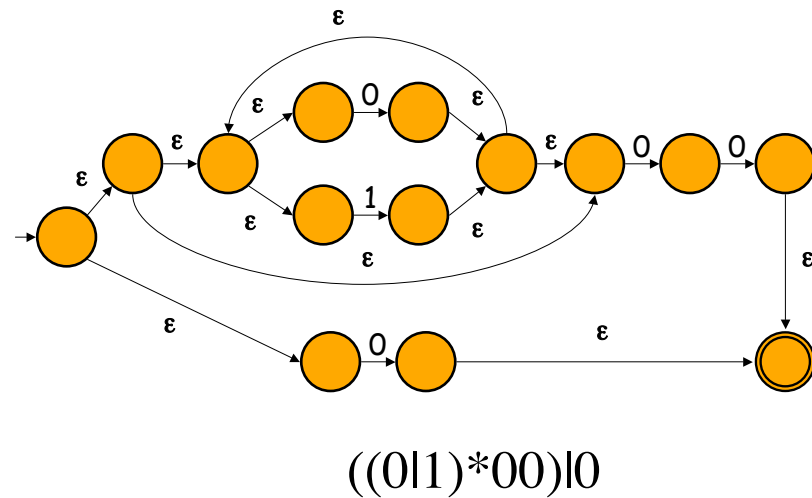
$(0|1)^*00$

9/8/05

26

9/8/05

28



9/8/05

29

Simulating NFAs

- Similar to DFA simulation
- But have to deal with ϵ transitions and multiple transitions on the same input
- Instead of one state, we have to consider *sets* of states
- Simulating NFAs is a problem that is closely linked to converting a given NFA to a DFA

9/8/05

30

NFA to DFA Conversion

- Subset construction
- Idea: subsets of set of all NFA states are *equivalent* and become one DFA state
- Algorithm simulates movement through NFA
- Key problem: how to treat ϵ -transitions?

9/8/05

31

ϵ -Closure

- Start state: q_0
- ϵ -closure(S): S is a set of states
initialize: $S \leftarrow \{q_0\}$
 $T \leftarrow S$
repeat $T' \leftarrow T$
 $T \leftarrow T' \cup [\cup_{s \in T'} \text{move}(s, \epsilon)]$
until $T = T'$

9/8/05

32

ϵ -Closure (T: set of states)

```
push all states in T onto stack
initialize  $\epsilon$ -closure(T) to T
while stack is not empty do begin
  pop t off stack
  for each state u with  $u \in \text{move}(t, \epsilon)$  do
    if  $u \notin \epsilon\text{-closure}(T)$  do begin
      add u to  $\epsilon\text{-closure}(T)$ 
      push u onto stack
    end
  end
end
```

9/8/05

33

NFA Simulation

- Start state: q_0
- Input: c_1, \dots, c_k
 $T \leftarrow \epsilon\text{-closure}(\{q_0\})$
for $i \leftarrow 1$ **to** k
 $T \leftarrow \text{DFAedge}(T, c_i)$

9/8/05

35

NFA Simulation

- After computing the ϵ -closure move, we get a set of states
- On some input extend all these states to get a new set of states

$$\text{DFAedge}(T, c) = \epsilon\text{-closure}(\cup_{q \in T} \text{move}(q, c))$$

9/8/05

34

Conversion from NFA to DFA

- Conversion method closely follows the NFA simulation algorithm
- Instead of simulating, we can collect those NFA states that behave identically on the same input
- Group this set of states to form one state in the DFA

