

## Pattern matching

---

What happens if, at a Unix/Linux shell prompt, you type

```
ls *
```

and press return?

Suppose the current directory contains files called `regfla.tex`, `regfla.aux`, `regfla.log`, `regfla.dvi`, and (strangely) `.aux`. What happens if you type

```
ls *.aux
```

and press return?

# Alphabets

---

An **alphabet** is specified by giving a finite set,  $\Sigma$ , whose elements are called **symbols**. For us, any set qualifies as a possible alphabet, so long as it is finite.

## Examples:

$\Sigma_1 = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$  — 10-element set of decimal digits.

$\Sigma_2 = \{a, b, c, \dots, x, y, z\}$  — 26-element set of lower-case characters of the English language.

$\Sigma_3 = \{S \mid S \subseteq \Sigma_1\}$  —  $2^{10}$ -element set of all subsets of the alphabet of decimal digits.

## Non-example:

$\mathbb{N} = \{0, 1, 2, 3, \dots\}$  — set of all non-negative whole numbers is not an alphabet, because it is infinite.

## Strings over an alphabet

---

A **string of length  $n$**  ( $\geq 0$ ) over an alphabet  $\Sigma$  is just an ordered  $n$ -tuple of elements of  $\Sigma$ , written without punctuation.

**Example:** if  $\Sigma = \{a, b, c\}$ , then  $a$ ,  $ab$ ,  $aac$ , and  $bbac$  are strings over  $\Sigma$  of lengths one, two, three and four respectively.

$\Sigma^*$  def set of all strings over  $\Sigma$  of any finite length.

N.B. there is a unique string of length zero over  $\Sigma$ , called the **null string** (or **empty string**) and denoted  $\epsilon$  (no matter which  $\Sigma$  we are talking about).

## Concatenation of strings

---

The **concatenation** of two strings  $u, v \in \Sigma^*$  is the string  $uv$  obtained by joining the strings end-to-end.

**Examples:** If  $u = ab$ ,  $v = ra$  and  $w = cad$ , then  $vu = raab$ ,  $uu = abab$  and  $wv = cadra$ .

This generalises to the concatenation of three or more strings.

E.g.  $uvwuv = abracadabra$ .

## Regular expressions over an alphabet $\Sigma$

---

- each symbol  $a \in \Sigma$  is a regular expression
- $\epsilon$  is a regular expression
- $\emptyset$  is a regular expression
- if  $r$  and  $s$  are regular expressions, then so is  $(r|s)$
- if  $r$  and  $s$  are regular expressions, then so is  $rs$
- if  $r$  is a regular expression, then so is  $(r)^*$

Every regular expression is built up inductively, by *finitely many* applications of the above rules.

(N.B. we assume  $\epsilon$ ,  $\emptyset$ ,  $($ ,  $)$ ,  $|$ , and  $*$  are **not** symbols in  $\Sigma$ .)

## Matching strings to regular expressions

---

- $u$  matches  $a \in \Sigma$  iff  $u = a$
- $u$  matches  $\varepsilon$  iff  $u = \varepsilon$
- no string matches  $\emptyset$
- $u$  matches  $r|s$  iff  $u$  matches either  $r$  or  $s$
- $u$  matches  $rs$  iff it can be expressed as the concatenation of two strings,  $u = vw$ , with  $v$  matching  $r$  and  $w$  matching  $s$
- $u$  matches  $r^*$  iff either  $u = \varepsilon$ , or  $u$  matches  $r$ , or  $u$  can be expressed as the concatenation of two or more strings, each of which matches  $r$

## Examples of matching, with $\Sigma = \{0, 1\}$

---

- $0|1$  is matched by each symbol in  $\Sigma$
- $1(0|1)^*$  is matched by any string in  $\Sigma^*$  that starts with a '1'
- $((0|1)(0|1))^*$  is matched by any string of even length in  $\Sigma^*$
- $(0|1)^*(0|1)^*$  is matched by any string in  $\Sigma^*$
- $(\epsilon|0)(\epsilon|1)|11$  is matched by just the strings  $\epsilon$ , **0**, **1**, **01**, and **11**
- $\emptyset 1|0$  is just matched by **0**

# Languages

---

A (formal) *language*  $L$  over an alphabet  $\Sigma$  is just a set of strings in  $\Sigma^*$ .

Thus any subset  $L \subseteq \Sigma^*$  determines a language over  $\Sigma$ .

The *language determined by a regular expression*  $r$  over  $\Sigma$  is

$$L(r) \stackrel{\text{def}}{=} \{u \in \Sigma^* \mid u \text{ matches } r\}.$$

Two regular expressions  $r$  and  $s$  (over the same alphabet) are *equivalent* iff  $L(r)$  and  $L(s)$  are equal sets (i.e. have exactly the same members).



