

# *Java II*

## *Finite Automata I*

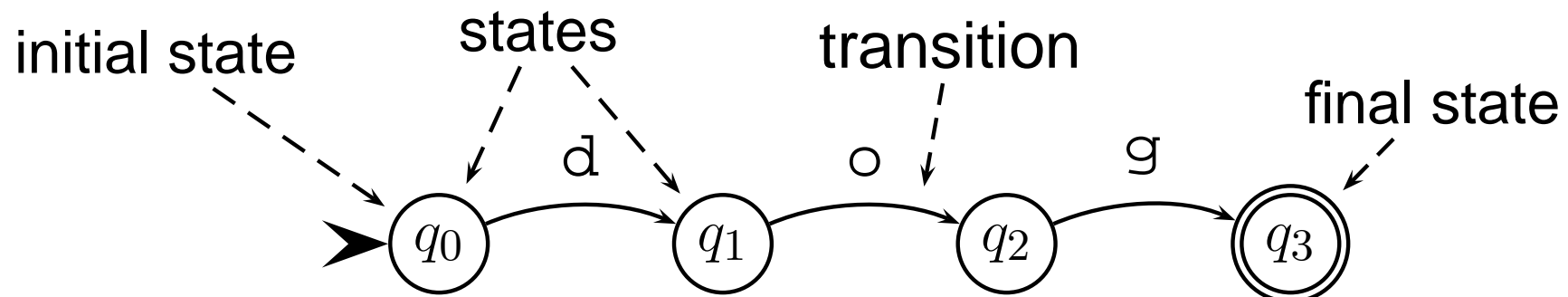
Bernd Kiefer

`Bernd.Kiefer@dfki.de`

Deutsches Forschungszentrum für künstliche Intelligenz

- We already learned about Java's regular expression functionality
- Now we get to know the machinery behind
  - `Pattern` and
  - `Matcher` classes
- Compiling a regular expression into a `Pattern` object produces a *Finite Automaton*
- This automaton is then used to perform the matching tasks
- We will see how to construct a finite automaton that *recognizes* an input string, i.e., tries to find a full match

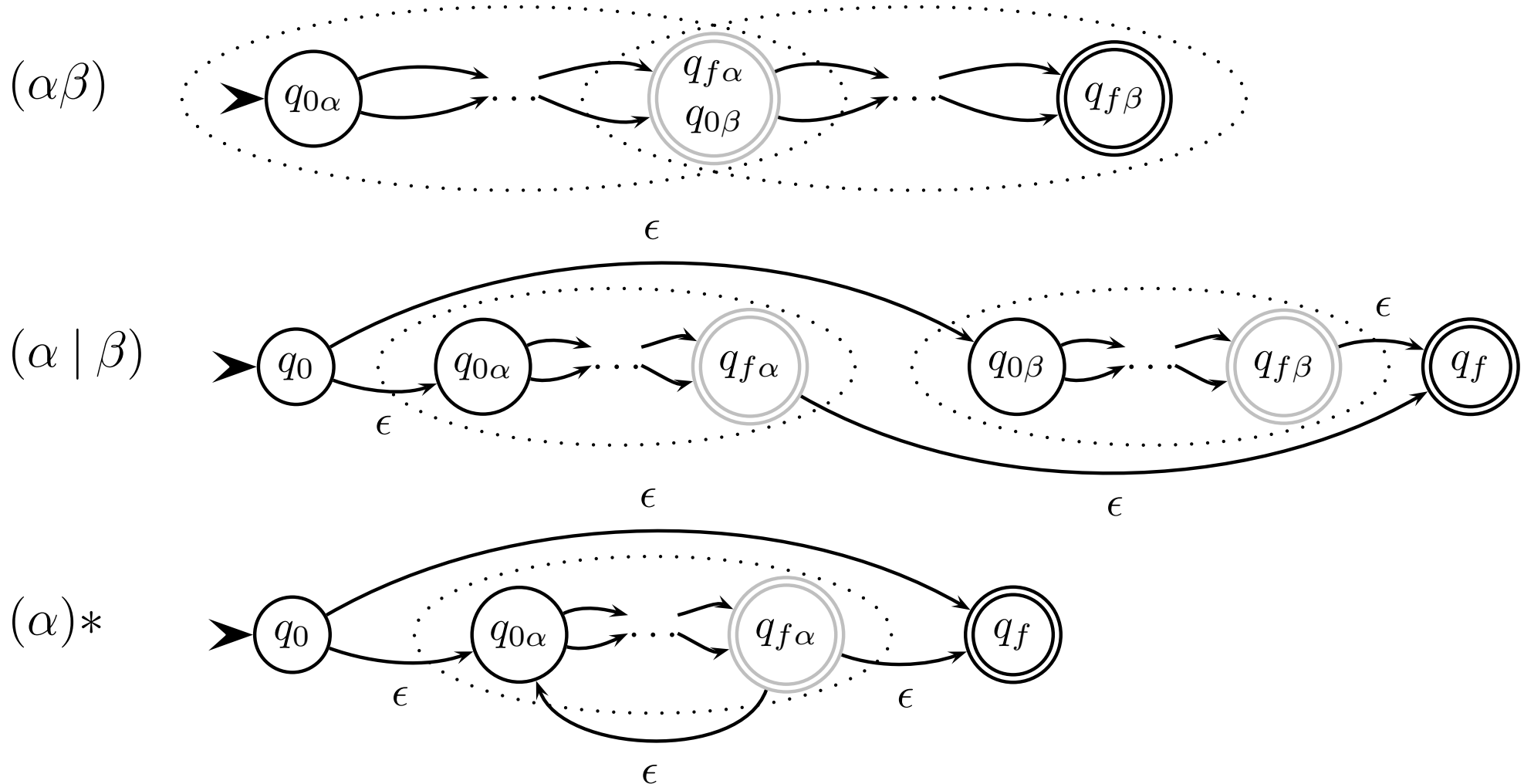
- A finite automaton (FA) is a tuple  $A = \langle Q, \Sigma, \delta, q_0, F \rangle$ 
  - $Q$  a finite non-empty set of states
  - $\Sigma$  a finite alphabet of input letters
  - $\delta$  a (total) transition function  $Q \times \Sigma \longrightarrow Q$
  - $q_0 \in Q$  the initial state
  - $F \subseteq Q$  the set of final (accepting) states
- Transition graphs (diagrams):



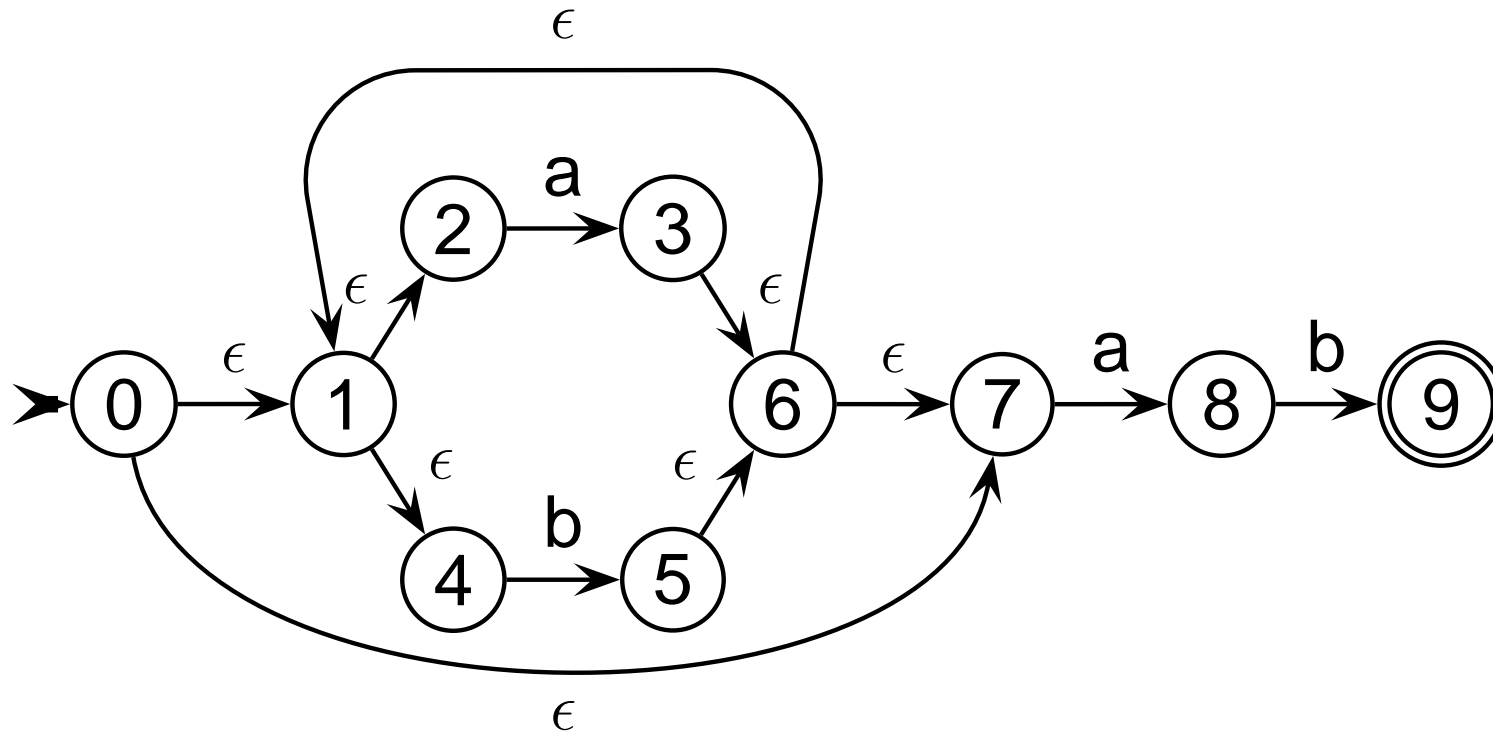
- A finite automaton *accepts* a given input string  $s$  if there is a sequence of states  $p_1, p_2, \dots, p_{|s|} \in Q$  such that
  1.  $p_1 = q_0$ , the start state
  2.  $\delta(p_i, s_i) = p_{i+1}$ , where  $s_i$  is the  $i$ -th character in  $s$
  3.  $p_{|s|} \in F$ , i.e., a final state
- A string is successfully *matched* if we have found the appropriate sequence of states
- Imagine the string on an input tape with a pointer that is advanced when using a  $\delta$  transition
- The set of strings accepted by an automaton is the accepted *language*, analogous to regular expressions

