# LANGAGES FORMELS

ENS Paris-Saclay

DER informatique

2ème semestre 2025-26

# Organisation du cours → https://samvangool.net/langform.html

- **Deux parties** de 6 semaines chacune :

  1. langages reguliers, automates, monoïdes — 22 janvier – 12 mars

     Sam van Gool (cours) et Guillaume Scerri (TD) → CCTD1

     → **Partiel:** 19 mars, 14h, salle 1Z61 (!)

  2. Langages algebriques, analyse syntaxique, langages d'arbres — 26 mars – 11 mai

     Stefan Schwoon (cours) et Luc Lapointe (TD) → CCTD2

     → **Examen:** Semaine du 18 mai

     → **Projet** 2ème partie

- **Conditions de validation** :

  - Session 1 : CCTD1 (1), CCTD2 (1), projet (1), partiel (2), examen (2)
  - Session 2: " , " , " , examen (4)

# 1. Words, languages, automata

The basic building blocks.

The rest of the slides are written in English, because this way, your lecturer will write less (fewer?!) grammar mistakes. We will speak French in class. If any of this causes (informal) language issues, don't hesitate to tell me. I hope you'll come to appreciate the mix as a feature, not a bug!

Let $\Sigma$ be a finite set, which we call the alphabet.

A word is a finite sequence of symbols from $\Sigma$. The set of words is denoted $\Sigma^*$.

Example. $-\Sigma = \{0,1\}$. Words are binary sequences, e.g., 1001, 01001, 000.

$\nabla$ $-\Sigma = \{a, b, ab\}$. The expression $abab$ does not define a word in this alphabet.

We would need to write $(a, b, ab)$ or $(ab, ab)$ or $(ab, a, b)$.

$\Rightarrow$ We usually avoid taking such sets as alphabets, and assume unique parsability without commas.

We usually write $\varepsilon$ for the empty word, i.e., the unique sequence of length 0. $\Sigma^+ := \Sigma^* \setminus \{\varepsilon\}$.

Given words $u, v \in \Sigma^*$, we can form their concatenation, $u \cdot v$. We often omit the $\cdot$.

We write $|u|$ for the length of $u$, defined inductively: $|\varepsilon| := 0$, $|ua| := |u| + 1$ for all $u \in \Sigma^*, a \in \Sigma$.

Fact. The structure $(\Sigma^*, \cdot, \varepsilon)$ is a monoid, which is free over $\Sigma$.

This means: for all $u, v, w \in \Sigma^*$, $(u \cdot v) \cdot w = u \cdot (v \cdot w)$, $u \cdot \varepsilon = u = \varepsilon \cdot u$, (monoid)

and any function $\Sigma \to M$, with $M$ a monoid, has a unique homomorphic extension $\Sigma^* \to M$ (free).

Let $w \in \Sigma^*$.

A **prefix** of $w$ is a word $u \in \Sigma^*$ such that
there exists $v \in \Sigma^*$ such that $w = uv$.

A **suffix** of $w$ is a word $v \in \Sigma^*$ such that
there exists $u \in \Sigma^*$ such that $w = uv$.

A **factor** of $w$ is a word $f \in \Sigma^*$ such that
there exist $u, v \in \Sigma^*$ such that $w = ufv$.

A **subword** of $w$ is a word $s = (s_1, \ldots, s_n) \in \Sigma^*$ $(n \geq 0)$ such that
there exist $u_0, u_1, \ldots, u_n \in \Sigma^*$ such that $w = u_0 s_1 u_1 \ldots s_n u_n$.

A $language$ over the alphabet $\Sigma$ is a subset of $\Sigma^*$.

$Formal\ language\ theory$ is the study of the set $\mathcal{P}(\Sigma^*)$ of languages.

$Fact$. There are uncountably many languages (if $\Sigma \neq \emptyset$, which we will always assume).

( Exercise ) $\rightarrow$ By this I mean: "If you don't know how to prove it, please try, or look it up

online, or in a book, or ask a friend, or ask me, or..."

I do NOT mean: "This is easy and you should feel bad if you find it difficult."

(We just don't have time to discuss everything.)

G. Cantor (1845-1918)

$\Rightarrow$ We need methods of describing some of these languages, at least.

$Automata$ and $regular\ expressions$ are such methods.

They are fundamental to syntactic analysis, logic, verification, and more.

photo credit $\rightleftharpoons$
http://www.math.uni-hamburg.de/home/grothkopf/fotos/math-ges/

An *automaton* over $\Sigma$ is a directed multigraph with $\Sigma$-labeled edges and two distinguished subsets.

Explicitly, an automaton is a quintuple $A = (Q, \Sigma, \delta, I, F)$, where

- $Q$ is a set, whose elements are called *states*;

- $\Sigma$ is an alphabet;

- $\delta \subseteq Q \times \Sigma \times Q$, its elements are called *transitions* or *edges*;

- $I \subseteq Q$ a set of *initial states*;

- $F \subseteq Q$ a set of *final states*.

We often denote an edge $(q, a, r)$ as $q \xrightarrow{a} r$. Its *source* is $q$, its *target* is $r$, its *label* is $a$.

We usually assume (without saying so) that $Q$ is finite.

Acronym: NFA ("AFN" en français)

An automaton $A = (Q, \Sigma, \delta, I, F)$ accepts a word $w \in \Sigma^*$ if there exists a path from an initial state to a final state that is labeled by $w$. Otherwise, $A$ rejects $w$.

Explicitly, $A$ accepts $w$ provided that there exists a sequence $\pi \in Q^*$ with $|\pi| = |w| + 1$ such that:

- $\pi_0 \in I$,

- $\pi_{|w|} \in F$,

- for each $0 \leq i < |w|$, $\pi_i \xrightarrow{w_i} \pi_{i+1}$ in $\delta$

$\left.\right\}$ $\pi$ is a **successful run** on $w$

The **language recognized by** $A$ is $\mathcal{L}(A) := \{ w \in \Sigma^* \mid A \text{ accepts } w \}$.

A language $L \subseteq \Sigma^*$ is **recognizable** if there exists an NFA that recognizes it, and we put

$$\text{Rec}(\Sigma^*) := \{ L \subseteq \Sigma^* \mid L \text{ recognizable} \}.$$

**Question.** When does $\varepsilon$ belong to $\mathcal{L}(A)$?

By this, I mean: "I think somebody will have a reasonable guess. Please answer me to avoid long awkward silences!"