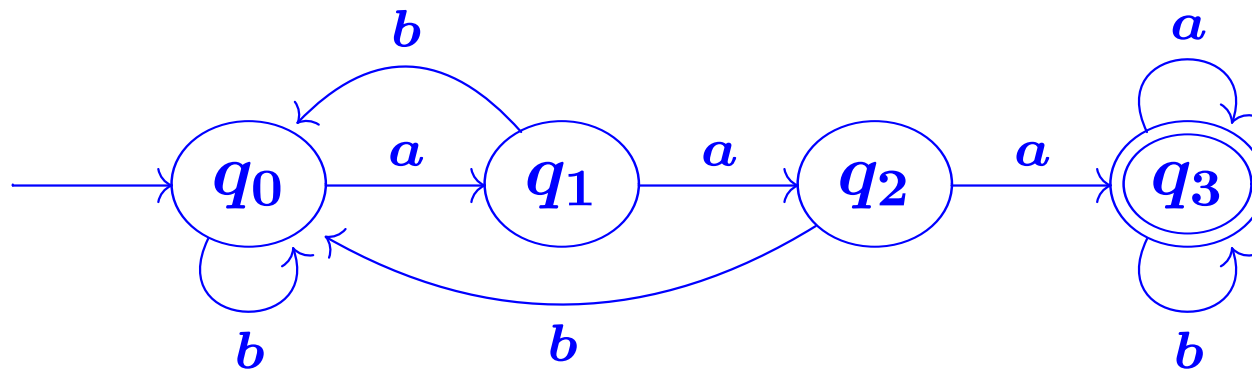
## Some questions

---

- (a) Is there an algorithm which, given a string  $u$  and a regular expression  $r$  (over the same alphabet), computes whether or not  $u$  matches  $r$ ?
- (b) In formulating the definition of regular expressions, have we missed out some practically useful notions of pattern?
- (c) Is there an algorithm which, given two regular expressions  $r$  and  $s$  (over the same alphabet), computes whether or not they are equivalent? (Cf. Slide 8.)
- (d) Is every language of the form  $L(r)$ ?

## Example of a finite automaton

---



States:  $q_0$ ,  $q_1$ ,  $q_2$ ,  $q_3$ .

Input symbols:  $a$ ,  $b$ .

Transitions: as indicated above.

Start state:  $q_0$ .

Accepting state(s):  $q_3$ .

## $L(M)$ , *language accepted* by a finite automaton $M$

---

consists of all strings  $u$  over its alphabet of input symbols satisfying  $q_0 \xrightarrow{u}^* q$  with  $q_0$  the start state and  $q$  some accepting state. Here

$$q_0 \xrightarrow{u}^* q$$

means, if  $u = a_1 a_2 \dots a_n$  say, that for some states

$q_1, q_2, \dots, q_n = q$  (not necessarily all distinct) there are transitions of the form

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \dots \xrightarrow{a_n} q_n = q.$$

**N.B.**

case  $n = 0$ :  $q \xrightarrow{\epsilon}^* q'$  iff  $q = q'$

case  $n = 1$ :  $q \xrightarrow{a}^* q'$  iff  $q \xrightarrow{a} q'$ .

A **non-deterministic finite automaton** (NFA),  $M$ , is specified by

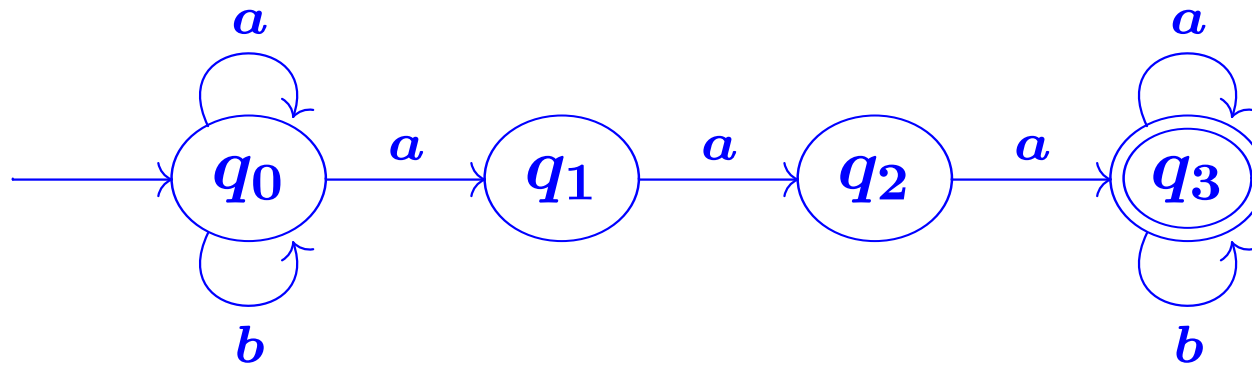
- a finite set  $States_M$  (of **states**)
- a finite set  $\Sigma_M$  (the alphabet of **input symbols**)
- for each  $q \in States_M$  and each  $a \in \Sigma_M$ , a subset  $\Delta_M(q, a) \subseteq States_M$  (the set of states that can be reached from  $q$  with a single **transition** labelled  $a$ )
- an element  $s_M \in States_M$  (the **start state**)
- a subset  $Accept_M \subseteq States_M$  (of **accepting states**)

## Example of a non-deterministic finite automaton

---

Input alphabet:  $\{a, b\}$ .