- in the definition of automata,  $\delta$  was a total function  $\Rightarrow$  given an input string, the path through the automaton is uniquely determined
- those automata are therefore called *deterministic*
- for *nondeterministic FA*,  $\delta$  is a *transition relation*
- $\delta : Q \times \Sigma \cup \{\epsilon\} \longrightarrow \mathcal{P}(Q)$ , where  $\mathcal{P}(Q)$  is the powerset of  $Q$
- allows transitions from one state into several states with the same input symbol
- need not be total
- can have transitions labeled  $\epsilon$  (not in  $\Sigma$ ), which represents the empty string

Construct nondeterministic automata from regular expressions

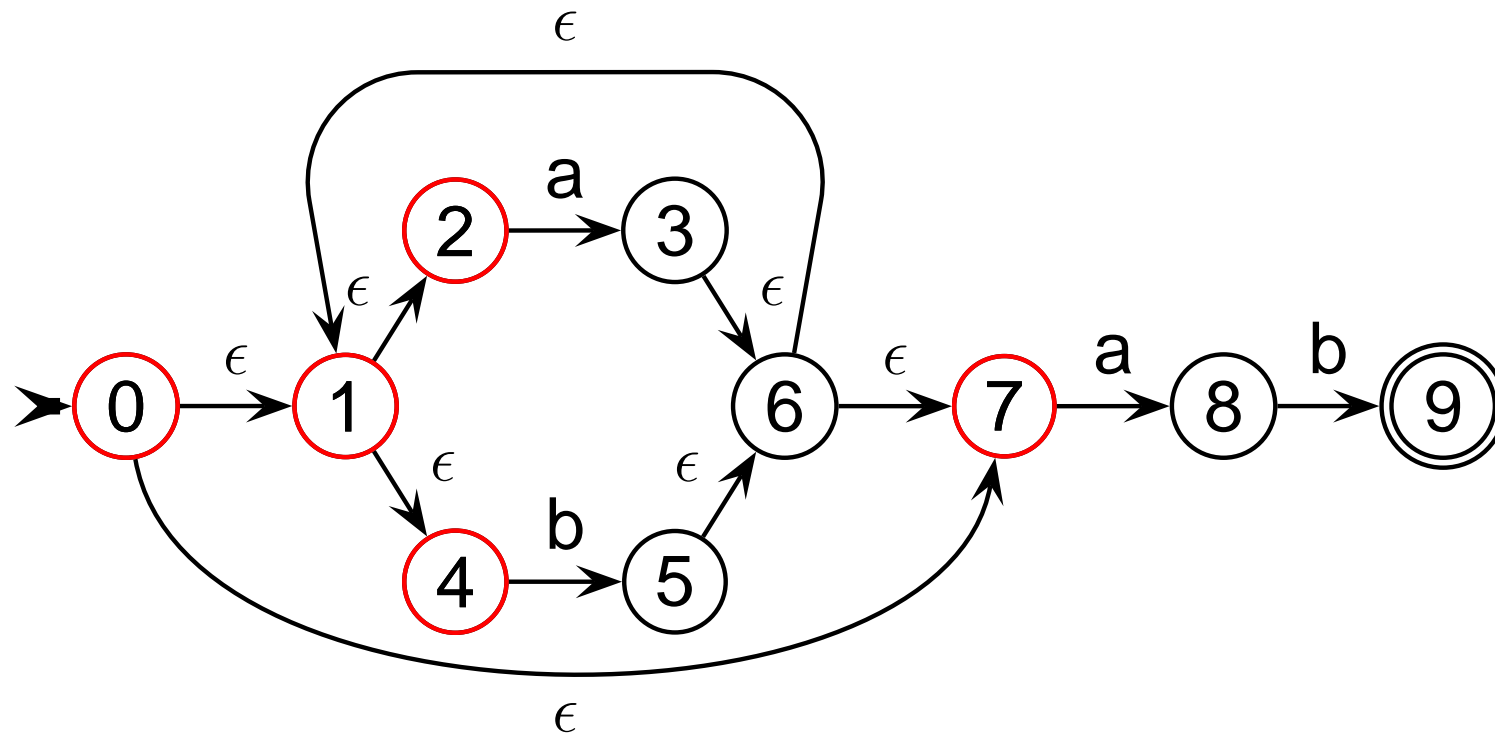


- Traversing a DFA is easy given the input string: the path is uniquely determined
- In contrast, traversing an NFA requires keeping track of a set of (current) states, starting with the set  $\{q_0\}$
- Processing the next input symbol means taking all possible outgoing transitions from this set and collecting the new set
- From every NFA, an equivalent DFA (one which does accept the same language), can be computed
- Basic Idea: track the subsets that can be reached for every possible input