9/8/05

36

Subset Construction

```

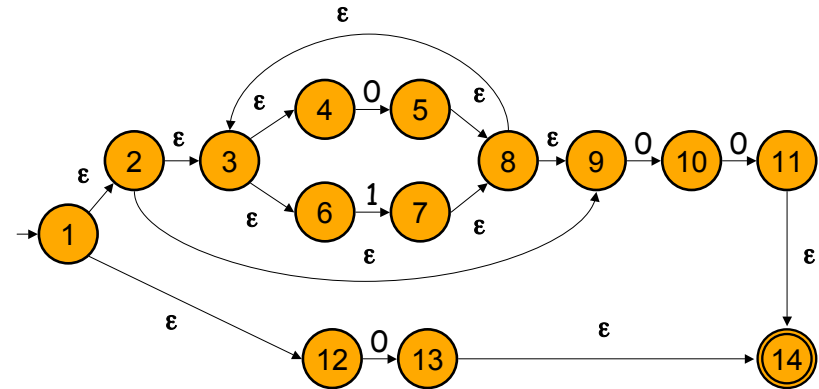
add  $\epsilon$ -closure( $q_0$ ) to  $Dstates$  unmarked
while  $\exists$  unmarked  $T \in Dstates$  do begin
  mark  $T$ ;
  for each symbol  $c$  do begin
     $U := \epsilon$ -closure(move( $T, c$ ));
    if  $U \notin Dstates$  then
      add  $U$  to  $Dstates$  unmarked
       $Dtrans[d, c] := U$ ;
    end
  end
end

```

9/8/05

37

Example: subset construction



9/8/05

39

Subset Construction

```

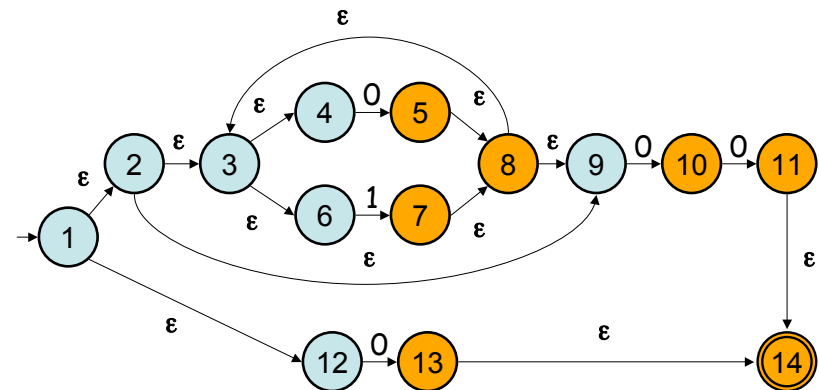
states[0] =  $\epsilon$ -closure( $\{q_0\}$ )
 $p = j = 0$ 
while  $j \leq p$  do begin
  for each symbol  $c$  do begin
     $e = DFAedge(states[j], c)$ 
    if  $e = states[i]$  for some  $i \leq p$ 
    then  $Dtrans[j, c] = i$ 
    else  $p = p + 1$ 
       $states[p] = e$ 
       $Dtrans[j, c] = p$ 
    end
   $j = j + 1$ 
end

```

9/8/05

38

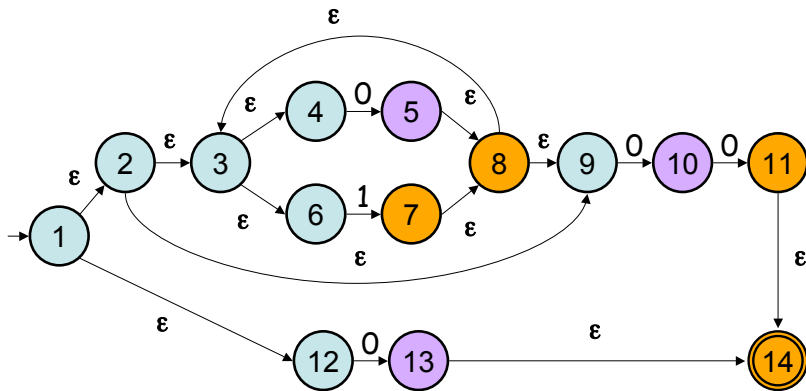
ϵ -closure(q_0)