States, transitions, start state, and accepting states as shown:



The language accepted by this automaton is the same as for the automaton on Slide 10, namely

$$\{u \in \{a, b\}^* \mid u \text{ contains three consecutive } a\text{'s}\}.$$

A **deterministic finite automaton** (DFA)

is an NFA  $M$  with the property that for each  $q \in States_M$  and  $a \in \Sigma_M$ , the finite set  $\Delta_M(q, a)$  contains exactly one element—call it  $\delta_M(q, a)$ .

Thus in this case transitions in  $M$  are essentially specified by a **next-state function**,  $\delta_M$ , mapping each (state, input symbol)-pair  $(q, a)$  to the unique state  $\delta_M(q, a)$  which can be reached from  $q$  by a transition labelled  $a$ :

$$q \xrightarrow{a} q' \quad \text{iff} \quad q' = \delta_M(q, a)$$

An **NFA with  $\epsilon$ -transitions** ( $\text{NFA}^\epsilon$ )

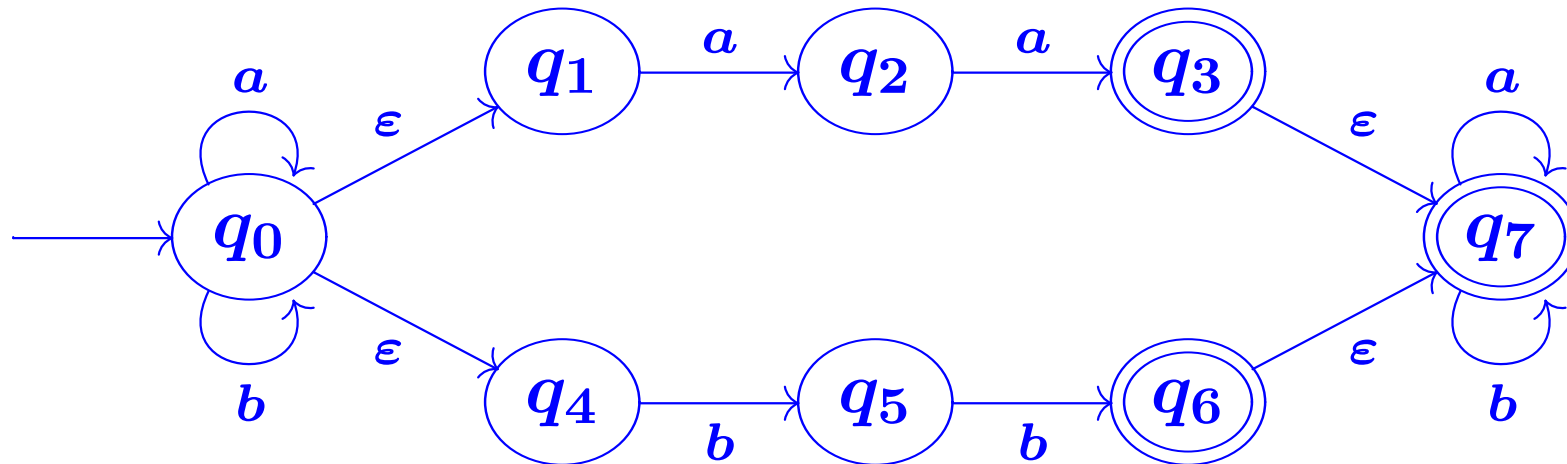
is specified by an NFA  $M$  together with a binary relation, called the  **$\epsilon$ -transition relation**, on the set  $\text{States}_M$ . We write

$$q \xrightarrow{\epsilon} q'$$

to indicate that the pair of states  $(q, q')$  is in this relation.

---

**Example** (with input alphabet =  $\{a, b\}$ ):



## $L(M)$ , *language accepted* by an NFA<sup>ε</sup> $M$

---

consists of all strings  $u$  over the alphabet  $\Sigma_M$  of input symbols satisfying  $q_0 \xRightarrow{u} q$  with  $q_0$  the initial state and  $q$  some accepting state.

