

Regular Expression

Beulah A.
AP/CSE

Introduction

- Regular expressions describe regular languages
- ie the language accepted by a finite automata are easily described by regular expression.
- Many programming languages provide regular expression capabilities,
 - Built-in → Perl, JavaScript, Ruby, AWK, Tcl,
 - Standard library → .NET, Java, Python C++
- REs are widely supported in programming languages, text processing programs (particular lexers, lex, yacc), advanced text editors

Introduction

- Let Σ be a finite set of symbols.
- Let L_1, L_2 be set of strings in Σ^* .
- The concatenation of L_1 and L_2 denoted by $L_1 L_2$ is the set of all strings of the form xy , where $x \in L_1$ and $y \in L_2$.
- $L^0 = \{\epsilon\}$
- $L^i = LL^{i-1}$ for $i \geq 1$.

Introduction

- Kleene Closure

$$L^* = \bigcup_{i=0}^{\infty} L^i = L^0 \cup L^1 \cup L^2 \cup \dots$$

- Positive Closure

$$L^+ = \bigcup_{i=1}^{\infty} L^i = L^1 \cup L^2 \cup \dots$$

Example

Let $L_1 = \{10, 01\}$, $L_2 = \{11, 00\}$

Then $L_1 L_2 = \{1011, 1000, 0111, 0100\}$

Let $L = \{10, 11\}$

Then $L^* = L_0 \cup L_1 \cup L_2 \cup \dots$

$= \{\epsilon\} \cup \{10, 11\} \cup \{1011, 1010, 1110, 1111\} \cup \dots$

$= \{\epsilon, 10, 11, 1011, 1010, 1110, 1111, \dots\}$

Operators of RE

$$* \rightarrow L^*$$

$$\cdot \rightarrow L_1 \cdot L_2, L_1 L_2$$

$$/ \rightarrow L_1 \cup L_2$$

Definition of Regular Expression

- Let Σ be an alphabet. The regular expressions over Σ and the sets that they denote are defined recursively as follows:
 1. φ is a regular expression and denotes the empty set $\{\}$.
 2. ε is a regular expression and denotes the set $\{\varepsilon\}$
 3. For each $a \in \Sigma$, 'a' is a regular expression and denotes the set $\{a\}$.
 4. If r and s are regular expressions denoting the languages R and S respectively then $(r + s)$, (rs) , $(r)^*$ are regular expressions that denotes the sets $R \cup S$, RS and R^* respectively.

Precedence of RE operators

* \rightarrow higher precedence

.

/ \rightarrow Lower precedence

Example

- $(0/1)^* = \{\epsilon, 0, 1, 00, 01, 10, 11, \dots\} = (0+1)^*$
(i.e.) all strings of 0 and 1
- $01^* = \{0, 01, 011, 0111, \dots\}$
- $0^* = \{\epsilon, 0, 00, 000, \dots\}$
- $1(1)^* = \{1, 11, 111, 1111, \dots\} = 1^+$