An automaton $A = (Q, \Sigma, \delta, I, F)$ is **deterministic** if, for each $a \in \Sigma$, the relation

$$\delta_a := \delta \cap (Q \times \{a\} \times Q) \text{ is functional and total, and } \#I = 1.$$

That is, for every $q \in Q$, $a \in \Sigma$, there exists a **unique** $r \in Q$ such that $q \xrightarrow{a} r$.

We sometimes write $q \cdot a$ for this unique state.

The transition relation of a DA can also be viewed as a function $\bullet : Q \times \Sigma \to Q$,

or as a function $\rho : \Sigma \longrightarrow Q^Q$.

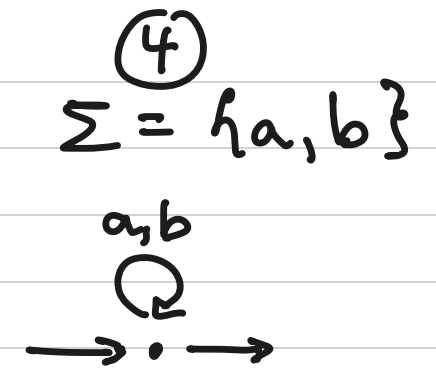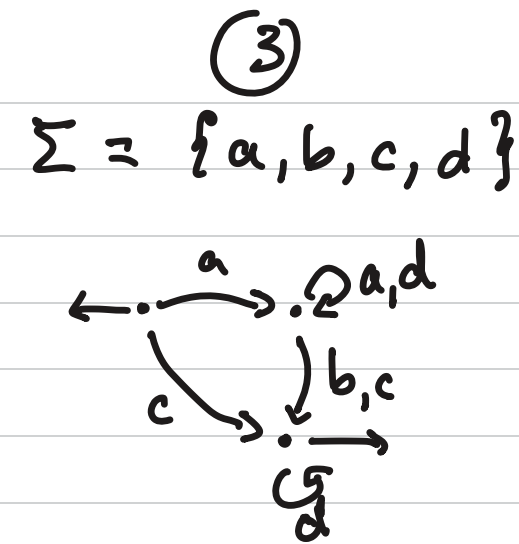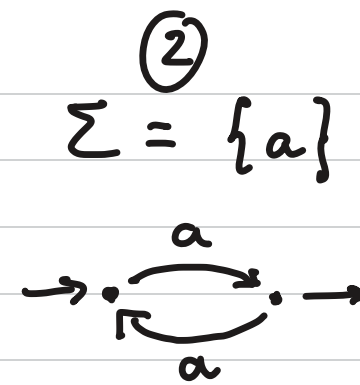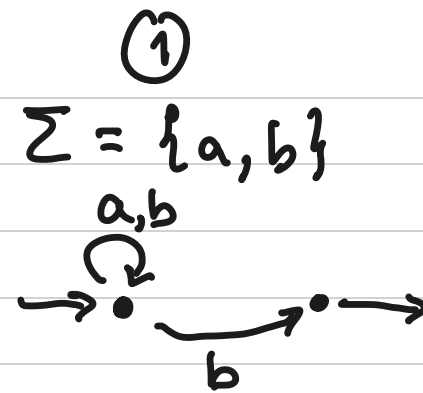**Fact.** For any set $Q$, the triple $(Q^Q, \circ, id_Q)$ is a monoid.

Thus, by the **free** property of $\Sigma^*$, there exists a unique homomorphic extension

$$\hat{\rho} : \Sigma^* \longrightarrow Q^Q \text{ of } \rho.$$

Explicitly, $\hat{\rho}$ can be defined by induction: $\hat{\rho}(\varepsilon) := id_Q$, and, for any $w \in \Sigma^*$, $a \in \Sigma$,

$$\hat{\rho}(wa) = \lambda q \cdot \rho(a)(\hat{\rho}(w)(q)). \quad (\text{"the function sending } q \in Q \text{ to: } \rho(a) \text{ applied to } \hat{\rho}(w)(q)\text{"}).$$

We also write $\tilde{\delta}(w) := q_0 \cdot w$, where $q_0$ is the initial state.

**Examples.**

① $\Sigma = \{a, b\}$



② $\Sigma = \{a\}$



③ $\Sigma = \{a, b, c, d\}$



④ $\Sigma = \{a, b\}$



$\longrightarrow \bullet$ : initial

$\bullet \longrightarrow$ : final

**Questions.** — What is $\mathcal{L}(A)$?

— Are there other automata recognizing the same language? How many? Bigger? Smaller?

— Is $A$ deterministic? If not, can we make it so, without changing $\mathcal{L}(A)$?

# 2. <u>Determinization</u>

Constructing an equivalent DFA out of an NFA using the power set.

**Theorem.** Any recognizable language can be recognized by a <span style="color:red">deterministic</span> automaton.

**Proof.** Let $A = (Q, \Sigma, \delta, I, F)$ be an automaton recognizing L.

Define the automaton $P(A) := (\mathcal{P}Q, \Sigma, \Delta, \{I\}, G)$, where:

- $\mathcal{P}Q$ is the set of subsets of $Q$ (the <span style="color:red">power set</span> of $Q$)

- $\Delta := \{(S, a, T) \in \mathcal{P}Q \times \Sigma \times \mathcal{P}Q \mid T = \{q \in Q \mid \exists s \in S, (s, a, q) \in \delta\}\}$

- $G := \{S \in \mathcal{P}Q \mid S \cap F \neq \emptyset\}$.

Note: $P(A)$ is deterministic.

We now claim that $L(A) = L(P(A))$.

**Proof of claim.** We will show by induction that, for any $w \in \Sigma^*$, $R(w) :=$

$$\tilde{\Delta}(w) = \{q \in Q \mid \text{there exists a } w\text{-path in } A \text{ such that } \pi_0 \in I \text{ and } \pi_{|w|} = q\}$$

**Base case.** $w = \varepsilon$ : both sets are equal to $I$.

**Inductive case.** $w = ua$ for $u \in \Sigma^*$, $a \in \Sigma$. By the induction hypothesis, $\tilde{\Delta}(u) = R(u)$.