Here  $\cdot \xRightarrow{\bar{\phantom{x}}} \cdot$  is defined by:

$q \xRightarrow{\varepsilon} q'$  iff  $q = q'$  or there is a sequence  $q \xrightarrow{\varepsilon} \dots q'$  of one or more  $\varepsilon$ -transitions in  $M$  from  $q$  to  $q'$

$q \xRightarrow{a} q'$  (for  $a \in \Sigma_M$ ) iff  $q \xRightarrow{\varepsilon} \cdot \xrightarrow{a} \cdot \xRightarrow{\varepsilon} q'$

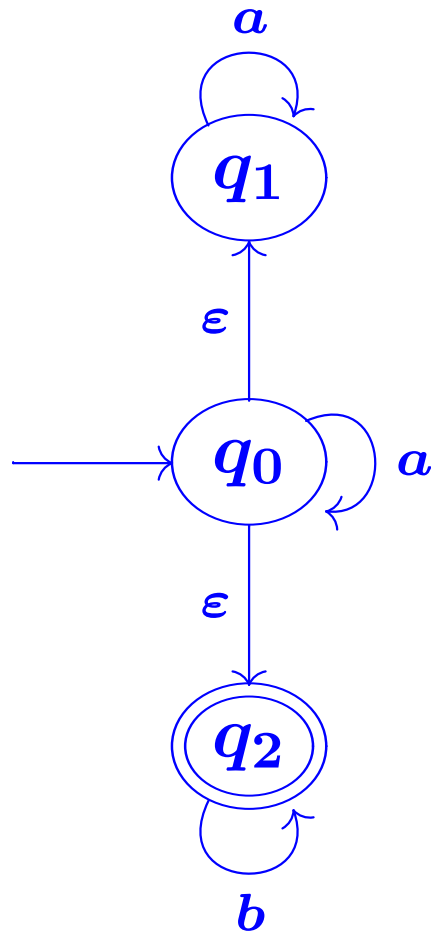
$q \xRightarrow{ab} q'$  (for  $a, b \in \Sigma_M$ ) iff  $q \xRightarrow{\varepsilon} \cdot \xrightarrow{a} \cdot \xRightarrow{\varepsilon} \cdot \xrightarrow{b} \cdot \xRightarrow{\varepsilon} q'$

and similarly for longer strings



# Example of the subset construction

$M :$



$\delta_{PM} :$

	$a$	$b$
$\emptyset$	$\emptyset$	$\emptyset$
$\{q_0\}$	$\{q_0, q_1, q_2\}$	$\{q_2\}$
$\{q_1\}$	$\{q_1\}$	$\emptyset$
$\{q_2\}$	$\emptyset$	$\{q_2\}$
$\{q_0, q_1\}$	$\{q_0, q_1, q_2\}$	$\{q_2\}$
$\{q_0, q_2\}$	$\{q_0, q_1, q_2\}$	$\{q_2\}$
$\{q_1, q_2\}$	$\{q_1\}$	$\{q_2\}$
$\{q_0, q_1, q_2\}$	$\{q_0, q_1, q_2\}$	$\{q_2\}$

**Theorem.** For each NFA<sup>ε</sup>  $M$  there is a DFA  $PM$  with the same alphabet of input symbols and accepting exactly the same strings as  $M$ , i.e. with  $L(PM) = L(M)$

Definition of  $PM$  (refer to Slides 12 and 14):

- $States_{PM} \stackrel{\text{def}}{=} \{S \mid S \subseteq States_M\}$
- $\Sigma_{PM} \stackrel{\text{def}}{=} \Sigma_M$
- $S \xrightarrow{a} S'$  in  $PM$  iff  $S' = \delta_{PM}(S, a)$ , where  
 $\delta_{PM}(S, a) \stackrel{\text{def}}{=} \{q' \mid \exists q \in S (q \xrightarrow{a} q' \text{ in } M)\}$
- $s_{PM} \stackrel{\text{def}}{=} \{q \mid s_M \xrightarrow{\epsilon} q\}$
- $Accept_{PM} \stackrel{\text{def}}{=} \{S \in States_{PM} \mid \exists q \in S (q \in Accept_M)\}$

## Definition

A language is *regular* iff it is the set of strings accepted by some deterministic finite automaton.

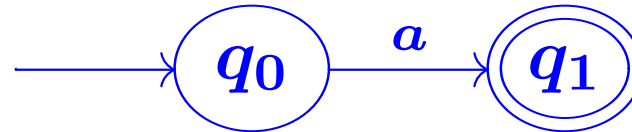
### Kleene's Theorem

(a) For any regular expression  $r$ ,  $L(r)$  is a regular language (cf. Slide 8).

(b) Conversely, every regular language is the form  $L(r)$  for some regular expression  $r$ .

