

# Finite Automata and Formal Languages

TMV027/DIT321– LP4 2013

Lecture 6  
Ana Bove

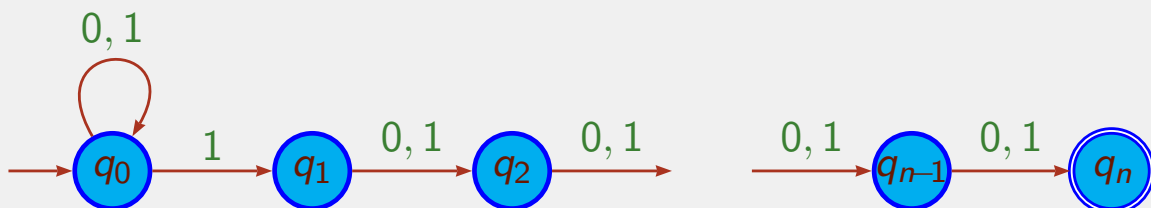
April 16th 2013

## Overview of today's lecture:

- More on NFA;
- NFA with  $\epsilon$ -Transitions;
- Equivalence between DFA and  $\epsilon$ -NFA;
- Regular expressions.

## A Bad Case for the Subset Construction

**Proposition:** Any DFA recognising the same language as the NFA below has at least  $2^n$  states:



This NFA recognises strings over  $\{0, 1\}$  such that the  $n$ th symbol from the end is a 1.

**Proof:** Let  $\mathcal{L}_n = \{x1u \mid x \in \Sigma^*, u \in \Sigma^{n-1}\}$  and  $D = (Q, \Sigma, \delta, q_0, F)$  a DFA.

We want to show that if  $|Q| < 2^n$  then  $\mathcal{L}(D) \neq \mathcal{L}_n$ .

## A Bad Case for the Subset Construction (Cont.)

**Lemma:** If  $\Sigma = \{0, 1\}$  and  $|Q| < 2^n$  then there exists  $x, y \in \Sigma^*$  and  $u, v \in \Sigma^{n-1}$  such that  $\hat{\delta}(q_0, x0u) = \hat{\delta}(q_0, y1v)$ .

**Proof:** Let us define a function  $h : \Sigma^n \rightarrow Q$  such that  $h(z) = \hat{\delta}(q_0, z)$ .

$h$  cannot be *injective* because  $|Q| < 2^n = |\Sigma^n|$ .

Hence, we have  $a_1 \dots a_n \neq b_1 \dots b_n$  such that

$$h(a_1 \dots a_n) = \hat{\delta}(q_0, a_1 \dots a_n) = \hat{\delta}(q_0, b_1 \dots b_n) = h(b_1 \dots b_n)$$

Let us assume that  $a_i = 0$  and  $b_i = 1$ .

Let  $x = a_1 \dots a_{i-1}$ ,  $y = b_1 \dots b_{i-1}$ ,  $u = a_{i+1} \dots a_n 0^{i-1}$ ,  $v = b_{i+1} \dots b_n 0^{i-1}$ .

Recall that for a DFA,  $\hat{\delta}(q, zw) = \hat{\delta}(\hat{\delta}(q, z), w)$  and hence:

$$\begin{aligned} \hat{\delta}(q_0, x0u) &= \hat{\delta}(q_0, a_1 \dots a_n 0^{i-1}) = \hat{\delta}(\hat{\delta}(q_0, a_1 \dots a_n), 0^{i-1}) = \\ &\hat{\delta}(\hat{\delta}(q_0, b_1 \dots b_n), 0^{i-1}) = \hat{\delta}(q_0, b_1 \dots b_n 0^{i-1}) = \hat{\delta}(q_0, y1v) \end{aligned}$$

## A Bad Case for the Subset Construction (Cont.)

**Lemma:** If  $|Q| < 2^n$  then  $\mathcal{L}(D) \neq \mathcal{L}_n$ .

**Proof:** Assume  $\mathcal{L}(D) = \mathcal{L}_n$ .

Let  $x, y \in \Sigma^*$  and  $u, v \in \Sigma^{n-1}$  as in previous lemma.

Then we must have that  $y1v \in \mathcal{L}(D)$  but  $x0u \notin \mathcal{L}(D)$ ,

That is,  $\hat{\delta}(q_0, y1v) \in F$  but  $\hat{\delta}(q_0, x0u) \notin F$ .

However, this contradicts the previous lemma that says that  $\hat{\delta}(q_0, x0u) = \hat{\delta}(q_0, y1v)$ .