Let $q \in Q$. $\quad q \in \tilde{\Delta}(w) \iff$ there exists $s \in \tilde{\Delta}(u)$ such that $(s, a, q) \in \delta \quad$ (def. of $\Delta$)

$\iff$ there exists a $u$-path $\pi$ s.t. $\pi_0 \in I$ and $\pi_{|u|} = s$ and $(s, a, q) \in \delta$ (def. of $R$)

$\iff q \in R(w)$ (def of $w$-path). $\square$ Now, $w \in L(P(A)) \iff \tilde{\Delta}(w) \in G$

$\iff R(w) \cap F \neq \emptyset \iff w \in L(A).$

The theorem gives a **determinization** construction.

Note that $\#\mathcal{P}(Q) = 2^{\#Q}$, so it has exponential cost, at worst.

**Question.** Is there a "better" determinization?

**Example.**

$$\xrightarrow{\quad} \overset{a,b}{\overset{\curvearrowright}{\cdot}} \xrightarrow{a} \cdot \xrightarrow{a,b} \cdot \xrightarrow{a,b} \cdot \xrightarrow{a,b} \cdot \xrightarrow{a,b} \cdot \xrightarrow{\quad} \text{recognizes} \quad \{a,b\}^* \cdot \{a\} \cdot \{a,b\}^4.$$

We get $2^6$ states if we determinize.

More generally, $L_n := \{a,b\}^* \cdot \{a\} \cdot \{a,b\}^n$ can be recognized by an NFA with $n+2$ states.

**However:** Any **deterministic** automaton recognizing $L_n$ requires at least $2^n$ states.

(We will see a proof of this later.)

The membership problem takes as input an automaton $A = (Q, \Sigma, \delta, I, F)$ and a word $w \in \Sigma^*$, and asks to determine whether or not $w \in \mathcal{L}(A)$.

Algorithm. We will keep track of a variable $W$ which contains the reachable states after reading $w$.

1) $W \leftarrow I$

2) while $w \neq \varepsilon$ :

3)      $new \leftarrow \emptyset$

4)      for all $q \in W$ :

5)          $new \leftarrow new \cup \{r \in Q \mid (q, head(w), r) \in \delta\}$

6)      $W \leftarrow new$

7)      $w \leftarrow tail(w)$        Runtime: $O(|w| \cdot (\#Q)^2)$.

8) return $(W \cap F \neq \emptyset)$

When $A$ is known to be deterministic, we can achieve $O(|w|)$ (exercise).

The **emptiness** problem asks, given an automaton $A$, whether or not $\mathcal{L}(A) = \emptyset$.

The **universality** problem " " " " " " " " $\mathcal{L}(A) = \Sigma^*$.

For a DFA $A$, first compute Reach $:= \{ q \in Q \mid$ there exist $i \in I$ and a path $i \to q \}$.

(how? complexity?)

- $\mathcal{L}(A) = \emptyset \iff F \cap \text{Reach} = \emptyset$

- $\mathcal{L}(A) = \Sigma^* \iff \text{Reach} \subseteq \bar{F}$.

For an NFA, emptiness can be done in the same way.

universality is PSPACE-complete.

# 3. Closure properties

Building recognizable languages via automata constructions

**Theorem.** Let $L \in \text{Rec}(\Sigma^*)$. The complement $\Sigma^* \setminus L$ is recognizable, too.

**Proof.** Pick $(Q, \Sigma, \delta, I, F)$ a DFA recognizing $L$. (This exists thanks to determinization!)

The DFA $(Q, \Sigma, \delta, I, Q \setminus F)$ recognizes $\Sigma^* \setminus L$.

Indeed, for $w \in \Sigma^*$, $w \in \Sigma^* - L \iff \tilde{\delta}(w) \notin F$. ⟥

**Theorem.** Let $L_1, L_2 \in \text{Rec}(\Sigma^*)$. The intersection $L_1 \cap L_2$ is recognizable, too.

**Proof.** Pick $A_i = (Q_i, \Sigma, \delta_i, I_i, F_i)$ an NFA recognizing $L_i$, for $i = 1, 2$.

Define the automaton $A_1 \times A_2 := (Q_1 \times Q_2, \Sigma, \delta, I_1 \times I_2, F_1 \times F_2)$

where $\delta := \{ ((q_1, q_2), a, (r_1, r_2)) \in (Q_1 \times Q_2) \times \Sigma \times (Q_1 \times Q_2) \mid$

$$q_1 \xrightarrow{a} r_1 \text{ in } A_1 \text{ and } q_2 \xrightarrow{a} r_2 \text{ in } A_2 \}.$$

Then (claim!) $L(A_1 \times A_2) = L(A_1) \cap L(A_2)$.

**Pf_of_claim** (idea). An induction on $w \in \Sigma^*$ shows that a $w$-path in $A_1 \times A_2$

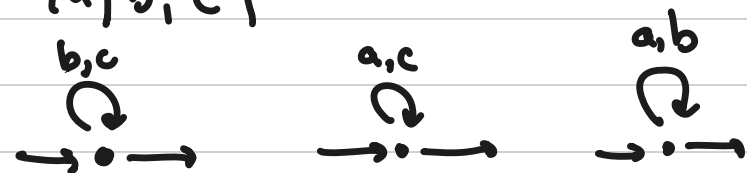is precisely given by a pair of $w$-paths in $A_1$ and $A_2$. ⬜ ⬜

**Example.**  For $\Sigma$ a finite alphabet, let $NA(\Sigma) := \{w \in \Sigma^* \mid \{a \in \Sigma \mid \exists_p, w_p = a\} \neq \Sigma\}$.