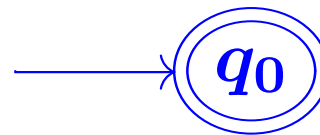
## NFAs for atomic regular expressions

---



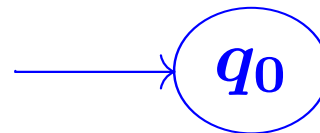
just accepts the one-symbol string  $a$

---



just accepts the null string,  $\epsilon$

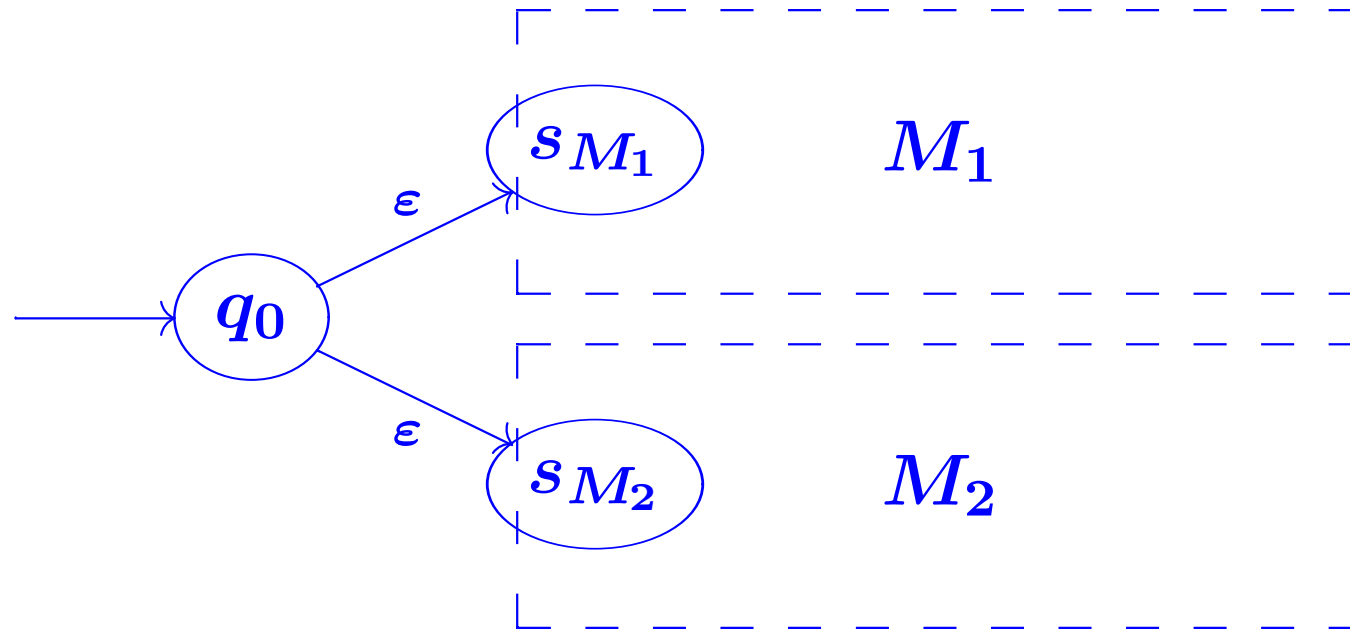
---



accepts no strings

## *Union( $M_1, M_2$ )*

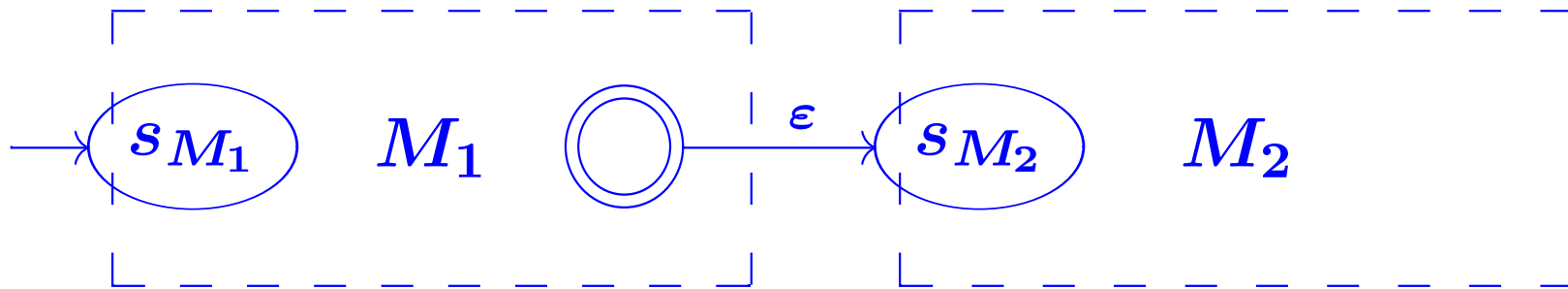
---



Set of accepting states is union of *Accept* $_{M_1}$  and *Accept* $_{M_2}$ .

