

# Pattern Matching

---

## Definitionen:

1. *Alphabet*: endliche Menge von Zeichen, die *Buchstaben* genannt werden.
2. *Wörter*: Sei  $\Sigma$  Alphabet. Ein  $\Sigma$ -*Wort* ist endliche Folge von Buchstaben aus  $\Sigma$ :

$$w = b_1 b_2 \cdots b_n \quad \text{mit } b_i \in \Sigma$$

ist  $\Sigma$ -Wort der Länge  $n$

$\Sigma^*$ : Menge der  $\Sigma$ -Wörter

3. Das *leere Wort* bezeichnen wir mit  $\varepsilon$ .
4. *Konkatenation* von Wörtern:  
Seien  $v = b_1 b_2 \cdots b_m \in \Sigma^*$  und  $w = c_1 c_2 \cdots c_n \in \Sigma^*$ .  
Wir definieren die Konkatenation  $vw$  von  $v$  und  $w$  als

$$vw := b_1 b_2 \cdots b_m c_1 c_2 \cdots c_n$$

5. *Sprache*:  
Eine beliebige Teilmenge  $\mathcal{L} \subseteq \Sigma^*$  heißt  $\Sigma$ -*Sprache*.

**Beispiel:** Sei  $\Sigma = \{a, b\}$ . Dann gilt

1.  $\Sigma^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$
2.  $\mathcal{L}_Q := \{w \in \Sigma^* \mid \exists v \in \Sigma^* : w = vv\}$   
ist eine  $\Sigma$ -Sprache

# Syntax–Analyse (Chomsky–Hierarchie)

## Aufgaben der *Syntax–Analyse* in der Informatik

1. Beschreibung von  $\Sigma$ –Sprachen
2. Bereitstellung von Algorithmen zum Erkennen von  $\Sigma$ –Wörtern

## Mögliche Ansätze

1. Beschreibung durch *reguläre Ausdrücke*  
Erkennung über *endliche Automaten*
2. Beschreibung durch *kontext–freie* Grammatik  
Erkennung durch *Keller–Automaten*
3. Beschreibung durch *kontext–sensitive* Grammatik  
algorithmisch entscheidbar, exponentieller Aufwand
4. Beschreibung durch beliebige Grammatik  
unentscheidbar

## Praktische Bedeutung in Informatik:

1. Reguläre Ausdrücke
  - (a) Scanner: Aufteilung eines Programms in *Identifier*, *Schlüsselwörter*, *Literale*, *Kommentare*, etc.
  - (b) Skript–Sprachen: *Tcl*, *Perl*, *Phyton*, ...
2. Kontext–freie Grammatik  
Parser: Strukturierung eines Programms in *Ausdrücke*, *Statements*, *Funktionen*, ...

# Reguläre Ausdrücke

**Gegeben:** Alphabet  $\Sigma$

**Induktive Definition** der

(a) Menge  $\text{RegExp}$  der *regulären Ausdrücke* und der

(b) Sprache  $\mathcal{L}(r)$  für  $r \in \text{RegExp}$

1.  $\emptyset \in \text{RegExp}$ :

$$\mathcal{L}(\emptyset) = \{\}$$

2.  $\varepsilon \in \text{RegExp}$ :

$$\mathcal{L}(\varepsilon) = \{\varepsilon\}$$

3.  $b \in \text{RegExp}$  für alle  $b \in \Sigma$

$$\mathcal{L}(b) = \{b\}$$

4. Konkatenation:  $r_1, r_2 \in \text{RegExp} \Rightarrow r_1 r_2 \in \text{RegExp}$

$$\mathcal{L}(r_1 r_2) = \{vw \in \Sigma^* \mid v \in \mathcal{L}(r_1) \wedge w \in \mathcal{L}(r_2)\}$$

5. Alternative:  $r_1, r_2 \in \text{RegExp} \Rightarrow (r_1 | r_2) \in \text{RegExp}$

$$\mathcal{L}((r_1 | r_2)) = \mathcal{L}(r_1) \cup \mathcal{L}(r_2)$$

6. Abschluß:  $r \in \text{RegExp} \Rightarrow (r)^* \in \text{RegExp}$

$$\mathcal{L}((r)^*) = \{\varepsilon\} \cup \{w_1 \cdots w_n \mid w_i \in \mathcal{L}(r)\}$$

# Vereinfachung von regulären Ausdrücken

## Bindungs-Regeln

1. Abschluß bindet stärker als Konkatination
2. Konkatination bindet stärker als Alternative

Dann können Klammern weggelassen werden:

1.  $ab^*$  steht für  $a(b)^*$
2.  $a|b|c$  steht für  $((a|b)|c)$

## Alternative Schreibweise:

$r_1 + r_2$  statt  $r_1|r_2$

**Definition:**  $r \simeq s \Leftrightarrow \mathcal{L}(r) = \mathcal{L}(s)$

## Eigenschaften regulärer Ausdrücke

1.  $\emptyset + r \simeq r, \quad \varepsilon r \simeq r \varepsilon = r, \quad \emptyset^* = \varepsilon, \quad \emptyset r = r \emptyset = \emptyset$
2.  $r + s \simeq s + r, \quad \varepsilon^* = \varepsilon$
3.  $r + r \simeq r, \quad (r^*)^* \simeq r^*$
4.  $(r + s) + t \simeq r + (s + t), \quad (rs)t \simeq r(st)$
5.  $(r + s)t \simeq rt + st, \quad t(r + s) \simeq tr + ts$
6.  $(\varepsilon + r)^* \simeq r^*, \quad rr^* = r^*r$
7.  $(r + s)^* \simeq (r^*s^*)^*$
8.  $(rs)^* \simeq \varepsilon + rs(rs)^*, \quad (rs)^*r \simeq r(sr)^*$

# Reguläre Ausdrücke: Erweiterungen & Beispiele

Abkürzungen: Sei  $\Sigma = \{a, b, c, \dots, x, y, z, 0, \dots, 9\}$

$[abc]$	$\hat{=}$	$a b c$
$[0-9]$	$\hat{=}$	$0 1 2 3 4 5 6 7 8 9$
$[\wedge 0-9]$	$\hat{=}$	$a b c \dots x y z$
$.$	$\hat{=}$	$a b c \dots x y z 0 \dots 9$
$r^+$	$\hat{=}$	$rr^*$
$r^?$	$\hat{=}$	$\varepsilon r$

## Beispiele

1.  $[0-9]^+$

nicht-leere Folge von Ziffern, z. B. "01"

2.  $0|[1-9][0-9]^*$

Zahl im Dezimal-System

3.  $\text{http://}[\wedge/]^+/\wedge[\_ ]^+$

(a) Wörtlich: "http://"

(b) nicht-leere Folge von Buchstaben, die keinen Slash  
"/" enthält

(c) Wörtlich: Slash "/"

