

Regular expressions: what for?

- They are used in any task that require *pattern matching*:
 - search engines,
 - protein analysis in bioinformatics,
 - search and replace functionalities in word processors and text editors,
 - text processing utilities like `sed` and `awk` in Unix and Linux.
- Many programming languages provide various regular expression capabilities.
- Lexical analysers in compilers also use regular expressions.
- ...

Finite representation of languages

Consider the following language:

$L =$ all strings of \square s and \diamond s that have two or three occurrences of \diamond , the first and second of which are not consecutive

Can we describe this **infinite language** with some kind of **finite ‘pattern’**?

- Strings in L can start with any number (possibly none) of \square s: \square^*
- Then comes the first \diamond .

It should be alone, so must be followed by at least one \square : $\square^*\diamond\square$

- Then any number (possibly none) of \square s: $\square^*\diamond\square\square^*$
- Then comes the second \diamond : $\square^*\diamond\square\square^*\diamond$
- Then

- **either** no more \diamond s but there can be any number (possibly none) \square s,
- **or**, after some (possibly none) \square s, there is a third \diamond somewhere, followed by any number (possibly none) \square s:

$\square^*\diamond\square\square^*\diamond(\square^* \cup \square^*\diamond\square^*)$

\leadsto

regular expression representing L

From language to regular expression: Example 2

L = all the strings consisting of some number (possibly none) of a 's, followed by some number (possibly none) of b 's
= $\{\varepsilon, a, b, aa, bb, ab, aaa, aab, bbb, abb, \dots, aaaabbbbbbb, \dots\}$

Let's try: a^*b^*

\leadsto L is represented by the regular expression a^*b^*

Notation:

$L = \text{Language_of}(a^*b^*)$

read: "the language of a^*b^* is L "

or "the language represented by a^*b^* is L "

From regular expression to language

Example 3: $\text{Language_of}(a(a^* \cup b^*)) = ?$

- = all words starting with a , followed by either a (possibly empty) word of a s, or a (possibly empty) word of b s
- = $\{a, aa, ab, aaa, abb, aaaa, abbb, aaaaaa, abbbbbb, \dots\}$

Example 4: $\text{Language_of}(a(a \cup b)^*) = ?$

- = all words starting with a , followed by any word over $\{a, b\}$
- = all words over the alphabet $\{a, b\}$ starting with a
- = $\{a, aa, ab, aaa, aab, aba, abb, aaaaaa, abbaabab, aabbababba, \dots\}$

The two examples are NOT the same!

From regular expression to language: Example 5

Language_of($(b \cup aaa^*)^*$) = ?

- Let's start with Language_of(aaa^*): all words of a s of length ≥ 2
 $= \{aa, aaa, aaaa, aaaaa, \dots\}$
- Language_of($b \cup aaa^*$): as above, plus the word b
 $= \{b, aa, aaa, aaaa, aaaaa, \dots\}$
- Language_of($(b \cup aaa^*)^*$): we can take any number (possibly none) of words from Language_of($b \cup aaa^*$) and concatenate them

Some questions checking our understanding of this language:

- Is there a word of a s and b s that is not in Language_of($(b \cup aaa^*)^*$) ?
- Do the words ε , $aabbaaabab$, and $bbbabbbaaabaa$ belong to the language
Language_of($(b \cup aaa^*)^*$) ?

Summary: describing regular expressions by induction

Each regular expression (over an alphabet S) is a string consisting of

- symbols from S ,
- plus symbols from $\cup, *, (,), \varepsilon, \emptyset$.

The set of all regular expressions over S is defined inductively:

- *Basis:* \emptyset, ε and each symbol in S are regular expressions.
- *Inductive step:* If R and Z are regular expressions,
then so are $(R \cup Z)$, (RZ) , and (R^*) .
- *Closure:* No other string is a regular expression.

