



CMP4060 Languages and Compilers

Lexical Analysis – Part2

Ayman AboElHassan, PhD
Assistant Professor

ayman.abo.elmaaty@eng.cu.edu.eg

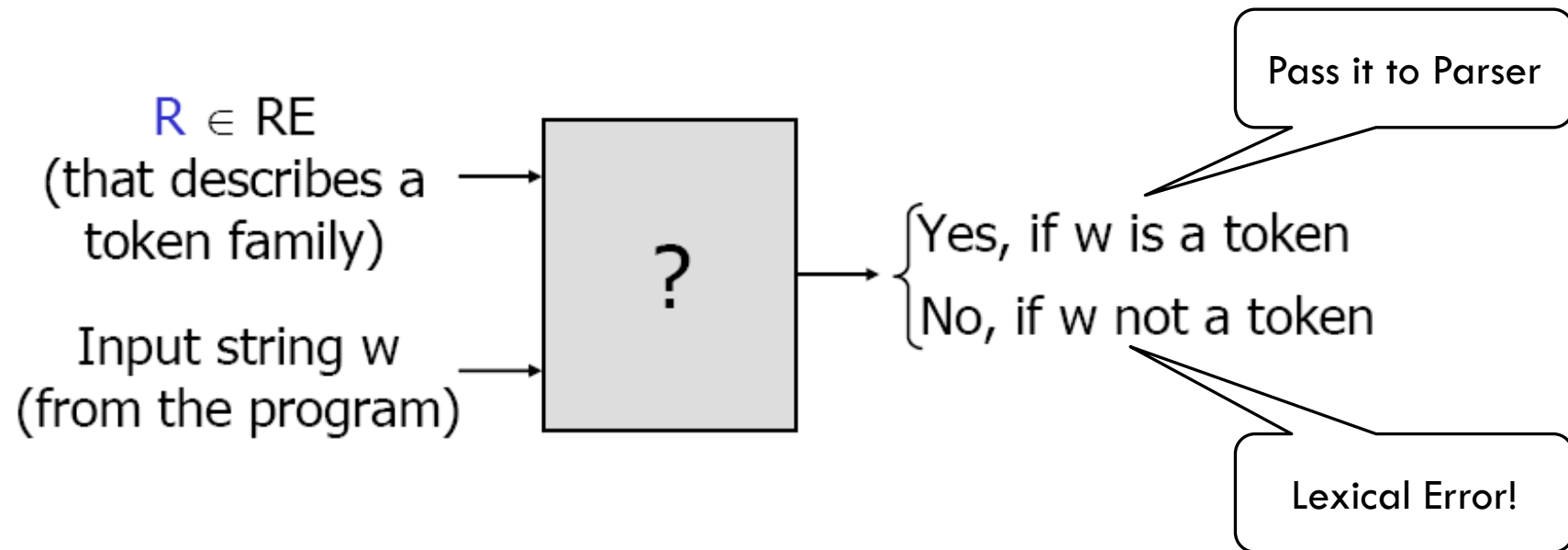


Course Outline

- Introduction to Compilers
- Lexical Analysis: Regular Grammars
- Lexical Implementation: Finite Automata
- Syntax Analysis: Context-Free Grammars
- Parser Implementation: Top-Down Parsers
- Parser Implementation: Bottom-Up Parsers
- Semantic Analysis
- Code Generation
- Code Optimization

How to use RE in Lexical Analyzer?

Given $R \in \text{RE}$ and input string w , need a mechanism to determine if $w \in L(R)$





Lexical Structure



Defining Lexical Structure

Step1: What is an acceptable token?

Write a regular expression for the lexemes of each token

- $\text{digit} = [0-9]^+$
- $\text{letter} = [a-zA-Z]^+$
- $\text{number} = \text{digit}^+$
- $\text{identifier} = \text{letter} (\text{letter} + \text{digit} + \text{'_'})^*$
- $\text{keyword} = \text{'if'} + \text{'else'} + \dots$
- $\text{openpar} = \text{'('}$
- \dots



Defining Lexical Structure

Step2: What regular expressions to use?

Construct R , matching all lexemes for all tokens

- $R = \text{keyword} + \text{identifier} + \text{number} + \dots$
 $= R1 + R2 + \dots$



Defining Lexical Structure

Step3: Check if a substring matches a regular expression?

For input string be $x_1 \dots x_n$

Check if substring $x_1 \dots x_i \in L(R)$ for $i \in \{1, 2, \dots, n\}$

If success:

- Mark substring $x_1 \dots x_i \in L(R)$
- Remove $x_1 \dots x_i$ from input string
- Repeat the operation on the remaining string

Lexical Implementation
performs the check



Defining Lexical Structure

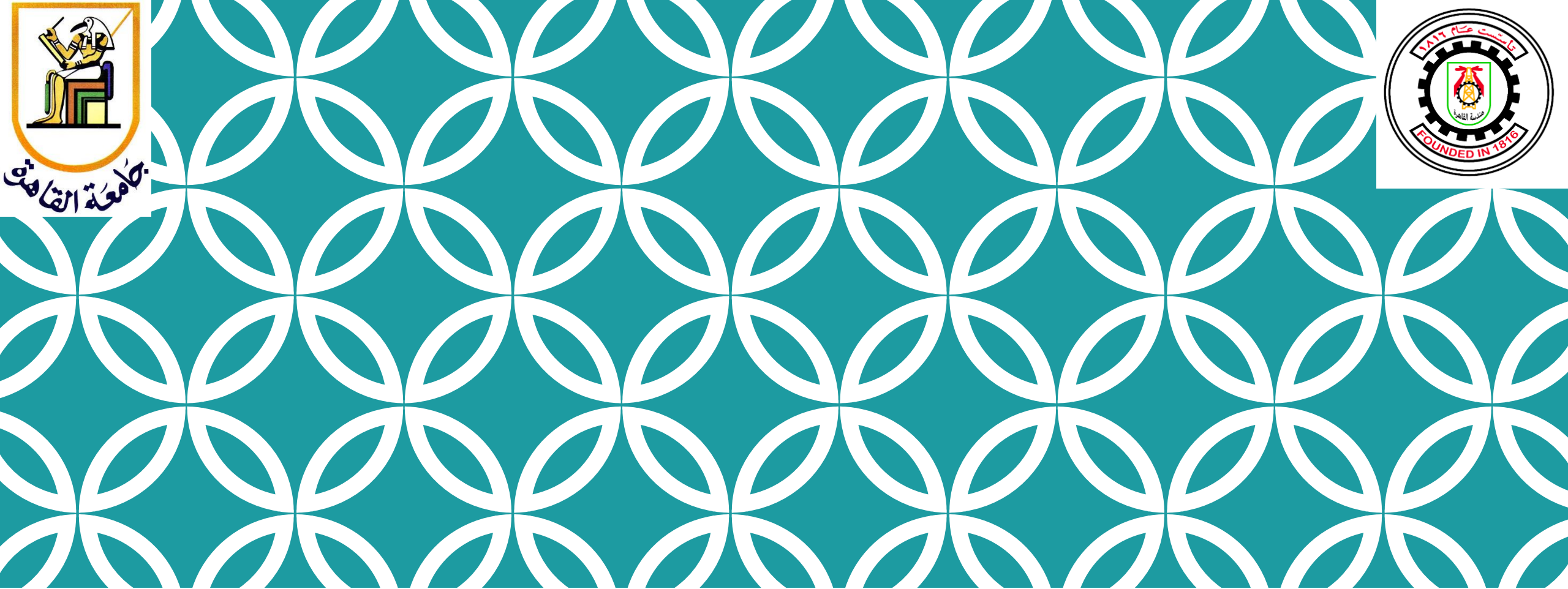
Lexical implementation needs to handle ambiguities

Token Length

- If 2 substrings $x_1 \dots x_i$ & $x_1 \dots x_k$ match the same regex $L(R)$
- Solution: Pick longest possible substring matching

Which Token

- If the same substring $x_1 \dots x_i$ matches 2 regex $L(R_m)$ & $L(R_n)$
- Solution: Pick the first regex in the list



Finite Automata



Finite Automata

Regular Expression → Specification “Rules”

Finite Automata → Implementation



Finite Automata

A finite automaton consists of

- An input alphabet Σ
- A set of states S
- A start state n
- A set of accepting states $F \subseteq S$
- A set of transitions: $state \xrightarrow{input} state$

Finite Automata (State Graph)

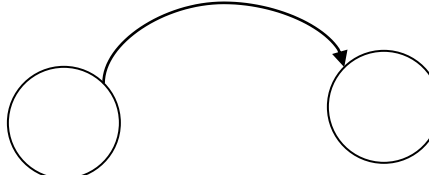
A finite automaton consists of

- An input alphabet Σ

- A set of states $S \rightarrow$ 

- A start state $n \rightarrow$ 

- A set of accepting states $F \subseteq S \rightarrow$ 

- A set of transitions: $state \xrightarrow{input} state \rightarrow$ 



Finite Automata (State Graph)

Example

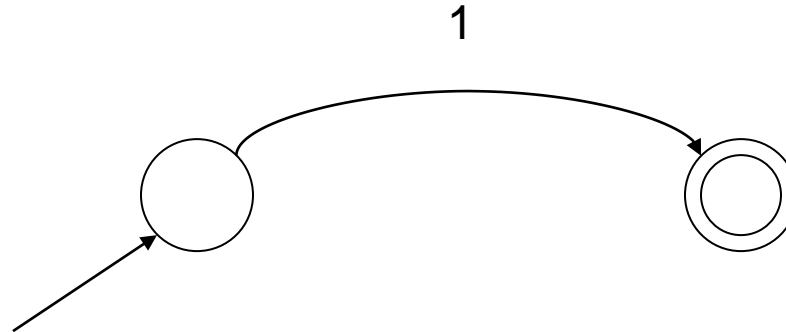
A finite automaton that accepts only “1”

Finite Automata (State Graph)

Example

A finite automaton that accepts only “1”

- $L(R) = '1'$





Finite Automata (State Graph)

