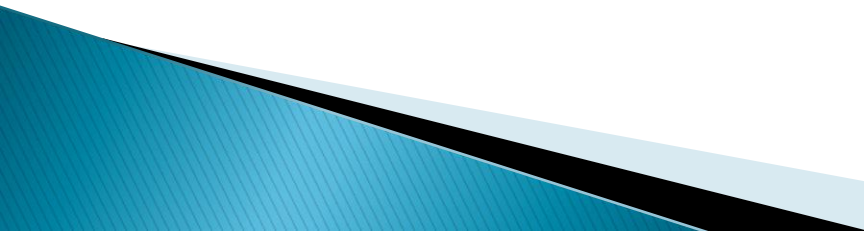


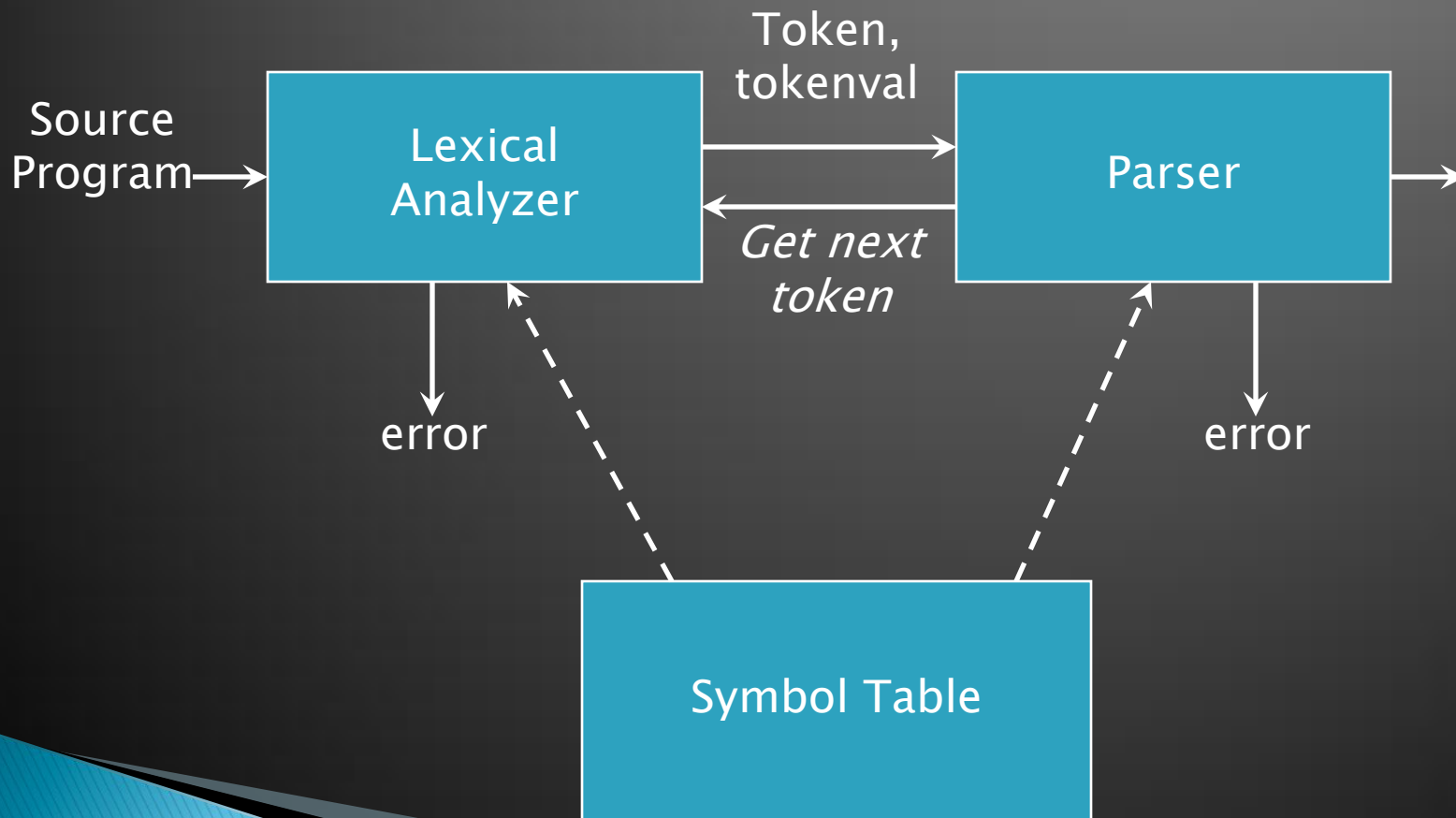
# Lecture # 5

Lexical Analysis

# Role of Lexical Analyzer

- ▶ It is the first phase of compiler
  - ▶ Its main task is to read the input characters and produce as output a sequence of tokens that the parser uses for syntax analysis
  - ▶ Reasons to make it a separate phase are:
    - Simplifies the design of the compiler
    - Provides efficient implementation(read the source code)
    - Improves portability
- 

# Interaction of the Lexical Analyzer with the Parser



# Tokens, Patterns, and Lexemes

- ▶ A *token* is a classification of lexical units
  - For example: `id` and `num`
- ▶ *Lexemes* are the specific character strings that make up a token
  - For example: `abc` and `123`
- ▶ *Patterns* are rules describing the set of lexemes belonging to a token
  - For example: “*letter followed by letters and digits*” and “*non-empty sequence of digits*”

# Diff b/w Token, Lexeme and Pattern

Token	Lexeme	Pattern
if	if	if
relation	<, <=,=,<>,>,>=	< or <= or = or <> or > or >=
id	y, x	Letter followed by letters and digits
num	31 , 28	Any numeric constant
operator	+ , * , - , /	Any arithmetic operator + or * or - or /

# Attributes of Tokens

`y := 31 + 28*x`

Lexical analyzer

`<id, "y"> <assign, :=> <num, 31> <operator, +> <num, 28> <operator, *>`

token

tokenval  
(token attribute)