FOR EXAMPLE:

- Over $S = \{a, b\}$: $\emptyset, \varepsilon, ((a^*)(b^*)), (((a((a \cup b)^*))b)(b^*))$
- Over $S = \{x, y\}$: $\emptyset, \varepsilon, (((x(y^*))y)(y \cup z)), (((x(y^*))y)(y \cup z))$
- Over $S = \{\diamond, \square\}$: $\emptyset, \varepsilon, ((\square \cup (\diamond^*))((\diamond \square)^*))$

Regular expressions: notational conventions

These brackets are pretty incomprehensible. So we introduce some conventions:

- We omit the outermost brackets.
- We omit the brackets when concatenate expressions, so e.g. we write $ababb$ instead of $((ab)(ab))b$.
- $*$ 'binds tighter' than concatenation, so e.g. we write aab^* instead of $(aa)(b^*)$.

For example: (cf. the previous slide)

$$a^*b^*, \quad a(a \cup b)^*bb^*, \quad xy^*y(y \cup z), \quad (\square \cup \diamond^*)(\diamond\square)^*$$

But take care: say, $(aa)^*b$ and aa^*b are NOT the same!

Always use brackets if in any doubt!

From regular expression to language: Example 6

$$\text{Language_of}((b \cup a)aa^*) = ?$$

Let's do it, 'from inside out':

- Language_of($b \cup a$) = $\{a, b\}$
- Language_of($(b \cup a)a$) = $\{aa, ba\}$
- Language_of($(b \cup a)aa^*$) = all words starting with aa followed by a (possibly empty) word of a s, and all words starting with ba followed by a (possibly empty) word of a s.

Compare this language with $\text{Language_of}(b \cup aaa^*)$ (see Example 5).

Brackets DO MATTER: the two languages are NOT the same!

Summary: regular expressions represent languages

Every regular expression over some alphabet S **represents** a language over S .

These are two different things:

a regular expression R

\longleftrightarrow

the language R represents: $\text{Language_of}(R)$

a string consisting of
symbols from S and

$\cup, *, (,), \varepsilon, \emptyset$

a set of words over S

FOR EXAMPLE:

$1(1 \cup 2)^*$

$\text{Language_of}(1(1 \cup 2)^*) =$ all words of 1s and 2s
starting with 1

$(a \cup b)bba$

$\text{Language_of}((a \cup b)bba) = \{abba, bbba\}$

finite 'pattern'

words following the pattern (can be infinitely many)

Summary: how regular expressions represent languages

Inductive definition of the language $\boxed{\text{Language_of}(R)}$ represented by the regular expression R :

We follow the inductive definition of regular expressions on slide 224.

Basis: $\text{Language_of}(\emptyset) = \emptyset$

$\text{Language_of}(\varepsilon) = \{\varepsilon\}$

$\text{Language_of}(a) = \{a\}$ for any symbol a in the alphabet S .

Inductive step: If R and Z are regular expressions, then

- $\text{Language_of}(R \cup Z) =$ all words that belong either to $\text{Language_of}(R)$
or to $\text{Language_of}(Z)$
- $\text{Language_of}(RZ) =$ any word from $\text{Language_of}(R)$
followed by any word from $\text{Language_of}(Z)$
- $\text{Language_of}(R^*) =$ any word from $\text{Language_of}(R)$
followed by any word from $\text{Language_of}(R)$
followed by any word from $\text{Language_of}(R) \dots$

the number of iterations is arbitrary (possibly none)

From regular expression to language: Example 7

$$\text{Language_of}(\varepsilon \cup c^*(a \cup bc^*)) = ?$$

all the strings that are either empty or start with some (possibly none) cs ,
followed by either an a , or a b followed by some (possibly none) cs

$$= \{\varepsilon, a, b, bc, bcc, ca, cb, cbcc, \dots\}$$

NOT in Language_of($\varepsilon \cup c^*(a \cup bc^*)$): $ab, aa, caac, \dots$

Regular languages

A regular language

is *any* language that can be represented by a regular expression.

In other words, a language L is regular if $L = \text{Language_of}(R)$
for some regular expression R .

For instance, all the languages in Examples 1–7 above are regular.

NOT every language is regular!

Regular languages and finite automata

- (1) **There is a general ‘mechanical’ procedure \mathcal{P}_1 that converts any regular expression R to an NFA A_R such that $L(A_R) = \text{Language_of}(R)$.**

So every regular language is accepted by some automaton.