Example

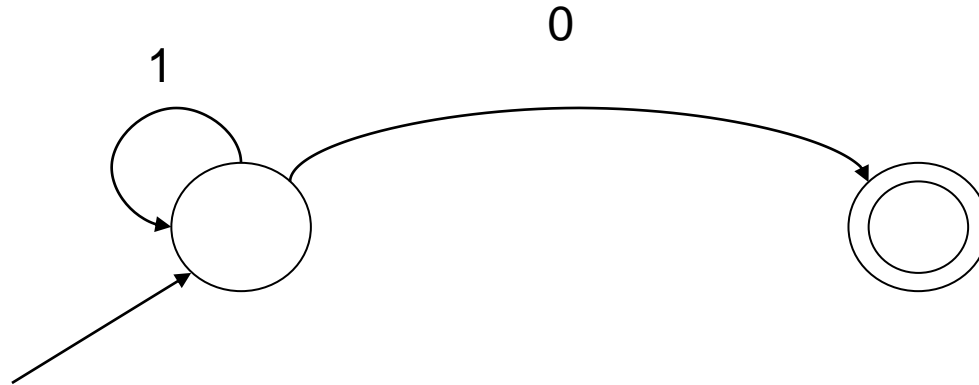
A finite automaton that any number of 1s followed by a single 0

Finite Automata (State Graph)

Example

A finite automaton that any number of 1s followed by a single 0

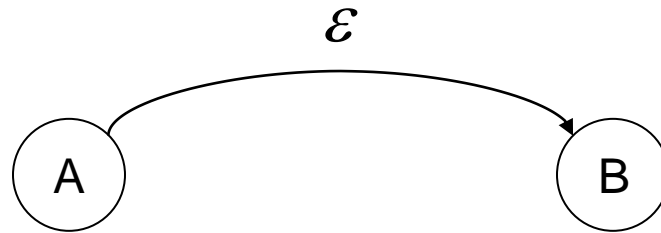
- $L(R) = (1)^*0$



Finite Automata (State Graph)

Epsilon Moves

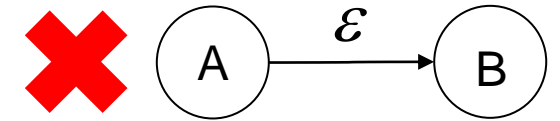
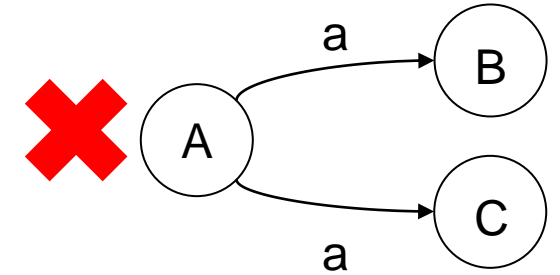
- Another kind of transition
- Machine can move from state A to state B without reading input



Finite Automata (State Graph)

Deterministic Finite Automata (DFA)

- One transition per input per state
- No epsilon-moves



Nondeterministic Finite Automata (NFA)

- Can have multiple transitions for one input in a given state
- Can have epsilon-moves

NFAs & DFAs recognize the same set of languages (regular languages)



NFA vs DFA

DFA can take only one path through the state graph

- Completely determined by input
- Faster to execute (There are no choices to consider)

NFAs can choose

- Whether to make epsilon-moves
- Which of multiple transitions for a single input to take
- Rule: NFA accepts if it can get to a final state



NFA vs DFA

For a given language NFA can be simpler than DFA

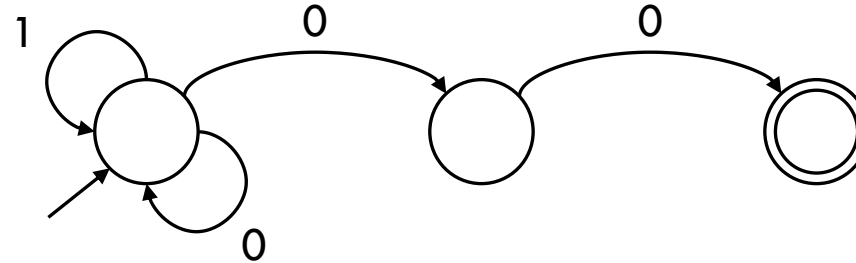
- Example: $L(R) = (1 | 0)^*00$

NFA vs DFA

For a given language NFA can be simpler than DFA

- Example: $L(R) = (1 \mid 0)^*00$

NFA

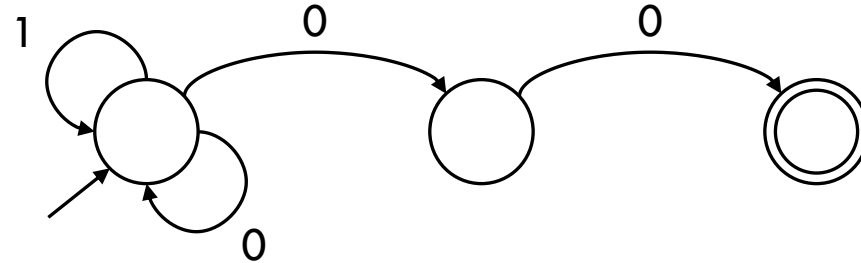


NFA vs DFA

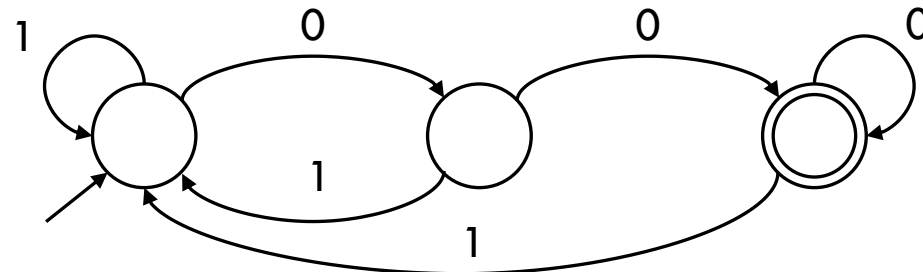
For a given language NFA can be simpler than DFA

- Example: $L(R) = (1 \mid 0)^*00$

NFA



DFA





NFA vs DFA

For a given language NFA can be simpler than DFA

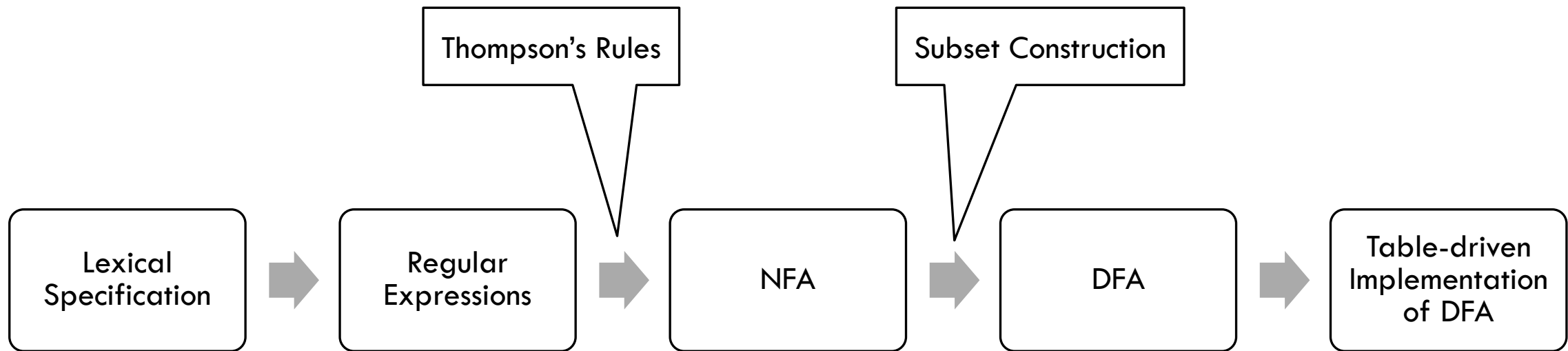
DFA can be exponentially larger than NFA

- It is possible to convert an NFA to a DFA

AUTOMATON	INITIAL	PER STRING
NFA	$O(r)$	$O(r \times x)$
DFA typical case	$O(r ^3)$	$O(x)$
DFA worst case	$O(r ^2 2^{ r })$	$O(x)$