Parser

# Specification of Tokens

- ▶ Alphabet: Finite, nonempty set of symbols

Example:  $\Sigma = \{0, 1\}$  binary alphabet

Example:  $\Sigma = \{a, b, c, \dots, z\}$  the set of all lower case letters

- ▶ Strings: Finite sequence of symbols from an alphabet e.g. 0011001

- ▶ Empty String: The string with zero occurrences of symbols from alphabet. The empty string is denoted by  $\epsilon$

# Continue...

- ▶ Length of String: Number of positions for symbols in the string.  $|w|$  denotes the length of string  $w$

Example  $|0110| = 4$ ;  $|\epsilon| = 0$

- ▶ Powers of an Alphabet:  $\Sigma^k$  = the set of strings of length  $k$  with symbols from  $\Sigma$

Example:

$$\Sigma = \{0, 1\}$$

$$\Sigma^1 = \{0, 1\}$$

$$\Sigma^2 = \{00, 01, 10, 11\}$$

$$\Sigma^0 = \{\epsilon\}$$



# Continue..

- ▶ The set of all strings over  $\Sigma$  is denoted  $\Sigma^*$

$$\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots$$

$$\Sigma^+ = \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup \dots$$

$$\Sigma^* = \Sigma^+ \cup \{\epsilon\}$$

# Continue..

- ▶ Language: is a specific set of strings over some fixed alphabet  $\Sigma$

Example:

The set of legal English words

The set of strings consisting of  $n$  0's followed by  $n$  1's

LP = the set of binary numbers whose value is prime  $\in \{\epsilon, 01, 0011, 000111, \dots\}$

$\{10, 11, 101, 111, 1011, \dots\}$

# Concatenation and Exponentiation

- ▶ The *concatenation* of two strings  $x$  and  $y$  is denoted by  $xy$
- ▶ The *exponentiation* of a string  $s$  is defined by

$$s^0 = \varepsilon$$
$$s^i = s^{i-1}s \quad \text{for } i > 0$$

note that  $s\varepsilon = \varepsilon s = s$

# Language Operations

- ▶ Union

$$L \cup M = \{s \mid s \in L \text{ or } s \in M\}$$

- ▶ Concatenation

$$LM = \{xy \mid x \in L \text{ and } y \in M\}$$

- ▶ Exponentiation

$$L^0 = \{\varepsilon\}; \quad L^i = L^{i-1}L$$

- ▶ Kleene closure

$$L^* = \cup_{i=0, \dots, \infty} L^i$$

- ▶ Positive closure

$$L^+ = \cup_{i=1, \dots, \infty} L^i$$

# Regular Expressions

- ▶ Basis symbols:
  - $\varepsilon$  is a regular expression denoting language  $\{\varepsilon\}$
  - $a \in \Sigma$  is a regular expression denoting  $\{a\}$
- ▶ If  $r$  and  $s$  are regular expressions denoting languages  $L(r)$  and  $M(s)$  respectively, then
  - $r|s$  is a regular expression denoting  $L(r) \cup M(s)$
  - $rs$  is a regular expression denoting  $L(r)M(s)$
  - $r^*$  is a regular expression denoting  $L(r)^*$
  - $(r)$  is a regular expression denoting  $L(r)$
- ▶ A language defined by a regular expression is called a *Regular set* or a *Regular Language*

# Regular Definitions

- ▶ Regular definitions introduce a naming convention:

$$d_1 \rightarrow r_1$$

$$d_2 \rightarrow r_2$$

...

$$d_n \rightarrow r_n$$

where each  $r_i$  is a regular expression over

$$\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$$

▶ Example:

**letter**  $\rightarrow$  **A** | **B** | ... | **Z** | **a** | **b** | ... | **z**

**digit**  $\rightarrow$  **0** | **1** | ... | **9**

**id**  $\rightarrow$  **letter** ( **letter** | **digit** )\*

▶ The following shorthands are often used:

$$r^+ = rr^*$$

$$r? = r | \epsilon$$

$$[\mathbf{a-z}] = \mathbf{a} | \mathbf{b} | \mathbf{c} | \dots | \mathbf{z}$$

▶ Examples:

**digit**  $\rightarrow$  **[0-9]**

**num**  $\rightarrow$  **digit**<sup>+</sup> ( **.** **digit**<sup>+</sup> )? ( **E** ( **+** | **-** )? **digit**<sup>+</sup> )?