## Product Construction for NFA

**Definition:** Given 2 NFA  $N_1 = (Q_1, \Sigma, \delta_1, q_1, F_1)$  and  $N_2 = (Q_2, \Sigma, \delta_2, q_2, F_2)$  over the same alphabet  $\Sigma$ , we define the product  $N_1 \times N_2 = (Q, \Sigma, \delta, q_0, F)$  as follows:

- $Q = Q_1 \times Q_2$ ;
- $\delta((p_1, p_2), a) = \delta_1(p_1, a) \times \delta_2(p_2, a)$ ;
- $q_0 = (q_1, q_2)$ ;
- $F = F_1 \times F_2$ .

**Lemma:**  $(t_1, t_2) \in \hat{\delta}((p_1, p_2), x)$  iff  $t_1 \in \hat{\delta}_1(p_1, x)$  and  $t_2 \in \hat{\delta}_2(p_2, x)$ .

**Proof:** By induction on  $x$ .

**Proposition:**  $\mathcal{L}(N_1 \times N_2) = \mathcal{L}(N_1) \cap \mathcal{L}(N_2)$ .

## Complement for NFA

**OBS:** Given NFA  $N = (Q, \Sigma, \delta, q, F)$  and  $N' = (Q, \Sigma, \delta, q, Q - F)$  we do *not* have in general that  $\mathcal{L}(N') = \Sigma^* - \mathcal{L}(N)$ .

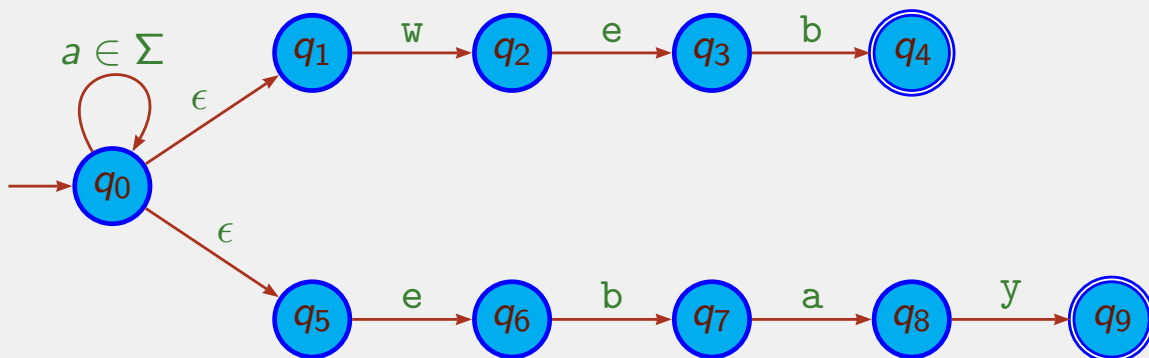
**Example:** Let  $\Sigma = \{a\}$  and  $N$  and  $N'$  as follows:



## NFA with $\epsilon$ -Transitions

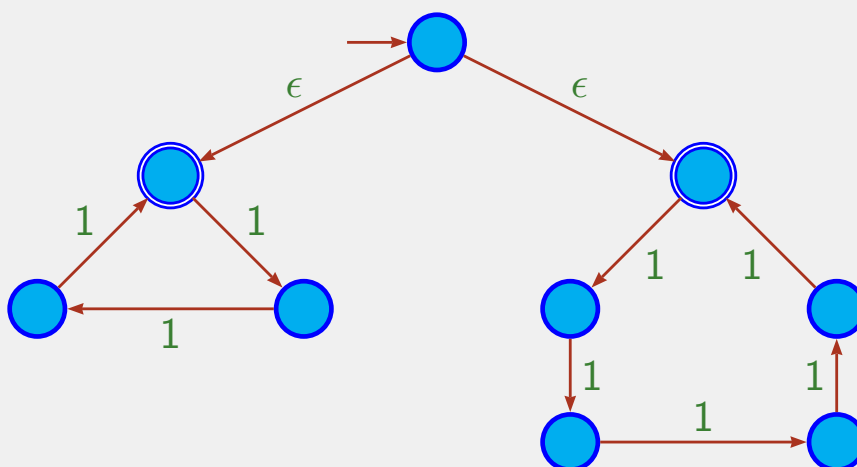
Another useful extension of automata that does not add more power is the possibility to allow  $\epsilon$ -transitions, that is, transitions from one state to another *without* reading any input symbol.