Moreover, there is a general way back:

- (2) **There is a general ‘mechanical’ procedure \mathcal{P}_2 that converts any NFA A to a regular expression R_A such that $\text{Language_of}(R_A) = L(A)$.**

So every language that is accepted by some automaton is regular.

In other words:

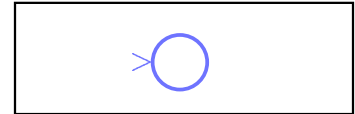
Regular languages are precisely those languages that are accepted by finite automata.

From regular expressions to NFA

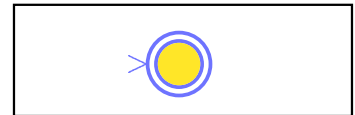
The procedure which
converts any regular expression R to an NFA A_R such that $L(A_R) = \text{Language_of}(R)$
operates along the inductive definition of the regular expression R
(again, see slide 224):

Basis cases:

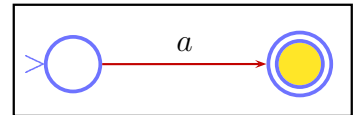
- If $R = \emptyset$. Then $\text{Language_of}(R) = \emptyset$. A_R :



- If $R = \varepsilon$. Then $\text{Language_of}(R) = \{\varepsilon\}$. A_R :

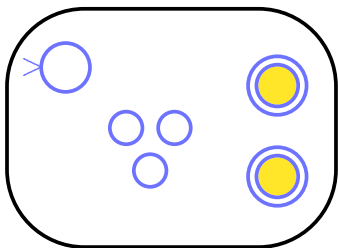


- If $R = a$. Then $\text{Language_of}(R) = \{a\}$. A_R :