9/8/05

40

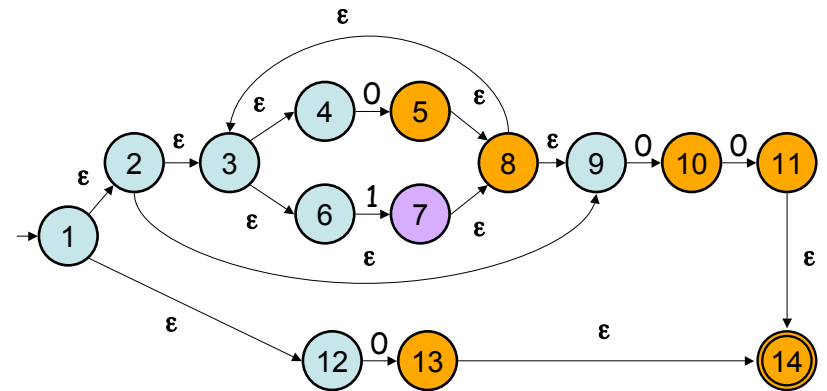
$\text{move}(\epsilon\text{-closure}(q_0), 0)$



9/8/05

41

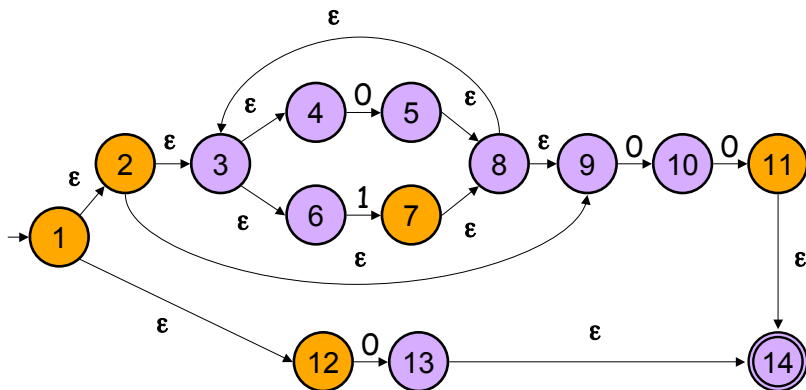
$\text{move}(\epsilon\text{-closure}(q_0), 1)$



9/8/05

43

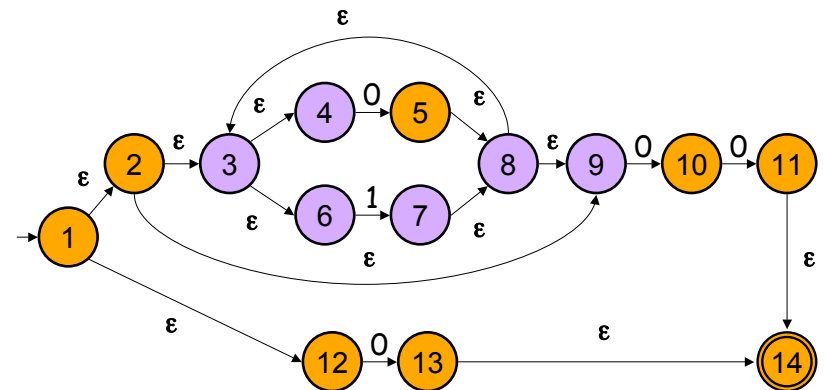
$\epsilon\text{-closure}(\text{move}(\epsilon\text{-closure}(q_0), 0))$



9/8/05

42

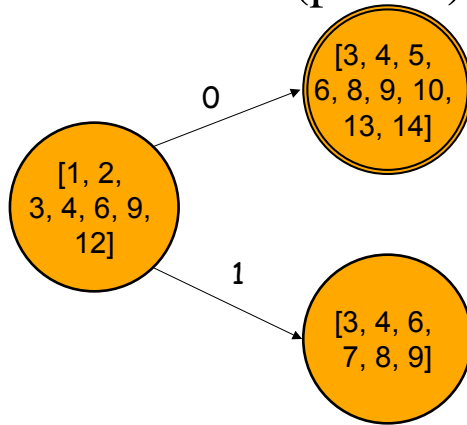
$\epsilon\text{-closure}(\text{move}(\epsilon\text{-closure}(q_0), 1))$