↳ "not all"

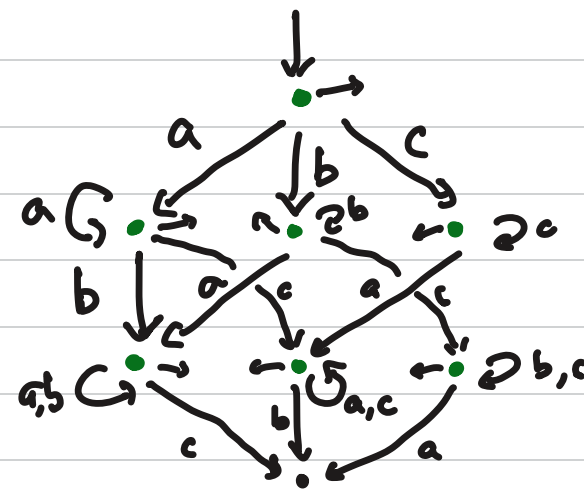The automaton $(\Sigma, \Sigma, \delta, \Sigma, \Sigma)$ with

$$\delta = \{(a,b,c) \in \Sigma \times \Sigma \times \Sigma \mid a \neq b, a = c\}$$
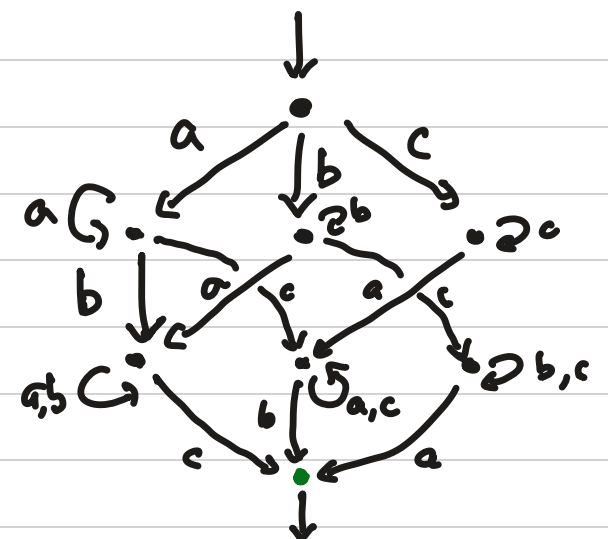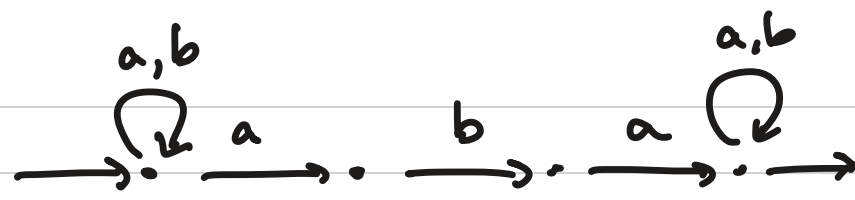
recognizes $NA(\Sigma)$.

$\Sigma = \{a, b, c\}$



**Q** Is there a smaller automaton than that one which recognizes $\Sigma^* - NA(\Sigma)$?

A smaller DFA? If not, how can we be sure?

**Example**

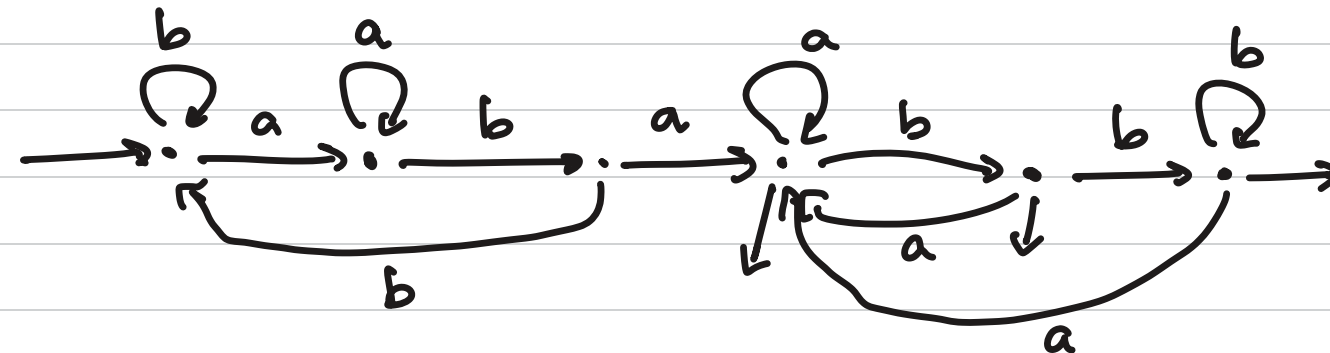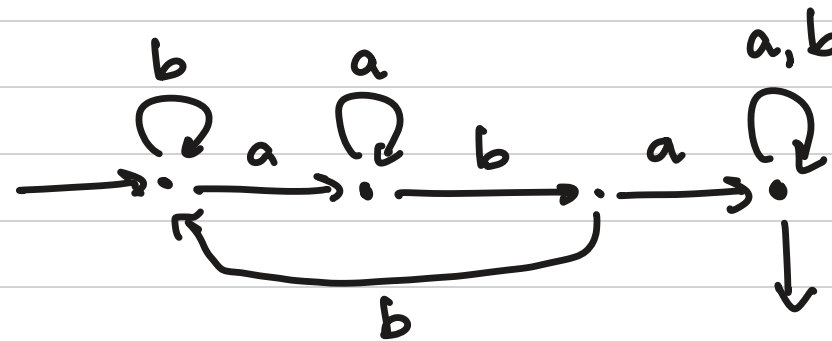$L = \{ w \in \Sigma^* \mid aba \text{ is a factor of } w \}$

*determinize*   $2^4 = 16$ states.



*delete inaccessible states*



*collapse last 3 states*

   only 4 states!

For $L \in \text{Rec}(\Sigma^*)$, a DFA $A$ such that $\mathcal{L}(A) = L$ is called **minimal**

   if, for every DFA $B$ such that $\mathcal{L}(B) = L$, $\#Q_B \geq \#Q_A$.

We will prove later that every recognizable $L$ has a unique minimal DFA.

**Corollary.** Let $L_1, L_2 \in Rec(\Sigma^*)$. Then $L_1 \cup L_2 \in Rec(\Sigma^*)$.

**Proof 1.** $L_1 \cup L_2 = \Sigma^* - ((\Sigma^* - L_1) \cap (\Sigma^* - L_2))$. $\square$

**Proof 2.** Let $A_i$ be an automaton recognizing $L_i$ $(i = 1, 2)$.

Define $A = A_1 \uplus A_2 = (Q_1 \uplus Q_2, \Sigma, \delta, I_1 \uplus I_2, F_1 \uplus F_2)$,

where $\delta := \{(q, a, r) \mid (q \in Q_1, r \in Q_1, \text{ and } (q, a, r) \in \delta_1) \text{ or } (q \in Q_2, r \in Q_2, \text{ and } (q, a, r) \in \delta_2\}$.