(d) nicht-leere Folge von Buchstaben, die kein Blank  
"\_" enthält

# Endliche Automaten (Finite State Machines)

**Definition:** Das 5–Tupel

$$\langle \Sigma, Z, A, s_0, \text{next} \rangle$$

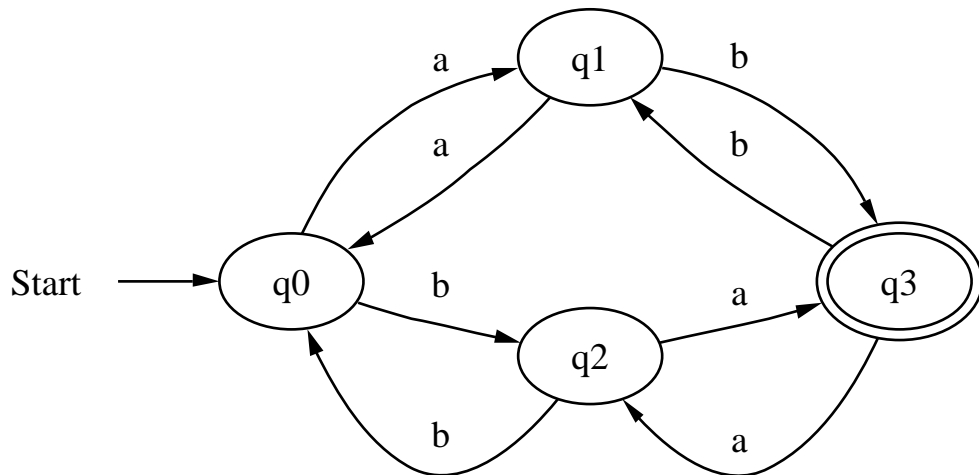
ist ein deterministischer endlicher Automat falls

1.  $\Sigma$ : Eingabe–Alphabet, endliche Menge
2.  $Z$ : Menge der Zustände, endlich
3.  $A$ : Menge der *akzeptierenden Zustände*,  $A \subseteq Z$
4.  $s_0$ : Start–Zustand,  $s_0 \in Z$
5.  $\text{next} : Z \times \Sigma \rightarrow Z$   
heißt *Zustands–Übergangs–Funktion*

Graphische Darstellung:

1. Zustände: Kreis um Namen des Zustands
2. Akzeptierende Zustände: doppelte Kreise
3. Zustands–Übergangs–Funktion:  
Gilt  $s_2 = \text{next}(s_1, a)$ , so verbinden wir  $s_1$  und  $s_2$  durch Pfeil, der mit  $a$  beschriftet ist.
4. Start–Zustand: steht in Diagramm entweder ganz links oder ganz oben, Pfeil zeigt auf Start–Zustand

## Ein Beispiel



Sei  $F_{ug} := \langle \Sigma, Z, A, s_0, \text{next} \rangle$  mit

1.  $\Sigma = \{a, b\}$ .
2.  $Z = \{q_0, q_1, q_2, q_3\}$ .
3.  $A = \{q_3\}$ .
4.  $s_0 = q_0$ .
5.  $\text{next}(q_0, a) = q_1, \text{next}(q_0, b) = q_2$ .
6.  $\text{next}(q_1, a) = q_0, \text{next}(q_1, b) = q_3$ .
7.  $\text{next}(q_2, a) = q_3, \text{next}(q_2, b) = q_0$ .
8.  $\text{next}(q_3, a) = q_2, \text{next}(q_3, b) = q_1$ .

# Berechnung eines endl. Automaten

**Gegeben:** Endlicher Automat

$$F = \langle \Sigma, Z, A, s_0, \text{next} \rangle$$

**Induktive Definition** der Berechnung eines endlichen Automaten:

1.  $q_1 \xrightarrow{c} q_2$  Berechnung mit Label  $c \in \Sigma$  falls  
 $q_2 = \text{next}(q_1, c)$
2.  $q_0 \xrightarrow{x_1} \dots \xrightarrow{x_n} q_n \xrightarrow{x_{n+1}} \dots \xrightarrow{x_m} q_m$   
Berechnung mit Label  $vw$  falls
  - (a)  $q_0 \xrightarrow{x_1} \dots \xrightarrow{x_n} q_n$  Berechnung mit Label  $v$
  - (b)  $q_n \xrightarrow{x_{n+1}} \dots \xrightarrow{x_m} q_m$  Berechnung mit Label  $w$  ist.

**Definition** der *akzeptierten Sprache*:

Ein Wort  $w \in \Sigma^*$  ist genau dann in der akzeptierten Sprache  $\mathcal{L}(F)$ , wenn es eine Berechnung

$$s_0 \xrightarrow{x_0} \dots \xrightarrow{x_n} q$$

mit Label  $w$  gibt und  $q \in A$  liegt.

**Beachte**, dass Berechnung im Start-Zustand  $s_0$  startet!

**Beispiel:** Betrachte letzten Automaten:

$$q_0 \xrightarrow{a} q_1 \xrightarrow{a} q_0 \xrightarrow{b} q_2 \xrightarrow{a} q_3$$

Berechnung mit Label  $aaba$ , also  $aaba \in \mathcal{L}(F)$ .



# Arbeitsweise eines endlichen Automaten

## Gegeben:

1. Endlicher Automat  $F = \langle \Sigma, Z, A, s_0, \text{next} \rangle$
2. Wort  $w \in \Sigma^*$

**Frage:** Wird  $w = w[1] \cdots w[n]$  von  $F$  akzeptiert?

1. Beginne im Start-Zustand  $s_0$ :

$$q = s_0$$

$q$  ist aktueller Zustand

2. Setze  $\text{idx} = 1$

3. Berechne Folge-Zustand

$$q = \text{next}(q, w[\text{idx}])$$

4. Inkrementiere  $\text{idx}$ .

5. Falls  $\text{idx} \leq n$ : Gehe zu 3.

6. Ist  $q \in A$ :

(a) Ja:  $F$  hat  $w$  akzeptiert

(b) Nein:  $F$  hat  $w$  nicht akzeptiert

## C–Implementierung einer FSM

```
typedef enum { Q0, Q1, Q2, Q3 } State;

bool accept(const char* w)
{
    State q = Q0;
    while (*w != 0) {
        switch (q) {
            case Q0: {
                switch (*w) {
                    case 'a': {
                        q = Q1;
                        break;
                    }
                    case 'b': {
                        q = Q2;
                        break;
                    }
                }
                break;
            }
            case Q1: { ... }
            case Q2: { ... }
            case Q3: { ... }
            ++w;
        }
        if (q == Q3)
            return true;
        return false;
    }
}
```

## Erweiterung der Zustands–Übergangs–Funktion

Wir erweitern die Zustands–Übergangs–Funktion

$$\text{next} : Z \times \Sigma \rightarrow Z$$

zu einer Funktion

$$\text{next}^* : Z \times \Sigma^* \rightarrow Z$$

Definition von  $\text{next}^*(w)$  induktiv über Länge von  $w$

1.  $\text{next}^*(q, \varepsilon) = q$
2.  $\text{next}^*(q, bw) = \text{next}^*(\text{next}(q, b), w)$

**Bemerkung:** Endlicher Automat

$$F = \langle \Sigma, Z, A, s_0, \text{next} \rangle$$

akzeptiert  $w \in \Sigma^*$  g.d.w.

$$\text{next}^*(s_0, w) \in A$$

**Bemerkung:** *akzeptierte Sprache*

$$\mathcal{L}(F) = \{w \in \Sigma^* \mid \text{next}^*(s_0, w) \in A\}$$

**Beispiel** von Folie 6:

Bezeichne  $\text{nr}(x, w)$  Anzahl der Buchstaben  $x$ , die in  $w$  auftreten. Dann gilt:

$$\mathcal{L}(F_{ug}) = \{w \mid \text{nr}(a, w) \% 2 = 1 \wedge \text{nr}(b, w) \% 2 = 1\}$$