Identities for Regular Expressions

$$\text{I1} \quad \varphi + R = R$$

$$\text{I2} \quad \varphi R = R\varphi = \varphi$$

$$\text{I3} \quad \lambda R = R\lambda = R$$

$$\text{I4} \quad \lambda^* = \lambda$$

$$\text{I5} \quad R + R = R$$

$$\text{I6} \quad R^*R^* = R^*$$

$$\text{I7} \quad RR^* = R^*R$$

$$\text{I8} \quad (R^*)^* = R^*$$

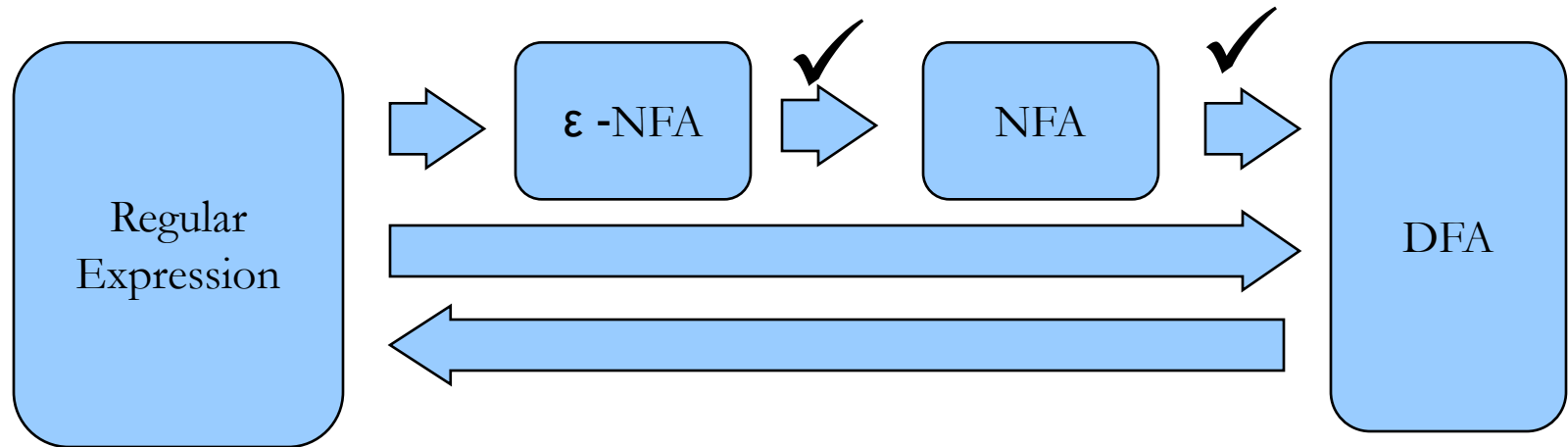
$$\text{I9} \quad \lambda + RR^* = R^* = \lambda + R^*R$$

$$\text{I10} \quad (PQ)^*P = P(QP)^*$$

$$\text{I11} \quad (P + Q)^* = (P^*Q^*)^* = (P^* + Q^*)^*$$

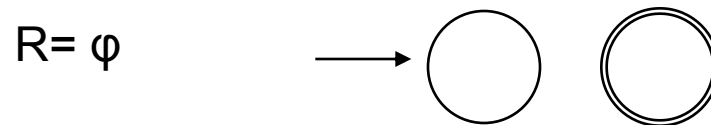
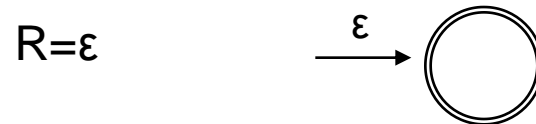
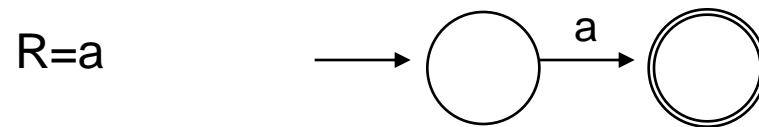
$$\text{I12} \quad (P + Q)R = PR + QR \text{ and} \\ R(P + Q) = RP + RQ$$