# Regular Definitions and Grammars

Grammar

$stmt \rightarrow \text{if } expr \text{ then } stmt$

|  $\text{if } expr \text{ then } stmt \text{ else } stmt$

|  $\epsilon$

$expr \rightarrow term \text{ relop } term$

|  $term$

$term \rightarrow \text{id}$

|  $\text{num}$

Regular definitions

$\text{if} \rightarrow \text{if}$

$\text{then} \rightarrow \text{then}$

$\text{else} \rightarrow \text{else}$

$\text{relop} \rightarrow < \mid <= \mid <> \mid > \mid >= \mid =$

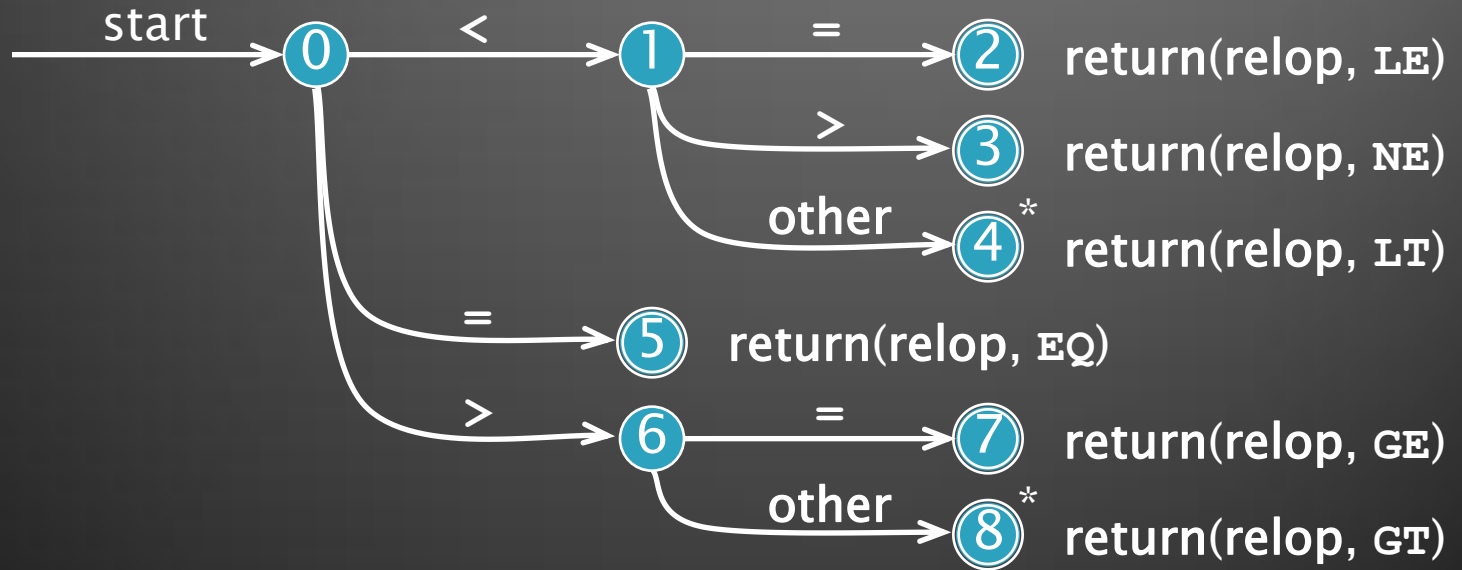
$\text{id} \rightarrow \text{letter} ( \text{letter} \mid \text{digit} )^*$

$\text{num} \rightarrow \text{digit}^+ ( . \text{digit}^+ )? ( \text{E} ( + \mid - )? \text{digit}^+ )?$

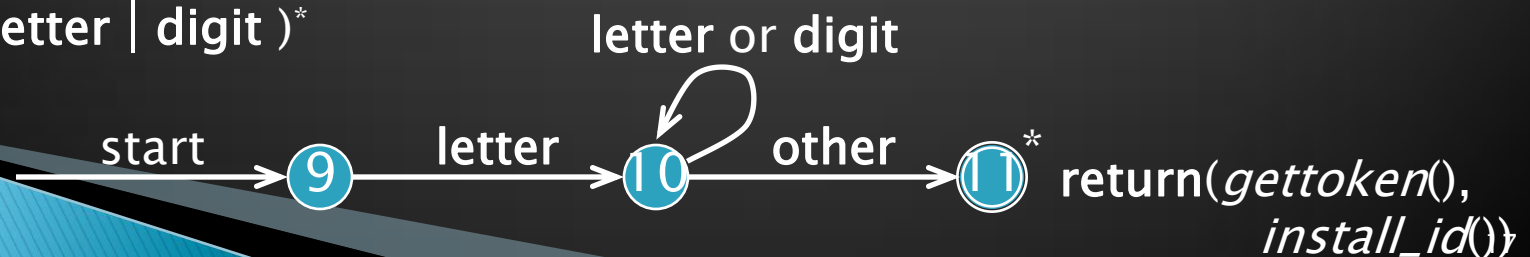


# Coding Regular Definitions in Transition Diagrams

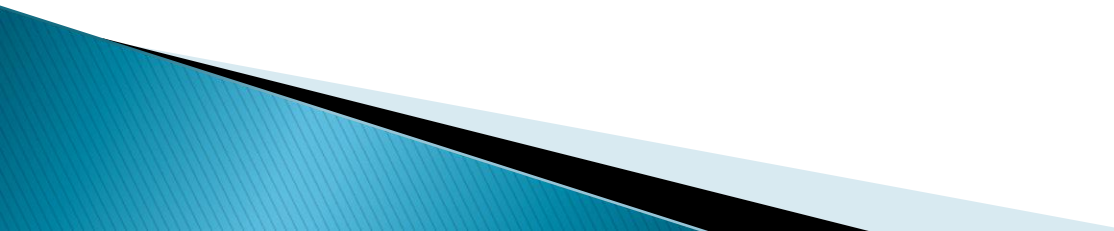
relop  $\rightarrow$  < | <= | <> | > | >= | =