9/8/05

44

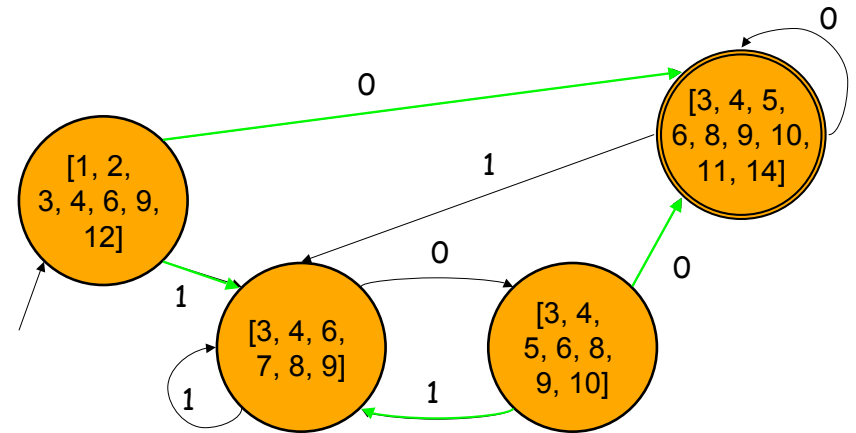
DFA (partial)



9/8/05

45

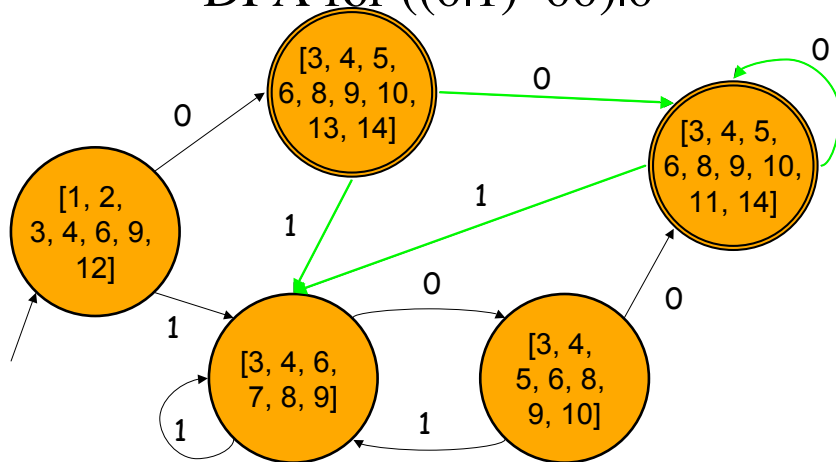
Minimization (I)



9/8/05

47

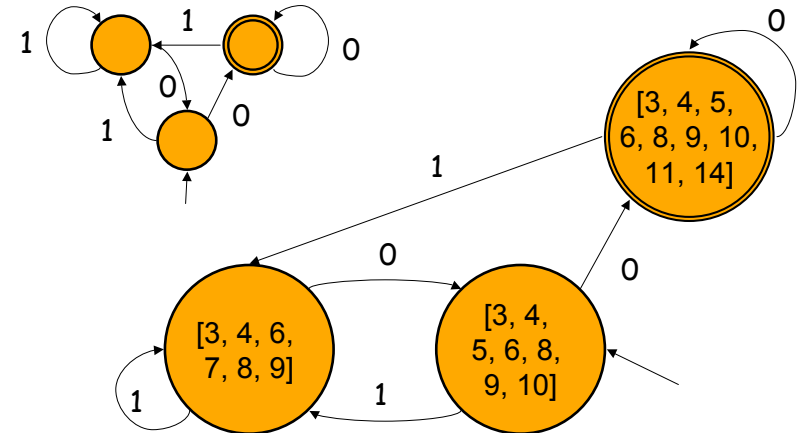
DFA for $((0|1)^*00)|0$



9/8/05

46

Minimization (II)

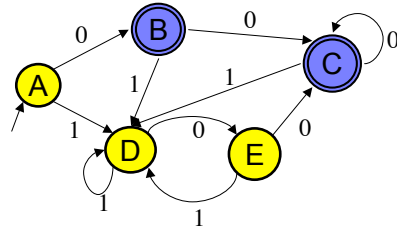


9/8/05

48

Minimization of DFAs

- Algorithm for minimizing the number of states in a DFA
- Step 1: partition states into 2 groups: accepting and non-accepting