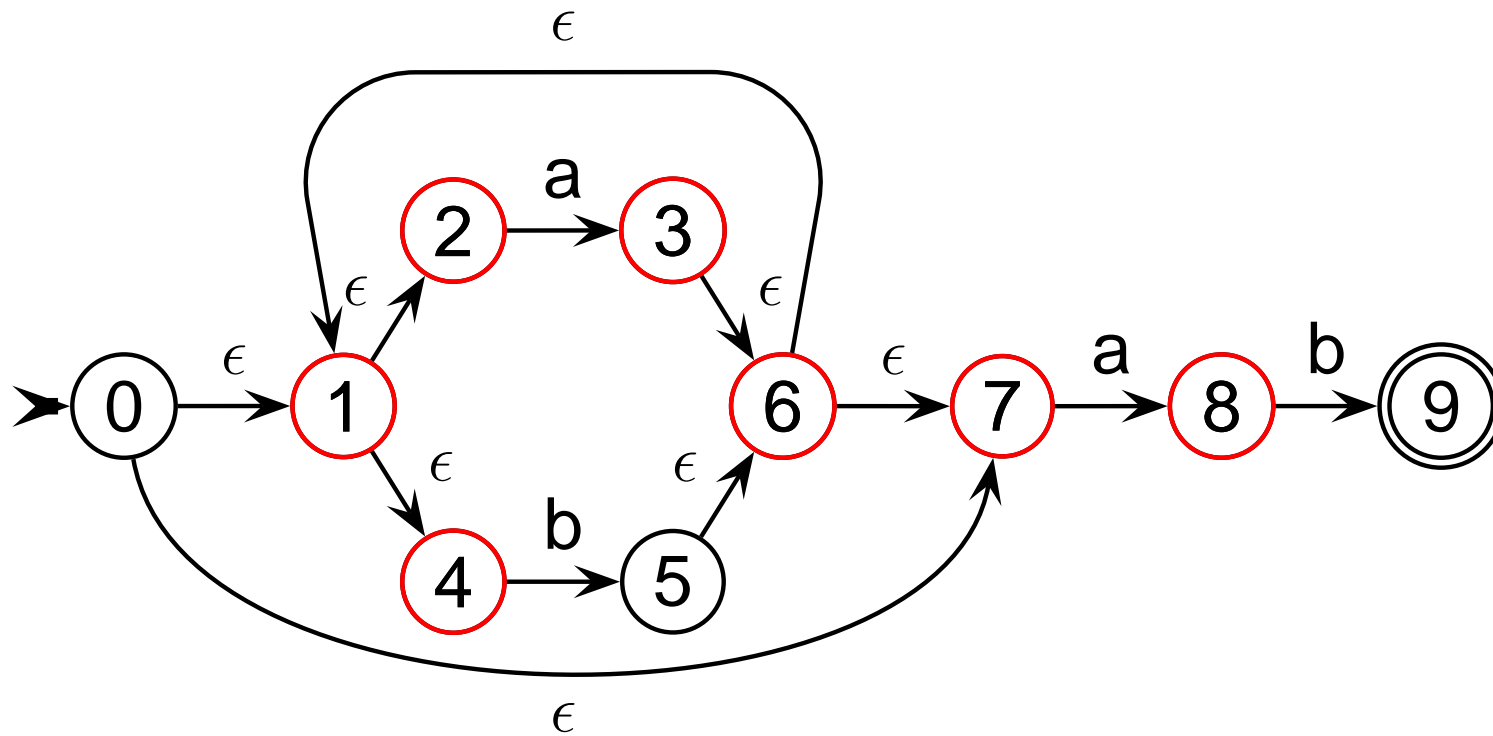


abab

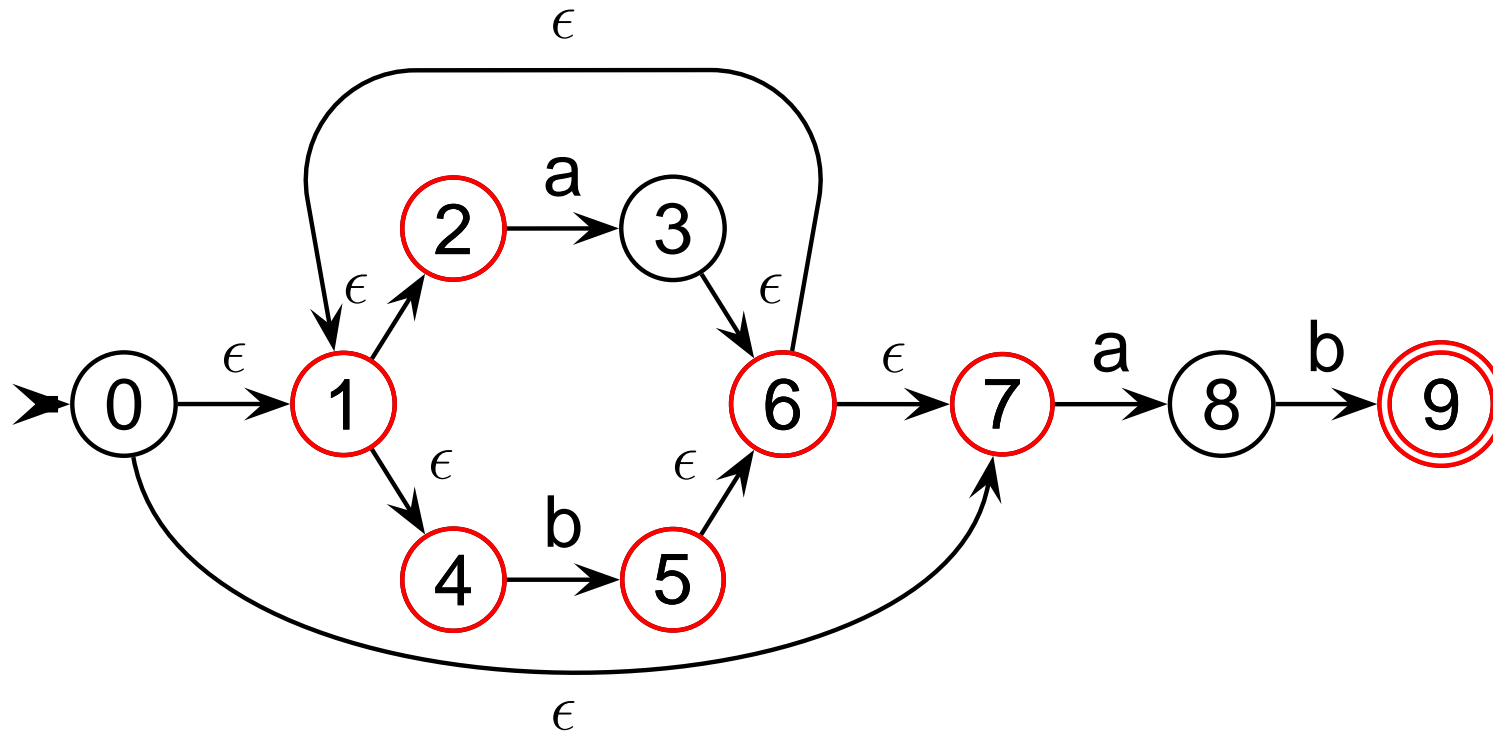




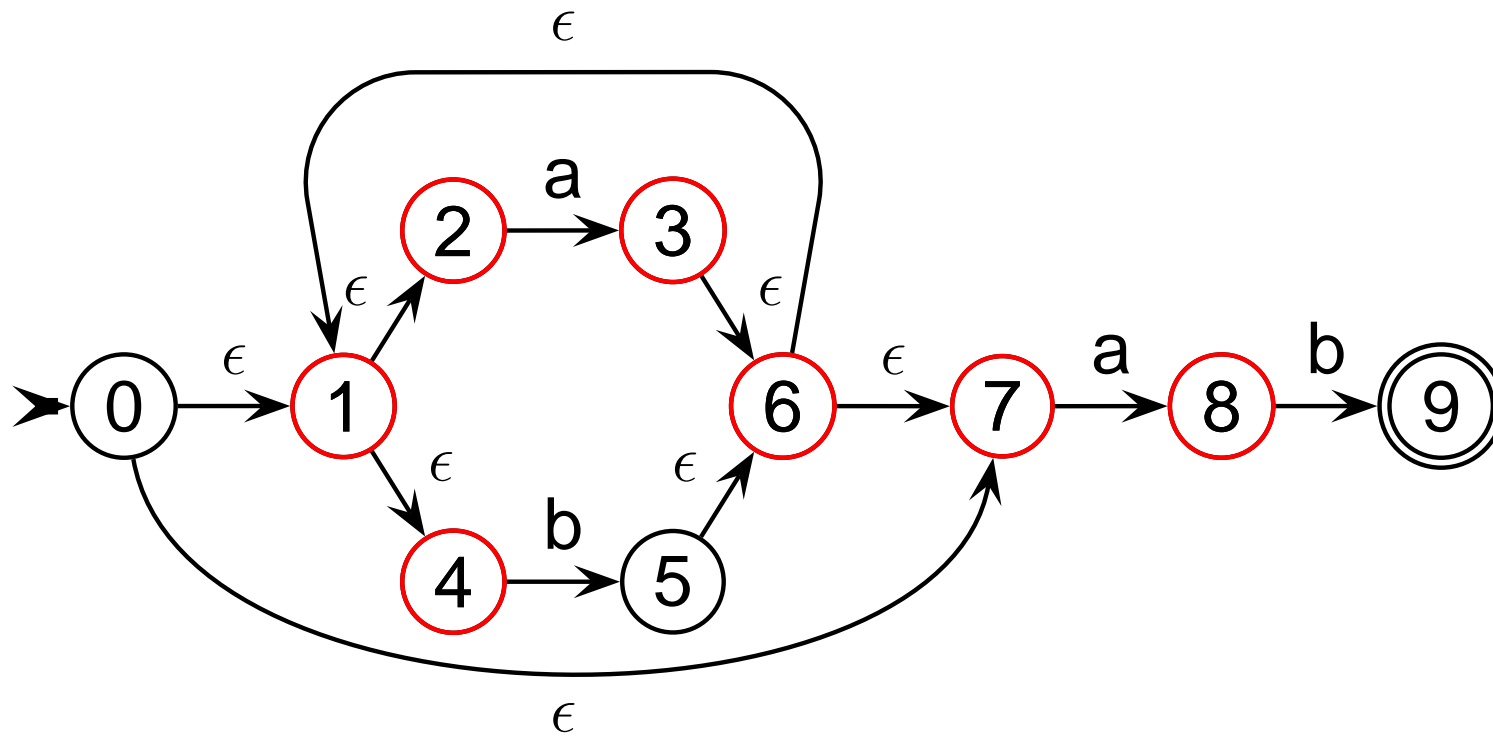
abab



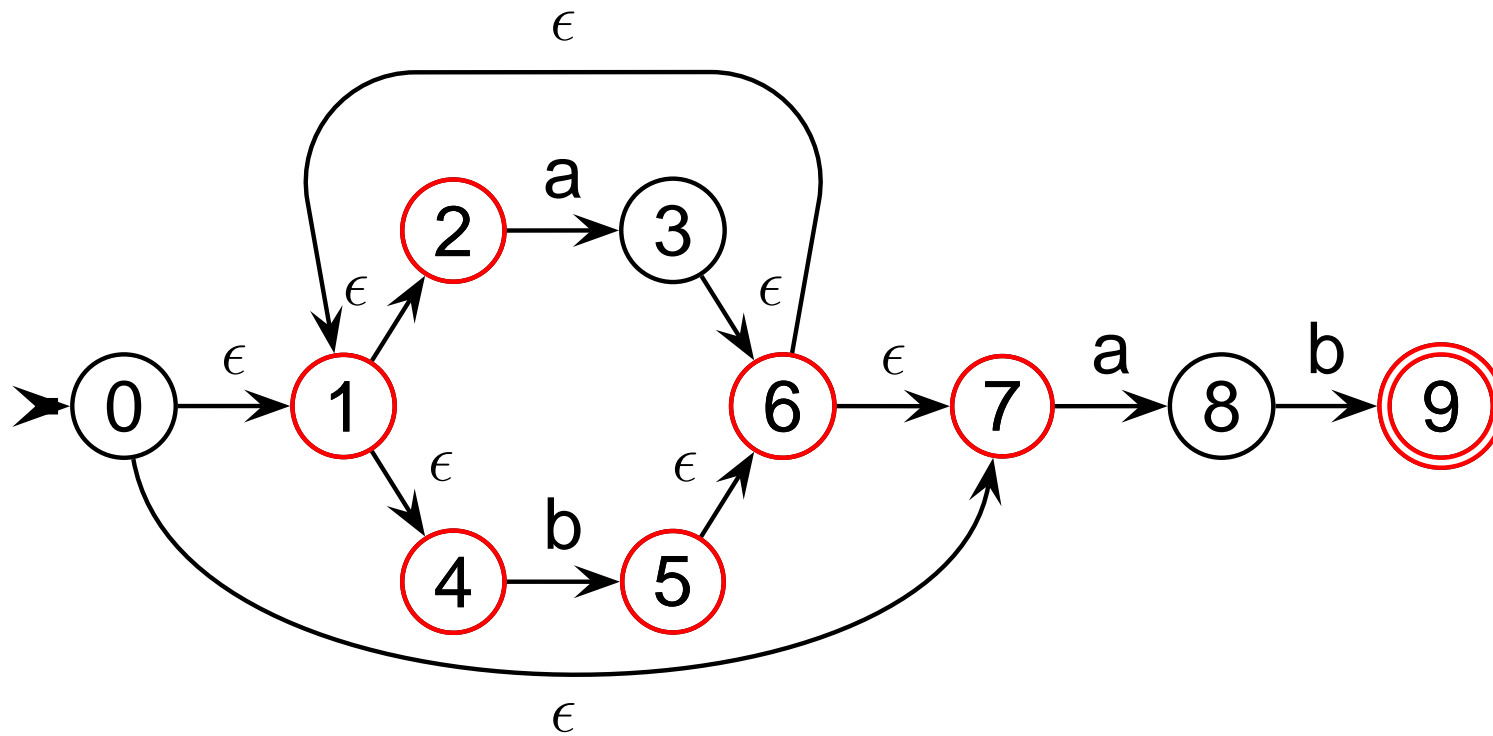
abab



abab



abab



abab

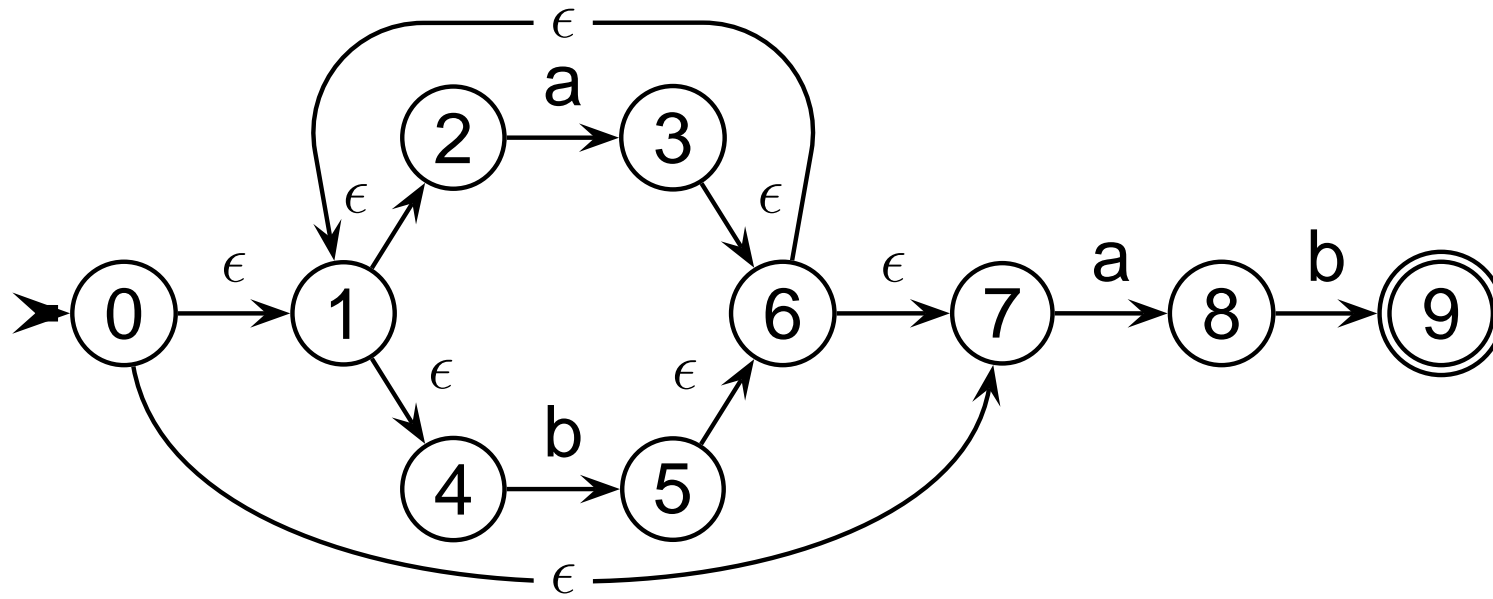
- Simulate “in parallel” all possible moves the automaton can make
- The states of the resulting DFA will represent sets of states of the NFA, i.e., elements of  $\mathcal{P}(Q)$