From Regular Expressions to Finite Automata



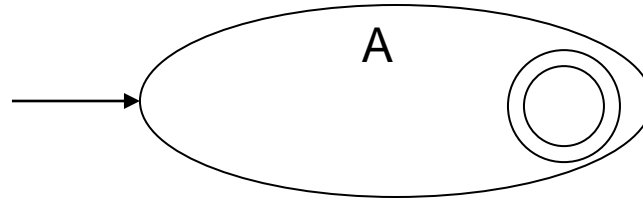


Thompson's Rules

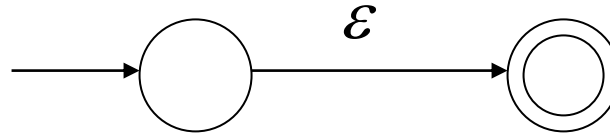
Thompson's Rules

For each kind of regular expression, define an NFA

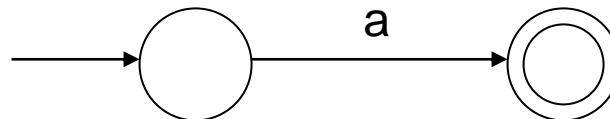
For A



For ε



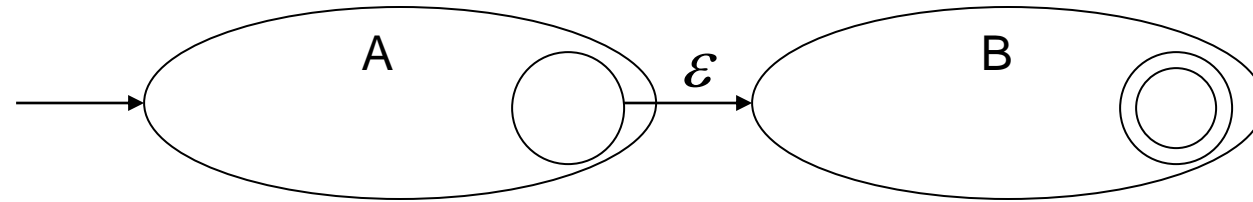
For input a



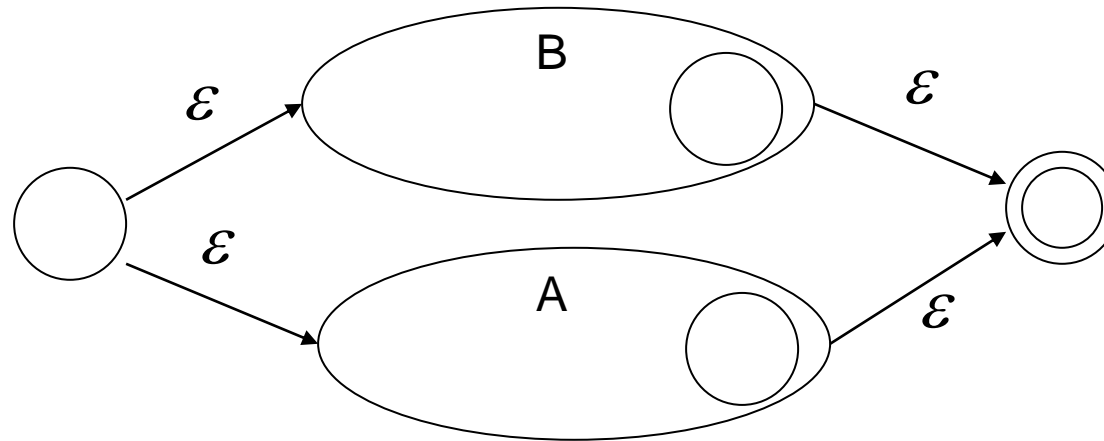
Thompson's Rules

For each kind of regular expression, define an NFA

For AB



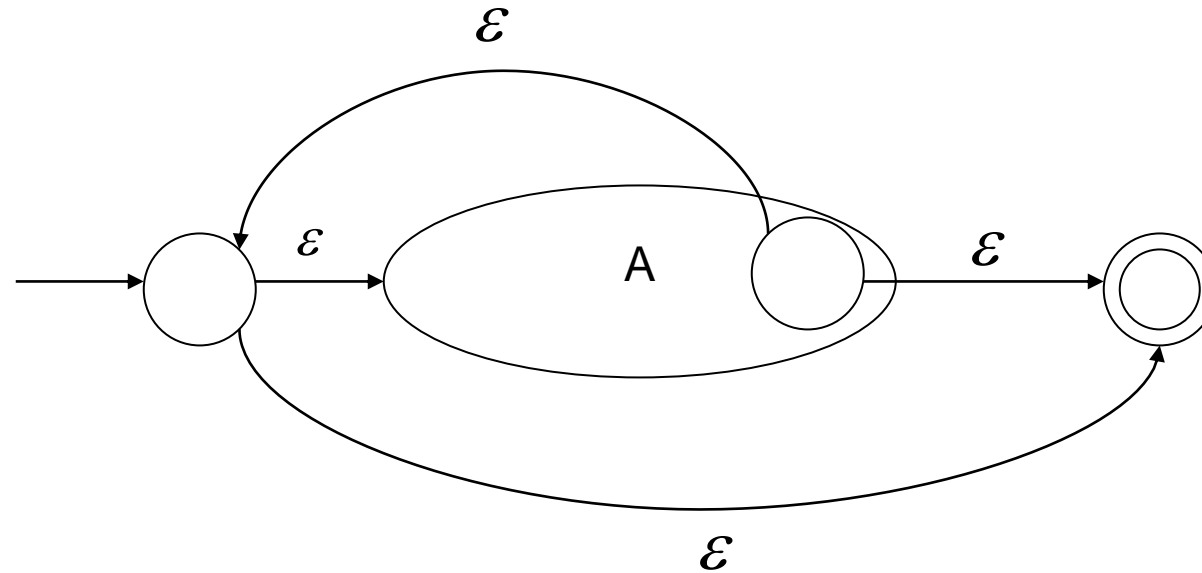
For $A+B$



Thompson's Rules

For each kind of regular expression, define an NFA

For A^*



Thompson's Rules



Example: $(1 + 0) * 1$



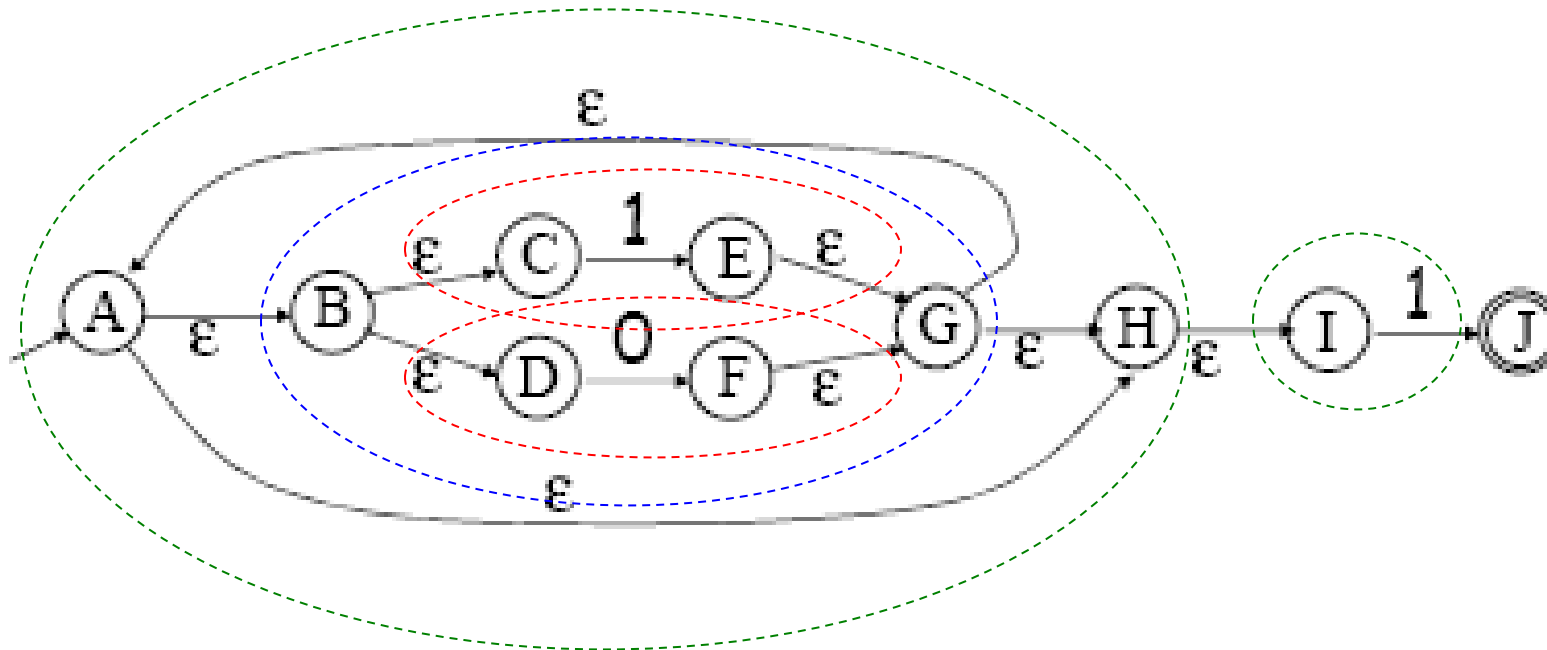
Thompson's Rules

Example: $(1 + 0)^* 1$

$$\begin{array}{c} (\underbrace{1}_{M_1} + \underbrace{0}_{M_2})^* \underbrace{1}_{M_3} \\ \underbrace{\hspace{10em}}_{M_4} \\ \underbrace{\hspace{15em}}_{M_5} \end{array}$$

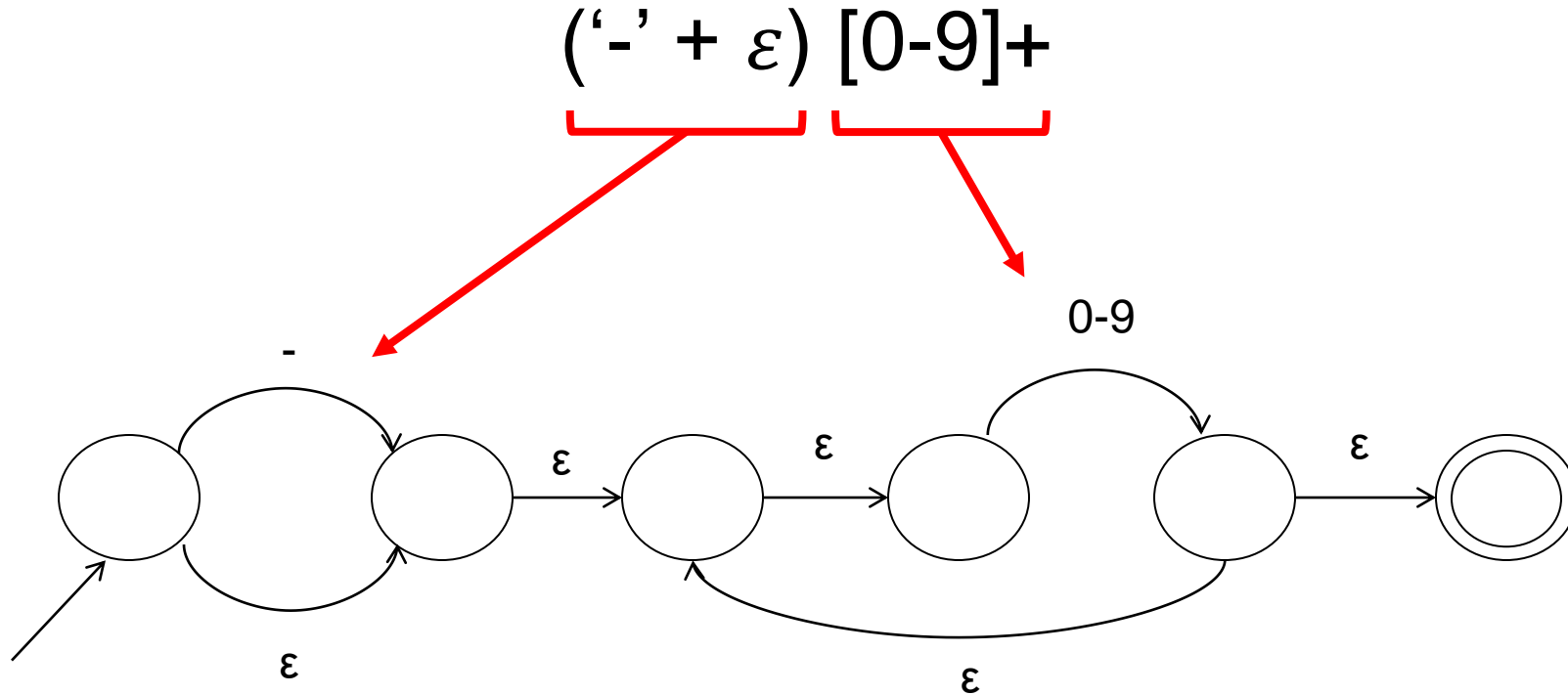
Thompson's Rules

Example: $(1+0)^*1$



Thompson's Rules

Example: $-?[0-9]^+$





Subset Construction Algorithm

Subset Construction Algorithm

Each state of DFA (D) is a set of NFA (N) states

Simulate “in parallel” all possible moves NFA can make on a given input string

- s is a single state of N
- T is a set of states of N

OPERATION	DESCRIPTION
$\epsilon\text{-closure}(s)$	Set of NFA states reachable from NFA state s on ϵ -transitions alone.
$\epsilon\text{-closure}(T)$	Set of NFA states reachable from some NFA state s in set T on ϵ -transitions alone; $= \cup_{s \in T} \epsilon\text{-closure}(s)$.
$\text{move}(T, a)$	Set of NFA states to which there is a transition on input symbol a from some state s in T .