Since $Q_1$ and $Q_2$ are disconnected, if $\pi$ is an accepting path in $A$ starting in $I_k$, then it must end in $F_k$.

Thus, $L(A) = L(A_1) \cup L(A_2)$. $\square$

**Question.** Which proof do you prefer? Why?

Let $K, L \subseteq \Sigma^*$. The concatenation of $K$ and $L$ is

$$K \cdot L := \{ u \cdot \sigma \mid u \in K, \sigma \in L \}.$$

Let $w \in \Sigma^*$. The left quotient of $L$ with respect to $w$ is

$$w^{-1} L := \{ u \in \Sigma^* \mid wu \in L \},$$

and the right quotient of $L$ w.r.t. $w$ is

$$L w^{-1} := \{ u \in \Sigma^* \mid uw \in L \}.$$

The left and right residuals of $L$ w.r.t. $K$ are

$$K \backslash L := \bigcap_{w \in K} w^{-1} L \qquad \text{and} \qquad L / K := \bigcap_{w \in K} L w^{-1}$$

$$= \{ u \in \Sigma^* \mid \text{for all } w \in K, wu \in L \} \qquad\qquad = \{ u \in \Sigma^* \mid \text{for all } w \in K, uw \in L \}.$$

For $n \in \mathbb{N}$, we define inductively the power by $L^0 := \{ \varepsilon \}$, and $L^{n+1} := L^n \cdot L$.

We define the Kleene star by $L^* := \bigcup_{n \geq 0} L^n$. Also $L^+ := \bigcup_{n \geq 1} L^n$ $(\neq L^* \smallsetminus \{\varepsilon\}$ in general!$)$

**Theorem.** Let $L \in \text{Rec}(\Sigma^*)$. For any $w \in \Sigma^*$, $w^{-1}L \in \text{Rec}(\Sigma^*)$ and $Lw^{-1} \in \text{Rec}(\Sigma^*)$.

**Proof.** Pick $A = (Q, \Sigma, \delta, I, F)$ an automaton for $L$. Define $A' = (Q, \Sigma, \delta, I', F)$, where

$$I' := \{ q \in Q \mid \text{there exist } q_0 \in I \text{ and a } w\text{-path from } q_0 \text{ to } q \}$$

Then, for any $u \in \Sigma^*$,

$A'$ accepts $u$ $\iff$ there exist $q \in I', r \in F$, and a $u$-path from $q$ to $r$

$\iff$ there exist $q_0 \in I, r \in F$, a $w$-path from $q_0$ to $q$, and

a $u$-path from $q$ to $r$

$\iff$ there exist $q_0 \in I, r \in F$, and a $(wu)$-path from $q_0$ to $r$

$\iff$ $A$ accepts $wu$.

So $L(A') = w^{-1}L(A) = w^{-1}L$.

The statement about $Lw^{-1}$ is proved similarly (exercise). $\qquad \square$

**Corollary.** Let $L \in \text{Rec}(\Sigma^*)$. The set $\{w^{-1}L : w \in \Sigma^*\}$ is finite.

**Proof.** Let $A$ be an automaton such that $\mathcal{L}(A) = L$.

We showed in the previous proof that, for any $w \in \Sigma^*$, $w^{-1}L$ is recognized by a variant of $A$ obtained by changing the inital states. Thus,

$$\{w^{-1}L : w \in \Sigma^*\} \subseteq \{\mathcal{L}(A') \mid A' \text{ a initial-state-variant of } A\}.$$

The second set has at most $2^{\#Q_A}$ elements. $\square$

**Exercise.** Let $L \in \text{Rec}(\Sigma^*)$. For any $K \subseteq \Sigma^*$, $K\backslash L$ and $L/K$ are recognizable.

**Theorem.** For any $K, L \in \text{Rec}(\Sigma^*)$, $K \cdot L$ is recognizable.

For the proof, we will use a variant of automata.

**Definition.** An **automaton with $\varepsilon$-transitions** is a tuple $A = (Q, \Sigma, \delta, I, F)$, where $Q, \Sigma, I, F$ are as in the definition of automaton, and $\delta \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times Q$.

$A$ **accepts** a word $w \in \Sigma$ if there exist $k_0, k_1, \ldots, k_{|w|} \in \mathbb{N}_{\geqslant 0}$ such that the NFA $A' := (Q, \Sigma \cup \{\varepsilon\}, \delta, I, F)$ accepts the word

$$\varepsilon^{k_0} w_1 \varepsilon^{k_1} \ldots \varepsilon^{k_{|w|-1}} w_{|w|} \varepsilon^{k_{|w|}}.$$

Equivalently: (exercise)

Let $w \in \Sigma^*$. A **$w$-path** in $A$ is a finite word of edges $\pi \in \delta^*$ such that: (1) for every $0 \leq i < |\pi| - 1$, the target of $\pi_i$ = the source of $\pi_{i+1}$,

(2) the concatenation of the labels of $\pi$ is equal to $w$

and **acceptance** is defined as for NFA's.

**Fact.** Any automaton with $\varepsilon$-transitions can be transformed into an automaton without $\varepsilon$-transitions that recognizes the same language.