- We use two operations on states/state-sets of the NFA

$\epsilon$ -closure( $T$ )	Set of states reachable from any state $s$ in $T$ on $\epsilon$ -transitions
$move(T, a)$	Set of states to which there is a transition from one state in $T$ on input symbol $a$

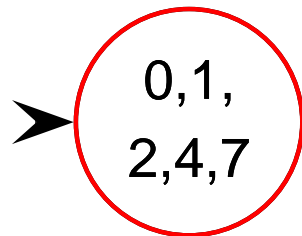
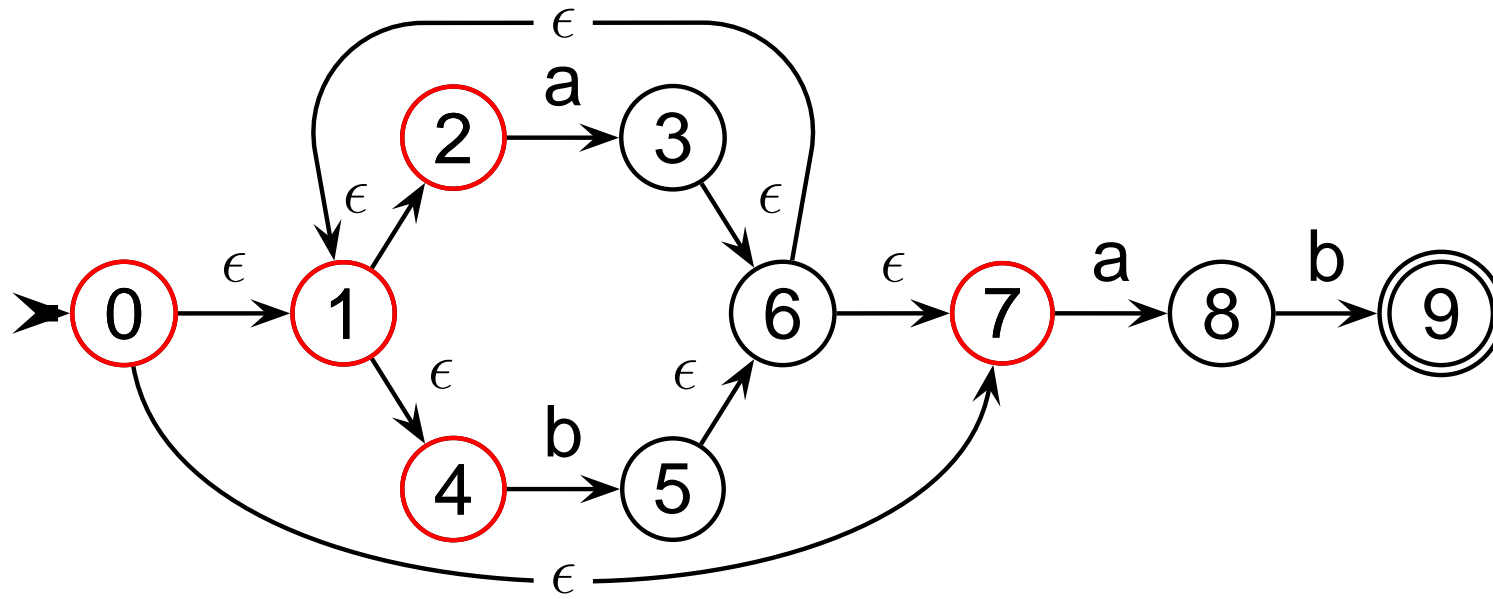
- The final states of the DFA are those where the corresponding NFA subset contains a final state