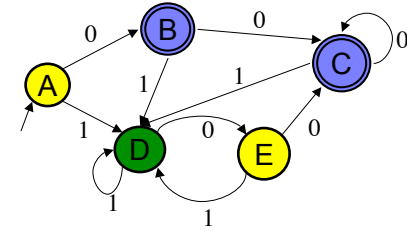
9/8/05

49

Minimization of DFAs

- Step 3: if a sub-group does not obey P split up the group into a separate group
- Go back to step 2. If no further sub-groups emerge then continue to step 4

A, 0: blue
A, 1: green
E, 0: blue
E, 1: green
D, 0: yellow
D, 1: green



B, 0: blue
B, 1: green
C, 0: blue
C, 1: green

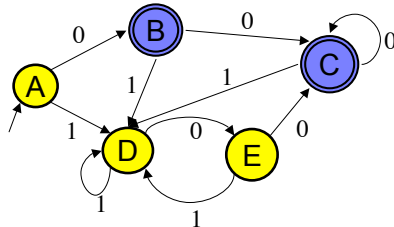
9/8/05

51

Minimization of DFAs

- Step 2: in each group, find a sub-group of states having property P
- P: The states have transitions on each symbol (in the alphabet) to the *same* group

A, 0: blue
A, 1: yellow
E, 0: blue
E, 1: yellow
D, 0: yellow
D, 1: yellow



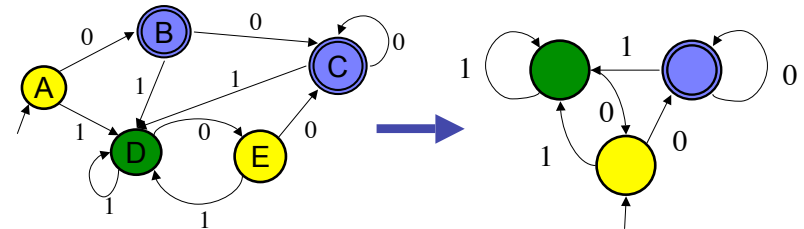
B, 0: blue
B, 1: yellow
C, 0: blue
C, 1: yellow

9/8/05

50

Minimization of DFAs

- Step 4: each group becomes a state in the minimized DFA
- Transitions to individual states are mapped to a single state representing the group of states



9/8/05

52

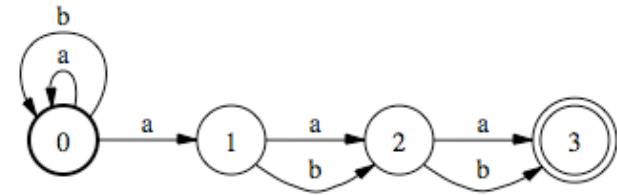
NFA to DFA

- Subset construction converts NFA to DFA
- Complexity:
 - in programs we measure time complexity in number of steps
 - For FSAs, we measure complexity in terms of the number of states

9/8/05

53

NFA to DFA



9/8/05

55

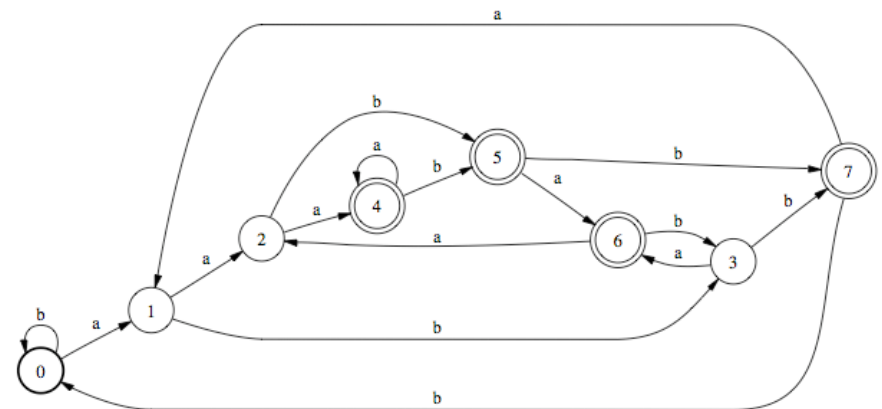
NFA to DFA

- Problem: An n state NFA can sometimes become a 2^n state DFA, an exponential increase in complexity
 - Try the subset construction on NFA built for the regexp A^*aA^{n-1} where A is the regexp (alb)
- Minimization can reduce the number of states
- But minimization requires determinization

