

# **Formální jazyky I**

Chomského hierarchie formálních jazyků. Regulární jazyky, jejich reprezentace a převody mezi nimi. Varianty konečných automatů. Nedeterminismus a determinizace automatů. Uzávěrové vlastnosti regulárních jazyků.

*IB102/IB005*

## Uvod

### Abeceda a jazyk

Abecedou rozumieme ľubovolnu konečnu množinu  $\Sigma$ . Jej prvky nazývame znaky (písmena alebo symboly).

*Priklad:*  $\{a,b\}$ ,  $\{0,1,\dots,9\}$ ,  $\emptyset$

Slovo (retazec) nad abecedou  $\Sigma$  je ľubovolna konečna postupnosť znakov z  $\Sigma$ . Dĺžka slova  $v$  je počet znakov v slove (znacime  $\#(v)$ ), počet vyskytu znaku  $a$  v slove znacime  $\#_a(v)$ . Špeciálny prípad je prázdne slovo s nulovou dĺžkou, znacime ho  $\varepsilon$ .

*Priklad:* pre  $\Sigma = \{a,b\}$  je slovo  $aabb$ ,  $aa$ ,  $b$ . Pre  $v = aabb$ ,  $\#(v) = 4$  a  $\#_a(v) = 2$

Množina všetkých slov nad  $\Sigma$  znacime  $\Sigma^*$ , množinu všetkých neprázdnych slov znacime  $\Sigma^+$ .

*Priklad:*  $\{a\}^* = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$ ,  $\{a\}^+ = \{a, aa, aaa, aaaa, \dots\}$

*Špeciálne:*  $\emptyset^* = \{\varepsilon\}$  a  $\emptyset^+ = \emptyset$ .

Jazyk nad abecedou  $\Sigma$  je ľubovolna množina slov nad  $\Sigma$  (i.e. podmnožiny  $\Sigma^*$ ). Možu byť konečné aj nekonečné.

*Priklad:*  $\Sigma = \{0, 1\}$ , potom  $L$  nad  $\Sigma$  je  $\{10, 1, 011101\}$ .  $L$  nad  $\{a,b\}$  definovaný ako  $\{w \in \{a,b\}^* \mid \#_a(w) = \#_b(w)\}$  je zas nekonečný.

*Špeciálne:*  $\emptyset$  je  $L$  nad ľubovolnou  $\Sigma$ .

### Operácie nad jazykmi

Kedže jazyku sú len množiny, môžeme aplikovať množinové operácie: zjednotenie, prienik a rozdiel.

*Priklad:* Ak  $L$  je nad  $\Sigma$ ,  $K$  je nad  $\Delta$ , tak  $L \cup K$  je nad  $\Sigma \cup \Delta$ .

Dalej definujeme (pre jazyky  $L$  je nad  $\Sigma$ ,  $K$  je nad  $\Delta$ )

- Zretazenie:  $K.L = \{uv \mid u \in K, v \in L\}$  nad  $\Sigma \cup \Delta$

- $i$ -ta mocnina jazyka  $L$ :  $L^0 = \{\varepsilon\}$  a  $L^{i+1} = L.L^i$

- Iterácia jazyka  $L$ :  $L^* = \bigcup_{i=0}^{\infty} L^i$ .

- Doplnok:  $\text{co-}L = \Sigma^* \setminus L$ .

- Substitúcie, homomorfizmus (prip. inverzni), zrkadlový obraz...

*Špeciálne:*  $\emptyset.L = L.\emptyset = \emptyset$  a  $\{\varepsilon\}.L = L.\{\varepsilon\} = L$ ,  $\emptyset^0 = \{\varepsilon\}$ ,  $\emptyset^i = \emptyset$  ( $i \in \mathbb{N}$ ),  $\{\varepsilon\}^j = \{\varepsilon\}$  ( $j \in \mathbb{N}_0$ ).