id  $\rightarrow$  letter ( letter | digit )<sup>\*</sup>

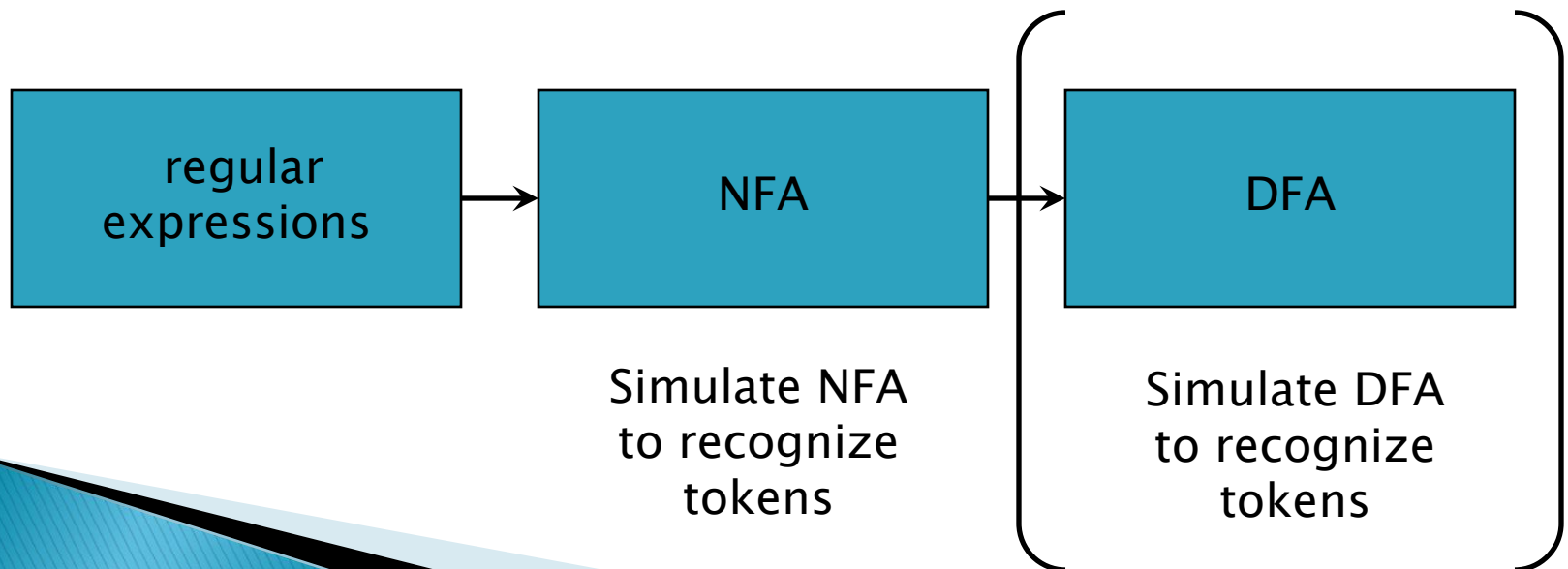


# Finite Automata

- ▶ Finite Automata are used as a model for:
    - Software for designing digital circuits
    - Lexical analyzer of a compiler
    - Searching for keywords in a file or on the web.
    - Software for verifying finite state systems, such as communication protocols.
- 

# Design of a Lexical Analyzer Generator

- ▶ Translate regular expressions to NFA
- ▶ Translate NFA to an efficient DFA



# Nondeterministic Finite Automata

- ▶ An NFA is a 5-tuple  $(S, \Sigma, \delta, s_0, F)$  where

$S$  is a finite set of *states*

$\Sigma$  is a finite set of symbols, the *alphabet*

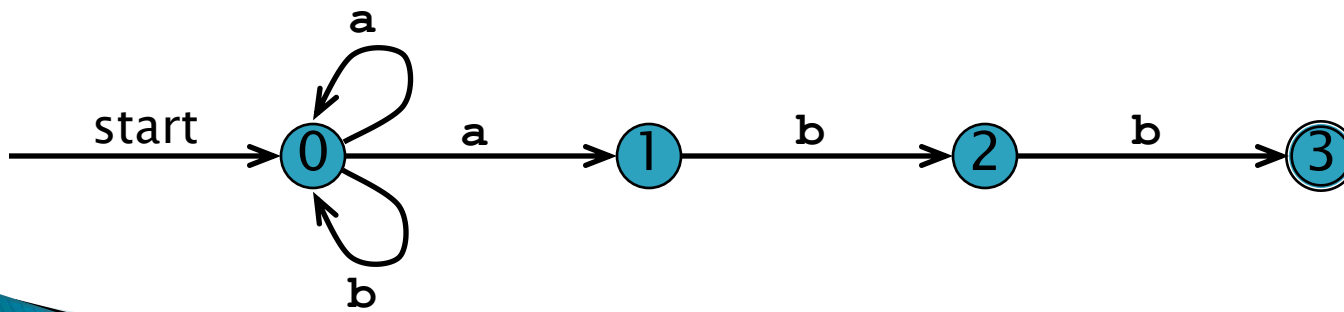
$\delta$  is a *mapping* from  $S \times \Sigma$  to a set of states

$s_0 \in S$  is the *start state*

$F \subseteq S$  is the set of *accepting* (or *final*) *states*

# Transition Graph

- ▶ An NFA can be diagrammatically represented by a labeled directed graph called a *transition graph*