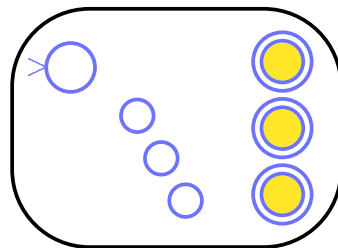


Inductive cases: Automaton accepting $\text{Language_of}(R \cup Z)$

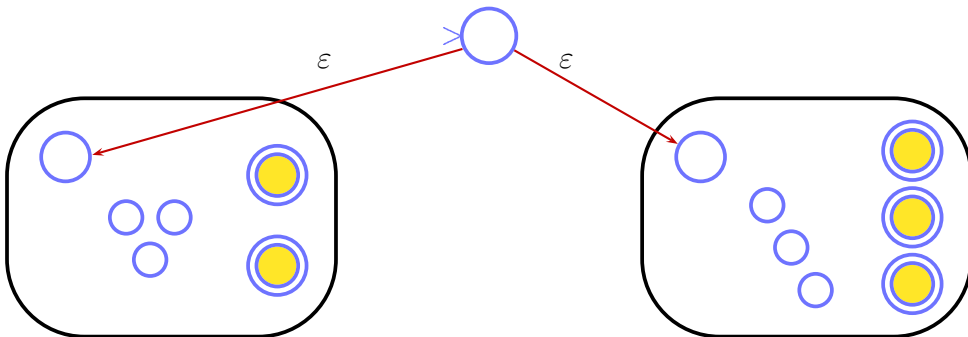
A_R : accepts $\text{Language_of}(R)$



A_Z : accepts $\text{Language_of}(Z)$

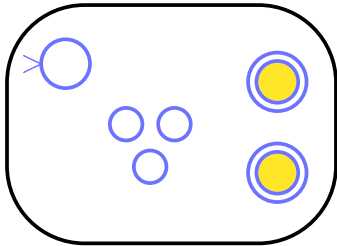


$A_{R \cup Z}$: accepts $\text{Language_of}(R \cup Z)$

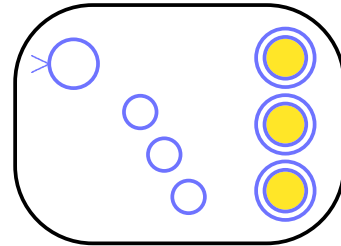


Inductive cases: Automaton accepting $\text{Language_of}(RZ)$

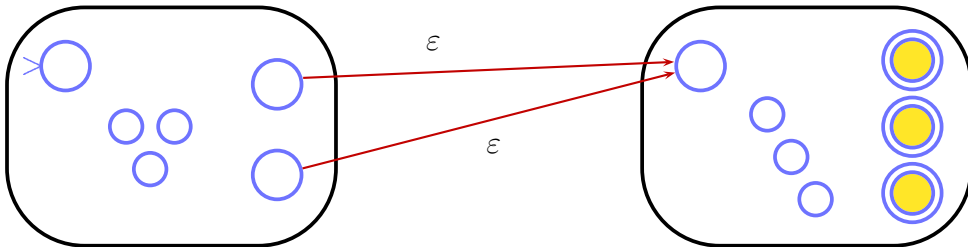
A_R : accepts $\text{Language_of}(R)$



A_Z : accepts $\text{Language_of}(Z)$

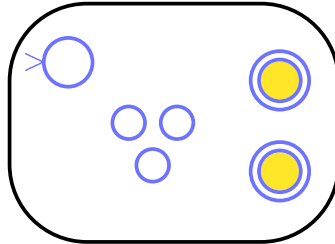


A_{RZ} : accepts $\text{Language_of}(RZ)$

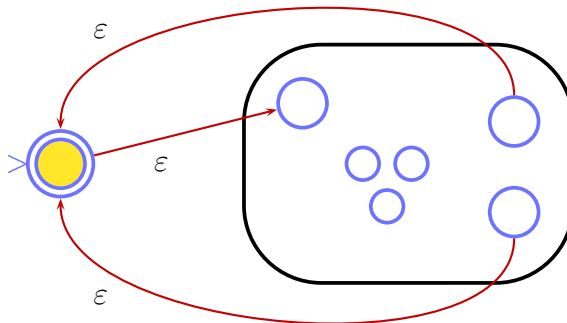


Inductive cases: Automaton accepting $\text{Language_of}(R^*)$

A_R accepts $\text{Language_of}(R)$



A_{R^*} accepts $\text{Language_of}(R^*)$



Example: how the procedure works

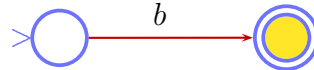
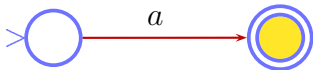
We apply the procedure to the regular expression

$$((a \cup ab)^*ba)^*$$

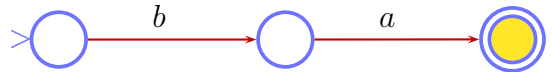
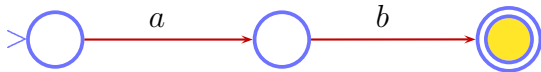
We construct step by step (going 'inside out') an NFA A such that

$$L(A) = \text{Language_of}(((a \cup ab)^*ba)^*)$$

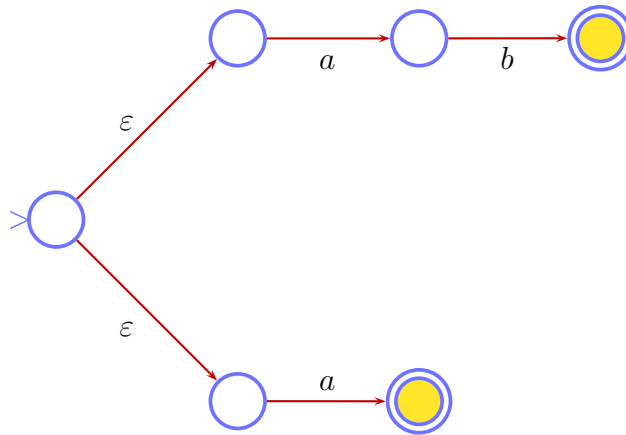
Step 0: automata accepting $\text{Language_of}(a)$ and $\text{Language_of}(b)$