**Definice 1.2.** Gramatika  $G$  je čtveřice  $(N, \Sigma, P, S)$ , kde

- $N$  je neprázdná konečná množina *neterminálních symbolů* (stručněji: *neterminálů*).
- $\Sigma$  je konečná množina *terminálních symbolů* (*terminálů*) taková, že  $N \cap \Sigma = \emptyset$ . Sjednocením  $N$  a  $\Sigma$  obdržíme množinu všech symbolů gramatiky, kterou obvykle označujeme symbolem  $V$ .
- $P \subseteq V^*NV^* \times V^*$  je konečná množina *pravidel*. Pravidlo  $(\alpha, \beta)$  obvykle zapisujeme ve tvaru  $\alpha \rightarrow \beta$  (a čteme jako " $\alpha$  přepiš na  $\beta$ ").
- $S \in N$  je speciální *počáteční neterminál* (nazývaný také *kořen gramatiky*).

Požadavky na pravidla  $(\alpha, \beta)$  :  $\alpha$  musí obsahovat aspon jeden neterminál.  $\beta$  může být i prázdná ( $\epsilon$ ).

Vetna forma gramatiky  $G$ : prvky množiny  $(N \cup \Sigma)^*$ , které lze odvodit z počátečního neterminálu za pomoci pravidel gramatiky.

Formálně:  $\alpha \in (N \cup \Sigma)^*$  je vetna forma  $\iff S \Rightarrow^* \alpha$ .

Vetna forma bez neterminálu je veta. Množina všech vet gramatiky je jazyk generovaný gramatikou:  $L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$

## Chomského hierarchie formálných jazykú.

- 4 skupiny na zaklade omedzenia na tvar pravidiel pre gramatiku

**typ 0** Libovolná gramatika je gramatikou typu 0; na tvar pravidiel se nekladou žádné omezující požadavky. Někdy též se takové gramatiky označují jako gramatiky bez omezení či frázové gramatiky (phrase grammars).

**typ 1** Gramatika je typu 1 (nebo též *kontextová*<sup>1</sup>, Context-Sensitive, CSG, méně často též *monotónní*), jestliže pro každé její pravidlo  $\alpha \rightarrow \beta$  platí  $|\alpha| \leq |\beta|$  s eventuelní výjimkou pravidla  $S \rightarrow \varepsilon$ , pokud se  $S$  nevyskytuje na pravé straně žádného pravidla.

**typ 2** Gramatika je typu 2 (též *bezkontextová*, Context-Free, CFG), jestliže každé její pravidlo je tvaru  $A \rightarrow \alpha$ ,  $|\alpha| \geq 0$ .

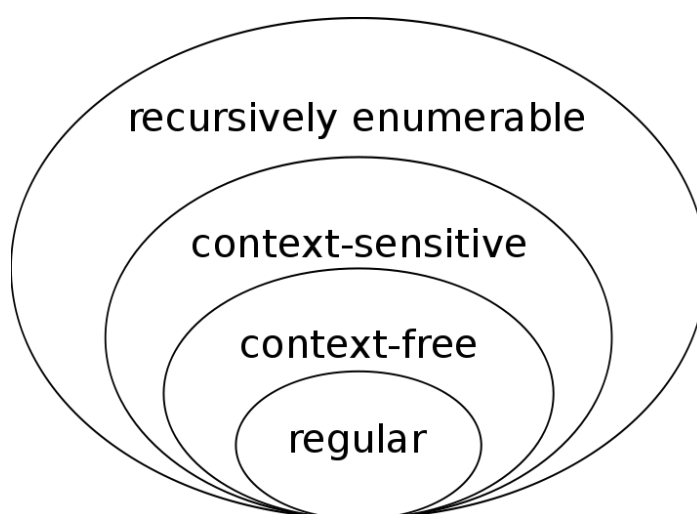
**typ 3** Gramatika je typu 3 (též *regulární* či *pravolineární*<sup>2</sup>), jestliže každé její pravidlo je tvaru  $A \rightarrow aB$  nebo  $A \rightarrow a$  s eventuelní výjimkou pravidla  $S \rightarrow \varepsilon$ , pokud se  $S$  nevyskytuje na pravé straně žádného pravidla.<sup>3</sup>

Mapovanie gramatika : jazyk:

Jazyk je:

- regularny, ak ho generuje gramatika typu 3
- bezkontextovy, ak ho generuje gramatika typu 2
- kontextovy, ak ho generuje gramatika typu 1
- typu 0, ak ho generuje gramatika typu 0

To znamena: Hierarchie gramatik urcuje prislusnu hierarchiu jazykov.