$S = \{0,1,2,3\}$   
 $\Sigma = \{a,b\}$   
 $s_0 = 0$   
 $F = \{3\}$

# Transition Table

- ▶ The mapping  $\delta$  of an NFA can be represented in a *transition table*

$$\delta(0, \mathbf{a}) = \{0, 1\}$$

$$\delta(0, \mathbf{b}) = \{0\}$$

$$\delta(1, \mathbf{b}) = \{2\}$$

$$\delta(2, \mathbf{b}) = \{3\}$$

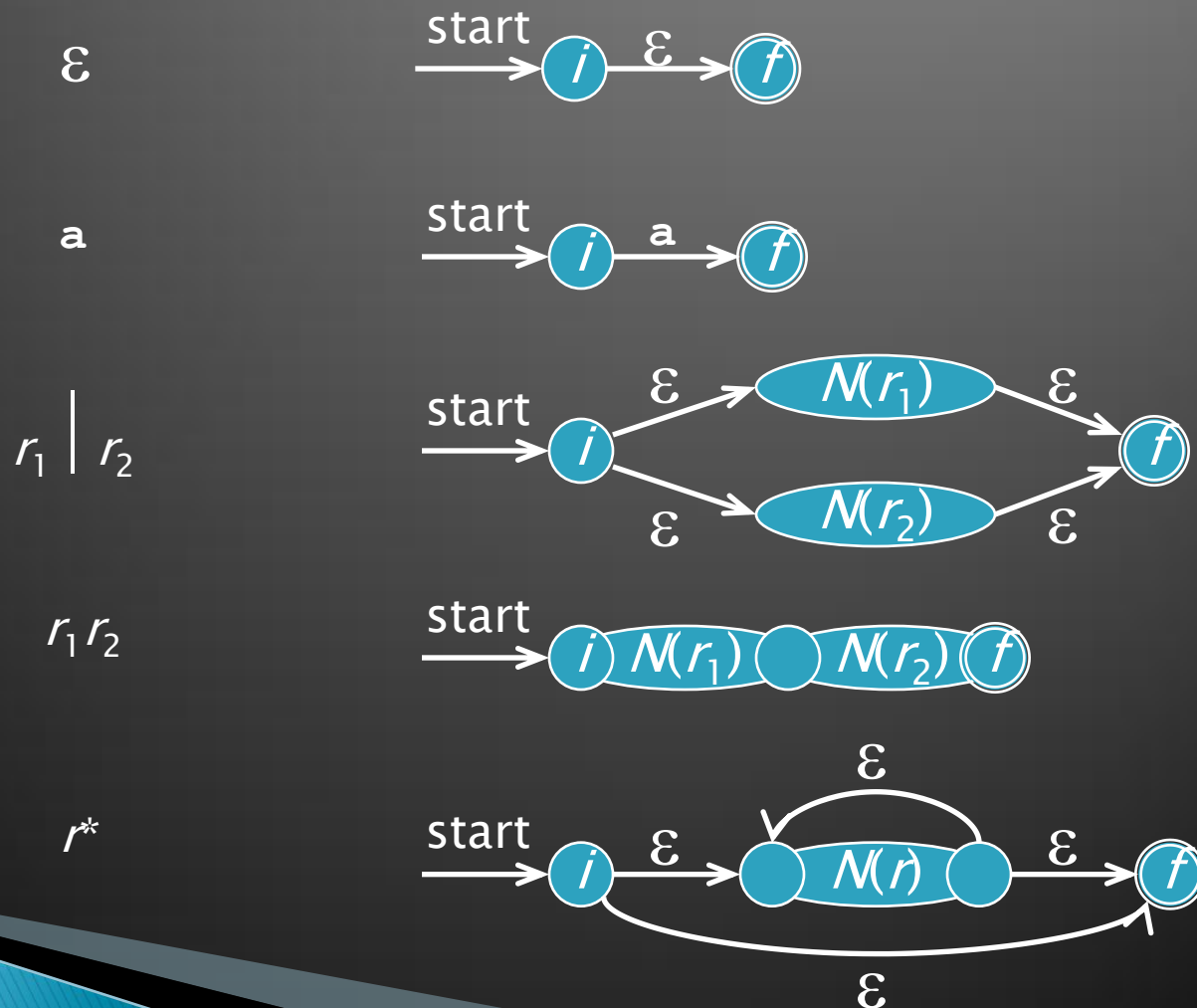


<i>State</i>	<i>Input a</i>	<i>Input b</i>
0	{0, 1}	{0}
1		{2}
2		{3}

# The Language Defined by an NFA

- ▶ An NFA *accepts* an input string  $x$  if and only if there is some path with edges labeled with symbols from  $x$  in sequence from the start state to some accepting state in the transition graph
- ▶ A state transition from one state to another on the path is called a *move*
- ▶ The *language defined by* an NFA is the set of input strings it accepts, such as  $(\mathbf{a} \mid \mathbf{b})^* \mathbf{abb}$  for the example NFA

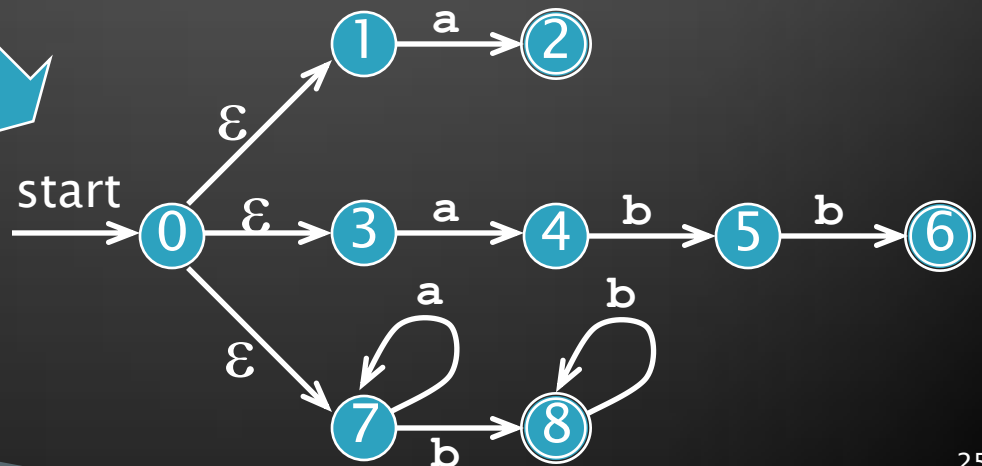
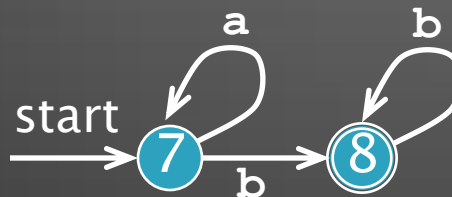
# From Regular Expression to $\epsilon$ -NFA (Thompson's Construction)





# Combining the NFAs of a Set of Regular Expressions

a        { *action*<sub>1</sub> }  
abb     { *action*<sub>2</sub> }  
a\*b+   { *action*<sub>3</sub> }

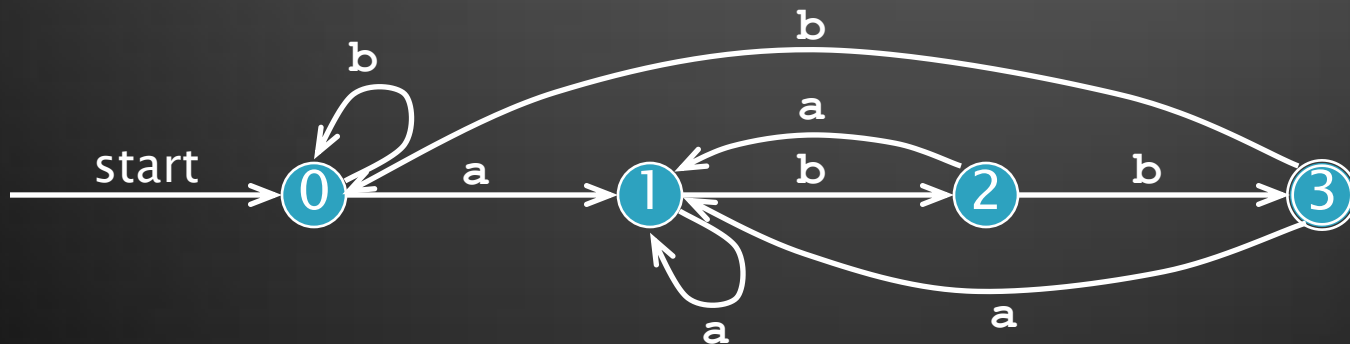


# Deterministic Finite Automata

- ▶ A *deterministic finite automaton* is a special case of an NFA
  - No state has an  $\varepsilon$ -transition
  - For each state  $s$  and input symbol  $a$  there is at most one edge labeled  $a$  leaving  $s$
- ▶ Each entry in the transition table is a single state
  - At most one path exists to accept a string
  - Simulation algorithm is simple