Road map

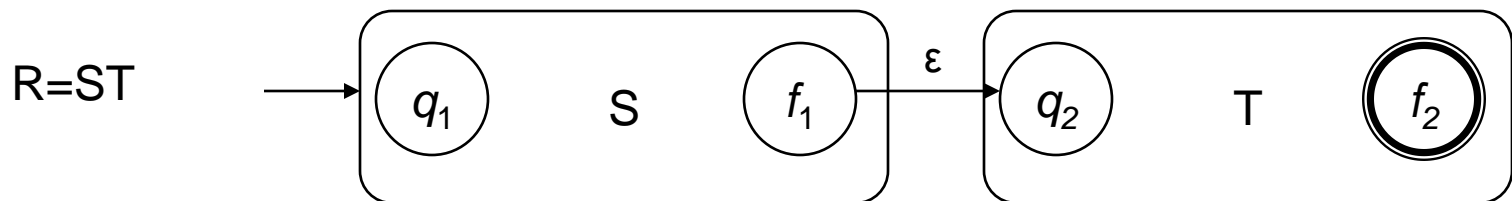
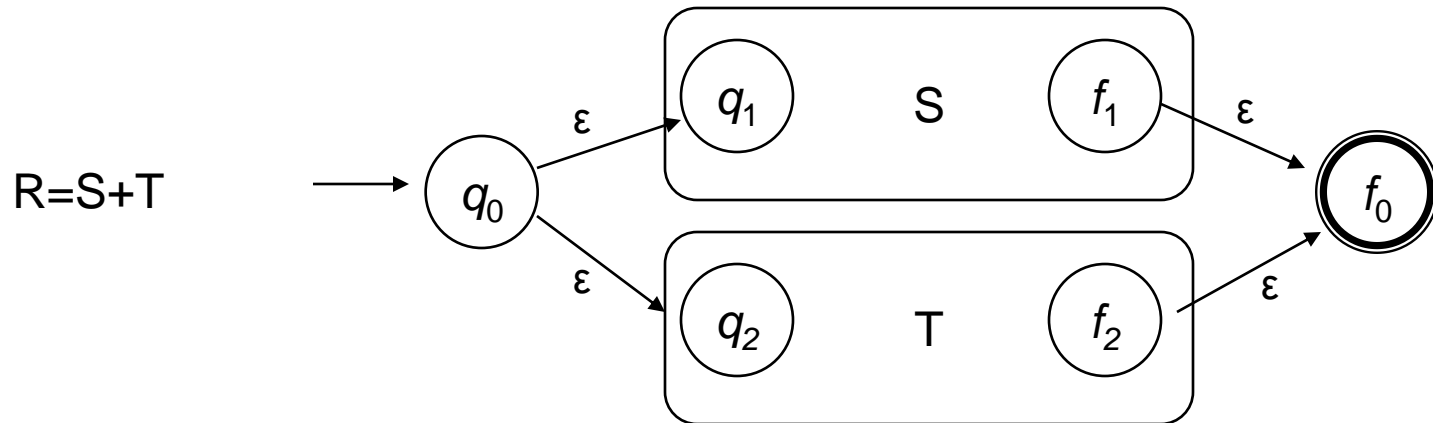