Pozor: obrazok ukazuje vlastnost: ak je jazyk regularny, potom je aj context-free.

Nehovori nic o velkosti jazyka (napr  $\Sigma^*$  je regularny,  $a^n b^n$  je "mensi", ale je context-free)

Prehľad všetkých rozhodnuteľných a nerozhodnuteľných problémov týkajúcich sa tried jazykov Chomského hierarchie:

	<i>R</i>	<i>DCF</i>	<i>CF</i>	<i>CS</i>	<i>Rec</i>	<i>RE</i>
Je $L(\mathcal{G})$ prázdny? konečný?	R	R	R	N	N	N
Je $L(\mathcal{G}) = \Sigma^*$ ?	R	R	N	N	N	N
Je $L(\mathcal{G}) = R$ ? ( $R$ je regulárna množina)	R	R	N	N	N	N
Je $L(\mathcal{G}_1) = L(\mathcal{G}_2)$ ?	R	R	N	N	N	N
Je $L(\mathcal{G}_1) \supseteq L(\mathcal{G}_2)$ ?	R	N	N	N	N	N
Je $L(\mathcal{G})$ regulárny jazyk?	ano	R	N	N	N	N
Je průnik dvou jazyků jazyk téhož typu?	ano	N	N	ano	ano	ano
Je sjednocení dvou jazyků jazyk téhož typu?	ano	N	ano	ano	ano	ano
Je komplement jazyka jazyk téhož typu?	ano	ano	N	ano	ano	N
Je zřetězení dvou jazyků jazyk téhož typu?	ano	N	ano	ano	ano	ano
Je gramatika $\mathcal{G}$ víceznačná?	R	N	N	N	N	N

## Regulární jazyky, jejich reprezentace a převody mezi nimi.

Jazyk  $L$  nazveme regulárním, ak je rozpoznatelný nějakým konečným automatem (vid. konečné automaty níže). Případně, ak je generován gramatikou typu 3 (vid. definice Chomského hierarchie).

*Pumping lemma:*

**Lemma 2.13** (o vkládání). *Nechť  $L$  je regulární jazyk. Pak existuje  $n \in \mathbb{N}$  takové, že libovolné slovo  $w \in L$ , jehož délka je alespoň  $n$ , lze psát ve tvaru  $w = xyz$ , kde  $|xy| \leq n$ ,  $y \neq \varepsilon$  a  $xy^iz \in L$  pro každé  $i \in \mathbb{N}_0$ . (Číslo  $n$  se neformálně nazývá pumpovací konstanta.)*

Pozor: Pumping lemma je v tvare implikace v tvare  $L \text{ je regulární} \Rightarrow Q$  (zvyšné tvrzenie), a teda platí, že existujú neregulárne jazyky, ktoré splňajú pumping lemma.

Pri dokazovaní, že  $L$  nie je regulárny používame kontrapozitívnu formu:  $\neg Q \Rightarrow L$  není regulární. Kucharka: dokaz sporom: Predpokladáme, že  $L$  je regulárny, ukážeme, že:

- pro lib.  $n \in \mathbb{N}$  (pumpovací konstanta)
- existuje slovo  $w \in L$  s délkou aspoň  $n$
- při lib. rozložení slova  $w$  na 3 části:  $x, y, z$ , tak, že  $|xy| \leq n$  a  $y \neq \varepsilon$ ,
- existuje  $i \in \mathbb{N}_0$ :  $xy^iz$  nepatří do  $L$ .