# Example DFA

A DFA that accepts  $(a | b)^*abb$



# Conversion of an NFA into a DFA

- ▶ The *subset construction algorithm* converts an NFA into a DFA using:

$$\varepsilon\text{-closure}(s) = \{s\} \cup \{t \mid s \xrightarrow{\varepsilon} \dots \xrightarrow{\varepsilon} t\}$$

$$\varepsilon\text{-closure}(T) = \bigcup_{s \in T} \varepsilon\text{-closure}(s)$$

$$\text{move}(T, a) = \{t \mid s \xrightarrow{a} t \text{ and } s \in T\}$$

- ▶ The algorithm produces:

*Dstates* is the set of states of the new DFA consisting of sets of states of the NFA

*Dtran* is the transition table of the new DFA

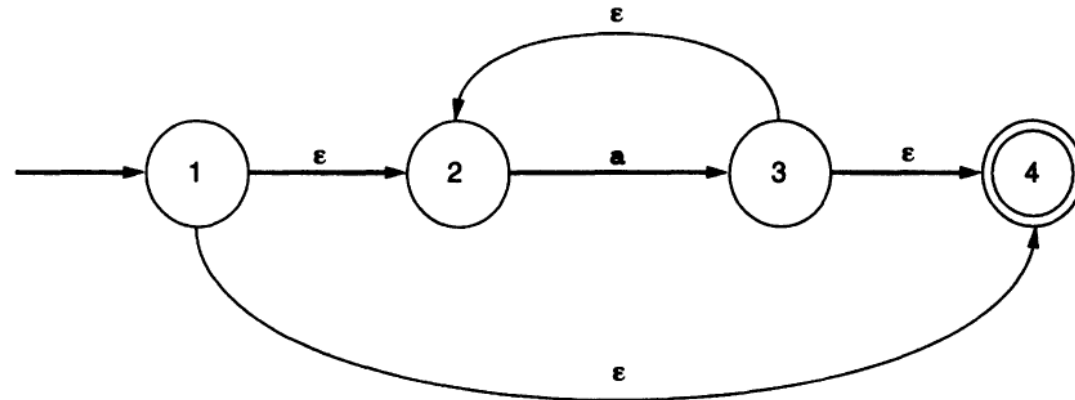
Thumbnails

Thus, this algorithm is called the **subset construction**. We first discuss the  $\epsilon$ -closure in a little more detail and then proceed to a description of the subset construction.

*The  $\epsilon$ -Closure of a Set of States* We define the  $\epsilon$ -closure of a single state  $s$  as the set of states reachable by a series of zero or more  $\epsilon$ -transitions, and we write this set as  $\bar{s}$ . We leave a more mathematical statement of this definition to an exercise and proceed directly to an example. Note, however, that the  $\epsilon$ -closure of a state always contains the state itself.

**Example 2.14**

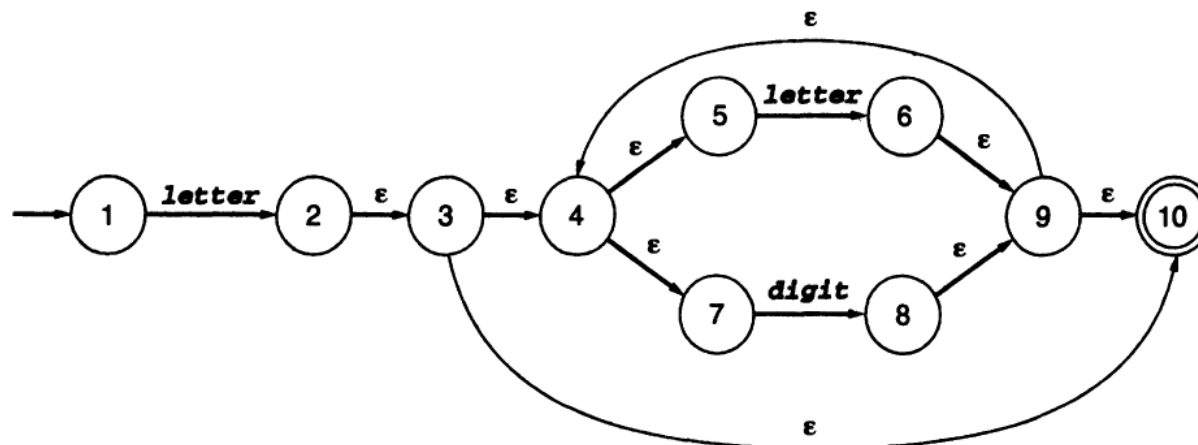
Consider the following NFA corresponding to the regular expression  $a^*$  under Thompson's construction:



In this NFA, we have  $\bar{1} = \{1, 2, 4\}$ ,  $\bar{2} = \{2\}$ ,  $\bar{3} = \{2, 3, 4\}$ , and  $\bar{4} = \{4\}$ . §