**Proof.** See TD1. □

**Proof of Theorem.** Let $A$ and $B$ be automata.

Construct the automaton with $\varepsilon$-transitions $C = (Q_C, \Sigma, \delta_C, I_A, F_B)$, where

$$Q_C := Q_A \uplus Q_B, \text{ and}$$

$$\delta_C := \delta_A \cup \delta_B \cup \{(q, \varepsilon, r) \mid q \in F_A, r \in I_B\}.$$

**Claim.** $L(C) = L(A) \cdot L(B)$

**Proof.** Let $w \in \Sigma$. $\quad w \in L(A) \cdot L(B) \overset{\text{def.}}{\iff}$ there exist $u \in L(A), v \in L(B)$ such that $w = uv$
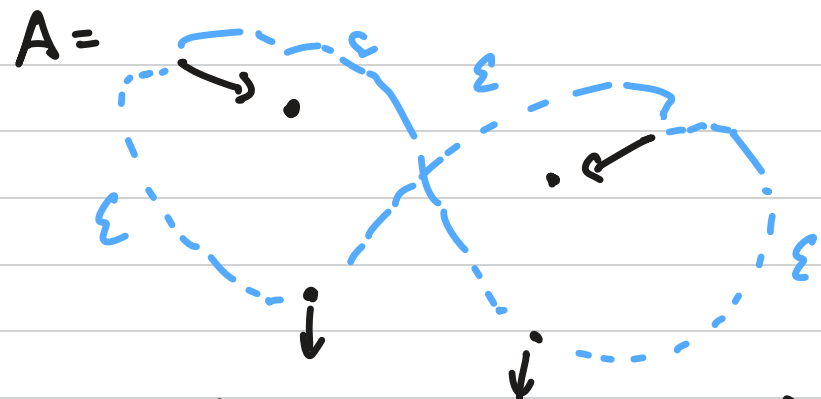
$\overset{\text{def.}}{\iff}$ there exist paths $\pi_1 : q_0 \overset{u}{\longrightarrow} q_1$ in $A$, $\pi_2 : r_0 \overset{v}{\longrightarrow} r_1$ in $B$ with $q_0 \in I_A, q_1 \in F_A$, $r_0 \in I_B, r_1 \in F_B$

**why?** $\iff$ there exists $w$-path $\pi$ in $C$ starting in $I_A$, ending in $F_B$ $\overset{\text{def.}}{\iff}$ $w \in L(C)$. □ □

**Theorem.** For any $L \in Rec(\Sigma^*)$, $L^* \in Rec(\Sigma^*)$.

Idea

↑
not a proof!

A =



Add $\varepsilon$-transition from any final to any initial state.

▽ does not work ▽

**Example.** A = • (one state, no edges). What does the construction above do?

An automaton is normalized if $\#I = \#F = 1$, the initial state is not a target of any edge, $I \cap F = \emptyset$, and the final " " " " source " " " .

**Lemma.** For any automaton A, there exists a normalized automaton A' such that

$$L(A') = L(A) - \{\varepsilon\}.$$

**Proof of Theorem.** Let A be normalized recognizing $L - \{\varepsilon\}$. Define B by adding to A an $\varepsilon$-transition from the final state to the initial state, and making the initial state also final. Then $L(B) = L^*$ (exercise). □

**Proof of lemma.** Given $A = (Q, \Sigma, \delta, I, F)$, define $A' = (Q', \Sigma, \delta', \{i_0\}, \{f_0\})$

where $\quad Q' := Q \uplus \{i_0, f_0\}$

and $\quad \delta' := \delta \cup \{(i_0, a, q) \mid \text{there exists } i \in I \text{ such that } i \xrightarrow{a} q \text{ in } A\}$

$\cup \{(q, a, f_0) \mid \text{ " " } f \in F \text{ " " } q \xrightarrow{a} f \text{ " " }\}$

$\cup \{(i_0, a, f_0) \mid \text{there exist } i \in I, f \in F \text{ such that } i \xrightarrow{a} f \text{ in } A\}.$

Then, for $w \in \Sigma^+$, we have $w \in L(A) \iff w \in L(A')$:

If $w \in L(A)$, let $\pi$ be a successful run on $w$, from $i \in I$ to $f \in F$.

Since $w \neq \varepsilon$, $|\pi| = |w| + 1 \geq 2$. If $|w| = 1$, then, since $\pi$ is successful, $i \xrightarrow{w} f$ in $\delta$,

so that $(i_0, w, f)$ is in $\delta$.

Suppose $|w| \geq 2$. Write $w = a w' b$, for $a, b \in \Sigma$ and $w' \in \Sigma^*$.

Let $\pi'$ be defined by replacing the first node in $\pi$ with $i_0$

and the last node in $\pi$ with $f_0$. $\pi'$ is a $w$-path in $A'$.

Conversely, if $w \in L(A')$, let $\pi$ be a successful run in $A'$.
If $|w| = 1$ then there exist $i \in I, f \in F$ such that $i \xrightarrow{w} f$.
Otherwise, write $w = aub$ as before. Since $\pi$ is successful, it begins with
$i_0 \xrightarrow{a} q$ for some $q \in Q$, so we can pick $i \in Q$ such that $i \xrightarrow{w_0} q$.
Similarly, $\pi$ ends with $r \xrightarrow{b} f_0$ for some $r \in Q$, so we can pick $f \in Q$
such that $r \xrightarrow{b} f$. Replacing $i_0$ with $i$ and $f_0$ with $f$ yields a
successful run in $A'$. $\square$

# 4. Regular expressions

*Describing recognizable languages syntactically.*

**Definition.** A *regular expression* over alphabet $\Sigma$ is an expression of one of the forms:

- $\phi$ ,      (r.e.)

- $\varepsilon$ ,

- for any $a \in \Sigma$ :   $a$ ,

or, for any regular expressions $r_1, r_2$ :

- $r_1 \cdot r_2$ ,

- $r_1 + r_2$ ,

- $(r_1)^*$ .

Same definition, written in Backus-Naur Form :

$$e ::= \phi \mid \varepsilon \mid a \mid e \cdot e \mid e + e \mid e^*$$

The *language*, $\mathcal{L}(e)$, of a r.e. $e$ is defined inductively:

$\mathcal{L}(\phi) := \phi$ ,   $\mathcal{L}(\varepsilon) := \{\varepsilon\}$,   $\mathcal{L}(a) := \{a\}$,   $\mathcal{L}(r_1 \cdot r_2) := \mathcal{L}(r_1) \cdot \mathcal{L}(r_2)$,   $\mathcal{L}(r_1 + r_2) := \mathcal{L}(r_1) \cup \mathcal{L}(r_2)$,   $\mathcal{L}(r_1^*) := \mathcal{L}(r_1)^*$.

A language $L$ is *regular* if there exists a regular expression $e$ such that $L = \mathcal{L}(e)$.

(Some people call this "rational". Terminology is difficult.)