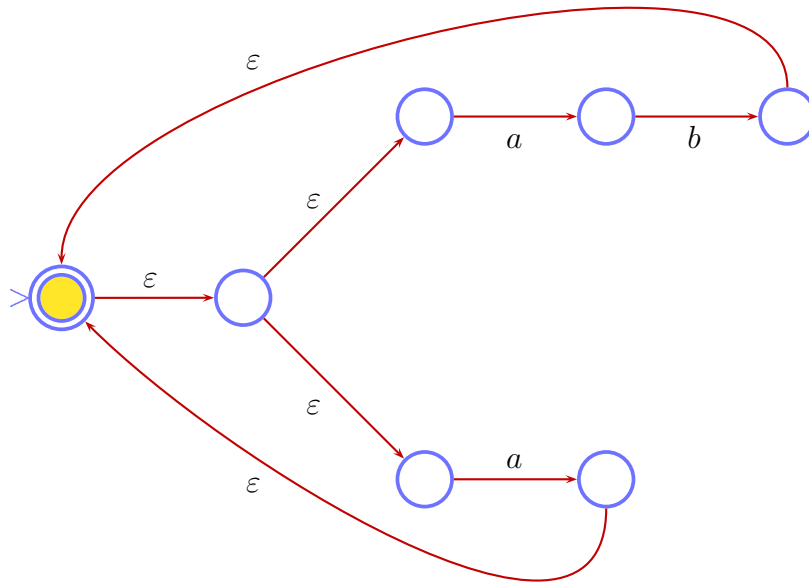
Step 1: automata accepting $\text{Language_of}(ab)$ and $\text{Language_of}(ba)$



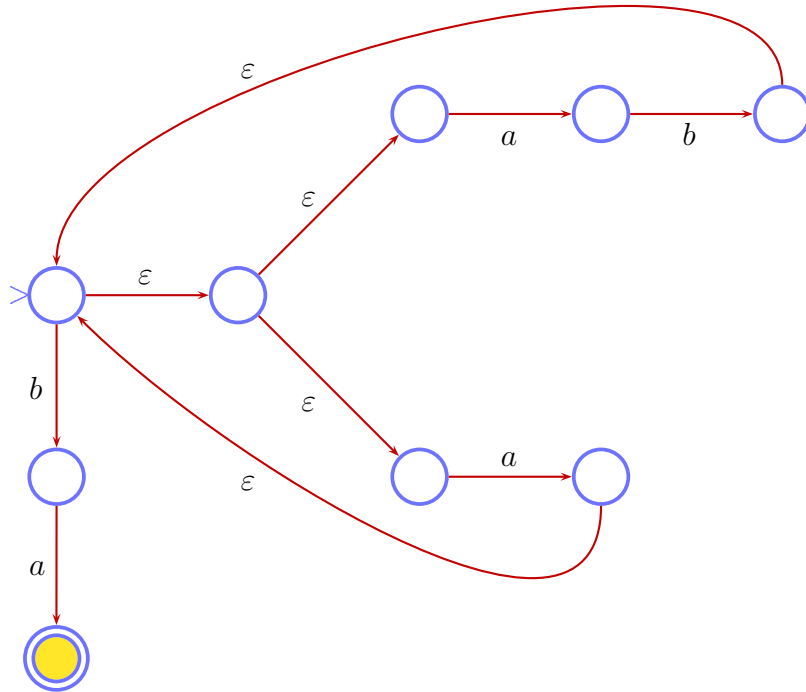
Step 2: automaton accepting Language_of($a \cup ab$)



Step 3: automaton accepting Language_of($(a \cup ab)^*$)