Subset Construction Algorithm

Each state of DFA \rightarrow a non-empty subset of states of the NFA

1. Start state = ε -closure(S_0)

- S_0 is its start state
- The set of NFA states reachable through ε -moves from NFA start state

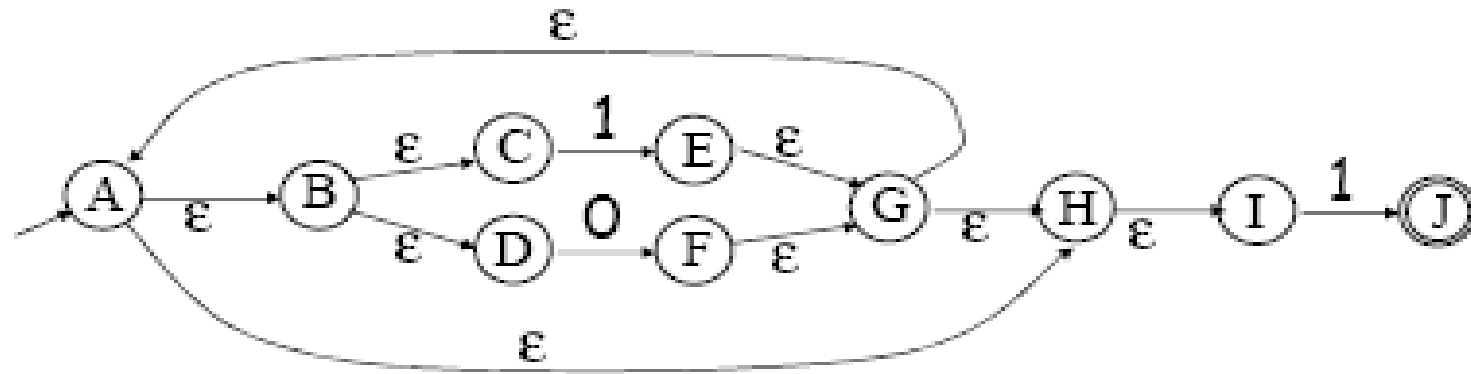
2. $T' = \varepsilon$ -closure(move(T, a))

- T' is the set of NFA states reachable from any state in T after reading the input a , considering ε -moves as well

3. Add a transition $T \xrightarrow{a} T'$

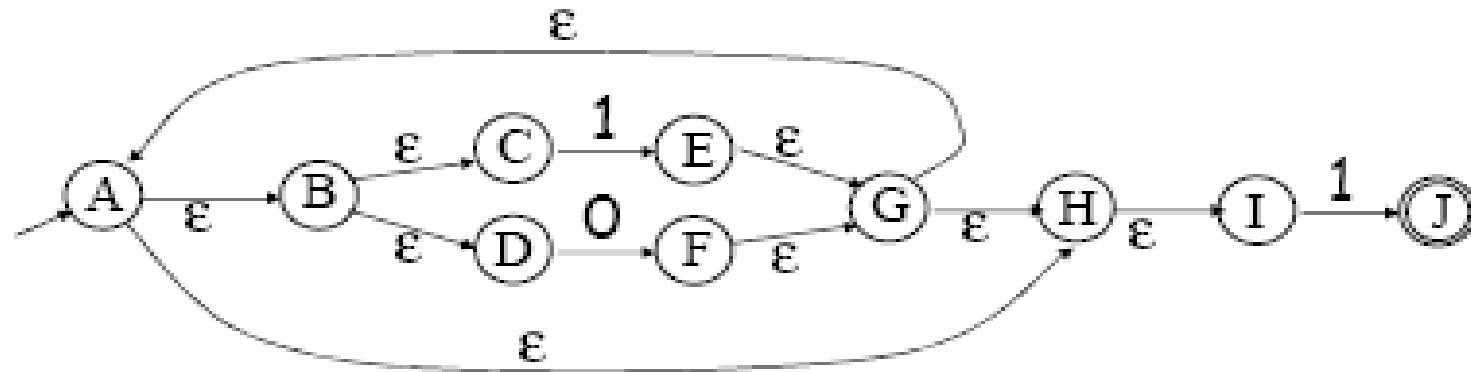
Subset Construction Algorithm

Example: $(1+0)^*1$



Subset Construction Algorithm

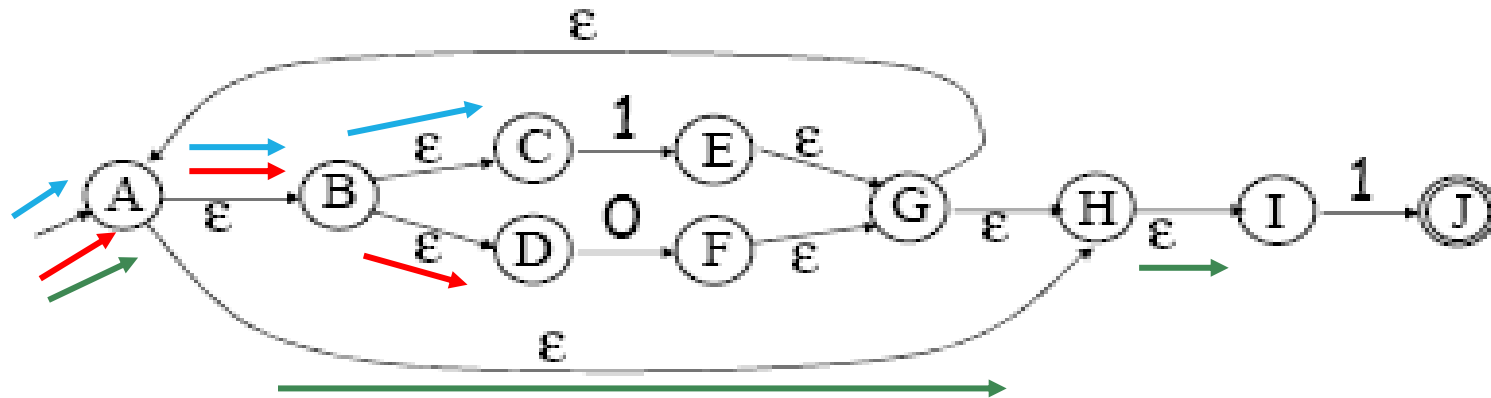
Example: $(1+0)^*1$



Step1: Get start state “reachable through ϵ -moves”

Subset Construction Algorithm

Example: $(1+0)^*1$

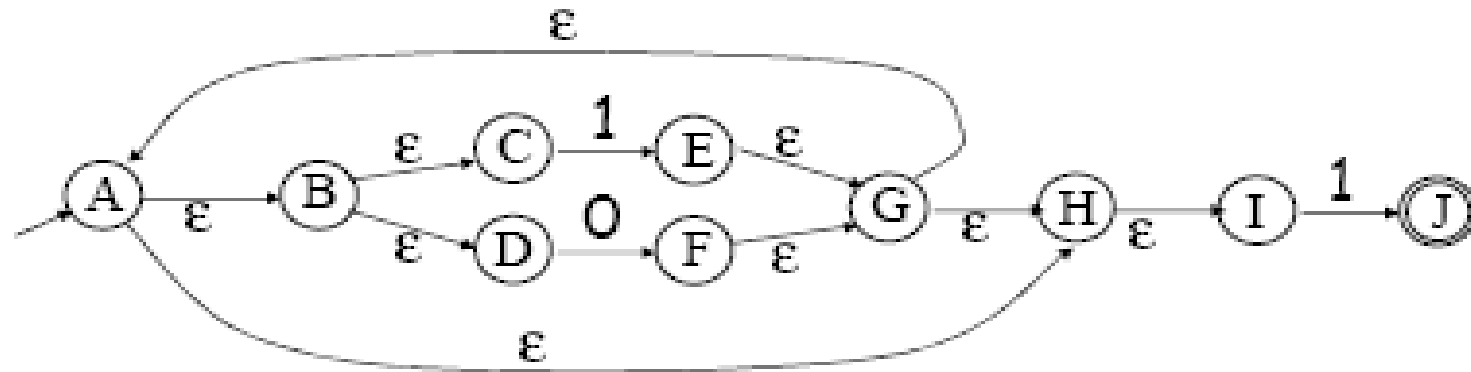


Step 1: Get start state “reachable through ϵ -moves”

ABCDHI

Subset Construction Algorithm

Example: $(1+0)^*1$

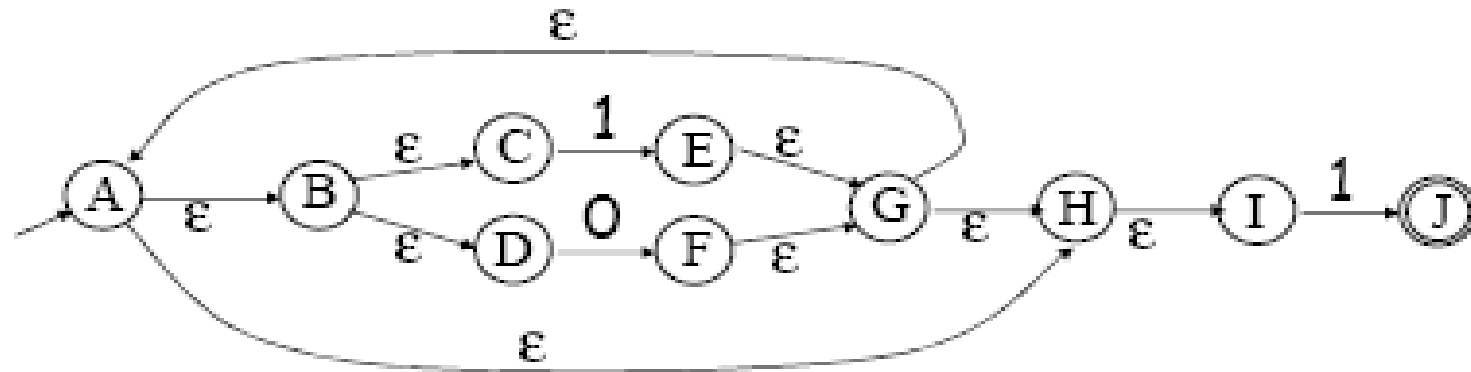


Step2: States T' reachable from T by a (& ϵ)