*Příklad:*

**Příklad 2.14.** *Ukážeme, že  $L = \{a^p \mid p \text{ je prvočíslo}\}$  nad abecedou  $\{a\}$  není regulární.*

*Důkaz.* Pro dosažení sporu předpokládejme, že  $L$  je regulární. Buď  $n \in \mathbb{N}$  libovolné (pumpovací konstanta z PL). Jelikož prvočísel je nekonečně mnoho, existuje prvočíslo  $p$ , které je větší nebo rovno  $n$ ; zvolme  $w = a^p$  patřící do  $L$ . Při jakémkoli rozdělení  $w$  na podslova  $x, y, z$  musí být  $y = a^k$ ,  $k \geq 1$ . Napumpujeme-li  $y$   $p + 1$ -krát, dostaneme:  $xy^{p+1}z = xyyp^pz = xyzyp^p = a^p a^{kp} = a^{p(k+1)}$ , což je jistě slovo, které nepatří do jazyka  $L$ , protože  $p(k+1)$  není prvočíslo – dostáváme tedy spor s naším předpokladem, že  $L$  je regulární. Podle PL tedy  $L$  regulární není.  $\square$

*Myhillova-Nerodova věta* (nebralo sa na lahsich automatoch):

**Věta 2.28** (Myhillova-Nerodova). *Nechť  $L$  je jazyk nad  $\Sigma$ , pak tato tvrzení jsou ekvivalentní:*

1.  $L$  je rozpoznatelný konečným automatem.
2.  $L$  je sjednocením některých tříd rozkladu určeného pravou kongruencí na  $\Sigma^*$  s konečným indexem.
3. Relace  $\sim_L$  má konečný index.

## Variety konečných automatů (minimalni, nedeterministicky, ε-kroky).

**Definice 2.1.** Konečný automat (Finite Automaton, FA)  $\mathcal{M}$  je pětice  $(Q, \Sigma, \delta, q_0, F)$ , kde

- $Q$  je neprázdná konečná množina stavů.
- $\Sigma$  je konečná množina vstupních symbolů, nazývaná také vstupní abeceda.
- $\delta : Q \times \Sigma \rightarrow Q$  je parciální přechodová funkce.
- $q_0 \in Q$  je počáteční stav.
- $F \subseteq Q$  je množina koncových stavů.

Abychom mohli definovat jazyk přijímaný daným FA  $\mathcal{M}$ , zavedeme rozšířenou přechodovou funkci  $\hat{\delta} : Q \times \Sigma^* \rightarrow Q$ , definovanou induktivně vzhledem k délce slova ze  $\Sigma^*$ :

- $\hat{\delta}(q, \varepsilon) = q$  pro každý stav  $q \in Q$ .
- $\hat{\delta}(q, wa) = \begin{cases} \delta(\hat{\delta}(q, w), a) & \text{je-li } \hat{\delta}(q, w) \text{ i } \delta(\hat{\delta}(q, w), a) \text{ definováno,} \\ \perp & \text{jinak.} \end{cases}$

Rozšířená přechodová funkce induktivně definuje přechod nad celým slovem za pomoci přechodu nad jedním znakem. Potom ( $M$  je konečný automat):

$$L(M) = \{ w \in \Sigma^* \mid \hat{\delta}(q_0, w) \in F \}$$

*Příklad:*

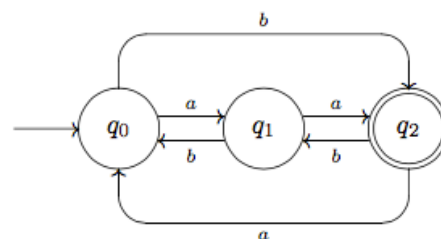
**Příklad 2.3.** Necht'  $\mathcal{M} = (\{q_0, q_1, q_2\}, \{a, b\}, \delta, q_0, \{q_2\})$  je FA, kde

$$\delta(q_0, a) = q_1 \quad \delta(q_0, b) = q_2$$

$$\delta(q_1, a) = q_2 \quad \delta(q_1, b) = q_0$$

$$\delta(q_2, a) = q_0 \quad \delta(q_2, b) = q_1$$

Pak  $L(\mathcal{M}) = \{w \in \{a, b\}^* \mid (\#_a(w) - \#_b(w)) \bmod 3 = 2\}$ .