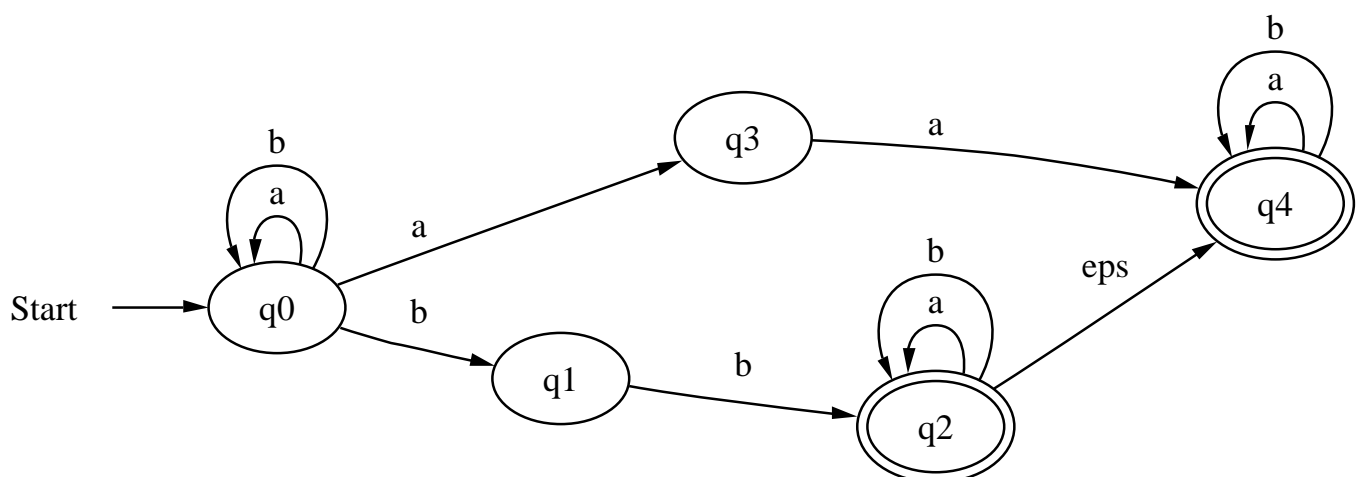
# Nicht-deterministische Endliche Automaten

**Definition:** Das 6-Tupel

$$\langle \Sigma, Z, A, s_0, \text{next}, \text{eps} \rangle$$

ist nicht-deterministischer endl. Automat mit  $\varepsilon$  Übergängen falls

1.  $\Sigma$ : Eingabe-Alphabet, endliche Menge
2.  $Z$ : Menge der Zustände, endlich
3.  $A$ : Menge der *akzeptierenden Zustände*,  $A \subseteq Z$
4.  $s_0$ : Start-Zustand,  $s_0 \in Z$
5.  $\text{next} : Z \times \Sigma \rightarrow 2^Z$   
heißt nicht-determ. *Zustands-Übergangs-Funktion*
6.  $\text{eps} : Z \rightarrow 2^Z$   
heißt *Epsilon-Übergangs-Funktion*



## Akzeptierte Sprache

**Gegeben:** Nicht-deterministischer endl. Automat

$$F = \langle \Sigma, Z, A, s_0, \text{next}, \text{eps} \rangle$$

**Induktive Definition** der Berechnung eines endlichen Automaten:

1.  $q_1 \xrightarrow{\varepsilon} q_2$  Berechnung mit Label  $\varepsilon$  falls  $q_2 \in \text{eps}(q_1)$
2.  $q_1 \xrightarrow{c} q_2$  Berechnung mit Label  $c \in \Sigma$  falls  $q_2 \in \text{next}(q_1, c)$
3.  $q_0 \xrightarrow{x_1} \dots \xrightarrow{x_n} q_n \xrightarrow{x_{n+1}} \dots \xrightarrow{x_m} q_m$   
Berechnung mit Label  $vw$  falls
  - (a)  $q_0 \xrightarrow{x_1} \dots \xrightarrow{x_n} q_n$  Berechnung mit Label  $v$
  - (b)  $q_n \xrightarrow{x_{n+1}} \dots \xrightarrow{x_m} q_m$  Berechnung mit Label  $w$  ist.

**Definition** der *akzeptierten Sprache*:

Ein Wort  $w \in \Sigma^*$  ist genau dann in der akzeptierten Sprache  $\mathcal{L}(F)$ , wenn es eine Berechnung

$$s_0 \xrightarrow{x_0} \dots \xrightarrow{x_n} q$$

mit Label  $w$  gibt, so dass  $q \in F$  ist.

**Beachte**, dass Berechnung im Start-Zustand  $s_0$  startet!

**Beispiel:** Betrachte letzten Automaten:

$$q_0 \xrightarrow{a} q_0 \xrightarrow{b} q_1 \xrightarrow{b} q_2 \xrightarrow{\varepsilon} q_4 \xrightarrow{a} q_4$$

Berechnung mit Label  $abba$ , also  $abba \in \mathcal{L}(F)$ .

## Berechnung der akzeptierten Sprache

Wir erweitern die Funktion `next` und `eps` zu Funktionen

$$\text{Next} : 2^Z \times \Sigma \rightarrow 2^Z \quad \text{und} \quad \text{Eps} : 2^Z \rightarrow 2^Z$$

$$\text{Next}(Q, b) := \{z \in Z \mid \exists q \in Q : z \in \text{next}(q, b)\}$$

$$\text{Eps}(Q) := Q \cup \{z \in Z \mid \exists q \in Q : z \in \text{eps}(q)\}$$

### Beispiel:

1.  $\text{Next}(\{q_0\}, a) = \{q_0, q_3\}$
2.  $\text{Next}(\{q_0, q_3\}, a) = \{q_0, q_3, q_4\}$
3.  $\text{Eps}(\{q_0\}) = \{q_0\}$
4.  $\text{Eps}(\{q_2\}) = \{q_2, q_4\}$

Wir erweitern die Funktion `Next` induktiv auf Worte  $w$  zu einer Funktion

$$\text{Next}^* : 2^Z \times \Sigma^* \rightarrow 2^Z$$

1.  $\text{Next}^*(Q, \varepsilon) := \text{Eps}(Q)$
2.  $\text{Next}^*(Q, bw) := \text{Next}^*\left(\text{Next}(\text{Eps}(Q), b), w\right)$

### Beispiel:

1.  $\text{Next}^*(\{q_0\}, aa) = \{q_0, q_3, q_4\}$
2.  $\text{Next}^*(\{q_0\}, abba) = \{q_0, q_2, q_3, q_4\}$

**Satz:** Ist  $F$  nicht-determ. Automat, so gilt:

$$\mathcal{L}(F) = \{w \in \Sigma^* \mid \text{Next}^*(\{s_0\}, w) \cap F \neq \emptyset\}$$

# Potenz–Mengen–Konstruktion

**Gegeben:** Nicht–deterministische FSM

$$F_{nd} = \langle \Sigma, Z, A, s_0, \text{next}, \text{eps} \rangle$$

**Gesucht:** Deterministische FSM

$$F_{det} = \langle \Sigma, \mathcal{Z}, \mathcal{A}, S_0, \mathcal{NS} \rangle$$

mit  $\mathcal{L}(F_{det}) = \mathcal{L}(F_{nd})$

**Definition:**

1.  $\mathcal{Z} := 2^Z$
2.  $\mathcal{A} := \{M \in 2^Z \mid M \cap F \neq \emptyset\}$
3.  $S_0 := \{s_0\}$
4.  $\mathcal{NS}(Q) := \text{Eps}\left(\text{Next}\left(\text{Eps}(Q), b\right)\right)$

Zahl der Zustände:  $|\mathcal{Z}| = 2^{|Z|}$ : sehr groß!