ABCDHI

Subset Construction Algorithm

Example: $(1+0)^*1$



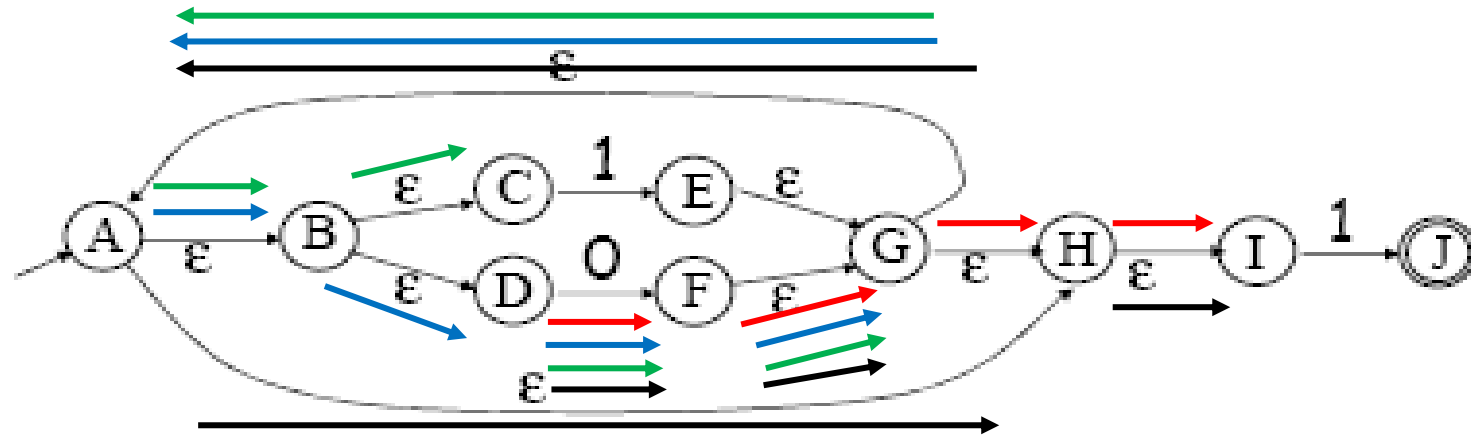
Step2: States T' reachable from T by a (& ϵ)

- $T = \{A, B, C, D, H, I\}$
- $a = 0$

ABCDHI

Subset Construction Algorithm

Example: $(1+0)^*1$



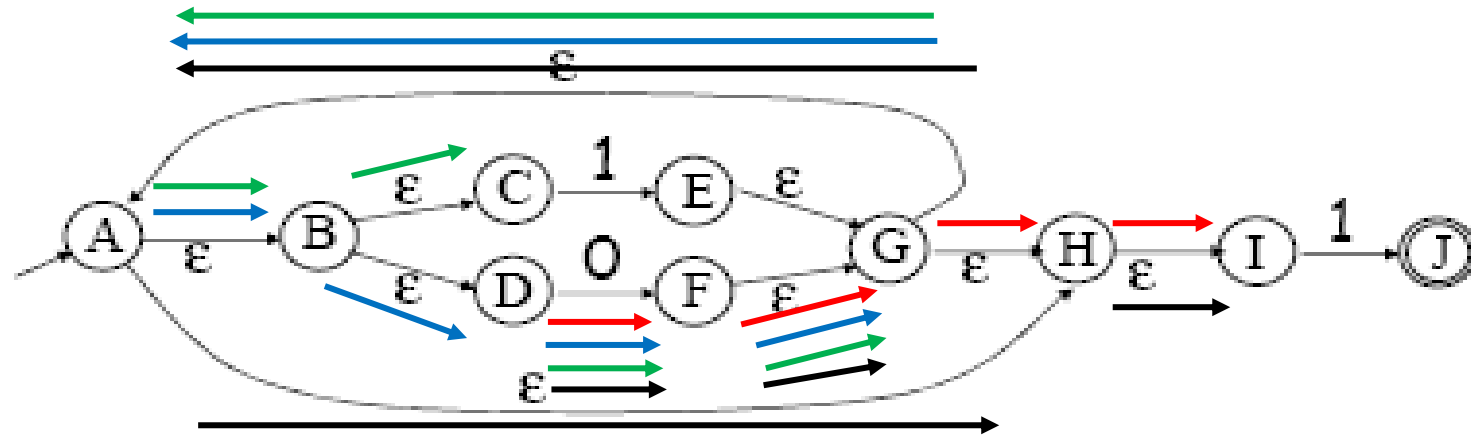
Step2: States T' reachable from T by a (& ϵ)

- $T = \{A, B, C, D, H, I\}$
- $a = 0$

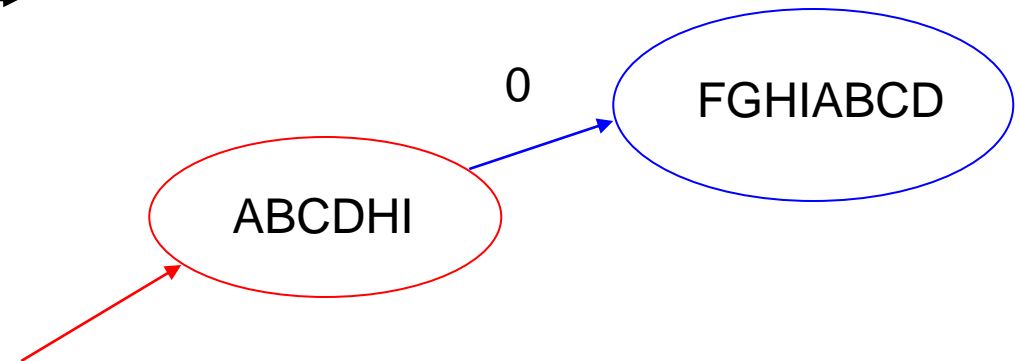
ABCDHI

Subset Construction Algorithm

Example: $(1+0)^*1$

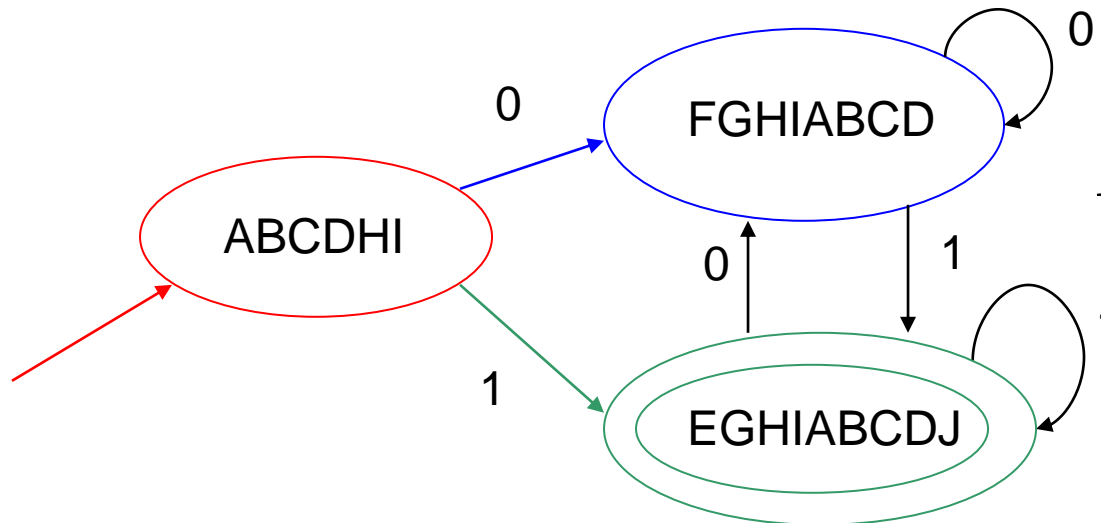
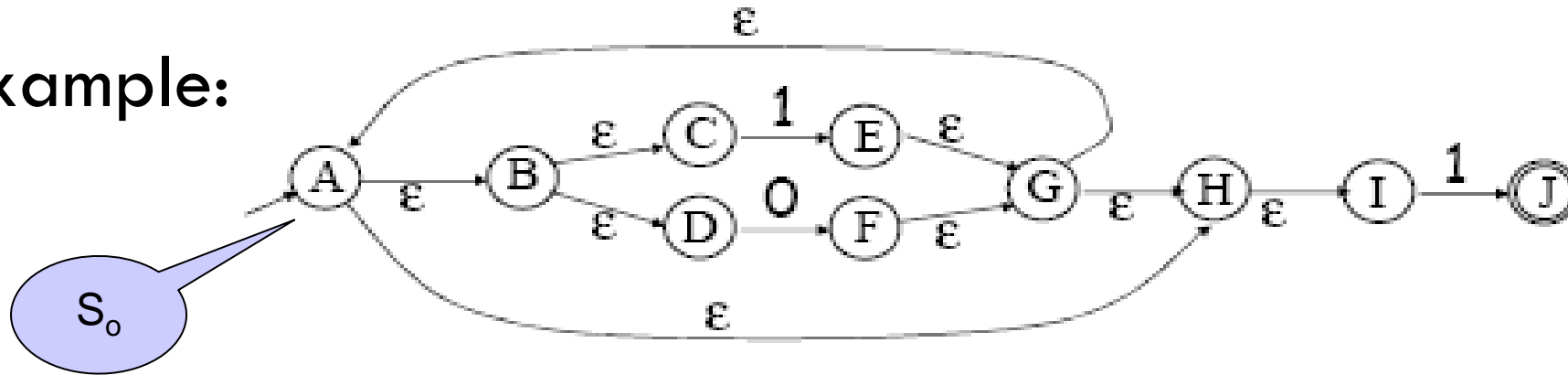