9/8/05

54

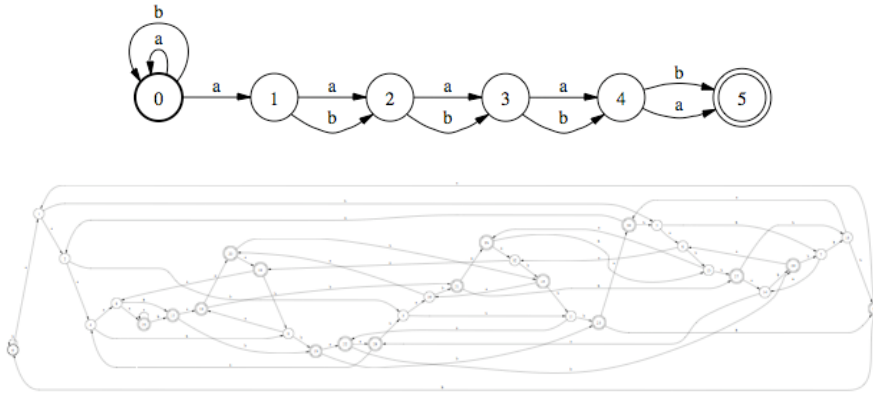
NFA to DFA



9/8/05

56

NFA to DFA

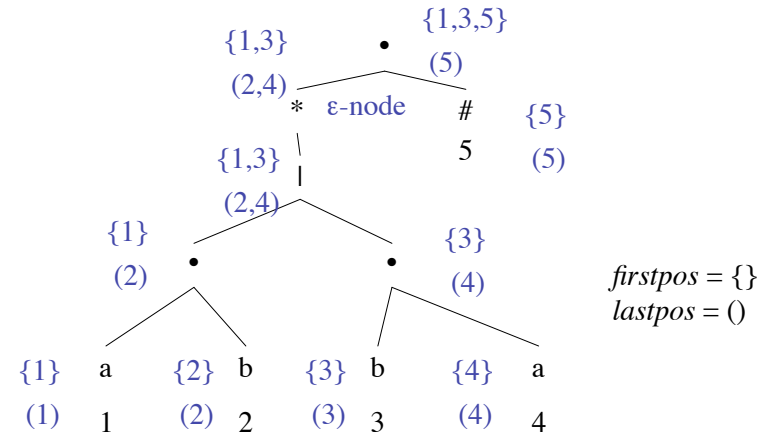


9/8/05

$2^5 = 32$ states

57

Regex to DFA: $(ab|ba)^*\#$



9/8/05

59

NFA vs. DFA in the wild

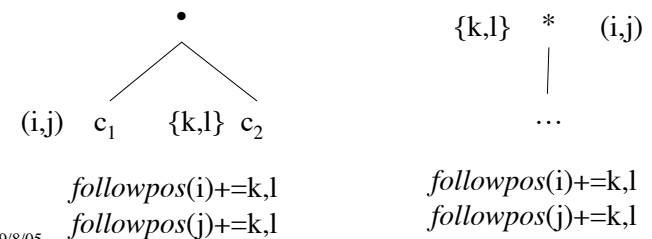
Engine Type	Programs
DFA	<i>awk</i> (most versions), <i>egrep</i> (most versions), <i>flex</i> , <i>lex</i> , MySQL, Procmail
Traditional NFA	GNU <i>Emacs</i> , Java, <i>grep</i> (most versions), <i>less</i> , <i>more</i> , .NET languages, PCRE library, Perl, PHP (pcre routines), Python, Ruby, <i>sed</i> (most versions), <i>vi</i>
POSIX NFA	<i>mawk</i> , MKS utilities, GNU <i>Emacs</i> (when requested)
Hybrid NFA/DFA	GNU <i>awk</i> , GNU <i>grep/egrep</i> , Tcl