Thompson's Construction

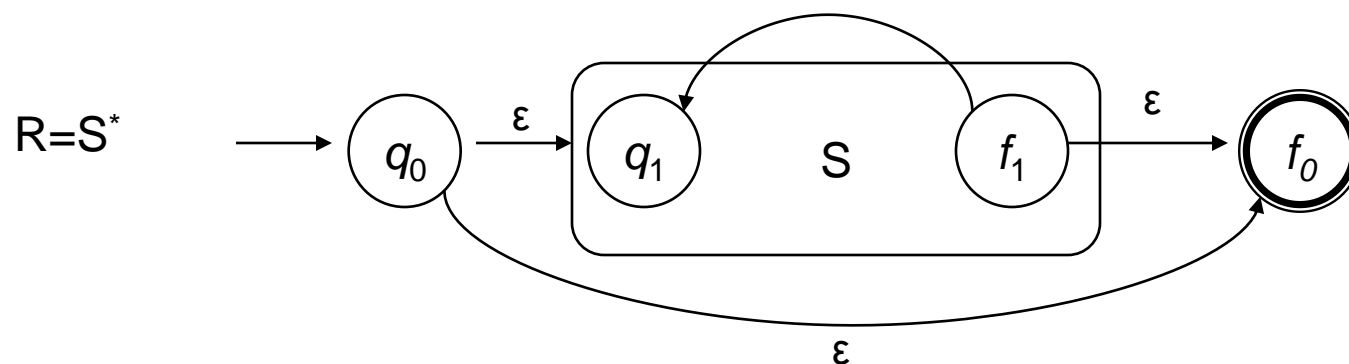
● Basis



Thompson's Construction



Thompson's Construction



Theorem

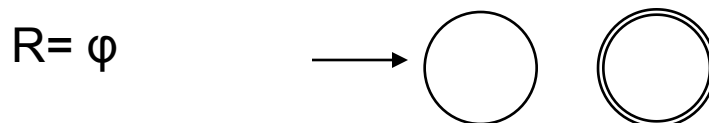
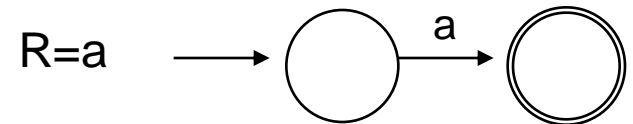
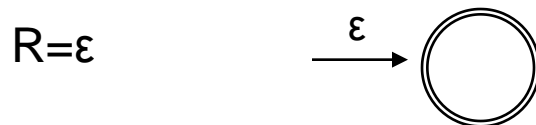
For every regular expression r there exists a NFA with ε -transitions that accepts $L(r)$

- Proof

- **Basis step (Zero operators)**

Suppose r is ε , φ or a for some $a \in \Sigma$.

Then the equivalent NFA's are:



Induction Case i

- $r = r_1 + r_2$
- $M_1 = (Q_1, \Sigma_1, \delta_1, q_1, \{f_1\})$ $L(M_1) = L(r_1)$
- $M_2 = (Q_2, \Sigma_2, \delta_2, q_2, \{f_2\})$ $L(M_2) = L(r_2)$.
- Assume Q_1 and Q_2 are disjoint.
- Let q_0, f_0 be a new initial and final state respectively.

Case i

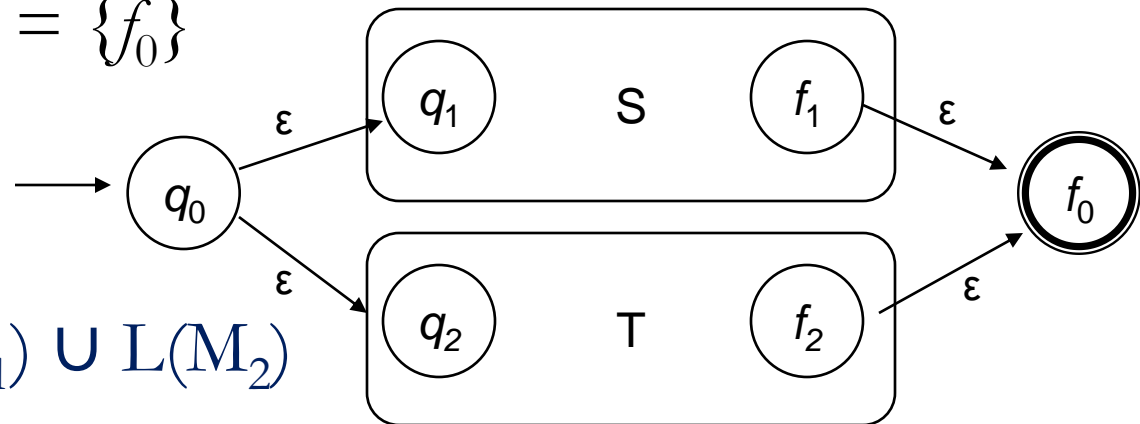
- $M = (Q_1 \cup Q_2 \cup \{q_0, f_0\}, \Sigma_1 \cup \Sigma_2, \delta, q_0, \{f_0\})$ where δ is defined by

$$\delta(q_0, \epsilon) = \{q_1, q_2\}$$

$$\delta(q, a) = \delta_1(q, a) \quad \text{if } q \in Q_1 - \{f_1\}, a \in \Sigma_1 \cup \{\epsilon\}$$

$$\delta(q, a) = \delta_2(q, a) \quad \text{if } q \in Q_2 - \{f_2\}, a \in \Sigma_2 \cup \{\epsilon\}$$

$$\delta_1(f_1, \epsilon) = \delta_2(f_2, \epsilon) = \{f_0\}$$



- $L(M) = L(M_1) \cup L(M_2)$

Case ii

- $r = r_1 \cdot r_2$
- $M_1 = (Q_1, \Sigma_1, \delta_1, q_1, \{f_1\})$ $L(M_1) = L(r_1)$
- $M_2 = (Q_2, \Sigma_2, \delta_2, q_2, \{f_2\})$ $L(M_2) = L(r_2)$
- $M = (Q_1 \cup Q_2, \Sigma_1 \cup \Sigma_2, \delta, \{q_1\}, \{f_2\})$, where δ is given by:

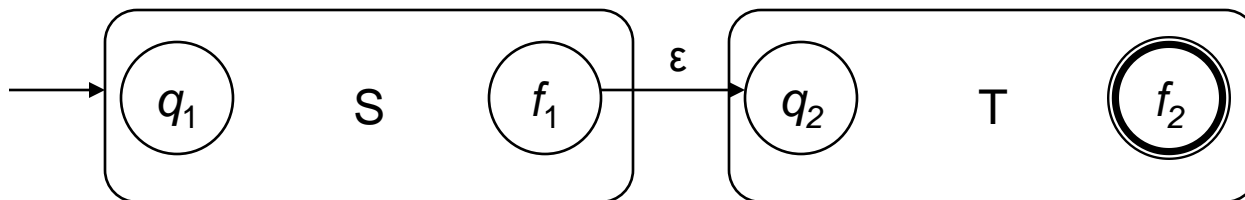
$\delta(q, a) = \delta_1(q, a)$ for q in $Q_1 - \{f_1\}$ and a in $\Sigma_1 \cup \{\epsilon\}$

$\delta(f_1, \epsilon) = \{q_2\}$

$\delta(q, a) = \delta_2(q, a)$ for q in Q_2 and a in $\Sigma_2 \cup \{\epsilon\}$

Case ii

- $L(M) = \{xy \mid x \text{ is in } L(M_1) \text{ and } y \text{ is in } L(M_2)\}$
- $L(M) = L(M_1) \cdot L(M_2).$