**Satz:**  $\mathcal{L}(F_{det}) = \mathcal{L}(F_{nd})$

**Beobachtung:** Nicht–deterministische Automaten sind nicht mächtiger als deterministische Automaten.

# Übersetzung von regulären Ausdrücken in nicht-deterministische FSMs

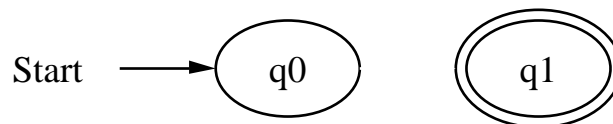
Wir definieren für jeden regulären Ausdruck  $r$  eine nicht-deterministische FSM durch Induktion über  $r$

Konstruktions-Invarianten:

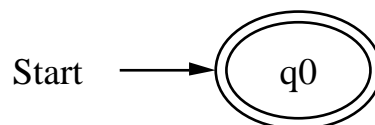
1. Menge der akzeptierenden Zustände ein-elementig
2.  $|\text{eps}(q)| \leq 2$
3.  $|\bigcup_{b \in \Sigma} \text{next}(q, b)| \leq 1$
4.  $\text{eps}(q) = \emptyset \vee \forall b \in \Sigma : \text{next}(q, b) = \emptyset$

## Induktive Definition:

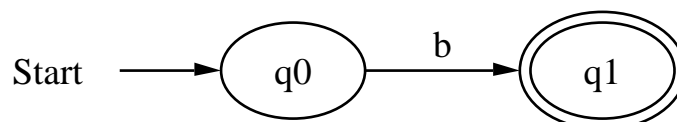
1.  $r = \emptyset$ :



2.  $r = \varepsilon$ :



3.  $r = b$  mit  $b \in \Sigma$ :

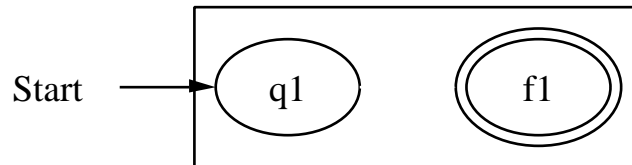




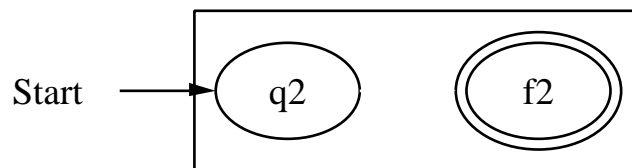
## RegExp $\mapsto$ FSM (Fortsetzung)

Nach IV seien FSMs für  $r_1$  und  $r_2$  wie folgt gegeben:

1.  $r_1$ :

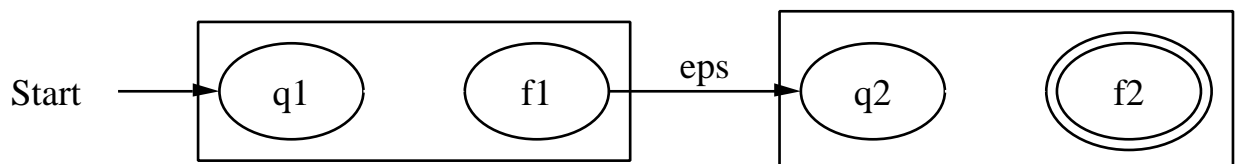


2.  $r_2$ :

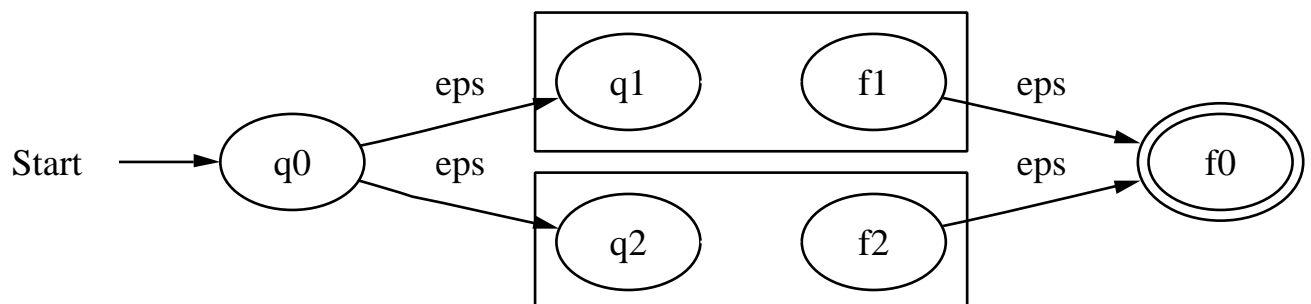


Fortsetzung der induktiven Definition:

4.  $r = r_1 r_2$

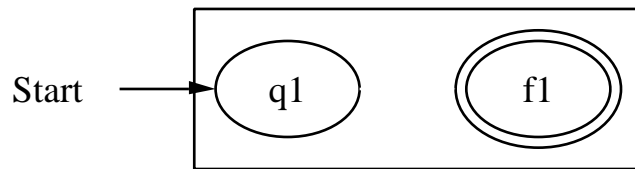


5.  $r = r_1 + r_2$



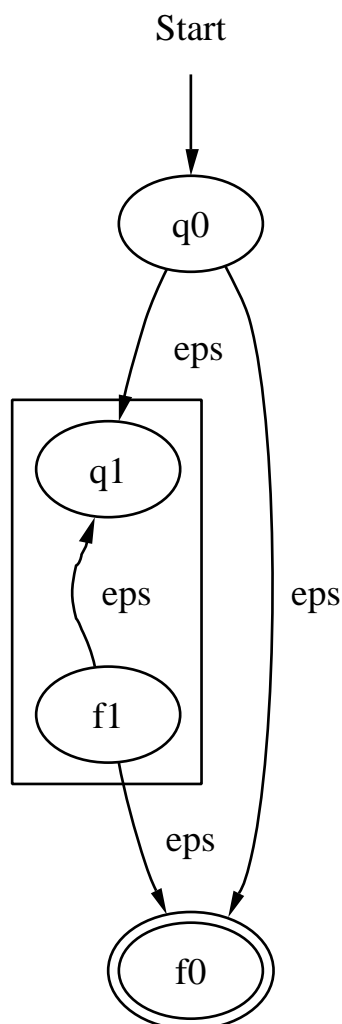
## Konstruktion der FSM für $r_1^*$

Nach IV sei FSM für  $r_1$  gegeben durch

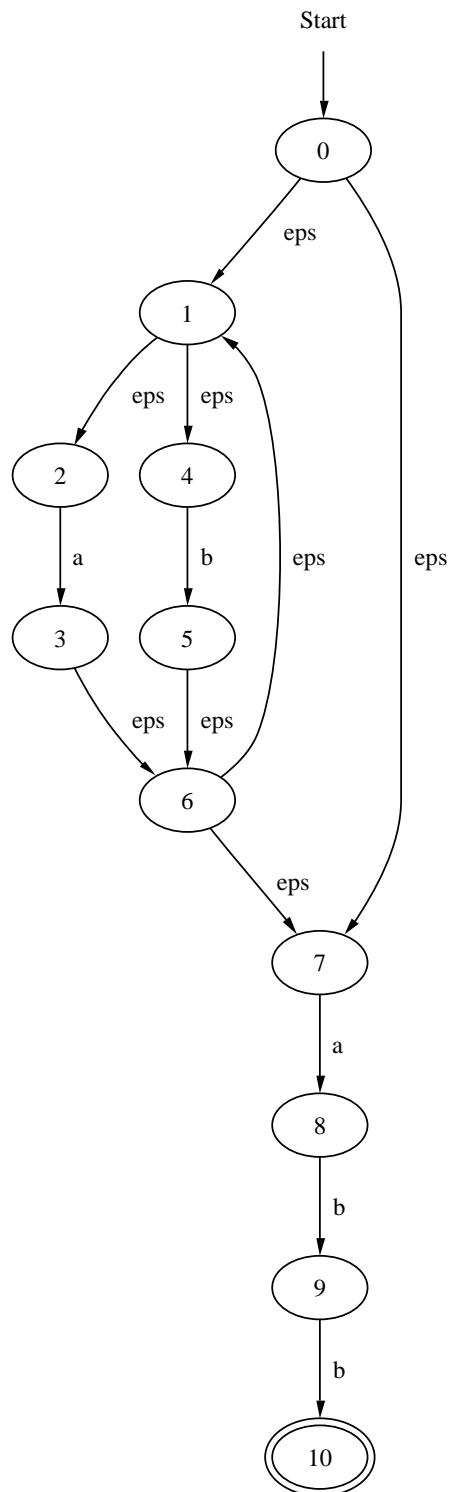


Fortsetzung der induktiven Definition

6.  $r = r_1^*$



Beispiel:  $(a|b)^*abb$



## Simulation nicht-deterministischer FSMs

**Problem:** Größe Potenz-Mengen-Konstruktion:  $2^{|Z|}$

**Lösung:** *Simulation* der Potenz-Mengen-Konstruktion  
Repräsentation einer FSM

1. Anzahl der Zustände: `numberStates`
2. Codierung der Zustände als Zahlen:  
 $0, \dots, \text{numberStates} - 1$
3. Feld von Symbolen: `symbol[numberStates]`
4. Feld von Folge-Zuständen: `next1[numberStates]`
5. Feld von Folge-Zuständen: `next2[numberStates]`

Bedeutung dieser Felder:

1.  $j \in \text{next}(i, c) \rightarrow \text{symbol}[i] = c$
2.  $(\forall c \in \Sigma : \text{next}(i, c) = \emptyset) \rightarrow \text{symbol}[i] = 0$
3.  $j \in \text{next}(i, c) \rightarrow \text{next1}[i] = j$
4.  $j \in \text{eps}(i) \rightarrow \text{next1}[i] = j \vee \text{next2}[i] = j$

C-Daten-Struktur

```
typedef struct {  
    unsigned    numberStates; // number of states  
    char*       symbol;       // array of characters  
    unsigned*   next1;        // possible next state  
    unsigned*   next2;        // possible next state  
} FSM;
```

## Representation der Beispiel-FSM

```
fsm->numberStates = 11;

// Numerierung der Zustände
//    0  1  2  3  4  5  6  7  8  9  10

fsm->symbol =
    { 0, 0, 'a', 0, 'b', 0, 0, 'a', 'b', 'b', 0 };

fsm->next1 =
    { 1, 2, 3, 6, 5, 6, 7, 8, 9, 10, 10 };

fsm->next2 =
    { 7, 4, 3, 6, 5, 6, 1, 8, 9, 10, 10 };
```

Berechnung des Epsilon-Abschluß

$$\text{epsClose} : 2^Z \rightarrow 2^Z$$

$\text{epsClose}(Q)$ : Menge der Zustände, die von Zuständen in  $Q$  durch Epsilon-Übergänge erreicht werden können.

Definition von  $\text{epsClose}(Q)$  iterativ:

1. Für alle  $Q \subseteq Z$  gilt:

$$Q \subseteq \text{epsClose}(Q)$$

2. Für alle  $Q \subseteq Z$  und alle  $q \in Z$  gilt:

$$q \in \text{epsClose}(Q) \Rightarrow \text{eps}(q) \subseteq \text{epsClose}(Q)$$

## Berechnung von $\text{epsClose}(Q)$

Repräsentation von Zustands-Mengen  $Q \subseteq Z$  durch Felder:

1. Für die Menge der Zustände gilt:

$$Z = \{q_0, q_1, \dots, q_n\}$$

mit  $n + 1 = \text{numberStates}$ .

2.  $Q \subseteq Z$  wird dargestellt durch

`bool state[numberStates]`

Dabei gilt

$$\text{state}[i] = \text{true} \leftrightarrow q_i \in Q$$

```
void epsClose(FSM* fsm, bool states[]) {
    bool change = true;
    while (change) {
        change = false;
        for (unsigned i = 0; i < fsm->numberStates; ++i)
        {
            if (states[i] && fsm->symbol[i] == 0) {
                if (!states[fsm->next1[i]]) {
                    change = true;
                    states[fsm->next1[i]] = true;
                }
                if (!states[fsm->next2[i]]) {
                    change = true;
                    states[fsm->next2[i]] = true;
                }
            }
        }
    }
}
```

# Simulation der Potenz–Mengen–Konstruktion

```
bool simulate(FSM* fsm, char* word) {
    bool currentStates[fsm->numberStates];
    bool nextStates[fsm->numberStates];
    currentStates[0] = true;
    for (unsigned i = 1; i < fsm->numberStates; ++i) {
        currentStates[i] = false;
    }
    epsClose(fsm, currentStates);
    while (*word != 0) {
        for (unsigned i = 0; i < fsm->numberStates; ++i)
            nextStates[i] = false;
        for (unsigned i = 0; i < fsm->numberStates; ++i) {
            if (    currentStates[i]
                && fsm->symbol[i] == *word)
            {
                nextStates[fsm->next1[i]] = true;
            }
        }
        for (unsigned i = 0; i < fsm->numberStates; ++i)
            currentStates[i] = nextStates[i];
        epsClose(fsm, currentStates);
        ++word;
    }
    return currentStates[fsm->numberStates - 1];
}
```

Zahl der Zustände:  $n$ , Länge des Wortes:  $m$

1. Komplexität `epsClose`:  $\mathcal{O}(n^2)$
2. Komplexität `simulate`:  $\mathcal{O}(m * n^2)$

# Konstruktion $\text{RegExp} \mapsto \text{FSM}$ : Implementierung

## 1. Speicherplatz reservieren

```
FSM* allocateFSM(unsigned n) {
    FSM* fsm = malloc( sizeof(FSM) );
    fsm->numberStates = n;
    fsm->symbol        = malloc( n * sizeof(int) );
    fsm->next1         = malloc( n * sizeof(int) );
    fsm->next2         = malloc( n * sizeof(int) );
    return fsm;
}
```