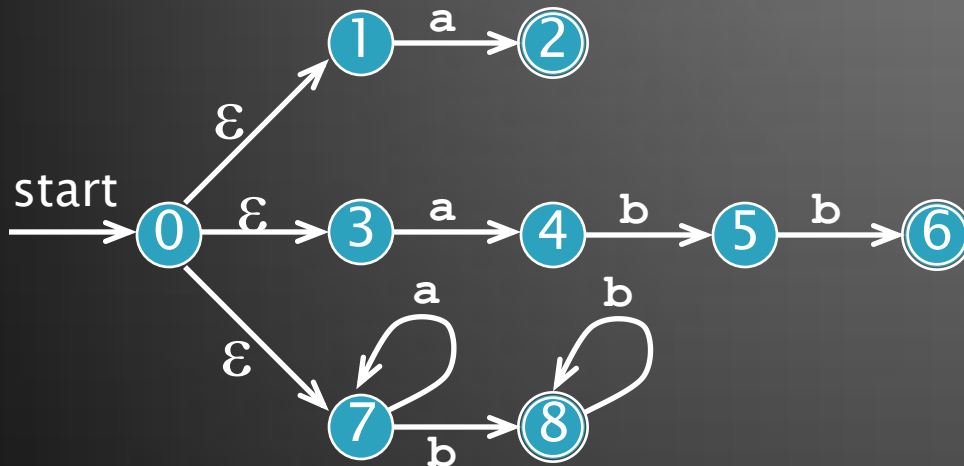
**Example 2.17**

Consider the NFA of Figure 2.9 (Thompson's construction for the regular expression **letter(letter|digit)\***):

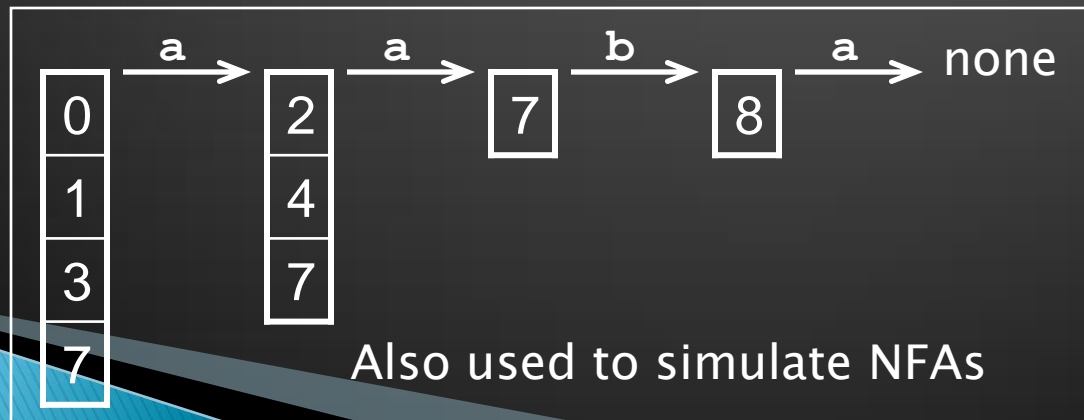


The subset construction proceeds as follows. The start state is  $\overline{\{1\}} = \{1\}$ . There is a transition on **letter** to  $\overline{\{2\}} = \{2, 3, 4, 5, 7, 10\}$ . From this state there is a transition on **letter** to  $\overline{\{6\}} = \{4, 5, 6, 7, 9, 10\}$  and a transition on **digit** to  $\overline{\{8\}} = \{4, 5, 7, 8, 9, 10\}$ . Finally, each of these states also has transitions on **letter** and **digit**, either to itself or to the other. The complete DFA is given in the following picture:

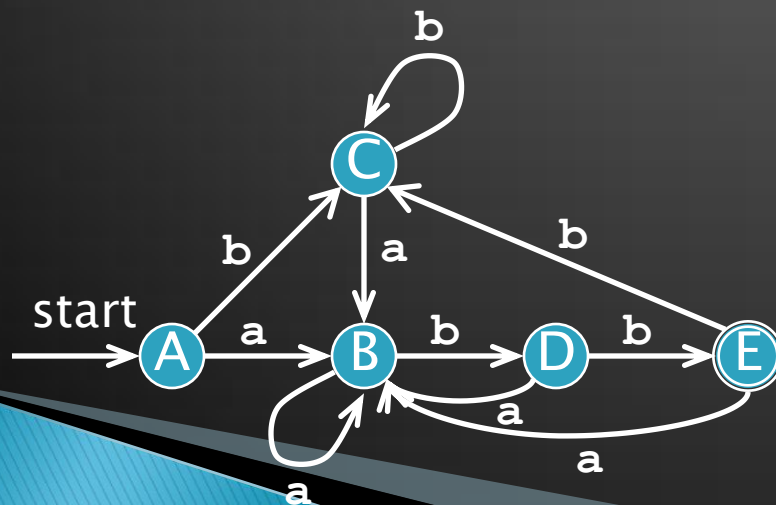
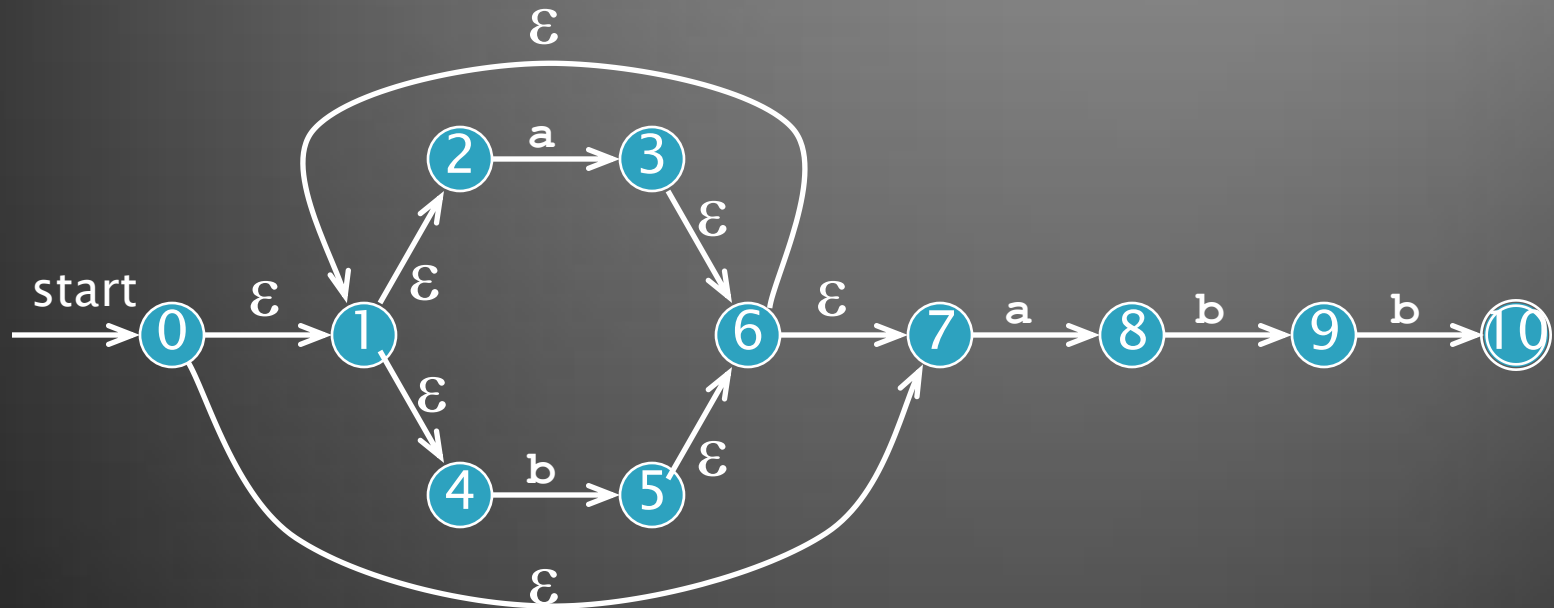
# $\epsilon$ -closure and move Examples