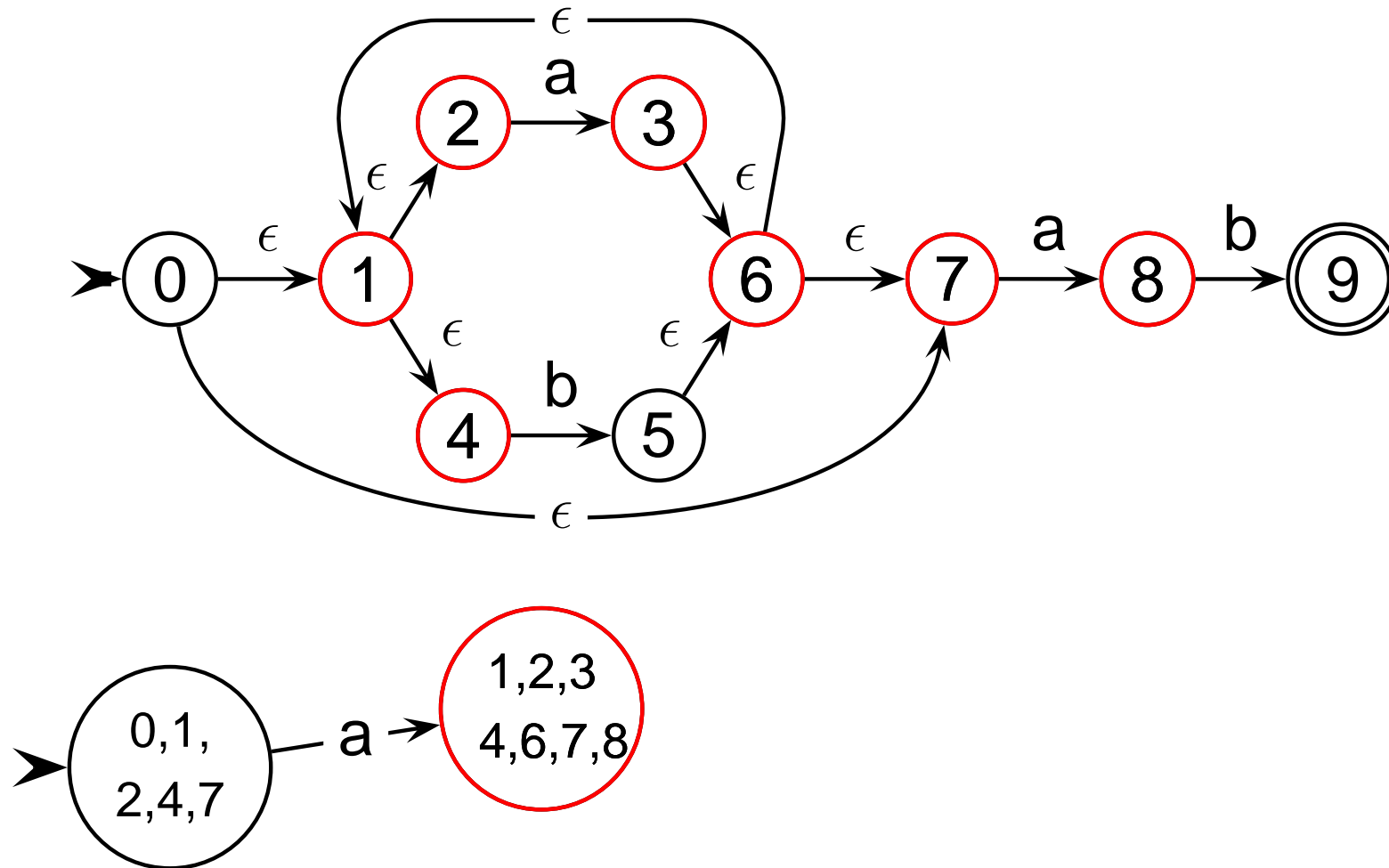
**proc** *SubsetConstruction*( $s_0$ )  $\equiv$   
   $DFASStates = \epsilon\text{-closure}(\{s_0\})$   
  **while** there is an unmarked state  $T$  in  $DFASStates$  **do**  
    mark  $T$   
    **for** each input symbol  $a$  **do**  
       $U := \epsilon\text{-closure}(\text{move}(T, a))$   
       $DFADelta[T, a] := U$   
      **if**  $U \notin DFASStates$  **then** add  $U$  as unmarked state to  $DFASStates$

**proc**  $\epsilon\text{-closure}(T)$   $\equiv$   
   $\epsilon\text{-closure} := T$ ;  $to\_check := T$   
  **while**  $to\_check$  not empty **do**  
    get some state  $t$  from  $to\_check$   
    **for** each state  $u$  with edge labeled  $\epsilon$  from  $t$  to  $u$   
      **if**  $u \notin \epsilon\text{-closure}$  **then** add  $u$  to  $\epsilon\text{-closure}$  and  $to\_check$

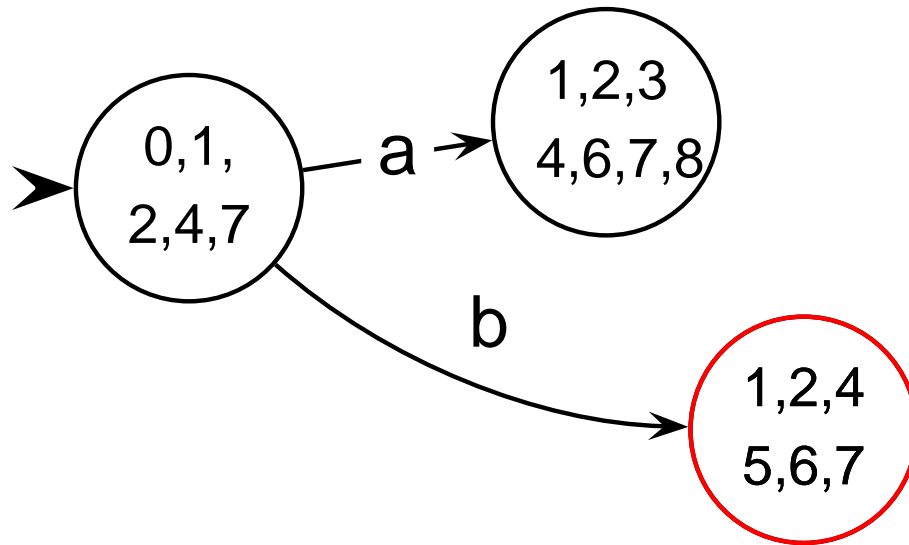
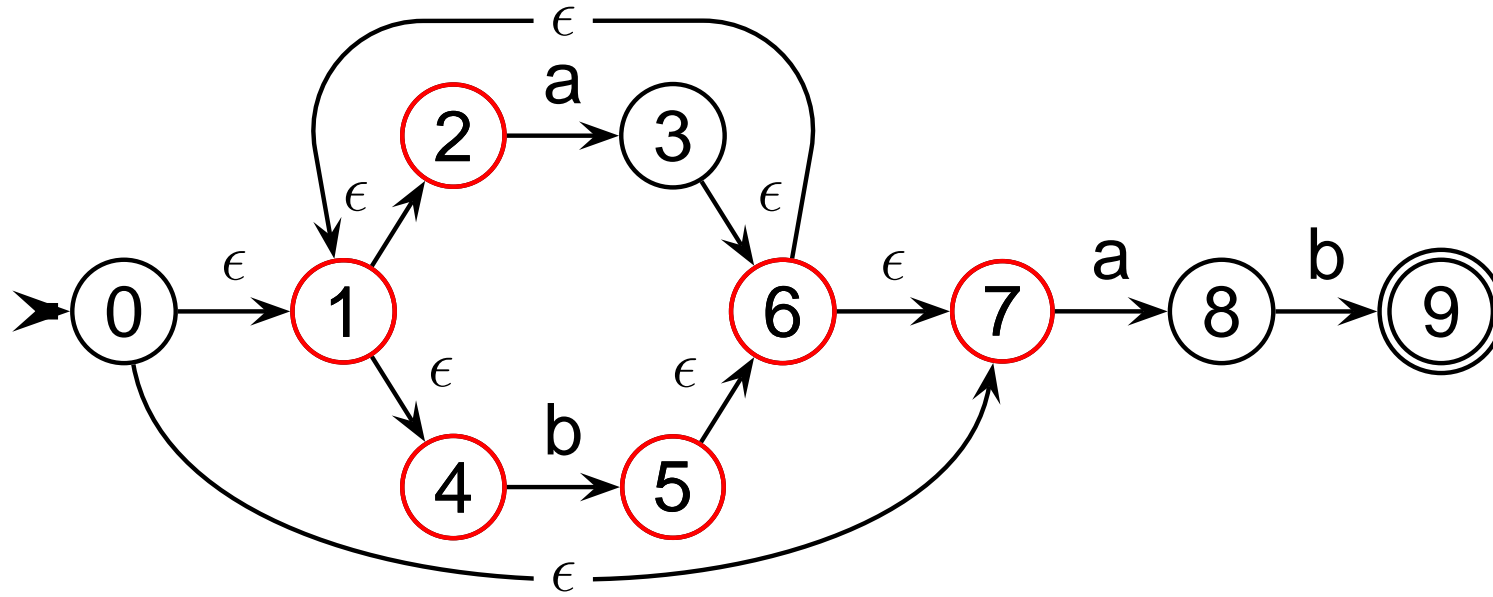