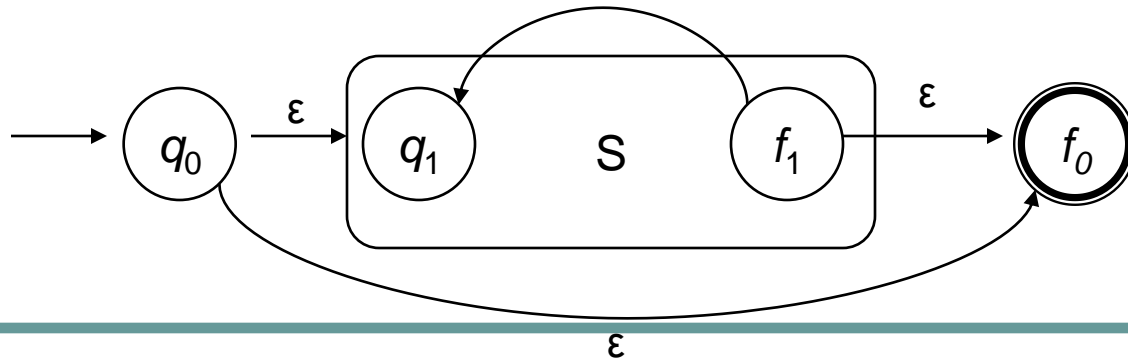


Case iii

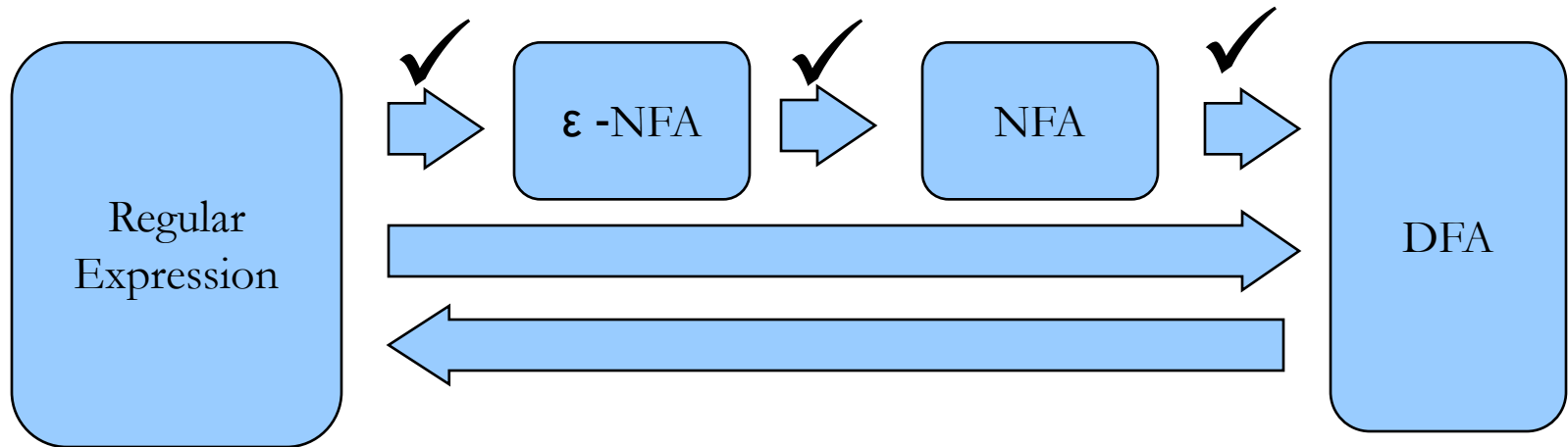
- $r = r_1^*$
- $M_1 = (Q_1, \Sigma_1, \delta_1, q_1, \{f_1\})$ $L(M_1) = r_1$
- $M = (Q_1 \cup \{q_0, f_0\}, \Sigma_1, \delta, q_0, \{f_0\})$, where δ is given by:

$$\delta(q, \varepsilon) = \delta(f_1, \varepsilon) = \{q_1, f_0\}$$

$$\delta(q, a) = \delta_1(q, a) \text{ for } q \text{ in } Q_1 - \{f_1\} \text{ and } a \text{ in } \Sigma_1 \cup \{\varepsilon\}$$