**Example:** The following  $\epsilon$ -NFA searches for the keyword web and ebay:



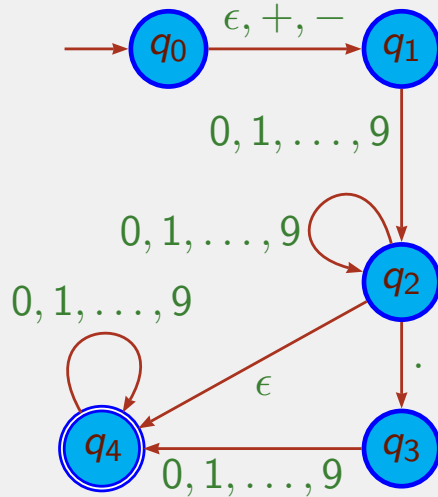
## $\epsilon$ -NFA Accepting Words of Length Divisible by 3 or by 5

**Example:** Let  $\Sigma = \{1\}$ .



## $\epsilon$ -NFA Accepting Decimal Numbers

**Exercise:** Define a NFA accepting number with an optional  $+/-$  symbol and an optional decimal part.



	$+,-$	$.$	$0,1,\dots,9$	$\epsilon$
$\rightarrow q_0$	$\{q_1\}$	$\emptyset$	$\emptyset$	$\{q_1\}$
$q_1$	$\emptyset$	$\emptyset$	$\{q_2\}$	$\emptyset$
$q_2$	$\emptyset$	$\{q_3\}$	$\{q_2\}$	$\{q_4\}$
$q_3$	$\emptyset$	$\emptyset$	$\{q_4\}$	$\emptyset$
$*q_4$	$\emptyset$	$\emptyset$	$\{q_4\}$	$\emptyset$

The uses of  $\epsilon$ -transitions represent the *optional* symbol  $+/-$  and the *optional* decimal part.

## NFA with $\epsilon$ -Transitions

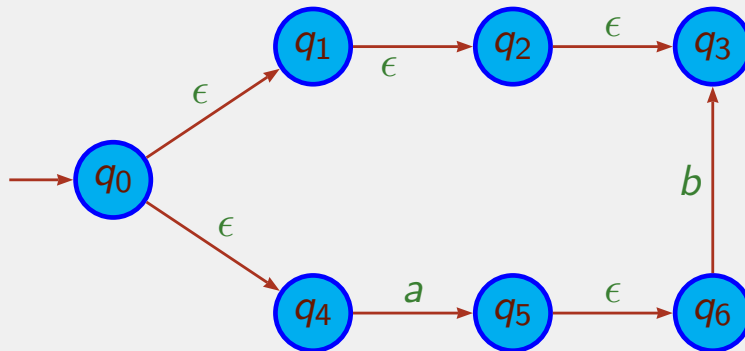
**Definition:** A *NFA with  $\epsilon$ -transitions* ( $\epsilon$ -NFA) is a 5-tuple  $(Q, \Sigma, \delta, q_0, F)$  consisting of:

- ① A finite set  $Q$  of *states*;
- ② A finite set  $\Sigma$  of *symbols* (alphabet);
- ③ A *transition function*  $\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow \mathcal{P}ow(Q)$  (“partial” function that takes as argument a state and a symbol or the  $\epsilon$ -transition, and returns a *set of states*);
- ④ A *start state*  $q_0 \in Q$ ;
- ⑤ A set  $F \subseteq Q$  of *final or accepting states*.

## $\epsilon$ -Closures

Informally, the  $\epsilon$ -closure of a state  $q$  is the set of states we can reach by doing nothing or by *only* following paths labelled with  $\epsilon$ .

**Example:** For the automaton



the  $\epsilon$ -closure of  $q_0$  is  $\{q_0, q_1, q_2, q_3, q_4\}$ .

## $\epsilon$ -Closures

**Definition:** Formally, we define the  $\epsilon$ -closure of a set of states as follows:

- If  $q \in S$  then  $q \in \text{ECLOSE}(S)$ ;
- If  $q \in \text{ECLOSE}(S)$  and  $p \in \delta(q, \epsilon)$  then  $p \in \text{ECLOSE}(S)$ .

**Note:** Alternative formulation

$$\frac{q \in S}{q \in \text{ECLOSE}(S)} \qquad \frac{q \in \text{ECLOSE}(S) \quad p \in \delta(q, \epsilon)}{p \in \text{ECLOSE}(S)}$$

**Definition:** We say that  $S$  is  $\epsilon$ -closed iff  $S = \text{ECLOSE}(S)$ .

## Remarks: $\epsilon$ -Closures

- Intuitively,  $p \in \text{ECLOSE}(S)$  iff there exists  $q \in S$  and a sequence of  $\epsilon$ -transitions such that



- The  $\epsilon$ -closure of a single state  $q$  can be computed as  $\text{ECLOSE}(\{q\})$ .
- $\text{ECLOSE}(\emptyset) = \emptyset$ .
- $S$  is  $\epsilon$ -closed iff  $q \in S$  and  $p \in \delta(q, \epsilon)$  implies  $p \in S$ .
- We can prove that  $\text{ECLOSE}(S)$  is the *smallest* subset of  $Q$  containing  $S$  which is  $\epsilon$ -closed.