# Theorem (Kleene). A language $L \subseteq \Sigma^*$ is regular if, and only if, it is recognizable.

**Proof.** "$\Rightarrow$" By induction, for every regular expression $e$, we construct an automaton $A_e$ such that $\mathcal{L}(e) = \mathcal{L}(A_e)$.

For the base cases, we have automata $A_e$:

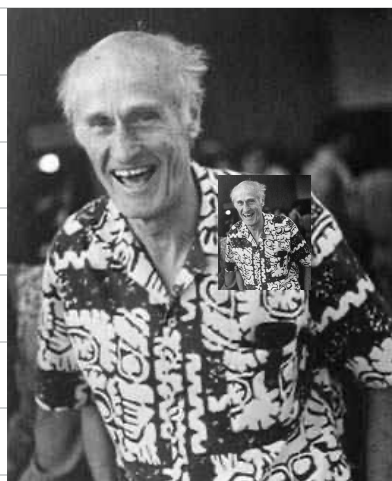| $e$ | $\emptyset$ | $\varepsilon$ | $a$ |
|---|---|---|---|
| $A_e$ | $\cdot$ | $\rightarrow \cdot \rightarrow$ | $\rightarrow \cdot \xrightarrow{a} \cdot \rightarrow$ |



S.C. Kleene
(1909-1994)

In each case, $\mathcal{L}(A_e) = \mathcal{L}(e)$.

For the inductive cases, we apply the closure properties that we proved before.

$e = r_1 \cdot r_2$: Given $A_{r_1}$ and $A_{r_2}$, construct $B$ such that $\mathcal{L}(B) = \mathcal{L}(A_{r_1}) \cdot \mathcal{L}(A_{r_2})$.

Define $A_e := B$. Then $\mathcal{L}(A_e) = \mathcal{L}(A_{r_1}) \cdot \mathcal{L}(A_{r_2}) \overset{IH}{=} \mathcal{L}(r_1) \cdot \mathcal{L}(r_2) = \mathcal{L}(e)$

The cases $e = r_1 + r_2$ and $e = r_1^*$ are similar, using that recognizable languages are closed under union and star.

**Proof** (continued). "$\Leftarrow$" Let $A$ be an automaton. We construct $r$ a regex such that $\mathcal{L}(r) = \mathcal{L}(A)$.

Let $n := \#Q$, and choose an (arbitrary) bijection $\{1, \dots, n\} \xrightarrow{\varphi} Q$. We write $q_i := \varphi(i)$, $1 \leq i \leq n$.

For each $0 \leq k \leq n$, $p, q \in Q$, define:

$$R^{(k)}_{p,q} := \{ w \in \Sigma^+ \mid \text{there exists a } w\text{-path } \pi \text{ from } p \text{ to } q$$

such that all states in $\pi$, except possibly the first and last,
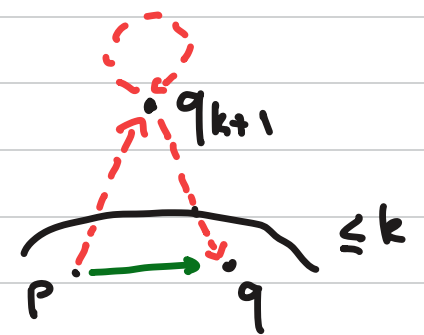
$$\text{belong to the set } \{q_1, \dots, q_k\} \}.$$

Let $p, q \in Q$. We construct regexes $r^{(k)}_{p,q}$ such that $\mathcal{L}(r^{(k)}_{p,q}) = R^{(k)}_{p,q}$, for each $0 \leq k \leq n$.

Induction on $k$: $\quad$ $k = 0$. The only possible paths have length 1.

$(0 \leq k \leq n)$

Let $r^{(k)}_{p,q} :=$ sum of $a \in \Sigma$ such that $(p, a, q) \in \delta$.

$k \Rightarrow k+1$. Note that $R^{(k+1)}_{p,q} = R^{(k)}_{p,q} + R^{(k)}_{p, q_{k+1}} \cdot \left( R^{(k)}_{q_{k+1}, q_{k+1}} \right)^* \cdot R^{(k)}_{q_{k+1}, q}$

So define $r^{(k+1)}_{p,q} := r^{(k)}_{p,q} + r^{(k)}_{p, q_{k+1}} \cdot \left( r^{(k)}_{q_{k+1}, q_{k+1}} \right)^* \cdot r^{(k)}_{q_{k+1}, q}$.

Finally, $r_{p,q} := \begin{cases} r^{(n)}_{p,q} & \text{if } p \neq q \\ r^{(n)}_{p,q} + \varepsilon & \text{if } p = q \end{cases}$, and $r :=$ sum of $r_{p,q}$ such that $p \in I$ and $q \in F$. $\quad \square$

The algorithm used for "⊆" in the proof is called McNaughton-Yamada. (+ Thompson, sometimes)
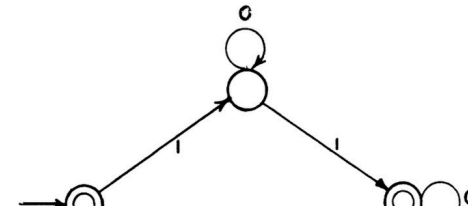
## Regular Expressions and State Graphs for Automata*

R. McNAUGHTON† AND H. YAMADA†

*Summary*—Algorithms are presented for 1) converting a state graph describing the behavior of an automaton to a regular expression describing the behavior of the same automaton (section 2), and 2) for converting a regular expression into a state graph (sections 3 and 4). These algorithms are justified by theorems, and examples are given. The first section contains a brief introduction to state graphs and the regular-expression language.