### Minimalni konečný automat

(k libovolnému konečnému automatu  $A$ ) je automat s nejmenším počtem stavů a totální přechodovou funkcí, který rozpoznává regulární jazyk  $L(A)$ .

Tvrzení o existenci minimálního konečného automatu plyne z Myhillovy-Nerodovy věty. Převod lze uskutečnit *algoritmicky*. Důkaz korektnosti tohoto postupu je pomocí kongruencí nad stavy. Hledáme jazykově ekvivalentní stavy: Stavy  $p, q$  nazveme jazykově ekvivalentní ( $p \equiv q$ ), ak  $L(p) = L(q)$ .

Postup: odstraníme nedosazitelné stavy. Potom “spojíme” stavy do množin podle relace  $\equiv$ . Na začátku je každý stav sam. Stavy  $p, q$  jsou v relaci  $\equiv$  ak jsou oba koncové (akceptují tedy rovnaké slova  $= \varepsilon$ ). Postupně hledáme  $\equiv$  až kým nenajdeme fixpoint.



Příklad:

	$\mathcal{M}$	$a$	$b$
→	1	2	—
	2	3	4
←	3	6	5
	4	3	2
←	5	6	3
←	6	2	—
	7	6	1

	$\mathcal{M}'$	$a$	$b$
→	1	2	$N$
	2	3	4
←	3	6	5
	4	3	2
←	5	6	3
←	6	2	$N$
	$N$	$N$	$N$

(odstranenie nedosiahnutelných)

	$\equiv_0$	$a$	$b$
$I$	1	$I$	$I$
	2	$II$	$I$
	4	$II$	$I$
	$N$	$I$	$I$
$II$	3	$II$	$II$
	5	$II$	$II$
	6	$I$	$I$

	$\equiv_1$	$a$	$b$
$I$	1	$II$	$I$
	$N$	$I$	$I$
$II$	2	$III$	$II$
	4	$III$	$II$
$III$	3	$IV$	$III$
	5	$IV$	$III$
$IV$	6	$II$	$I$

	$\equiv_2$	$a$	$b$
$I$	1	$III$	$II$
$II$	$N$	$II$	$II$
$III$	2	$IV$	$III$
	4	$IV$	$III$
$IV$	3	$V$	$IV$
	5	$V$	$IV$
$V$	6	$III$	$II$

(rozklad podľa  $\equiv_i$ )

Relace  $\equiv$  je tedy v tomto případě rovna relaci  $\equiv_2$ . Minimální automat pro jazyk  $L(\mathcal{M})$  vypadá takto:

	$\mathcal{M}/\equiv$	$a$	$b$
→	$I$	$III$	$II$
	$II$	$II$	$II$
	$III$	$IV$	$III$
←	$IV$	$V$	$IV$
←	$V$	$III$	$II$

## Nedeterministické konečné automaty

$\delta : Q \times \Sigma \rightarrow Q$  zmeníme na  $\delta : Q \times \Sigma \rightarrow 2^Q$ . A teda z jedného stavu môžeme pod jedným symbolom prejsť do viacerých stavov, a teda do určitej podmnožiny  $Q$ .

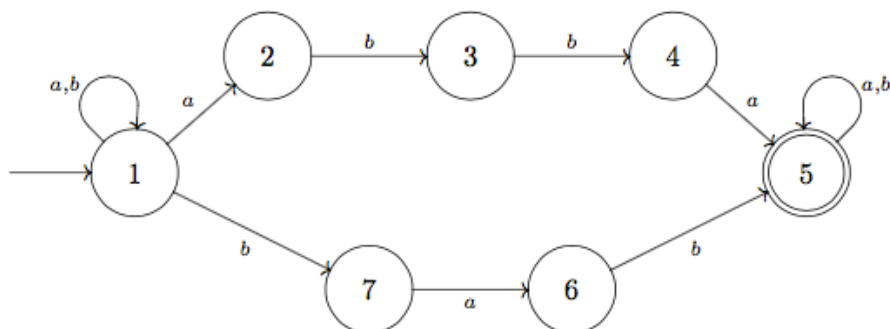
Podobne ako v prípade deterministických automatů zavedeme *rozšířenou* přechodovou funkci  $\hat{\delta} : Q \times \Sigma^* \rightarrow 2^Q$ :