9/8/05

58

Regex to DFA: *followpos*

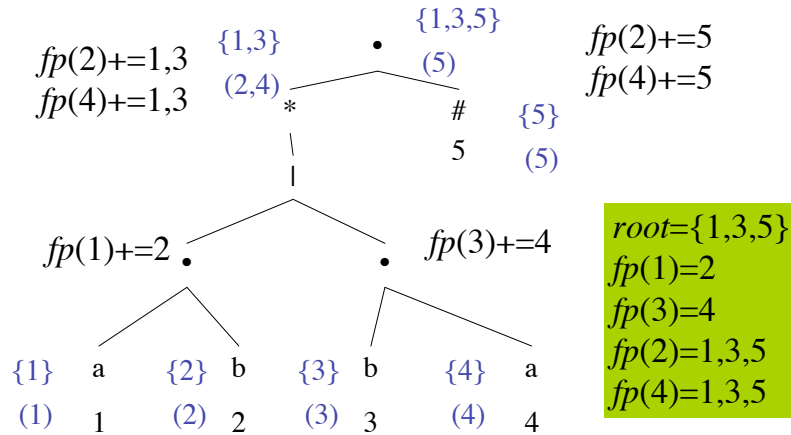
- $followpos(p)$ tells us which positions can follow a position p
- There are two rules that use the $firstpos \{ \}$ and $lastpos ()$ information



9/8/05

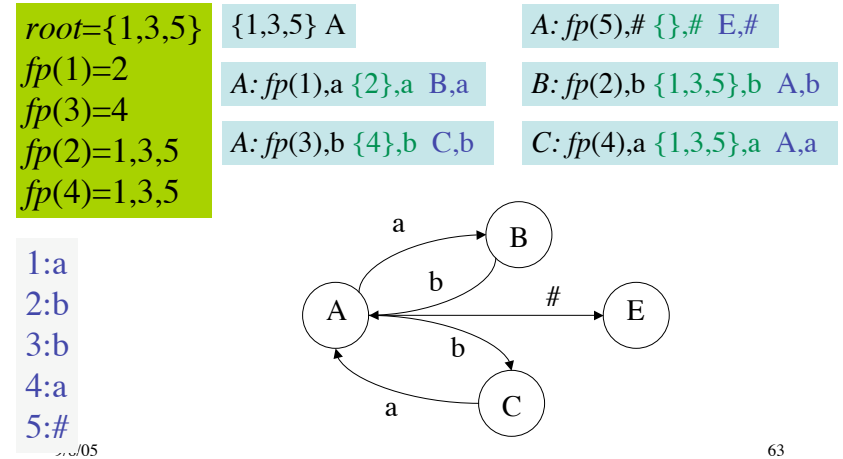
60

Regexp to DFA: $(ab|ba)^* \#$



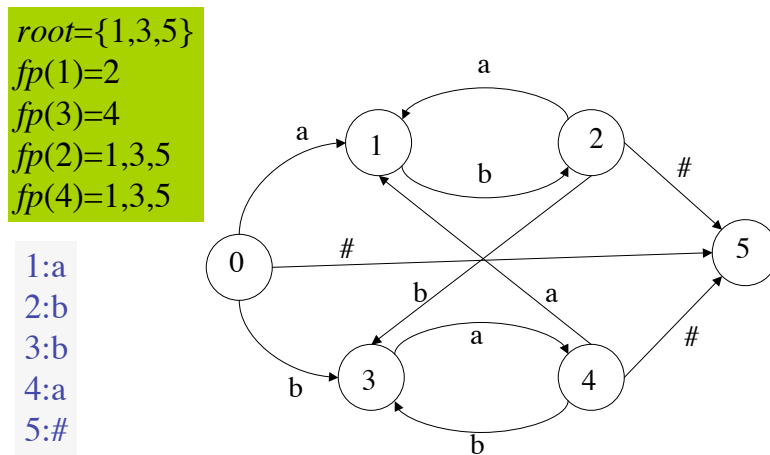
61

Regex to DFA: $(ab|ba)^*\#$



63

Regexp to DFA: $(ab|ba)^* \#$



62

Equivalence of Regexps

- $(R|S)|T == R|(S|T) == R|S|T$
- $(RS)T == R(ST)$
- $(R|S) == (S|R)$
- $R^*R^* == (R^*)^* == R^* == RR^*|\epsilon$
- $R^{**} == R^*$
- $(R|S)T = RT|ST$
- $R(S|T) == RS|RT$
- $(R|S)^* == (R^*S^*)^* == (R^*S)^*R^* == (R^*|S^*)^*$
- $RR^* == R^*R$
- $(RS)^*R == R(SR)^*$
- $R = R|R = R|\epsilon$