## Extending the Transition Function to Strings

**Definition:** Given an  $\epsilon$ -NFA  $E = (Q, \Sigma, \delta, q_0, F)$  we define

$$\hat{\delta} : Q \times \Sigma^* \rightarrow [Q]$$

$$\hat{\delta}(q, \epsilon) = \text{ECLOSE}(\{q\})$$

$$\hat{\delta}(q, ax) = \bigcup_{p \in \Delta(\text{ECLOSE}(\{q\}), a)} \hat{\delta}(p, x)$$

$$\text{where } \Delta(S, a) = \bigcup_{p \in S} \delta(p, a)$$

**Remark:** By definition,  $\hat{\delta}(q, a) = \text{ECLOSE}(\Delta(\text{ECLOSE}(\{q\}), a))$ .

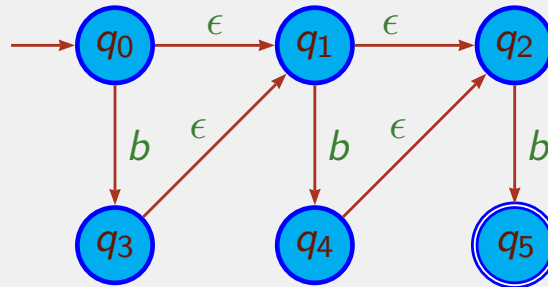
**Remark:** We can prove by induction on  $x$  that all sets  $\hat{\delta}(q, x)$  are  $\epsilon$ -closed.

This result uses that the union of  $\epsilon$ -closed sets is also a  $\epsilon$ -closed set.

## Language Accepted by a $\epsilon$ -NFA

**Definition:** The *language* accepted by the  $\epsilon$ -NFA  $(Q, \Sigma, \delta, q_0, F)$  is the set  $\mathcal{L} = \{x \in \Sigma^* \mid \hat{\delta}(q_0, x) \cap F \neq \emptyset\}$ .

**Example:** Let  $\Sigma = \{b\}$ .



The automaton accepts the language  $\{b, bb, bbb\}$ .

**Note:** Yet again, we could write a program that simulates a  $\epsilon$ -NFA and let the program tell us whether a certain string is accepted or not.

**Exercise:** Do it!

## Eliminating $\epsilon$ -Transitions

**Definition:** Given an  $\epsilon$ -NFA  $E = (Q_E, \Sigma, \delta_E, q_E, F_E)$  we define a DFA  $D = (Q_D, \Sigma, \delta_D, q_D, F_D)$  as follows:

- $Q_D = \{\text{ECLOSE}(S) \mid S \in \mathcal{P}ow(Q_E)\}$ ;
- $\delta_D(S, a) = \text{ECLOSE}(\Delta(S, a))$  with  $\Delta(S, a) = \cup_{p \in S} \delta(p, a)$ ;
- $q_D = \text{ECLOSE}(\{q_E\})$ ;
- $F_D = \{S \in Q_D \mid S \cap F_E \neq \emptyset\}$ .

**Note:** This construction is similar to the subset construction but now we need to  $\epsilon$ -close after each step.



# Eliminating $\epsilon$ -Transitions

Let  $E$  be an  $\epsilon$ -NFA and  $D$  the corresponding DFA.

**Lemma:**  $\forall x \in \Sigma^*. \hat{\delta}_E(q_E, x) = \hat{\delta}_D(q_D, x)$ .

**Proof:** By induction on  $x$ .

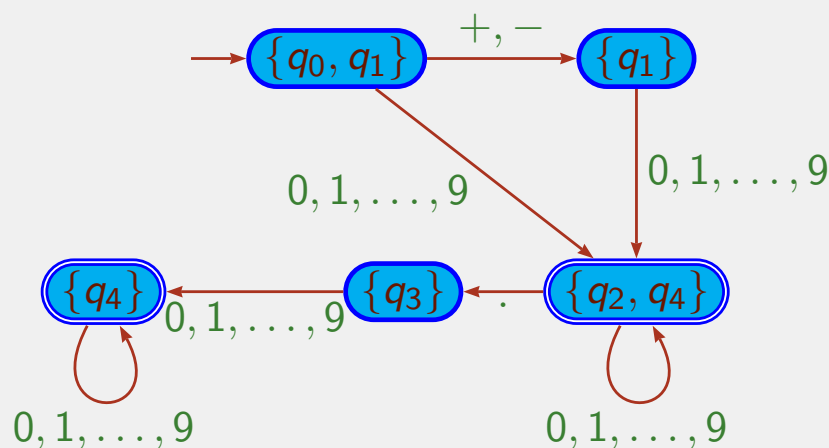
**Proposition:**  $\mathcal{L}(E) = \mathcal{L}(D)$ .