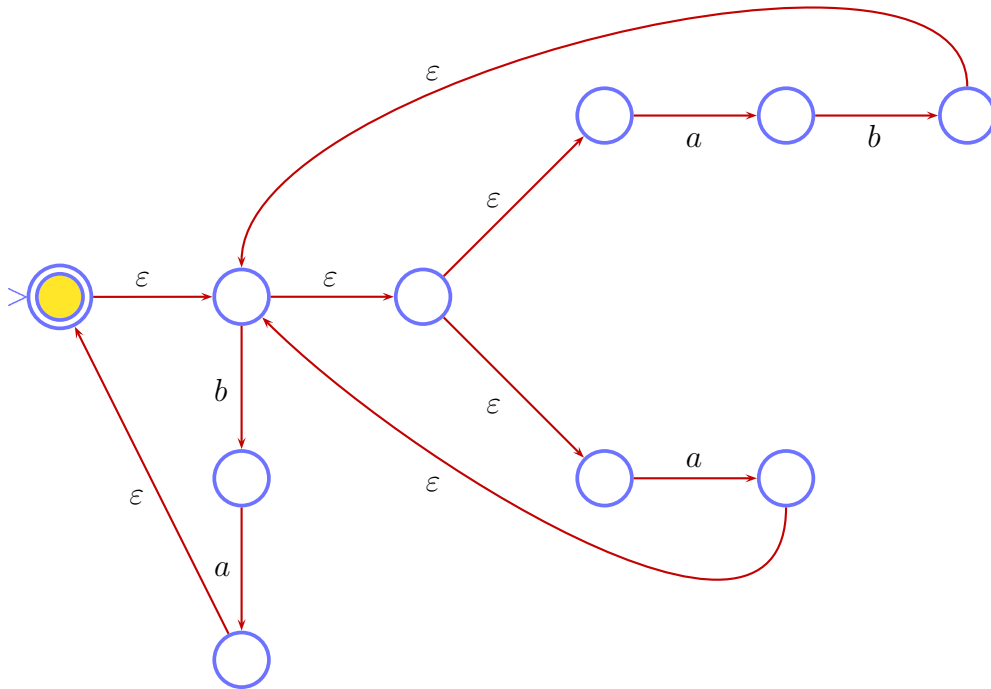


Step 4: automaton accepting Language of $((a \cup ab)^*ba)$



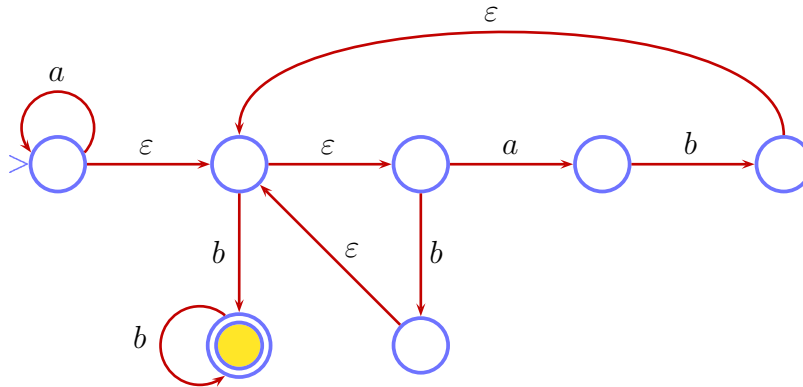
Final step of the procedure

Automaton accepting Language.of $\left((a \cup ab)^* ba \right)^*$:



Another example

Automaton accepting Language of $(a^*(b \cup ab)^*bb^*)$:



Finite automata and regular languages: summary

- Finite automata may be regarded as programs that use **fixed amounts of memory** (represented by states) regardless of the input.
- Finite automata can be used as **recognition devices**:
they accept certain inputs and reject others.
- **Nondeterminism** does not increase the computational power of finite automata, but nondeterministic automata are easier to design than **deterministic** ones.
- The languages accepted by finite automata are precisely the **regular** ones.

Nonregular languages

Since finite automata are theoretical models for programs using a constant amount of memory regardless of the input, one ought to expect that the power of such programs should be quite limited.

Since regular languages are all accepted by some finite automaton, there must be languages which are **not regular**, they cannot be given with the help of *any* regular expression.

Nonregular languages: an example

Consider the following language over the alphabet $\{a, b\}$:

$$\begin{aligned} L &= \text{all words starting with a word of } a\text{'s} \\ &\quad \text{followed by an equal-length word of } b\text{'s} \\ &= \{a^n b^n \mid n = 0, 1, 2, \dots\} \end{aligned}$$

Suppose that a finite automaton A tries to recognise words in this language. Then A must store the entire prefix a^n before the first b shows up. Otherwise A will not be able to compare the length of the coming word of b 's with the length of the prefix of a 's.

As each automaton is capable for a fixed finite amount of storage with the help of its states, no automaton A exists such that $L(A) = L$.

There is a precise mathematical proof justifying this *informal* argument:

this language L is indeed **not regular**.