64

Equivalence of Regexps

- $0(10)^*1|(01)^*$
- $(01)(01)^*(01)^*$
- $(01)(01)^*(01)(01)^*\epsilon$
- $(01)(01)^*\epsilon$
- $(01)^*$
- $(RS)^*R == R(SR)^*$
- $RS == (RS)$
- $R^* == RR^*\epsilon$
- $R == R|R$
- $R^* == RR^*|\epsilon$

9/8/05

65

Lexical Analyzer using DFAs

- The DFA recognizer has to find the *longest match* for a token
 - e.g. `<print>` and not `<pr>`, `<int>`
- If two patterns match the same token, pick the one that was listed earlier in R
 - e.g. prefer final state (in the original NFA) of r_2 over r_3

9/8/05

67

Lexical Analyzer using DFAs

- Each token is defined using a regexp r_i
- Merge all regexps into one big regexp
 - $R = (r_1 | r_2 | \dots | r_n)$
- Convert R to an NFA, then DFA, then minimize
 - remember orig NFA final states with each DFA state

9/8/05

66

Lexical Analyzer using DFAs

- Alternative method:
 - Organize all the DFAs for each token in an ordered list
 - For input i_1, i_2, \dots, i_n run all DFAs until some reach a final state (pick the longest match for each DFA)
 - Pick the token for which some DFA could read the longest match in the input,
 - e.g. prefer DFA #8 over all others because it read the input until i_{30} and none of the other DFAs reached i_{30}
 - If two DFAs reach the same input character then pick the one that is listed first in the ordered list

9/8/05

68

Implementing DFAs

- 2D array storing the transition table
- Adjacency list, more space efficient but slower
- Merge two ideas: array structures used for sparse tables like DFA transition tables
 - base & next arrays: Tarjan and Yao, 1979
 - Dragon book (default+base & next+check)

9/8/05

69

Implementing DFAs

	a	b	c	d
0	-	1	-	2
1	1	-	1	-
2	1	2	1	-

base

0	2
1	4
2	0

		-	1	-	2		
				1	-	1	-
1	2	1	-				
1	2	1	1	1	2	1	-
0	1	2	3	4	5	6	7
2	2	2	0	1	0	1	-

next

check

nextstate(s, x) :

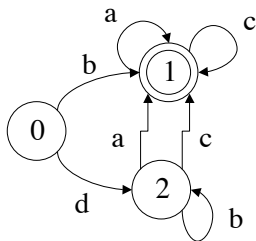
$L := \text{base}[s] + x$

return next[L] **if** check[L] eq s

9/8/05

71

Implementing DFAs



	a	b	c	d
0	-	1	-	2
1	1	-	1	-
2	1	2	1	-

9/8/05

70

Implementing DFAs

	a	b	c	d
0	-	1	-	2
1	1	-	1	-
2	1	2	1	-

base

0	2
1	3
2	0

default

		-	1	-	2		
				1	-	1	-
-	2	-	-				
-	2	1	1	2	1	-	
0	1	2	3	4	5	6	
-	2	0	1	0	1	-	

next

check

nextstate(s, x) :

$L := \text{base}[s] + x$

return next[L] **if** check[L] eq s

else return *nextstate*(default[s], x)

9/8/05

72

Summary

- Token \Rightarrow Pattern
- Pattern \Rightarrow Regular Expression
- Regular Expression \Rightarrow NFA
 - Thompson's Rules
- NFA \Rightarrow DFA
 - Subset construction
- DFA \Rightarrow minimal DFA
 - Minimization

\Rightarrow Lexical Analyzer (multiple patterns)