## *Concat*( $M_1, M_2$ )

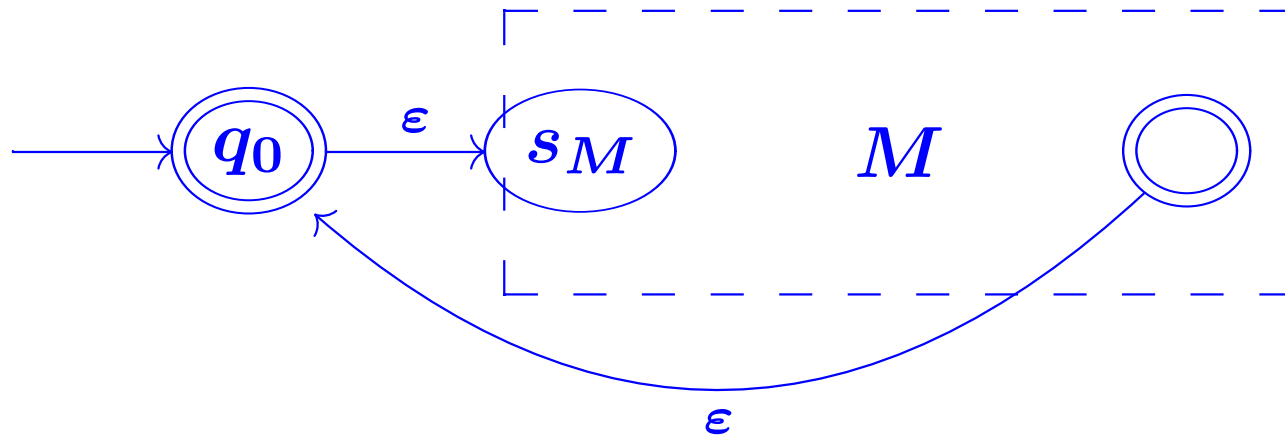
---



Set of accepting states is *Accept* $_{M_2}$ .

## *Star(M)*

---



The only accepting state of *Star(M)* is  $q_0$ .

**Lemma** Given an NFA  $M$ , for each subset  $Q \subseteq \text{States}_M$  and each pair of states  $q, q' \in \text{States}_M$ , there is a regular expression  $r_{q,q'}^Q$  satisfying

$$L(r_{q,q'}^Q) = \{u \in (\Sigma_M)^* \mid q \xrightarrow{u}^* q' \text{ in } M \text{ with all intermediate states of the sequence in } Q\}.$$

Hence  $L(M) = L(r)$ , where  $r = r_1 | \cdots | r_k$  and

$k =$  number of accepting states,

$r_i = r_{s,q_i}^Q$  with  $Q = \text{States}_M$ ,

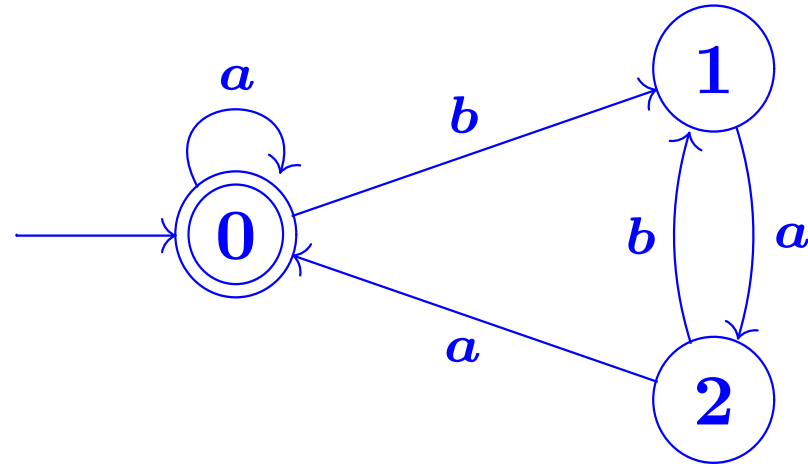
$s =$  start state,

$q_i = i$ th accepting state.

(In case  $k = 0$ , take  $r$  to be the regular expression  $\emptyset$ .)