Step3: Add a transition $T \xrightarrow{a} T'$



Subset Construction Algorithm

Example:

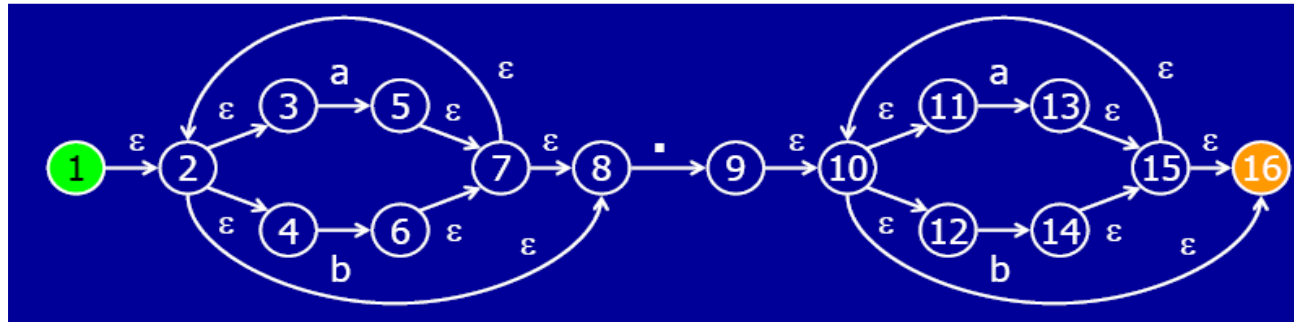


DFA Graph:

- No epsilon moves
- Single path per input

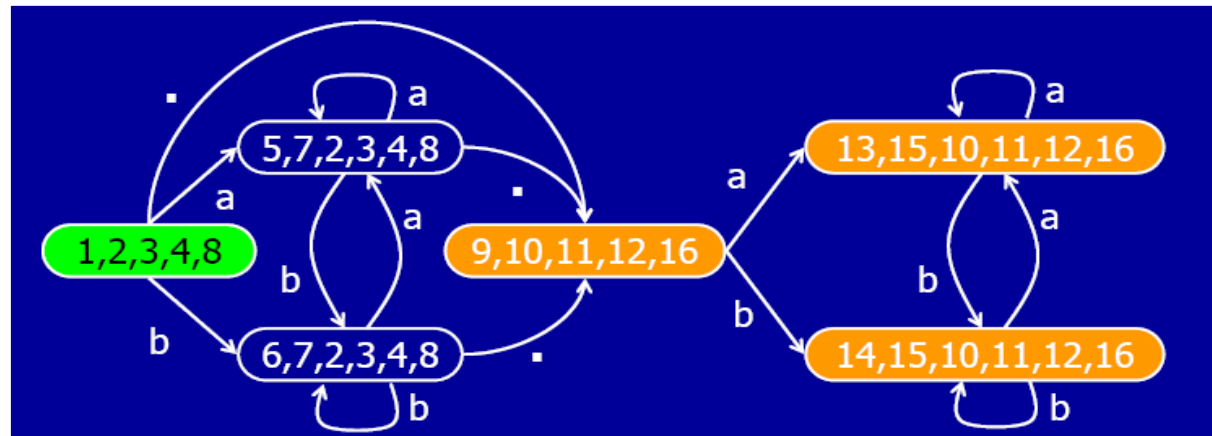
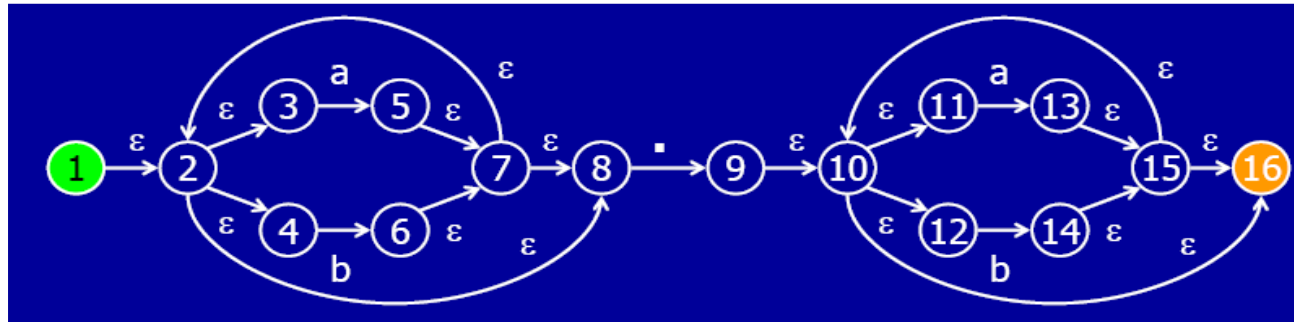
Subset Construction Algorithm

Example $(a+b)^*. (a+b)$



Subset Construction Algorithm

Example $(a+b)^*. (a+b)$





DFA Table Implementation



DFA Table Implementation

A DFA can be implemented by a 2D table T

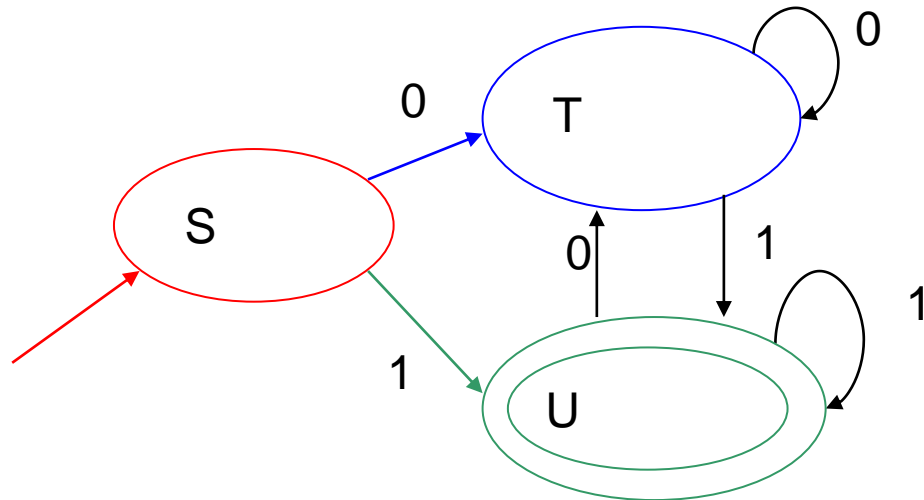
- One dimension is **states**
- Other dimension is **input symbol**
- For every transition $S_i \xrightarrow{a} S_k$ define $T[i,a] = k$

DFA execution

- If in state S_i and input a , read $T[i,a] = k$ and skip to state S_k
- Very efficient