**Proof:**  $x \in \mathcal{L}(E)$  iff  $\hat{\delta}_E(q_E, x) \cap F_E \neq \emptyset$  iff  $\hat{\delta}_E(q_E, x) \in F_D$  iff (by previous lemma)  $\hat{\delta}_D(q_D, x) \in F_D$  iff  $x \in \mathcal{L}(D)$ .

## Example: Eliminating $\epsilon$ -Transitions

Let us eliminate the  $\epsilon$ -transitions in  $\epsilon$ -NFA that recognises numbers in slide 8.

We obtain the following DFA:



# Finite Automata and Regular Languages

We have shown that DFA, NFA and  $\epsilon$ -NFA are equivalent in the sense that we can transform one to the other.

Hence, a language is *regular* iff there exists a finite automaton (DFA, NFA or  $\epsilon$ -NFA) that accepts the language.

## Regular Expressions

*Regular expressions* (RE) are an “algebraic” way to denote languages.

RE are a simple way to express the strings we want to accept.

They serve as input language for certain systems.

**Example:** grep command in UNIX (K. Thompson) is given a (variation) of a RE as input

We will show that RE are as expressive as DFA and hence, they define all and only the *regular languages*.

## Inductive Definition of Regular Expressions

**Definition:** Given an alphabet  $\Sigma$ , we inductively define the *regular expressions* over  $\Sigma$  as follows:

- Base cases:**
- The constants  $\emptyset$  and  $\epsilon$  are RE;
  - If  $a \in \Sigma$  then  $a$  is a RE.

**Inductive steps:** Given the RE  $R$  and  $S$ , we define the following RE:

- $R + S$  and  $RS$  are RE;
- $R^*$  is RE.

The precedence of the operands is the following:

- The closure operator  $*$  has the highest precedence;
- Next comes concatenation;
- Finally, comes the operator  $+$ ;
- We use parentheses  $(, )$  to change the precedence.

## Another Way to Define the Regular Expressions

A nicer way to define the regular expressions is by giving the following BNF (Backus-Naur Form), for  $a \in \Sigma$ :

$$R ::= \emptyset \mid \epsilon \mid a \mid R + R \mid RR \mid R^*$$

alternatively

$$R, S ::= \emptyset \mid \epsilon \mid a \mid R + S \mid RS \mid R^*$$

**Note:** BNF is a way to declare the syntax of a language.

It is very useful when describing *context-free grammars* and in particular the syntax of most programming languages.

# Functional Representation of Regular Expressions

```
data RExp a = Empty | Epsilon | Atom a |
            Plus (RExp a) (RExp a) |
            Concat (RExp a) (RExp a) |
            Star (RExp a)
```

For example the expression  $b + (bc)^*$  is given as

```
Plus (Atom "b") (Star (Concat (Atom "b") (Atom "c")))
```

## Language Defined by the Regular Expressions

Given a RE  $R$ , it defines the language  $\mathcal{L}(R)$ .

**Definition:** The *language* defined by a regular expression is defined by recursion on the expression:

- Base cases:
- $\mathcal{L}(\emptyset) = \emptyset$ ;
  - $\mathcal{L}(\epsilon) = \{\epsilon\}$ ;
  - Given  $a \in \Sigma$ ,  $\mathcal{L}(a) = \{a\}$ .

- Recursive cases:
- $\mathcal{L}(R + S) = \mathcal{L}(R) \cup \mathcal{L}(S)$ ;
  - $\mathcal{L}(RS) = \mathcal{L}(R)\mathcal{L}(S)$ ;
  - $\mathcal{L}(R^*) = \mathcal{L}(R)^*$ .

**Note:**  $x \in \mathcal{L}(R)$  iff  $x$  is generated/accepted by  $R$ .

**Notation:** We write  $x \in R$  or  $x \in \mathcal{L}(R)$  indistinctly.

# Overview of Next Lecture

Sections 3.2, 3.4:

- More on RE;
- Equivalence between FA and RE;
- Algebraic laws for regular expressions.