# Example



Direct inspection yields:

$r_{i,j}^{\{0\}}$	0	1	2
0			
1	$\emptyset$	$\epsilon$	$a$
2	$aa^*$	$a^*b$	$\epsilon$

$r_{i,j}^{\{0,2\}}$	0	1	2
0	$a^*$	$a^*b$	
1			
2			

## $Not(M)$

---

- $States_{Not(M)} \stackrel{\text{def}}{=} States_M$
- $\Sigma_{Not(M)} \stackrel{\text{def}}{=} \Sigma_M$
- transitions of  $Not(M)$  = transitions of  $M$
- start state of  $Not(M)$  = start state of  $M$
- $Accept_{Not(M)} = \{q \in States_M \mid q \notin Accept_M\}$ .

Provided  $M$  is a *deterministic* finite automaton, then  $u$  is accepted by  $Not(M)$  iff it is not accepted by  $M$ :

$$L(Not(M)) = \{u \in \Sigma^* \mid u \notin L(M)\}.$$

## $And(M_1, M_2)$

---

- states of  $And(M_1, M_2)$  are all ordered pairs  $(q_1, q_2)$  with  $q_1 \in States_{M_1}$  and  $q_2 \in States_{M_2}$
- alphabet of  $And(M_1, M_2)$  is the common alphabet of  $M_1$  and  $M_2$
- $(q_1, q_2) \xrightarrow{a} (q'_1, q'_2)$  in  $And(M_1, M_2)$  iff  $q_1 \xrightarrow{a} q'_1$  in  $M_1$  and  $q_2 \xrightarrow{a} q'_2$  in  $M_2$
- start state of  $And(M_1, M_2)$  is  $(s_{M_1}, s_{M_2})$
- $(q_1, q_2)$  accepting in  $And(M_1, M_2)$  iff  $q_1$  accepting in  $M_1$  and  $q_2$  accepting in  $M_2$ .

## Examples of non-regular languages

---

- The set of strings over  $\{(, ), a, b, \dots, z\}$  in which the parentheses '(' and ')' occur well-nested.
- The set of strings over  $\{a, b, \dots, z\}$  which are **palindromes**, i.e. which read the same backwards as forwards.
- $\{a^n b^n \mid n \geq 0\}$

# The Pumping Lemma

---

For every regular language  $L$ , there is a number  $\ell \geq 1$  satisfying the **pumping lemma property**:

all  $w \in L$  with  $\text{length}(w) \geq \ell$  can be expressed as a concatenation of three strings,  $w = u_1vu_2$ , where  $u_1$ ,  $v$  and  $u_2$  satisfy:

- $\text{length}(v) \geq 1$   
(i.e.  $v \neq \varepsilon$ )
- $\text{length}(u_1v) \leq \ell$
- for all  $n \geq 0$ ,  $u_1v^n u_2 \in L$   
(i.e.  $u_1u_2 \in L$ ,  $u_1vu_2 \in L$  [but we knew that anyway],  
 $u_1v^2u_2 \in L$ ,  $u_1v^3u_2 \in L$ , etc).

If  $n \geq \ell =$  number of states of  $M$ , then in

$$s_M = \underbrace{q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \cdots \xrightarrow{a_\ell} q_\ell}_{\ell+1 \text{ states}} \cdots \xrightarrow{a_n} q_n \in \text{Accept}_M$$

$q_0, \dots, q_\ell$  can't all be distinct states. So  $q_i = q_j$  for some  $0 \leq i < j \leq \ell$ . So the above transition sequence looks like

$$s_M = q_0 \xrightarrow{u_1^*} q_i \overset{v}{\curvearrowright} q_j^* \xrightarrow{u_2^*} q_n \in \text{Accept}_M$$

where

$$u_1 \stackrel{\text{def}}{=} a_1 \dots a_i \quad v \stackrel{\text{def}}{=} a_{i+1} \dots a_j \quad u_2 \stackrel{\text{def}}{=} a_{j+1} \dots a_n.$$

## How to use the Pumping Lemma to prove that a language $L$ is *not* regular

---

For each  $\ell \geq 1$ , find some  $w \in L$  of length  $\geq \ell$  so that

(†)  $\left\{ \begin{array}{l} \text{no matter how } w \text{ is split into three, } w = u_1 v u_2, \\ \text{with } \textit{length}(u_1 v) \leq \ell \text{ and } \textit{length}(v) \geq 1, \\ \text{there is some } n \geq 0 \text{ for which } u_1 v^n u_2 \text{ is not in } L. \end{array} \right.$

## Examples

---

(i)  $L_1 \stackrel{\text{def}}{=} \{a^n b^n \mid n \geq 0\}$  is not regular.

[For each  $\ell \geq 1$ ,  $a^\ell b^\ell \in L_1$  is of length  $\geq \ell$  and has property ( $\dagger$ ) on Slide 31.]

(ii)  $L_2 \stackrel{\text{def}}{=} \{w \in \{a, b\}^* \mid w \text{ a palindrome}\}$  is not regular.