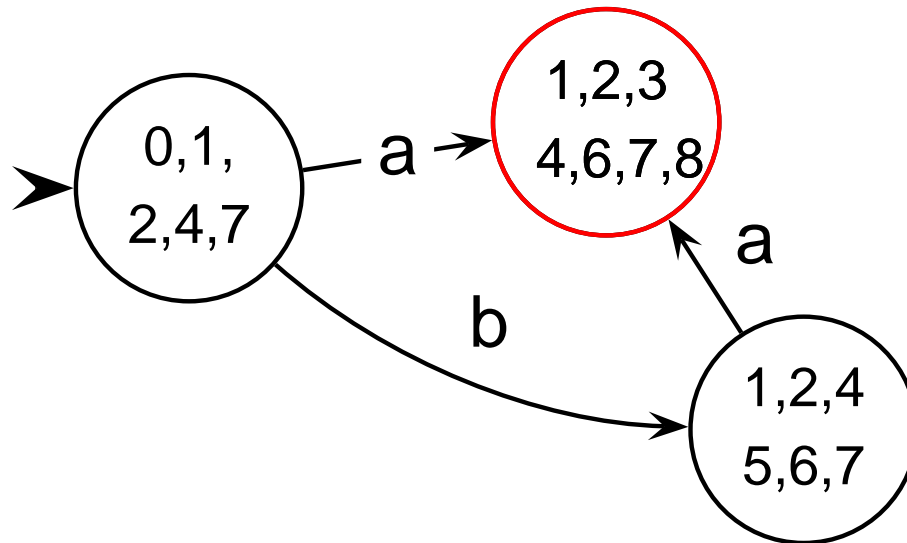
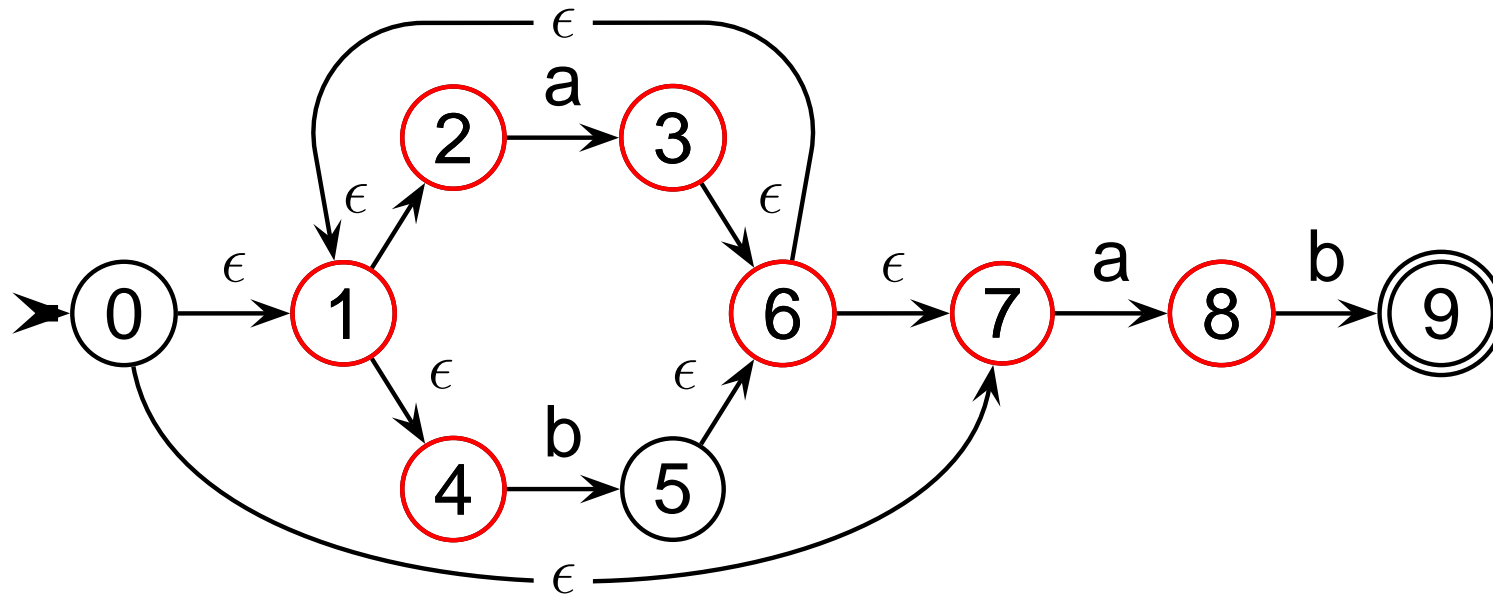




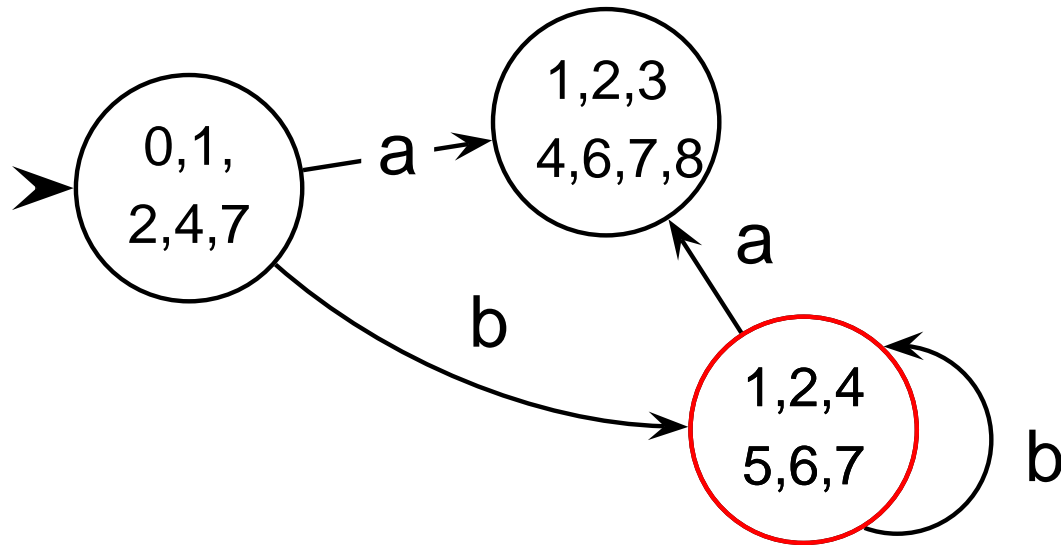
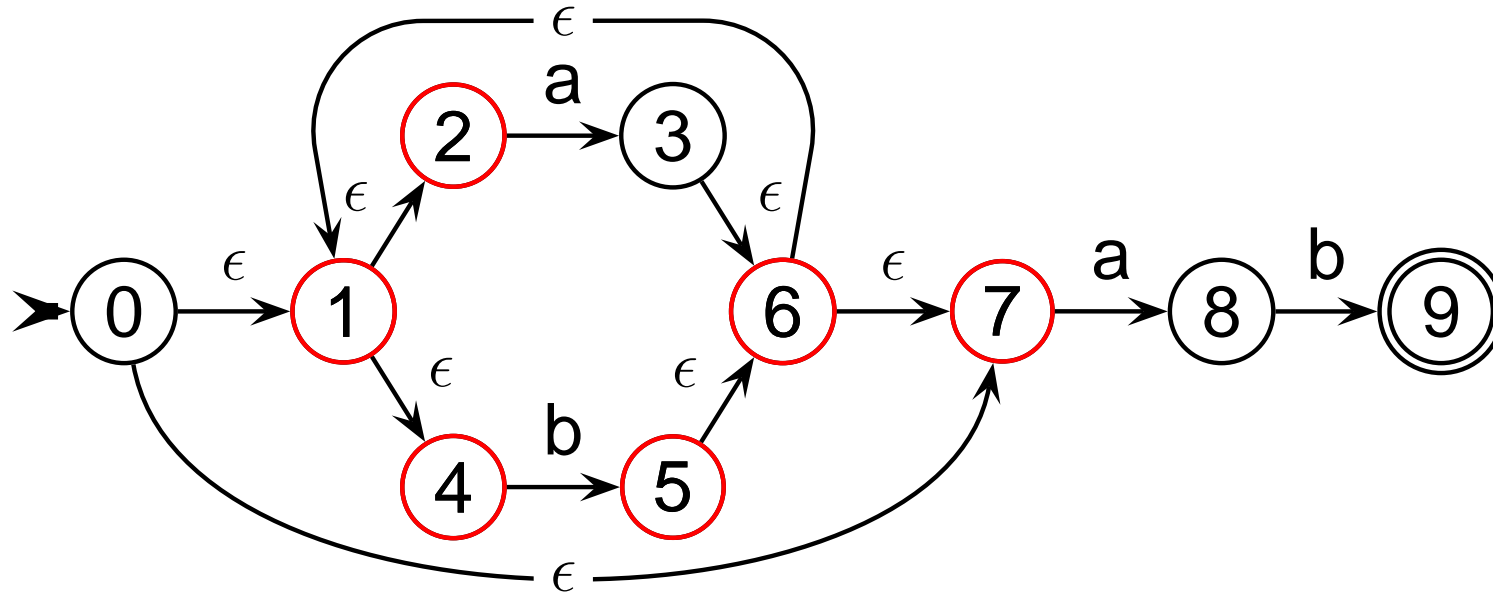
# Example: Subset Construction