$\epsilon$ -closure( $\{0\}$ ) =  $\{0,1,3,7\}$   
 $move(\{0,1,3,7\},a) = \{2,4,7\}$   
 $\epsilon$ -closure( $\{2,4,7\}$ ) =  $\{2,4,7\}$   
 $move(\{2,4,7\},a) = \{7\}$   
 $\epsilon$ -closure( $\{7\}$ ) =  $\{7\}$   
 $move(\{7\},b) = \{8\}$   
 $\epsilon$ -closure( $\{8\}$ ) =  $\{8\}$   
 $move(\{8\},a) = \emptyset$



# Subset Construction Example 1



*Dstates*

A = {0,1,2,4,7}

B = {1,2,3,4,6,7,8}

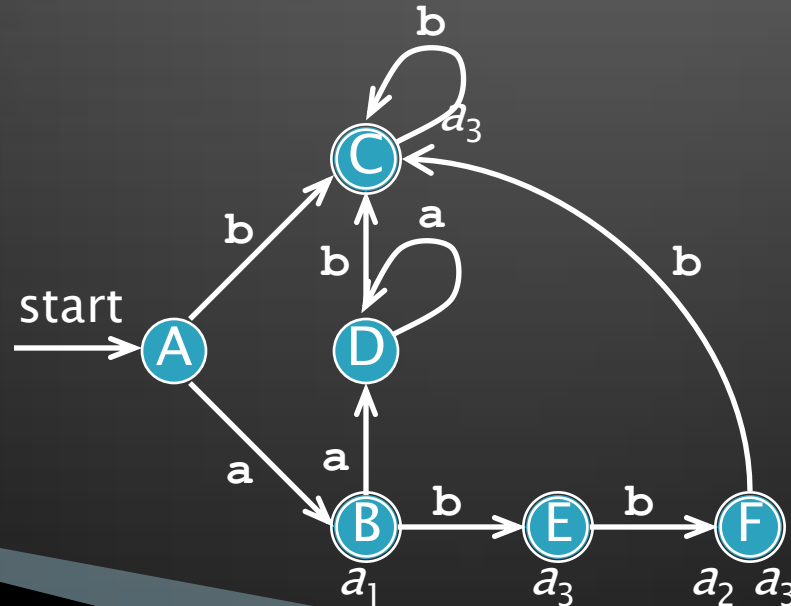
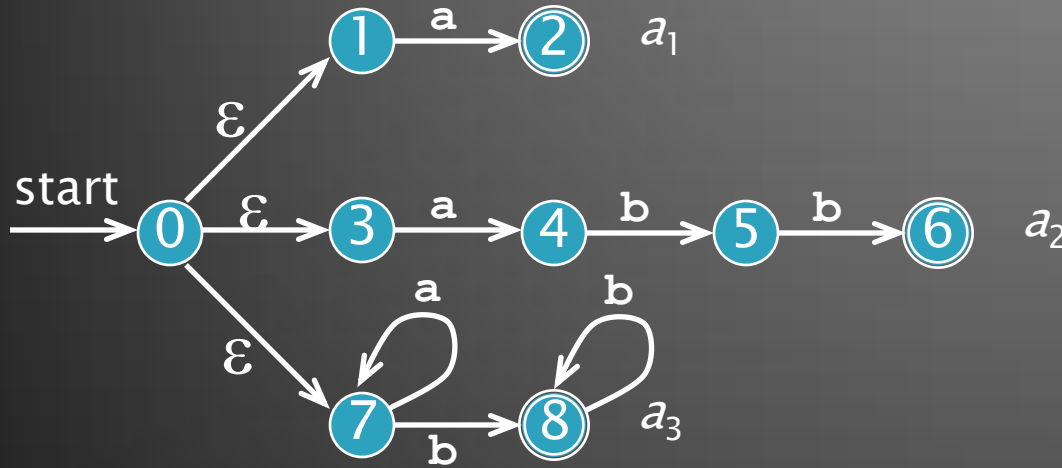
C = {1,2,4,5,6,7}

D = {1,2,4,5,6,7,9}

E = {1,2,4,5,6,7,10}



# Subset Construction Example 2



*Dstates*

A = {0, 1, 3, 7}

B = {2, 4, 7}

C = {8}

D = {7}

E = {5, 8}

F = {6, 8}