[For each  $\ell \geq 1$ ,  $a^\ell b a^\ell \in L_2$  is of length  $\geq \ell$  and has property ( $\dagger$ ).]

(iii)  $L_3 \stackrel{\text{def}}{=} \{a^p \mid p \text{ prime}\}$  is not regular.

[For each  $\ell \geq 1$ , we can find a prime  $p$  with  $p > 2\ell$  and then  $a^p \in L_3$  has length  $\geq \ell$  and has property ( $\dagger$ ).]



## Example of a non-regular language that satisfies the ‘pumping lemma property’

---

$$L \stackrel{\text{def}}{=} \{c^m a^n b^n \mid m \geq 1 \text{ and } n \geq 0\} \\ \cup \\ \{a^m b^n \mid m, n \geq 0\}$$

satisfies the pumping lemma property on Slide 29 with  $\ell = 1$ .

[For any  $w \in L$  of length  $\geq 1$ , can take  $u_1 = \varepsilon$ ,  $v =$  first letter of  $w$ ,  
 $u_2 =$  rest of  $w$ .]

But  $L$  is not regular. [See Exercise ??.]

**Lemma** *If a DFA  $M$  accepts any string at all, it accepts one whose length is less than the number of states in  $M$ .*

*Proof.* Suppose  $M$  has  $\ell$  states (so  $\ell \geq 1$ ). If  $L(M)$  is not empty, then we can find an element of it of shortest length,  $a_1 a_2 \dots a_n$  say (where  $n \geq 0$ ). Thus there is a transition sequence

$$s_M = q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \cdots \xrightarrow{a_n} q_n \in \text{Accept}_M.$$

If  $n \geq \ell$ , then not all the  $n + 1$  states in this sequence can be distinct and we can shorten it as on Slide 30. But then we would obtain a strictly shorter string in  $L(M)$  contradicting the choice of  $a_1 a_2 \dots a_n$ . So we must have  $n < \ell$ . □

# Some production rules for 'English' sentences

---

SENTENCE → SUBJECT VERB OBJECT

SUBJECT → ARTICLE NOUNPHRASE

OBJECT → ARTICLE NOUNPHRASE

ARTICLE → a

ARTICLE → the

NOUNPHRASE → NOUN

NOUNPHRASE → ADJECTIVE NOUN

ADJECTIVE → big

ADJECTIVE → small

NOUN → cat

NOUN → dog

VERB → eats

## A derivation

---

SENTENCE → SUBJECT VERB OBJECT  
→ ARTICLE NOUNPHRASE VERB OBJECT  
→ the NOUNPHRASE VERB OBJECT  
→ the NOUNPHRASE eats OBJECT  
→ the ADJECTIVE NOUN eats OBJECT  
→ the big NOUN eats OBJECT  
→ the big cat eats OBJECT  
→ the big cat eats ARTICLE NOUNPHRASE  
→ the big cat eats a NOUNPHRASE  
→ the big cat eats a ADJECTIVE NOUN  
→ the big cat eats a small NOUN  
→ the big cat eats a small dog

## Example of Backus-Naur Form (BNF)

---

Terminals:

x ' + - \* ( )

Non-terminals:

id op exp

Start symbol:

exp

Productions:

id ::= x | id'

op ::= + | - | \*

exp ::= id | exp op exp | (exp)

## A context-free grammar for the language

$$\{a^n b^n \mid n \geq 0\}$$

---

Terminals:

$a$     $b$

Non-terminal:

$I$

Start symbol:

$I$

Productions:

$$I ::= \varepsilon \mid aIb$$

## Every regular language is context-free

---

Given a DFA  $M$ , the set  $L(M)$  of strings accepted by  $M$  can be generated by the following context-free grammar:

set of terminals =  $\Sigma_M$

set of non-terminals =  $States_M$

start symbol = start state of  $M$

productions of two kinds:

$$q \rightarrow aq'$$

whenever  $q \xrightarrow{a} q'$  in  $M$

$$q \rightarrow \epsilon$$

whenever  $q \in Accept_M$

**Definition** A context-free grammar is **regular** iff all its productions are of the form

$$x \rightarrow uy$$

or

$$x \rightarrow u$$

where  $u$  is a string of terminals and  $x$  and  $y$  are non-terminals.

### Theorem

- (a) *Every language generated by a regular grammar is a regular language (i.e. is the set of strings accepted by some DFA).*
- (b) *Every regular language can be generated by a regular grammar.*



## Example of the construction used in the proof of the Theorem on Slide 40

regular grammar:

$$S \rightarrow abX$$

$$X \rightarrow bbY$$

$$Y \rightarrow X$$

$$X \rightarrow a$$

$$Y \rightarrow \epsilon$$

(start symbol =  $S$ )

$\rightsquigarrow$  NFA $^\epsilon$ :