## 2. Fsm f2 an Stelle o in f1 kopieren

```
void move(FSM* f1, FSM* f2, unsigned o) {
    for (unsigned i = 0; i < f2->numberStates; ++i)
    {
        f1->symbol[i + o] = f2->symbol[i];
        f1->next1[i + o] = f2->next1[i] + o;
        f1->next2[i + o] = f2->next2[i] + o;
    }
}
```

## 3. $i \xrightarrow{\varepsilon} j, i \xrightarrow{\varepsilon} k$

```
void epsTransition(FSM* fsm, unsigned i,
                  unsigned j, unsigned k) {
    fsm->symbol[i] = 0;
    fsm->next1[i] = j;  fsm->next2[i] = k;
}
```

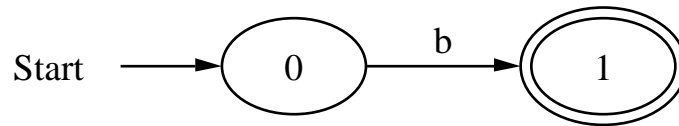
## 4. $i \xrightarrow{c} j$

```
void charTransition(FSM* fsm, unsigned i,
                   unsigned j, char c) {
    fsm->symbol[i] = c;
    fsm->next1[i] = j;  fsm->next2[i] = j;
}
```



# FSM zur Erkennung von Buchstaben

Akzeptieren des Buchstaben b



1. 2 Zustände

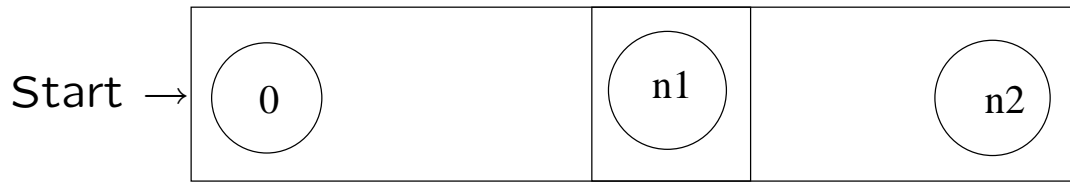
2.  $0 \xrightarrow{b} 1$

3.  $1 \xrightarrow{\varepsilon} 1$

Implementierung:

```
FSM* createCharacter(char c) {  
    FSM* fsm = allocateFSM(2);  
    charTransition(fsm, 0, 1, c);  
    epsTransition (fsm, 1, 1, 1);  
    return fsm;  
}
```

## Implementierung der Konkatenation

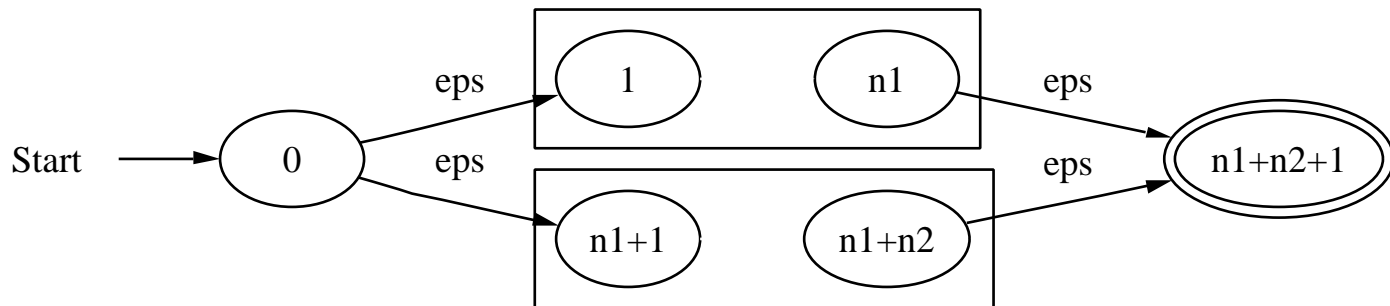


### Implementierung

1. Größe Fsm f1:  $n1$ .
2. Größe Fsm f2:  $n2$ .
3. Größe neue Fsm:  $n1 + n2 - 1$
4. Kopiere f1 an Offsett 0 in neuer Fsm.
5. Kopiere f2 an Offsett  $n1-1$  in neuer Fsm.

```
FSM* concat(FSM* f1, FSM* f2) {  
    unsigned n1 = f1 ->numberStates;  
    unsigned n2 = f2 ->numberStates;  
    unsigned n  = n1 + n2 - 1;  
    FSM* fsm = allocateFSM(n);  
    move(fsm, f1, 0);  
    move(fsm, f2, n1 - 1);  
    freeFsm(f1);  
    freeFsm(f2);  
    return fsm;  
}
```

## Implementierung der Alternative

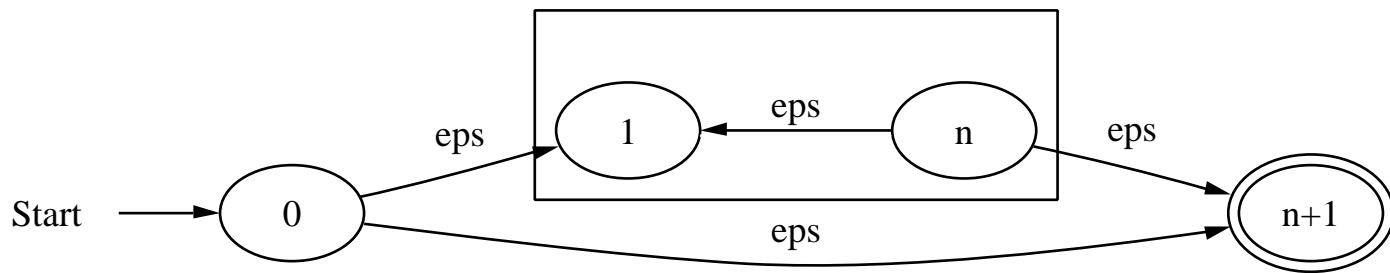


### Implementierung

1. Größe Fsm f1: n1.
2. Größe Fsm f2: n2.
3. Größe neue Fsm:  $n1 + n2 + 2$
4. Kopiere f1 an Offsett 1 in neuer Fsm.
5. Kopiere f2 an Offsett n1+1 in neuer Fsm.

```
FSM* alternative(FSM* f1, FSM* f2) {
    unsigned n1 = f1 ->numberStates;
    unsigned n2 = f2 ->numberStates;
    unsigned n = n1 + n2 + 2;
    FSM* fsm = allocateFSM(n);
    epsTransition(fsm, 0, 1, n1 + 1);
    move(fsm, f1, 1);
    move(fsm, f2, n1 + 1);
    epsTransition(fsm, n1, n - 1, n - 1);
    epsTransition(fsm, n - 2, n - 1, n - 1);
    epsTransition(fsm, n - 1, n - 1, n - 1);
    freeFsm(f1);    freeFsm(f2);
    return fsm;
}
```

## Implementierung des Abschlusses



1. Größe Fsm f: n.
2. Größe neue Fsm:  $n + 2$
3. Kopiere f an Offsett 1 in neuer Fsm.

### Implementierung

```
FSM* closure(FSM* f) {  
    unsigned n = f->numberStates;  
    FSM* fsm = allocateFSM(n+2);  
    epsTransition(fsm, 0, 1, n + 1);  
    move(fsm, f, 1);  
    epsTransition(fsm, n, n + 1, 1);  
    epsTransition(fsm, n + 1, n + 1, n + 1);  
    freeFsm(f);  
    return fsm;  
}
```