- $\hat{\delta}(q, \varepsilon) = \{q\}$
- $\hat{\delta}(q, wa) = \bigcup_{p \in \hat{\delta}(q, w)} \delta(p, a)$



A potom plati (M je nedeterministický automat):  $L(M) = \{ w \in \Sigma^* \mid \delta^*(q_0, w) \cap F \neq \emptyset \}$   
 Nedeterminizmus nám pomáha konstruovať “krajšie” automaty:

$L = \{ w \in \{a, b\}^* \mid w \text{ obsahuje podslovo abba alebo bab} \}$



Plati ale: Pre každý NFA  $M = (Q, \Sigma, \delta, q_0, F)$  existuje ekvivalentný DFA. (sekcia Nedeterminizmus a determinizácia automatu).

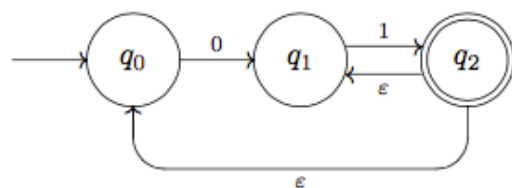
### Automaty s $\epsilon$ -kroky

Automat takto môže svoj stav zmeniť aj bez precitania vstupného symbolu (a teda úplne samovoľne). Ako značenie používame prechod pod  $\epsilon$ .

Prechodovú funkciu teda redefinujeme:

$$\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q.$$

Pozor: rozšírená prechodová funkcia už potom nie je tak intuitívna. V každom stave musíme zahrnúť ešte aj stavy, do ktorých sa dostaneme pomocou  $\epsilon$  prechodu (predtým než načítame ďalší symbol).

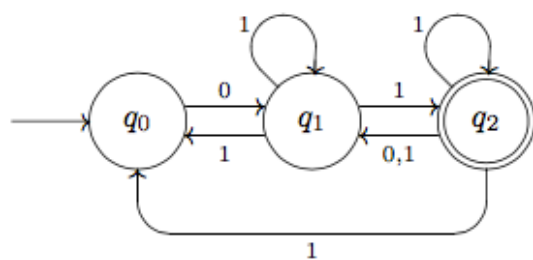


Plati: ku každému automatu M s  $\epsilon$ -krokmi existuje ekvivalentný nedeterministický automat M.

Hľadáme cesty tvaru:

$$p_1 \xrightarrow{\epsilon} \dots \xrightarrow{\epsilon} p_m \xrightarrow{a} q_1 \xrightarrow{\epsilon} \dots \xrightarrow{\epsilon} q_n$$

, pre ktoré pridáme pravidlo  $p_1 \xrightarrow{a} q_n$ .