Road map



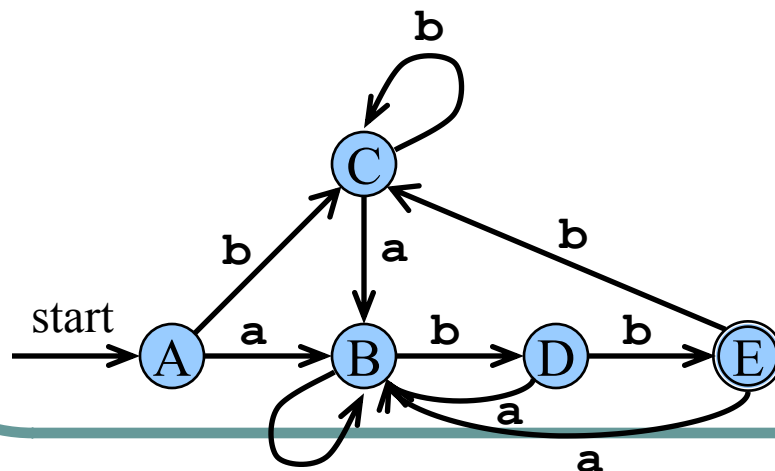
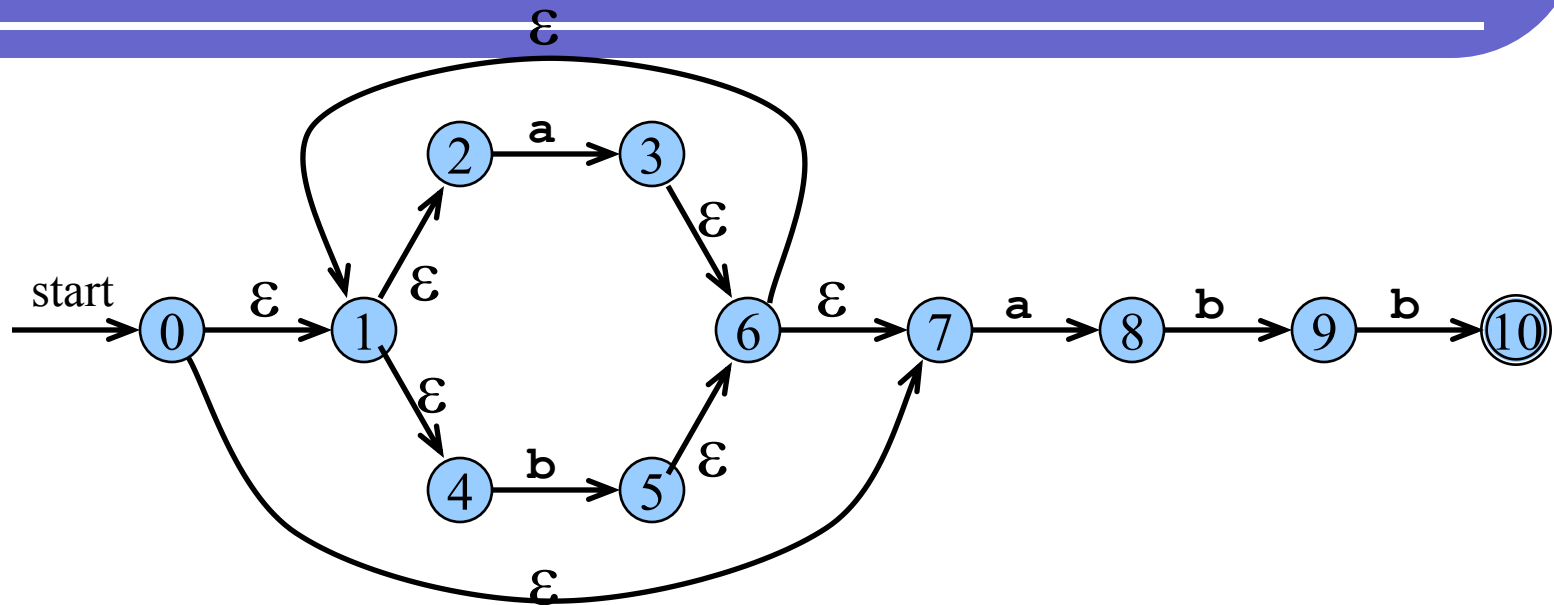
Conversion of ε -NFA to DFA

1. Find the ε -CLOSURE of the state q_0 from the constructed ε -NFA (i.e) from state q_0 , ε transition to other states are identified as well as ε transitions from other states are also identified and combined as one set (new state).

Conversion of ϵ -NFA to DFA

2. Perform the following steps until there are no more new states as been constructed.
 - i. Find the transition of the given regular expression symbols over Σ from the new state (i.e) move (new state, symbol)
 - ii. Find the ϵ -CLOSURE of move (new state, symbol).

Example



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}

Summary

- Definition of RE
- Precedence, identities, properties of RE.
- Thomson's construction to convert RE to NFA and then to DFA

Test Your Knowledge

- Which of the following does not represent the given language?

Language: $\{0,01\}$

- a) $0+01$
- b) $\{0\} \cup \{01\}$
- c) $\{0\} \cup \{0\}\{1\}$
- d) $\{0\} \wedge \{01\}$

Test Your Knowledge

- Regular Expression R and the language it describes can be represented as:
 - a) $R, R(L)$
 - b) $L(R), R(L)$
 - c) $R, L(R)$
 - d) All of the mentioned

Reference

- Hopcroft J.E., Motwani R. and Ullman J.D, “Introduction to Automata Theory, Languages and Computations”, Second Edition, Pearson Education, 2008