# Konstruktion einer RegExp aus einer FSM

**Definitionen:** Sei  $\Sigma$  ein endliches Alphabet.

1. Es seien  $\mathcal{L}_1, \mathcal{L}_2 \subseteq \Sigma^*$ . Wir definieren die *Konkatenation* von  $\mathcal{L}_1$  und  $\mathcal{L}_2$  als:

$$\mathcal{L}_1 \mathcal{L}_2 := \{vw \mid v \in \mathcal{L}_1 \wedge w \in \mathcal{L}_2\}$$

2. Sei  $\mathcal{L} \subseteq \Sigma^*$ . Für  $n \in \mathbb{N}$  definieren wir die *Potenz*  $\mathcal{L}^n$  durch Induktion über  $n$ :

$$\text{I.A. } n \mapsto 0: \quad \mathcal{L}^0 := \{\varepsilon\}$$

$$\text{I.S. } n \mapsto n + 1: \quad \mathcal{L}^{n+1} := \mathcal{L}^n \mathcal{L}$$

3. Sei  $\mathcal{L} \subseteq \Sigma^*$ . Wir definieren den *Abschluß* von  $\mathcal{L}$  als

$$\mathcal{L}^* = \bigcup_{n=0}^{\infty} \mathcal{L}^n$$

4. Seien  $u, w \in \Sigma^*$ .  $u$  ist *Präfix* von  $w$ , g.d.w. es gibt  $v \in \Sigma^*$  mit  $v \neq \varepsilon$  und  $uv = w$ . Schreibweise:  $u \prec w$

$$u \prec w \Leftrightarrow \exists v \in \Sigma^* : v \neq \varepsilon \wedge uv = w.$$

**Beispiele:** Sei  $\Sigma = \{a, b\}$ ,  $\mathcal{L}_1 = \{ba, b\}$  und  $\mathcal{L}_2 = \{abb, bb\}$ .

$$1. \quad \mathcal{L}_1 \mathcal{L}_2 = \{baabb, babb, babb, bbb\}$$

$$2. \quad \mathcal{L}_1^3 = \{bbb, bbba, bbab, babb, babab, babba, bbaba, bababa\}$$

$$3. \quad abbabb \prec abbabbabb$$

## Weitere Definitionen

**Definition:** Gegeben sei eine FSM

$$F = \langle \Sigma, \{q_1, \dots, q_n\}, A, q_1, \text{next} \rangle$$

Wir definieren eine partielle Funktion

$$\text{next}^{(k)} : Z \times \Sigma^* \rightarrow Z$$

1. Fall:  $\forall v \in \Sigma^* \setminus \{\varepsilon\} : v \prec w \rightarrow \text{next}(q, v) \in \{q_1, \dots, q_k\}$

Dann setzen wir

$$\text{next}^{(k)}(q, w) = \text{next}^*(q, w).$$

2. Fall:  $\exists v \in \Sigma^* \setminus \{\varepsilon\} : v \prec w \wedge \text{next}(q, v) \notin \{q_1, \dots, q_k\}$

Dann sei  $\text{next}^{(k)}(q, w)$  undefiniert, wir setzen

$$\text{next}^{(k)}(q, w) = \uparrow.$$

Falls  $\text{next}^{(k)}(q, w) = p$  ist, so sagen wir, dass die Berechnung von  $w$  im Zustand  $q$ , die Zustände aus  $\{q_{k+1}, \dots, q_n\}$  *vermeidet*.

Für  $k \in \{0, 1, \dots, n\}$  und  $i, j \in \{1, \dots, n\}$  definieren wir Sprachen  $\mathcal{R}_{i,j}^{(k)}$

$$\mathcal{R}_{i,j}^{(k)} := \{w \in \Sigma^* \mid \text{next}^{(k)}(q_i, w) = q_j\}$$

**Bemerkung:** Es gilt

$$\mathcal{R}_{i,j}^{(n)} := \{w \in \Sigma^* \mid \text{next}^*(q_i, w) = q_j\}$$

## Konstruktion einer RegExp aus einer FSM

**Gegeben:** FSM  $F = \langle \Sigma, \{q_1, \dots, q_n\}, A, q_1, \text{next} \rangle$

**Gesucht:**  $r \in \text{RegExp}$  mit  $\mathcal{L}(r) = \mathcal{L}(F)$

Wir geben induktive Konstruktionen der Sprachen  $\mathcal{R}_{i,j}^{(k)}$ :

1. Induktions-Anfang:  $k = 0$

(a)  $\exists a \in \Sigma : \text{next}(q_i, a) = q_j$  und  $i \neq j$

$$\mathcal{R}_{i,j}^{(0)} := \{a\}$$

(b)  $\exists a \in \Sigma : \text{next}(q_i, a) = q_i$

$$\mathcal{R}_{i,i}^{(0)} := \{a, \varepsilon\}$$

(c)  $\forall a \in \Sigma : \text{next}(q_i, a) \neq q_j$  und  $i \neq j$

$$\mathcal{R}_{i,j}^{(0)} := \{\}$$

(d)  $\forall a \in \Sigma : \text{next}(q_i, a) \neq q_i$  und  $i = j$

$$\mathcal{R}_{i,i}^{(0)} := \{\varepsilon\}$$

2. Induktions-Schritt:  $k \mapsto k + 1$

$$\mathcal{R}_{i,j}^{(k+1)} := \mathcal{R}_{i,j}^{(k)} \cup \mathcal{R}_{i,k+1}^{(k)} \left( \mathcal{R}_{k+1,k+1}^{(k)} \right)^* \mathcal{R}_{k+1,j}^{(k)}$$

Begründung: Es ist  $w \in \mathcal{R}_{i,j}^{(k+1)}$  g.d.w.

(a)  $w \in \mathcal{R}_{i,j}^{(k)}$  oder

(b)  $w = xv_1v_2 \cdots v_ny$  mit  $n \geq 0$  und

- $x \in \mathcal{R}_{i,k+1}^{(k)}, \quad y \in \mathcal{R}_{k+1,j}^{(k)}$
- $v_i \in \mathcal{R}_{k+1,k+1}^{(k)}$  für alle  $i = 1, \dots, n$

## Abschluß der Konstruktion

Wir definieren reguläre Ausdrücke  $r_{i,j}^{(k)}$  mit

$$\mathcal{L}(r_{i,j}^{(k)}) = \mathcal{R}_{i,j}^{(k)}$$

1. Induktions–Anfang:  $k = 0$

(a)  $\exists a \in \Sigma : \text{next}(q_i, a) = q_j$  und  $i \neq j$

$$r_{i,j}^{(0)} := a$$

(b)  $\exists a \in \Sigma : \text{next}(q_i, a) = q_i$

$$r_{i,i}^{(0)} := a + \varepsilon$$

(c)  $\forall a \in \Sigma : \text{next}(q_i, a) \neq q_j$  und  $i \neq j$

$$r_{i,j}^{(0)} := \emptyset$$

(d)  $\forall a \in \Sigma : \text{next}(q_i, a) \neq q_i$

$$r_{i,i}^{(0)} := \varepsilon$$

2. Induktions–Schritt:  $k \mapsto k + 1$

$$r_{i,j}^{(k+1)} := r_{i,j}^{(k)} + r_{i,k+1}^{(k)} \left( r_{k+1,k+1}^{(k)} \right)^* r_{k+1,j}^{(k)}$$

Setze  $r_{i,j} := r_{i,j}^{(n)}$ .