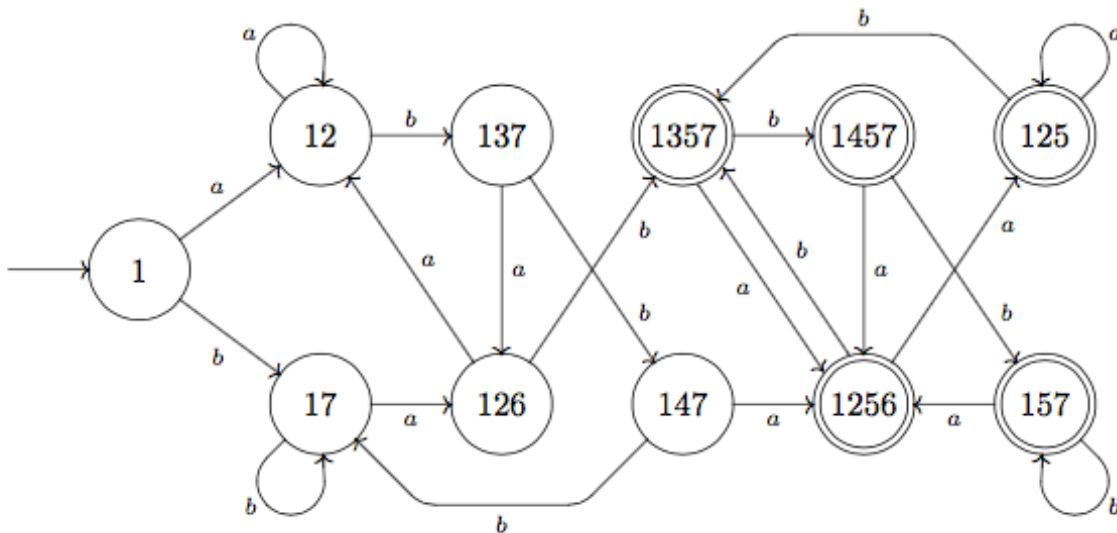
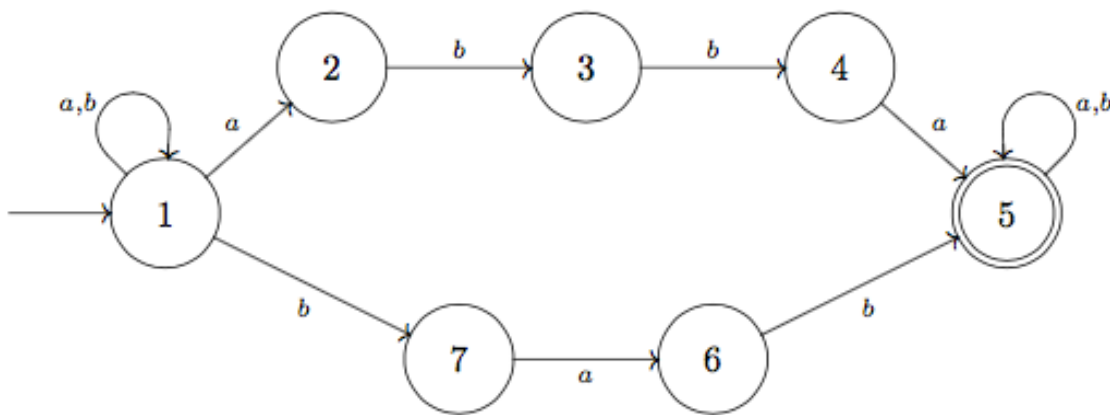
## Nedeterminismus a determinizace automatů.

Pro každý NFA  $M = (Q, \Sigma, \delta, q_0, F)$  existuje ekvivalentní DFA.

*Důkaz.* Nechť  $\mathcal{M}' = (Q', \Sigma, \delta', \{q_0\}, F')$  je deterministický konečný automat, kde:

- $Q' = 2^Q$ , tj. stavy automatu  $\mathcal{M}'$  jsou všechny podmnožiny  $Q$ .
- Přejchodová funkce  $\delta'$  je definována předpisem  $\delta'(P, a) = \bigcup_{q \in P} \delta(q, a)$
- Množina koncových stavů  $F'$  je tvořena právě těmi podmnožinami  $Q$ , které obsahují alespoň jeden prvek množiny  $F$ .

Zřejmě  $\mathcal{M}'$  je deterministický konečný automat (dokonce s totální přechodovou funkcí).



To, že akceptují ten istý jazyk, se dá ukázat vzhledem k délce odvozeného slova  $w$  (indukcí), ale je asi nad rámec statnic (aj časovo). Pozor na nárůst stavů (exponenciální). Automat se dá následně ještě minimalizovat.