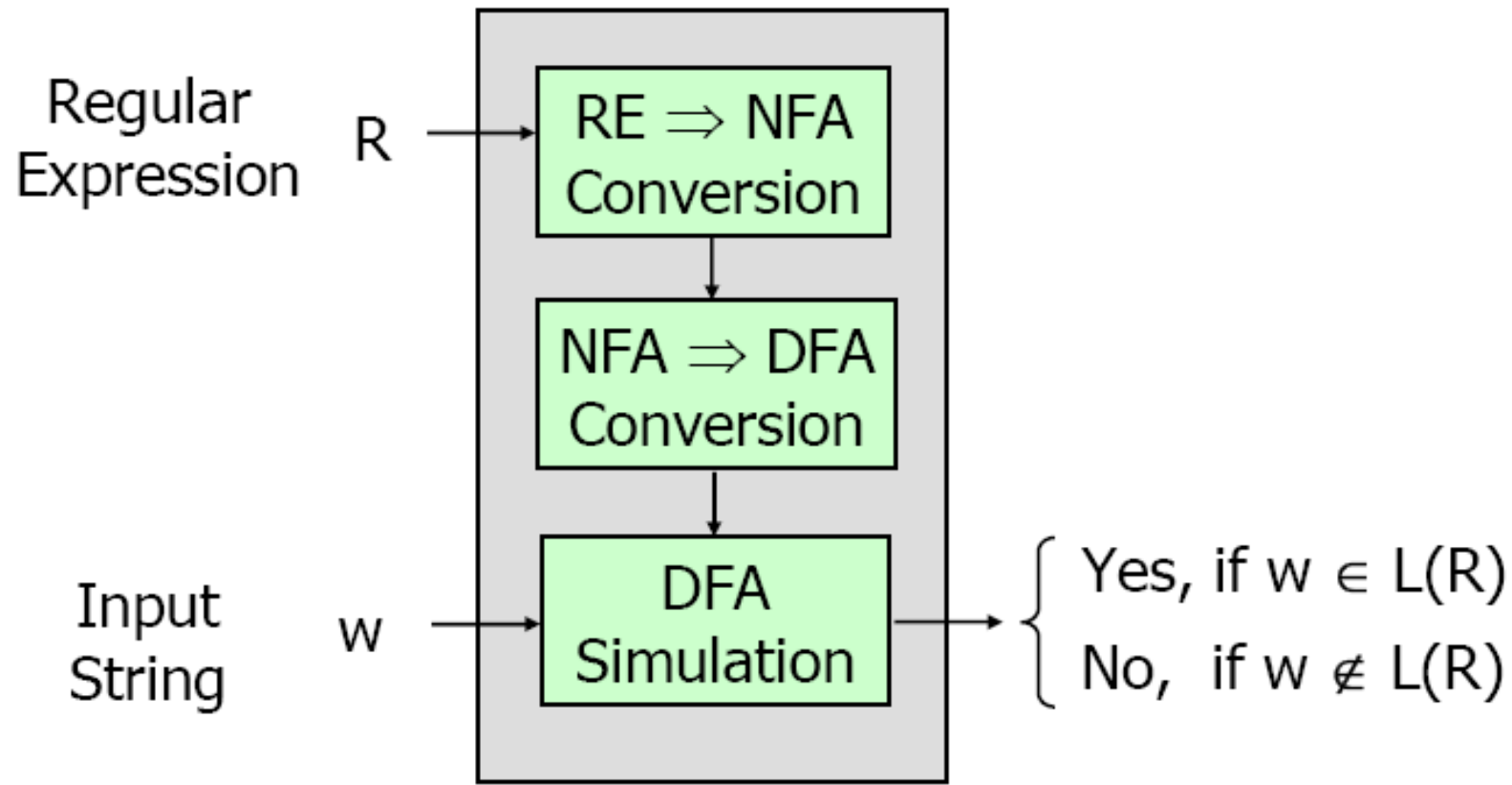
DFA Table Implementation

Example



	0	1
S	T	U
T	T	U
U	T	U

Putting the Pieces Together



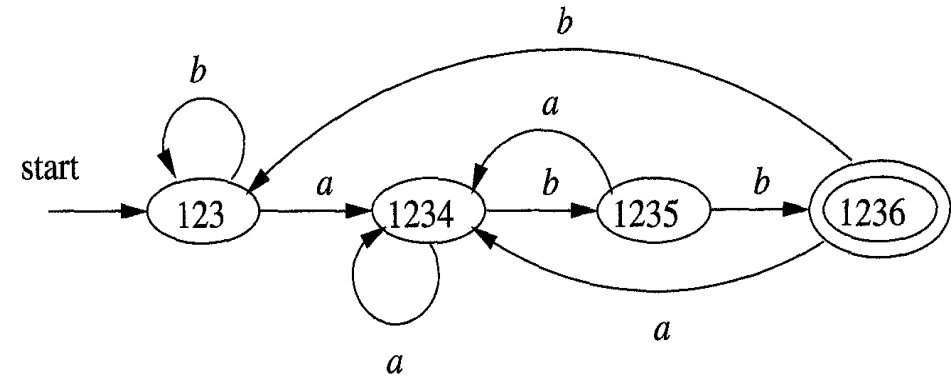
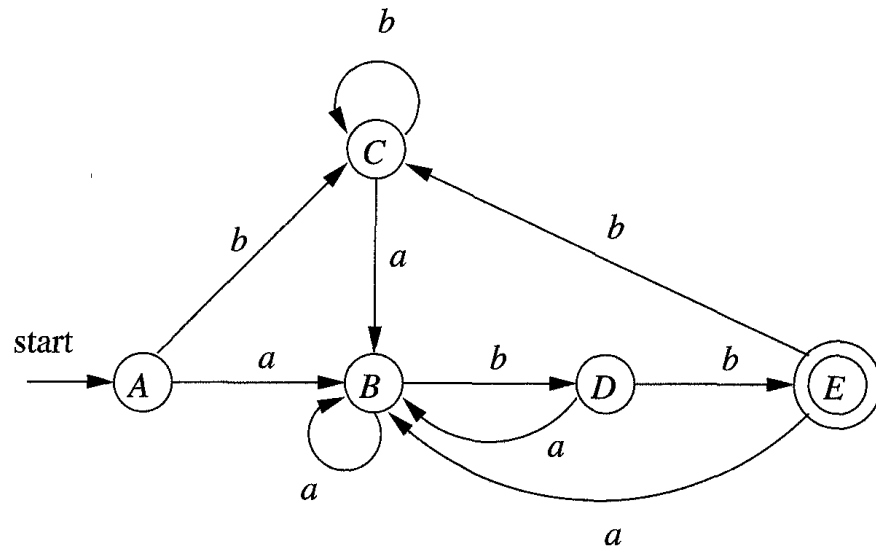


DFA Minimization

DFA Minimization

Multiple DFAs can be used to recognize the same language

Example: $(a+b)^*abb$





DFA Minimization

These automata don't have the same number of states

When implementing a lexical analyzer as a DFA

- Generally, prefer a DFA with as few states as possible
- Each state requires an entry in the DFA table implementation

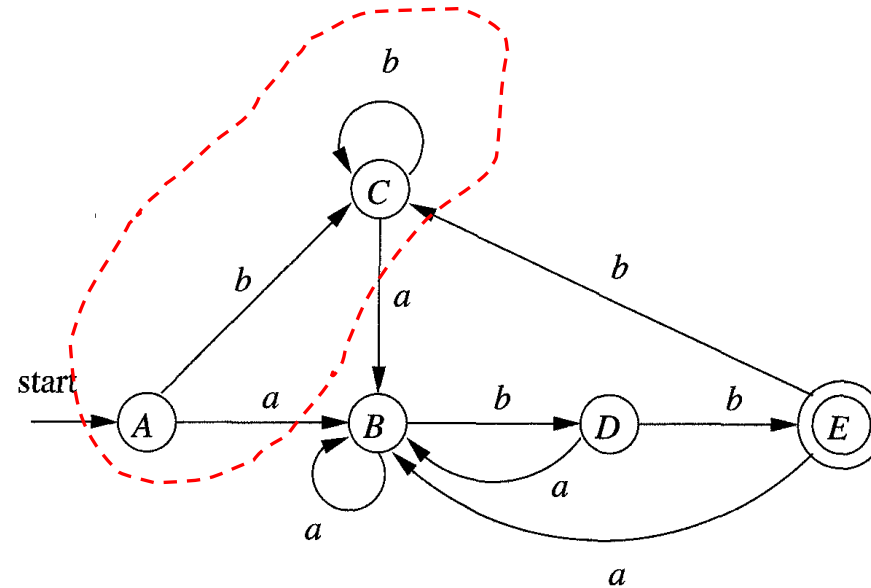
DFA Minimization

State A and C are equivalent

- On any input they transfer to the same state

- $A \xrightarrow{a} B, C \xrightarrow{a} B$

- $A \xrightarrow{b} C, C \xrightarrow{b} C$

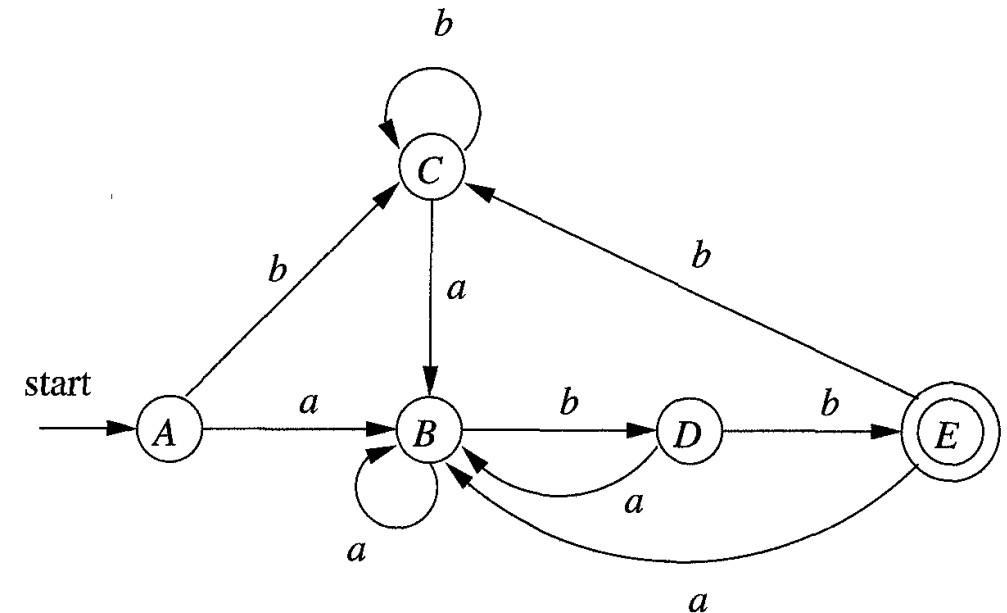


DFA Minimization

Step 1: Partition graph π into 2 groups (accepting, non-accepting)

Example

- $\{A, B, C, D\} \rightarrow$ non-accepting state
- $\{E\} \rightarrow$ accepting state



DFA Minimization

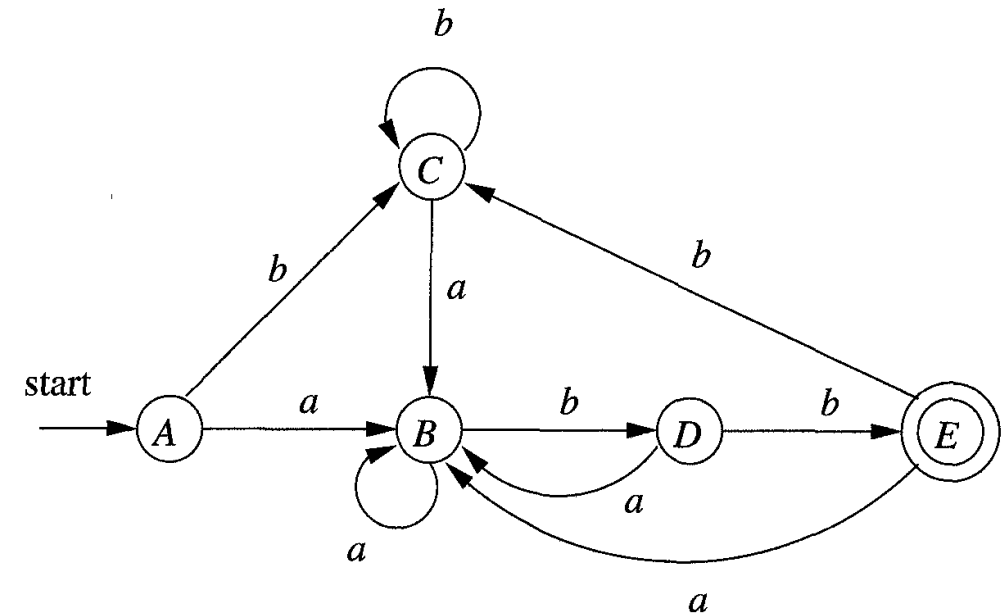
Step 1: Partition graph π into 2 groups (accepting, non-accepting)

Example

- $\{A, B, C, D\} \rightarrow$ non-accepting state
- $\{E\} \rightarrow$ accepting state