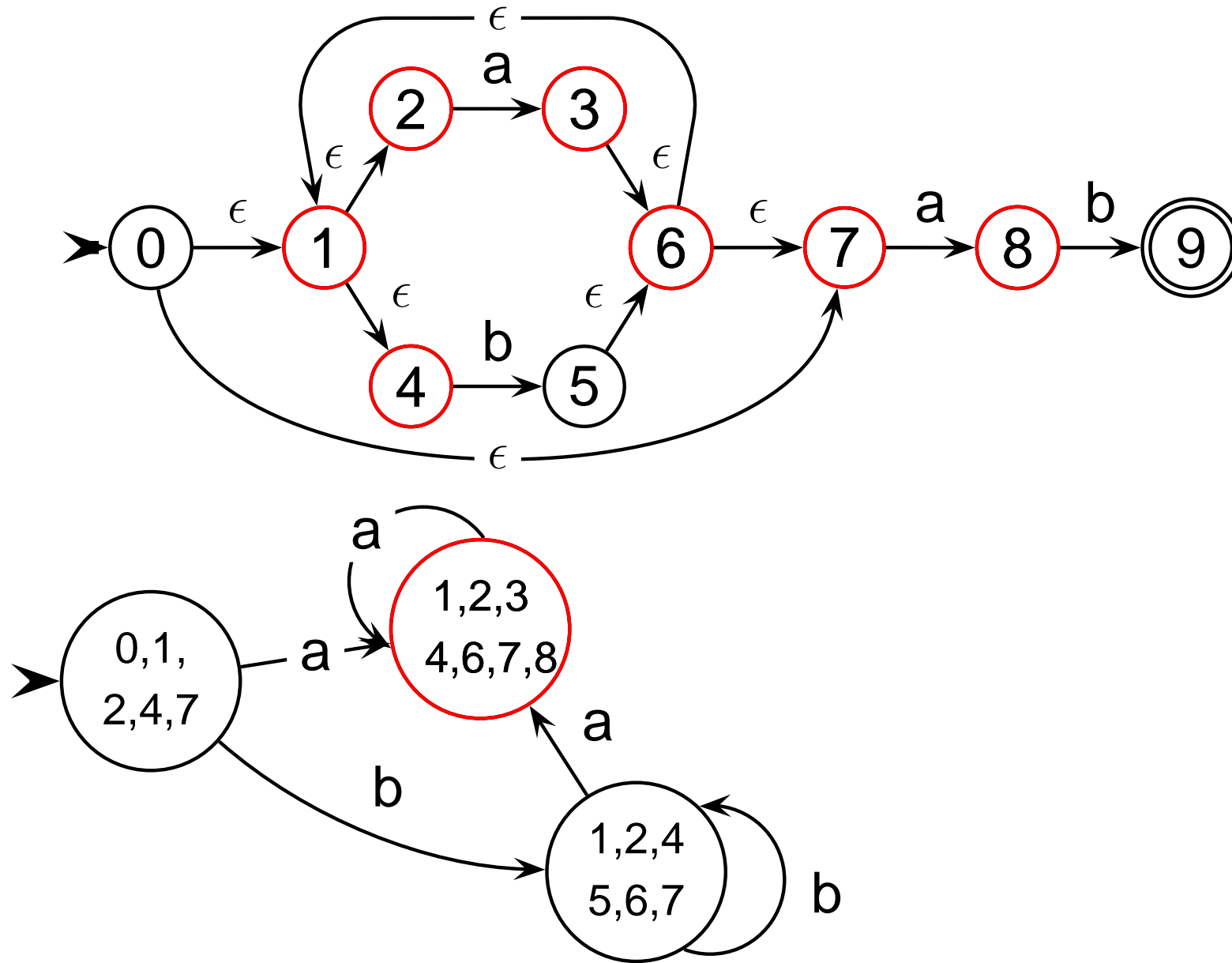
# Example: Subset Construction



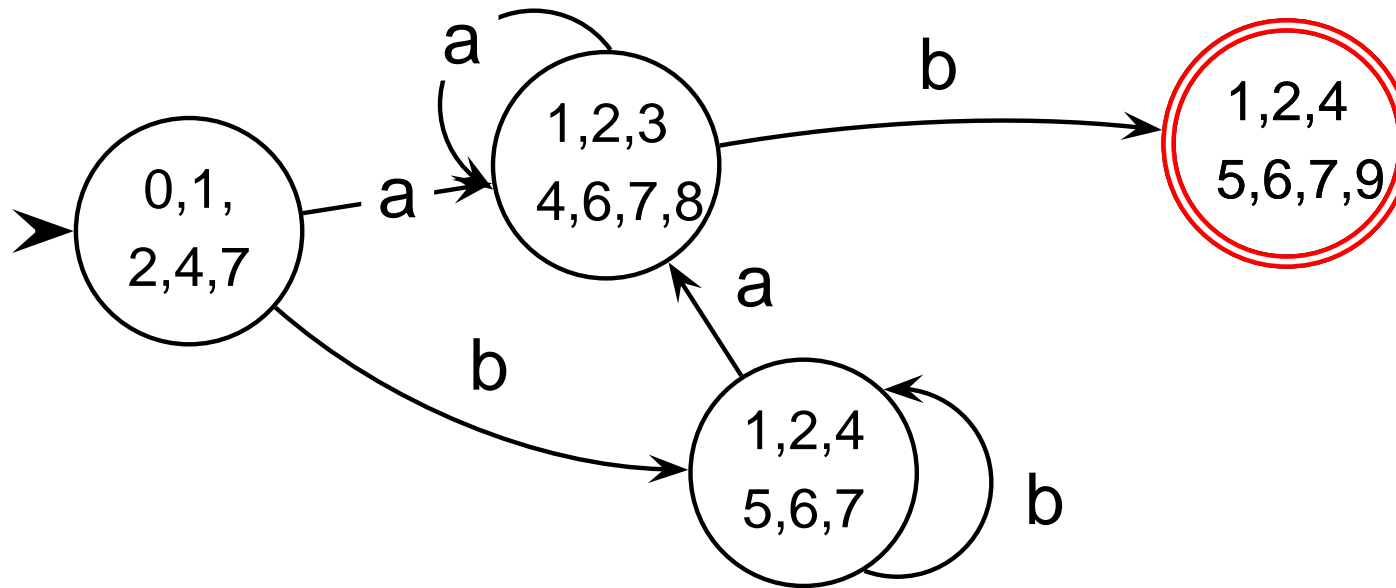
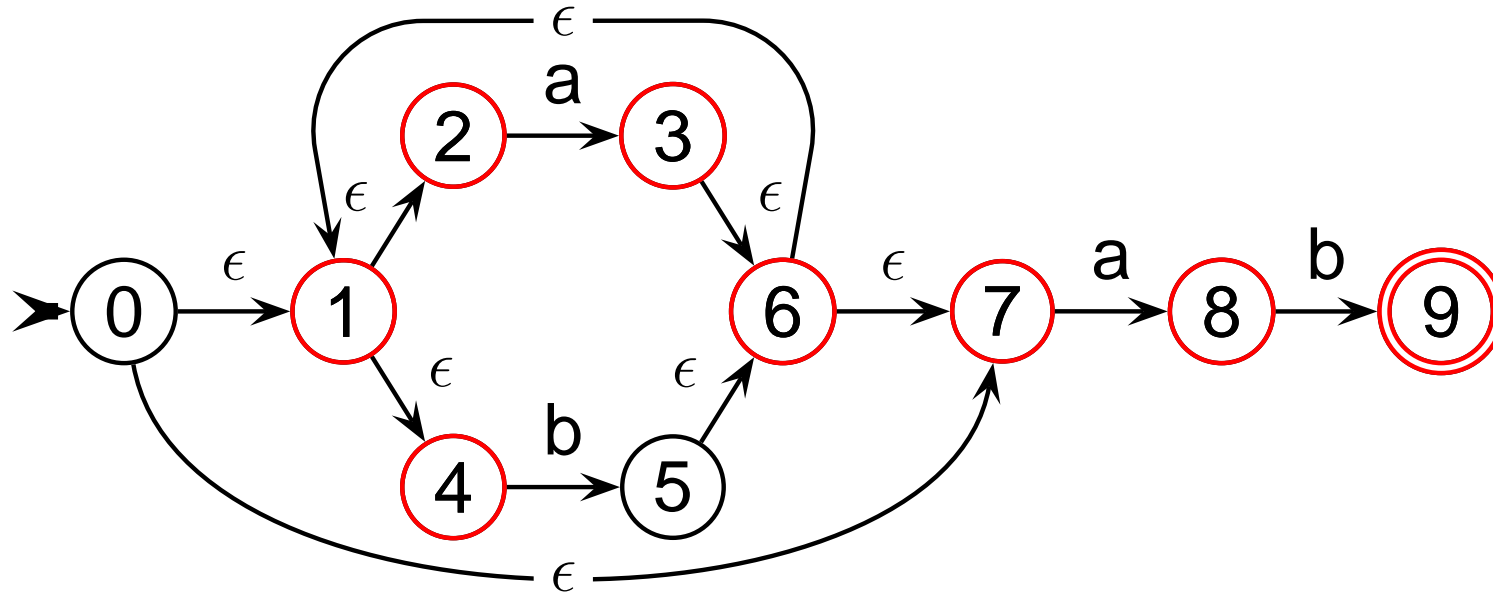
# Example: Subset Construction