## Uzávěrové vlastnosti regulárních jazyků.

### 1. Třída regulárních jazyků je uzavřena na zjednotění, průnik a rozdíl.

$L_1 \cup L_2$ : Majme  $A_1, A_2$ , platí  $L_1 \cup L_2 = A_1 \cup A_2$ , který je definovaný ako buď paralelná kompozícia, kde akceptujeme práve vtedy keď aspoň jeden automat dosiahne akceptujúci stav, alebo urobíme kartezský súčin stavov, prechodová relácia je dvojica, a akceptujeme stavmi, kde aspoň jeden z pôvodných bol akceptujúci.

$L_1 \cap L_2$ : Majme  $A_1, A_2$ , platí  $L_1 \cap L_2 = A_1 \cap A_2$ , podobne ako zjednotenie, ale musia akceptovať oba.

$L_1 - L_2$ : Majme  $A_1, A_2$ , platí  $L_1 - L_2 = A_1 - A_2$ , podobne, ale akceptujeme len ak prvý akceptoval a druhý zamietol, pri súčine sú akceptujúce stavy dvojica  $(p, q)$ ,  $p$  patrí  $A_1$  a je akceptujúci (patrí do  $F_1$ ) a  $q$  patrí  $A_2$  a je neakceptujúci (nepatrí do  $F_2$ ).

### 2. Třída regulárních jazyků je uzavřena na komplement.

Plati:  $\text{co-}L = \Sigma^* - L$

Postup: vymeníme koncové a nekonečné stavy:  $\text{co-}A = (Q, \Sigma, \delta, q_0, Q - F)$  k pôvodnému  $A = (Q, \Sigma, \delta, q_0, F)$ .

### 3. Třída regulárních jazyků je uzavřena na zretazenie.

Pridáme  $\epsilon$ -krok pre koncové stavy z  $A_1$  do počiatočných z  $A_2$ .

### 4. Třída regulárních jazyků je uzavřena na iteráciu.

Pridáme nový stav  $q$ , ktorý bude po novom jediný konečný. Zavedieme  $\epsilon$ -krok pre koncové stavy z  $A_1$  do  $q$ , a taktiež  $\epsilon$ -krok z  $q$  do počiatočného stavu.

### 5. Třída regulárních jazyků je uzavřena na substitúciu, homomorfizmus... (uz prilis specificke, no time for that)

### Příklady:

5.1 Rozhodněte, zda platí: jsou-li jazyky  $L_1, L_2, L_3, \dots$  regulární, pak i jazyk

$$\bigcup_{i=1}^{\infty} L_i$$

je regulární jazyk.

Neplatí. Jazyky  $L_i = \{a^i b^i\}$  pro každé  $i > 0$  jsou konečné a tudíž regulární, ale  $\bigcup_{i=1}^{\infty} L_i = \{a^n b^n \mid n > 0\}$  není regulární.

- a)  $L_1$  je regulární,  $L_2$  je neregulární  $\Rightarrow L_1 \cap L_2$  je neregulární
- b)  $L_1$  je regulární,  $L_2$  je neregulární  $\Rightarrow L_1 \cap L_2$  je regulární
- c)  $L_1$  je regulární,  $L_2$  je neregulární  $\Rightarrow L_1 \setminus L_2$  je neregulární
- d)  $L_1$  je regulární,  $L_2$  je neregulární  $\Rightarrow L_1 \setminus L_2$  je regulární
- e)  $L_1$  je regulární,  $L_2$  je neregulární  $\Rightarrow L_2 \setminus L_1$  je neregulární
- f)  $L_1$  je regulární,  $L_2$  je neregulární  $\Rightarrow L_2 \setminus L_1$  je regulární

- a)  $\{ab\}$  je regulární jazyk,  $a^n b^n$  nie je. Ich prienik je  $\{a,b\}$
- b)  $\Sigma^*$  je regulární jazyk,  $a^n b^n$  nie je. Ich prienik je  $a^n b^n$
- c)  $\{ab\}$  je regulární jazyk,  $a^n b^n$  nie je. výsledok je  $\emptyset$
- d)  $\Sigma^*$  je regulární jazyk,  $a^n b^n$  nie je. výsledok je  $\text{co-}a^n b^n$
- e) ....