Other algorithms, sometimes with better outcomes, exist, notably,

Brzozowski - McCluskey, which you will see in the TD, and the

state elimination method, which relies on Arden's Lemma:

**Lemma.** Let $K, L \subseteq \Sigma^*$. The equation $X = K \cdot X + L$ has a smallest solution in $\mathcal{P}(\Sigma^*)$, namely, $K^*L$. If $\varepsilon \notin K$, the solution is unique.
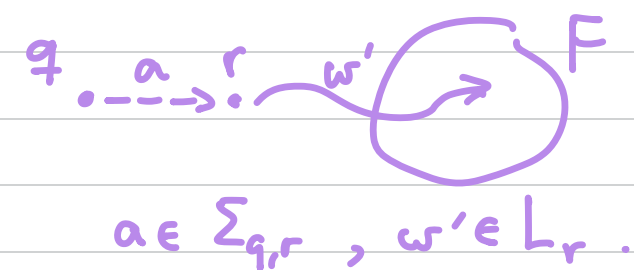
**Proof.** See TD. □

## Proof of "⇐" in Kleene's theorem, using state elimination method.

Let $A$ be an automaton. For every $q \in Q$, define $L_q := \{w \in \Sigma^* \mid \text{there exist } f \in F \text{ and } q \xrightarrow{w} f\}$.

Note that $L_q = \bigcup_{r \in Q} \Sigma_{q,r} \cdot L_r \cup 1_{q,F}$, where

$$\Sigma_{q,r} := \{a \in \Sigma \mid (q,a,r) \in \delta\} \quad \text{and} \quad 1_{q,F} := \begin{cases} \varepsilon & \text{if } q \in F \\ \emptyset & \text{if } q \notin F \end{cases}.$$

$q \dashrightarrow{a} r \xrightarrow{w'} F$
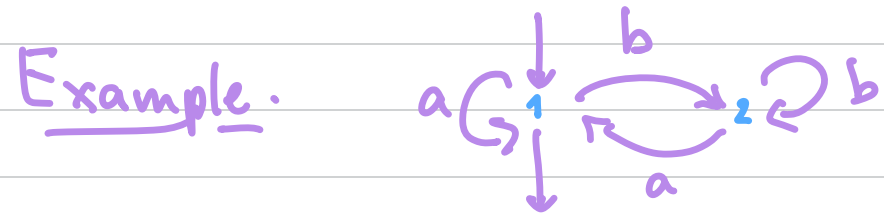
$a \in \Sigma_{q,r}, \ w' \in L_r.$

To obtain a regular expression, we successively solve the system of equations in $\#Q$ unknowns:

$$\begin{cases} R_1 = \sum_{i=1}^{n} \Sigma_{q_1,q_i} \cdot R_i + 1_{q_1,F} \\ \quad \vdots \\ R_n = \sum_{i=1}^{n} \Sigma_{q_n,q_i} \cdot R_i + 1_{q_n,F} \end{cases},$$

using Arden's Lemma, in the same way as Gaussian elimination for systems of linear equations.

For instance, starting from the last equation, it has the form $X = KX + L$, where

$$X = R_n, \quad K = \Sigma_{q_n,q_n}, \quad L = \sum_{i=1}^{n-1} \Sigma_{q_n,q_i} \cdot R_i + 1_{q_n,F}. \quad \text{So } R_n = K^* L.$$

**Example.**



$$\text{McNaughton - Yamada}: \quad r_{p,q}^{(k+1)} := r_{p,q}^{(k)} + r_{p,q_{k+1}}^{(k)} \cdot \left( r_{q_{k+1}, q_{k+1}}^{(k)} \right)^* \cdot r_{q_{k+1}, q}^{(k)}$$

$$r^{(0)} = \begin{pmatrix} a & b \\ a & b \end{pmatrix} \quad \rightsquigarrow \quad r^{(1)} = \begin{pmatrix} a + a \cdot a^* \cdot a & a + b \cdot a^* \cdot b \\ a + a \cdot a^* \cdot a & a + b \cdot a^* \cdot b \end{pmatrix} = \begin{pmatrix} a^+ & a^* b \\ a^+ & a^* b \end{pmatrix}$$

($r_{1,2}^{(0)}$ labeled)

$$\rightsquigarrow \quad r^{(2)} = \cdots = \begin{pmatrix} (a^* b)^* a^+ & (a^* b)^+ \\ (a^* b)^* a^+ & (a^* b)^+ \end{pmatrix}. \qquad r = r_{1,1}^{(2)} + \varepsilon = (a^* b)^* a^+ + \varepsilon.$$

$$a + b \cdot b^* \cdot a$$
$$= a + b^+ a$$
$$= (\varepsilon + b^+) a$$

**Elimination:**
$$\begin{cases} r_1 = a \cdot r_1 + b \cdot r_2 + \varepsilon \\ r_2 = a \cdot r_1 + b \cdot r_2 \end{cases} \quad \rightsquigarrow \quad \begin{cases} r_1 = a \cdot r_1 + b \cdot b^* \cdot a \cdot r_1 + \varepsilon = b^* a \cdot r_1 + \varepsilon \\ r_2 = b^* \cdot a \cdot r_1 \end{cases}$$

$$\rightsquigarrow \quad \begin{cases} r_1 = (b^* a)^* \\ r_2 = b^* \cdot a \cdot (b^* a)^* = (b^* a)^+ \end{cases}$$

**Another way to write this:** $(a+b)^* a + \varepsilon.$

A    Kleene algebra is a tuple $(K, 0, 1, +, \cdot, (\ )^*)$ where:

- $(K, 0, 1, +, \cdot)$ is a unital semiring, and

  (+ and $\cdot$ are associative, + is commutative, $\cdot$ distributes over +,

  0 is neutral for + and 1 is neutral for $\cdot$.)

for all $a, b \in K$:
- $1 + a \cdot a^* = a^* = 1 + a^* a$

John H. Conway (1937-2020)

- If $b + a \cdot c \leq c$ then $a^* b \leq c$, and if $b + ca \leq c$ then $ba^* \leq c$

  where "$x \leq y$" means: "$x + y = y$".

Example. $(\text{Rec}(\Sigma^*), \emptyset, \{\varepsilon\}, \cup, \cdot, (\ )^*)$ is a Kleene algebra.

Theorem.    Let $r, s$ be regular expressions. Then $\mathcal{L}(r) = \mathcal{L}(s)$ if, and only if,    Dexter Kozen

(Conway, Kozen) for every Kleene algebra $K$ and $(k_a)_{a \in \Sigma} \in K^\Sigma$, $r[a \mapsto k_a] = s[a \mapsto k_a]$ in $K$.

"Kleene algebras are a sound and complete axiomatization of regular languages."

Proof. Omitted. See: D. Kozen, "A completeness theorem for Kleene algebras..." (1994).