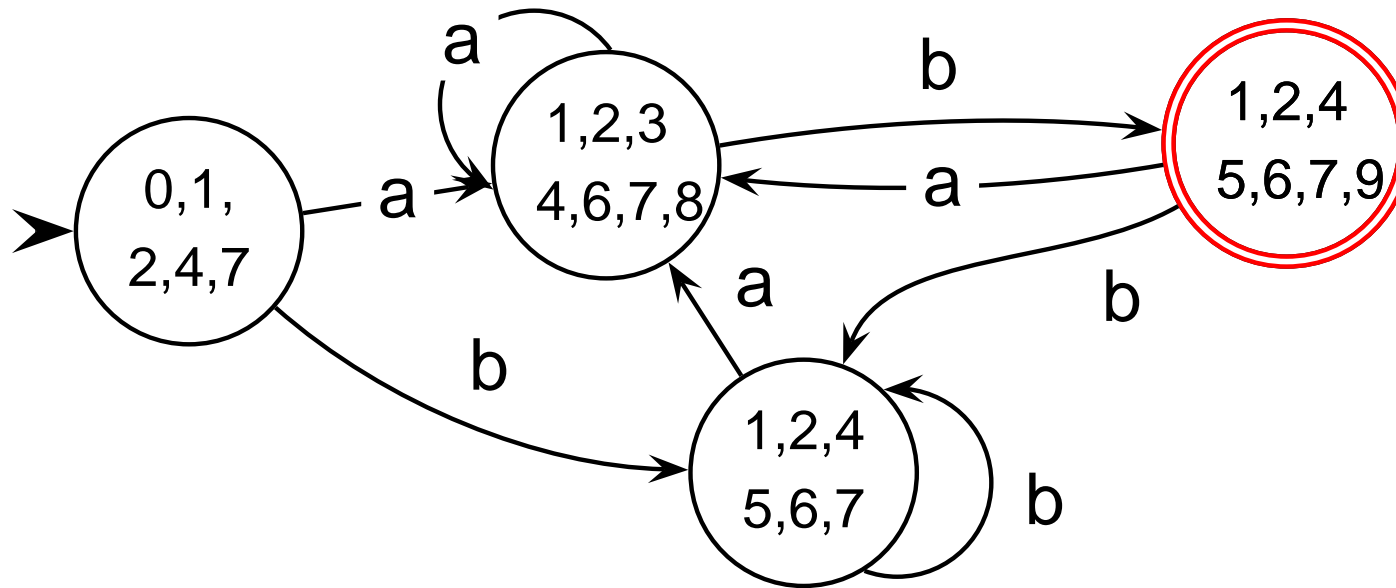
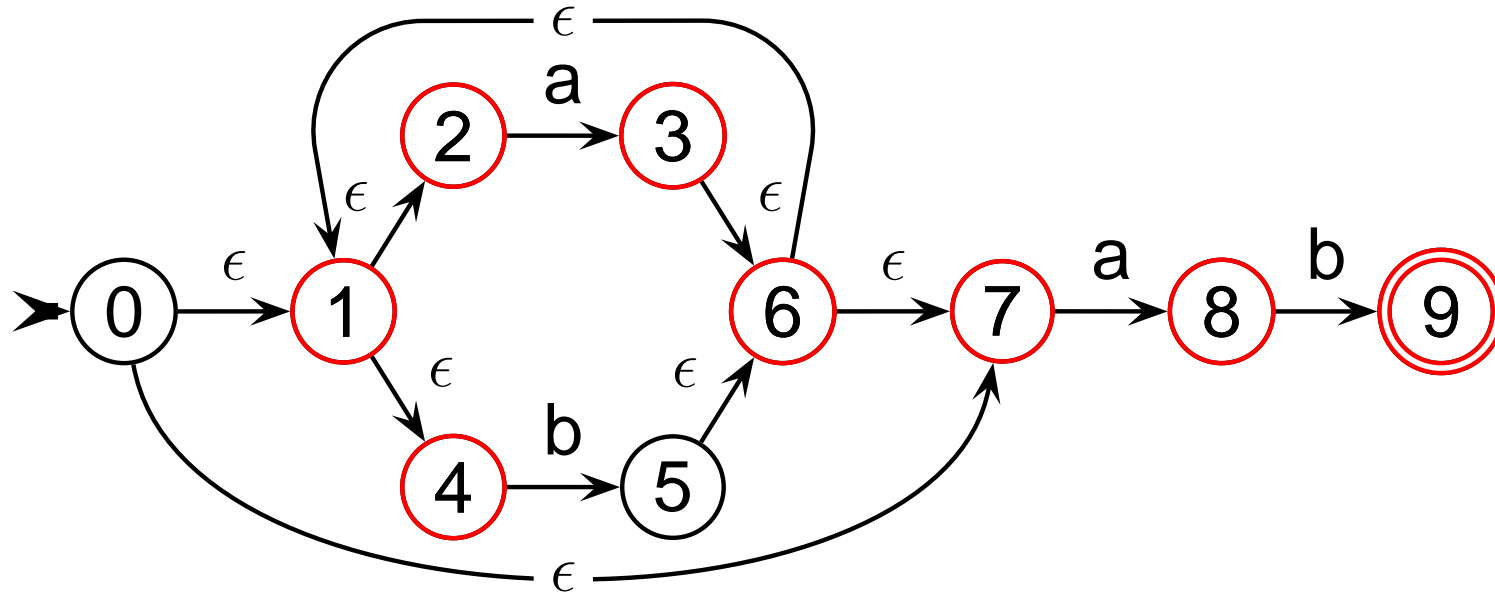
# Example: Subset Construction



# Example: Subset Construction



# Example: Subset Construction





- DFA traversal is linear to the length of input string  $x$
- NFA needs  $\mathcal{O}(n)$  space (states+transitions), where  $n$  is the length of the regular expression
- NFA traversal may need time  $n \times |x|$ , so why use NFAs?

- DFA traversal is linear to the length of input string  $x$
- NFA needs  $\mathcal{O}(n)$  space (states+transitions), where  $n$  is the length of the regular expression
- NFA traversal may need time  $n \times |x|$ , so why use NFAs?
- There are DFA that have at least  $2^n$  states!

- DFA traversal is linear to the length of input string  $x$
- NFA needs  $\mathcal{O}(n)$  space (states+transitions), where  $n$  is the length of the regular expression
- NFA traversal may need time  $n \times |x|$ , so why use NFAs?
- There are DFA that have at least  $2^n$  states!
- Solution 1: “Lazy” construction of the DFA: construct DFA states on the fly up to a certain amount and cache them

- DFA traversal is linear to the length of input string  $x$
- NFA needs  $\mathcal{O}(n)$  space (states+transitions), where  $n$  is the length of the regular expression
- NFA traversal may need time  $n \times |x|$ , so why use NFAs?
- There are DFA that have at least  $2^n$  states!
- Solution 1: “Lazy” construction of the DFA: construct DFA states on the fly up to a certain amount and cache them
- Solution 2: Try to minimize the DFA:  
There is a unique (modulo state names) minimal automaton for a regular language!

**proc** *Minimize*()  $\equiv$   
   $B_1 = F; B_2 = Q F$   
   $E = \{B_1, B_2\}$   
   $k = 3$   
  **for**  $a \in \Sigma$  **do**  
     $a(i) = \{s \in Q \mid s \in B_i \wedge \exists t : \delta(t, a) = s\}$   
     $L =$  the smaller of the  $a(i)$   
    **while**  $L \neq \emptyset$  **do**  
      take some  $i \in L$  and delete it  
      **for**  $j < k$  s.th.  $\exists t \in B_j$