Construct a new graph

- $\pi_{\text{new}} = \{ \{A, B, C, D\} \{E\} \}$



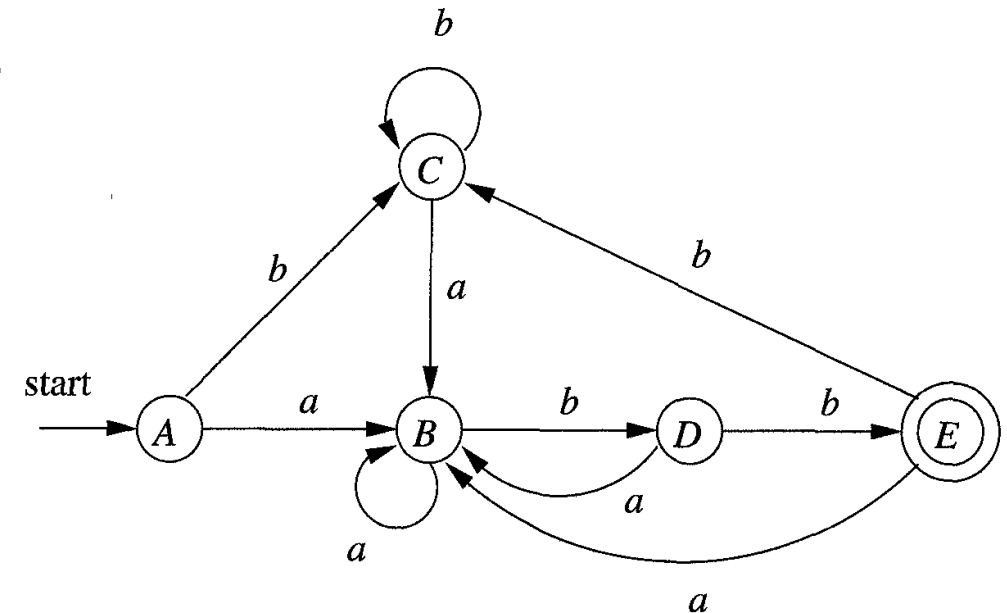
DFA Minimization

Step2: For each group G in π_{new}

- Split group G into subgroups
 - 2 states s & t in same subgroup iff for all input symbol a , state s & t have transitions on a to states in the same group

Example

- For $\{A, B, C, D\}$
 - On input $a \rightarrow \{A, B, C, D\}$ go to B
 - On input $b \rightarrow \{A, B, C\}$ go to $\{A, B, C, D\}$
 $\rightarrow \{D\}$ go to $\{E\}$



$$\pi_{\text{new}} = \{ \{A, B, C\} \{D\} \{E\} \}$$

DFA Minimization

Step3: Repeat Step2 till converge

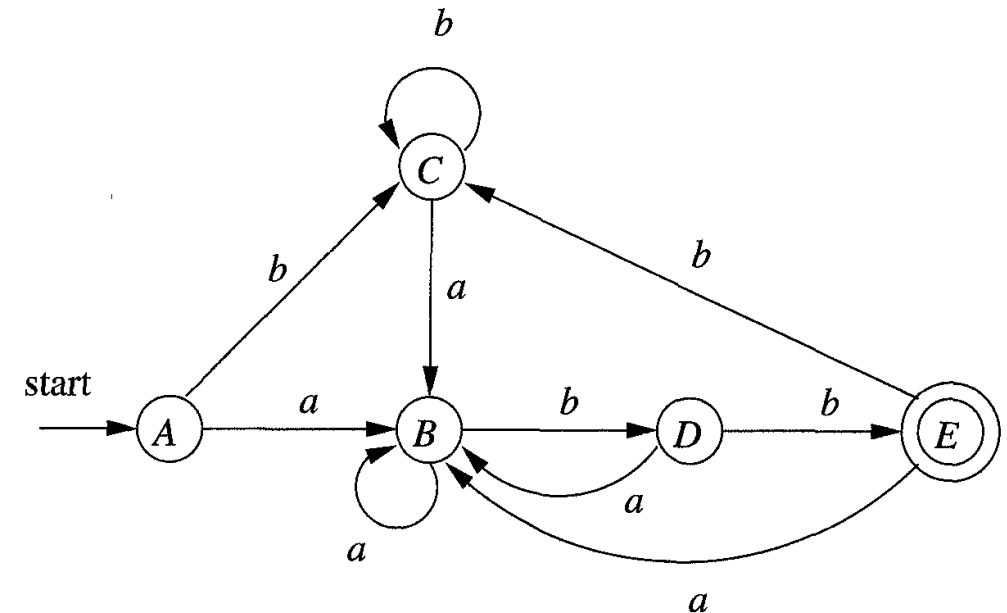
Example

- For $\{A, B, C\}$
 - On input $b \rightarrow \{A, C\}$ go to $\{A, B, C\}$
 $\rightarrow \{B\}$ go to $\{D\}$

$$\pi_{\text{new}} = \{ \{A, C\} \{B\} \{D\} \{E\} \}$$

- For $\{A, C\}$
 - On input $a \rightarrow \{A, C\}$ go to $\{A, C\}$
 - On input $b \rightarrow \{A, C\}$ go to $\{A, C\}$

$$\pi_{\text{new}} = \{ \{A, C\} \{B\} \{D\} \{E\} \} = \pi_{\text{final}}$$



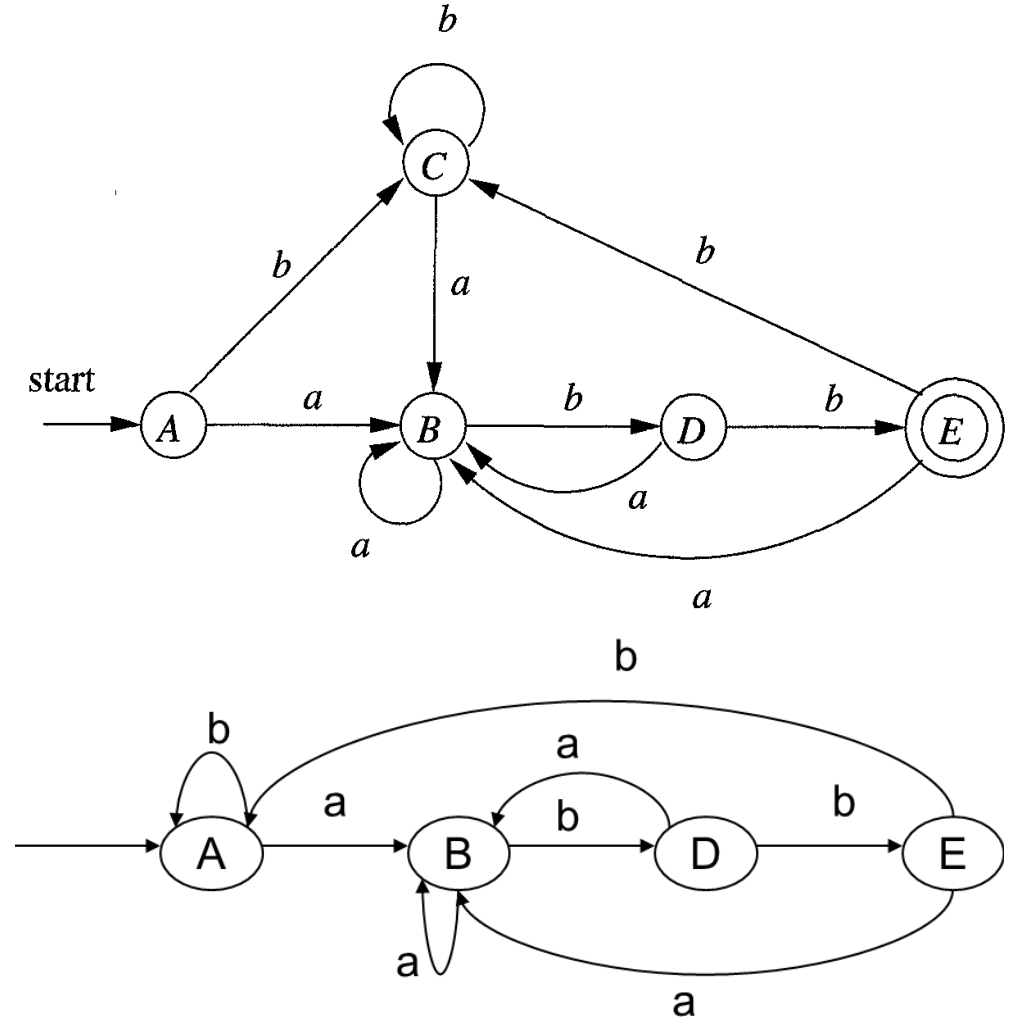
DFA Minimization

Step4: Choose 1 state in each group
in π_{final}

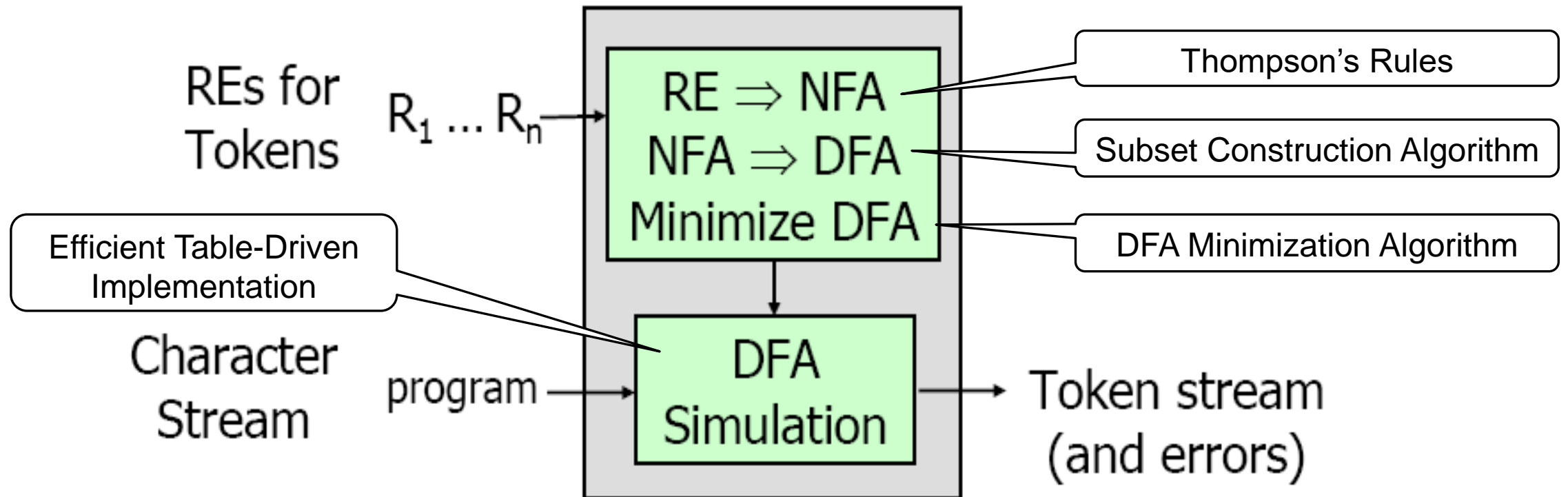
- Start State
- Accepting State
- Transition from that state to another outside on input a

Example

- $\pi_{\text{final}} = \{ \{A, C\} \{B\} \{D\} \{E\} \}$
- Pick A



Summary





Thank you