Sei  $A = \{q_l, \dots, q_n\}$  Menge der akzeptierenden Zustände.

$$r_F := r_{1,l} + r_{1,l+1} + \dots + r_{1,n}.$$

Dann gilt

$$\mathcal{L}(F) = \mathcal{L}(r_F)$$



## Ein Beispiel

**Definition:** Eine deterministische FSM

$$F = \langle \Sigma, Z, A, s_0, \text{next} \rangle$$

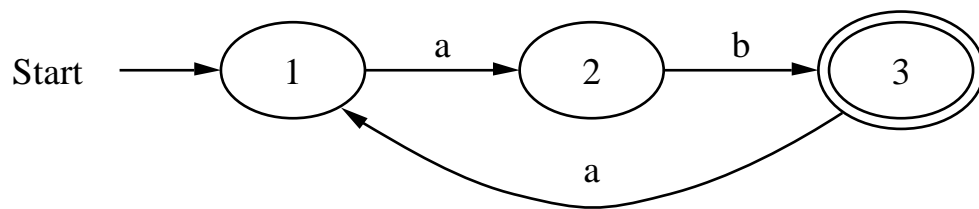
ist eine *partielle* FSM, wenn die Funktion `next` teilweise undefiniert ist.

**Bemerkung:** Die Konstruktion eines regulären Ausdrucks  $r$  mit

$$\mathcal{L}(r) = \mathcal{L}(F)$$

funktioniert auch für partielle FSMs.

**Aufgabe:** Berechnen Sie für folgende FSM einen äquivalenten regulären Ausdruck.



### Nützliche Vereinfachungs-Regeln

1.  $(\varepsilon + r)^* \simeq r^*$
2.  $(\varepsilon + r)(\varepsilon + r)^* \simeq r^*$
3.  $r_1 + r_1 r_2^* \simeq r_1 r_2^*$
4.  $r_1 + r_1 r_2 r_2^* \simeq r_1 r_2^*$
5.  $r_1 (r_2 r_1)^* r_2 = r_1 r_2 (r_1 r_2)^*$

## Berechnung des regulären Ausdrucks

	$k = 0$	$k = 1$	$k = 2$	$k = 3$
$r_{1,1}^{(k)}$	$\varepsilon$	$\varepsilon$	$\varepsilon$	$(aba)^*$
$r_{1,2}^{(k)}$	$a$	$a$	$a$	$a(baa)^*$
$r_{1,3}^{(k)}$	$\emptyset$	$\emptyset$	$ab$	$ab(aab)^*$
$r_{2,1}^{(k)}$	$\emptyset$	$\emptyset$	$\emptyset$	$ba(aba)^*$
$r_{2,2}^{(k)}$	$\varepsilon$	$\varepsilon$	$\varepsilon$	$(baa)^*$
$r_{2,3}^{(k)}$	$b$	$b$	$b$	$b(baa)^*$
$r_{3,1}^{(k)}$	$a$	$a$	$a$	$a(aba)^*$
$r_{3,2}^{(k)}$	$\emptyset$	$aa$	$aa$	$aa(baa)^*$
$r_{3,3}^{(k)}$	$\varepsilon$	$\varepsilon$	$\varepsilon + aab$	$(aab)^*$

1. Spalte: folgt unmittelbar aus Diagramm
2. Spalte: einzige Änderung dort, wo Pfade der Länge 2 sind, deren mittlerer Knoten 1 ist:

$$3 \xrightarrow{a} 1 \xrightarrow{a} 2, \text{ also } r_{3,2}^{(1)} = aa$$

3. Spalte: einzige Änderung dort, wo Pfade mit Knoten 2 hinzukommen:

$$3 \xrightarrow{a} 1 \xrightarrow{a} 2 \xrightarrow{b} 3, \text{ also } r_{3,3}^{(2)} = \varepsilon + aab$$

$$1 \xrightarrow{a} 2 \xrightarrow{b} 3, \text{ also } r_{1,3}^{(2)} = ab$$

4. Spalte: Sei  $i \xrightarrow{u} j$  und  $j \xrightarrow{v} j$

$$r_{i,j}^{(3)} = uv^*$$

Also:  $\mathcal{L}(F) = \mathcal{L}(ab(aab)^*)$

# Pumping Lemma

**Definition:** Für ein Wort  $w \in \Sigma^*$  bezeichnet

$$|w|$$

die *Länge* (Zahl der Buchstaben) von  $w$ .

Für  $n \in \mathbb{N}$  definieren wir  $w^n$  induktiv:

I.A.  $n \mapsto 0: \quad w^0 := \varepsilon$

I.S.  $n \mapsto n + 1: \quad w^{n+1} := ww^n$

**Definition:** Sei  $\Sigma$  Alphabet. Eine Sprache  $\mathcal{L} \subseteq \Sigma^*$  ist *regulär der Ordnung  $n$*  g.d.w. es einen endl. Automaten

$$F = \langle \Sigma, Z, A, s_0, \text{next} \rangle$$

gibt, so daß gilt:  $\mathcal{L} = \mathcal{L}(F)$  und  $\text{card}(Z) = n$ .

Hier bezeichnet  $\text{card}(Z)$  die Zahl der Zustände.

**Satz** (*Pumping Lemma*)

**Vor.:**  $\mathcal{L}$  ist reguläre Sprache der Ordnung  $n$ .

**Beh.:** Für alle  $w \in \mathcal{L}$  mit  $|w| \geq n$  gibt es  $x, y, z \in \Sigma^*$  mit

1.  $w = xyz$
2.  $|xy| \leq n$
3.  $|y| \geq 1$
4.  $\forall k \in \mathbb{N} : xy^kz \in \Sigma^*$

## Beweis des Pumping Lemma

Sei  $\mathcal{L} = \mathcal{L}(F)$  mit  $F = \langle \Sigma, Z, q_0, \text{next} \rangle$  und  $n = \text{card}(Z)$ .

Sei  $w = a_1 a_2 \cdots a_m \in \mathcal{L}(F)$  mit  $m \geq n$ . Dann gibt es Berechnung der FSM  $F$  mit Label  $w$ :

$$q_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \xrightarrow{a_3} \cdots \xrightarrow{a_m} q_m \quad \text{und } q_m \in A$$

Es können nicht alle Zustände in der Menge

$$\{q_0, q_1, q_2, \dots, q_n\}$$

verschieden sein, es gibt also  $i, j \in \{0, \dots, n\}$  mit

$$i < j \quad \text{und} \quad q_i = q_j.$$

Wir definieren

$$x := a_1 \cdots a_i, \quad y := a_{i+1} \cdots a_j, \quad z := a_{j+1} \cdots a_m$$

Dann haben wir die folgenden Berechnungen der FSM  $F$

1.  $q_0 \xrightarrow{a_1} \cdots \xrightarrow{a_i} q_i$ , also  $q_0 \xrightarrow{x} q_i$
2.  $q_i \xrightarrow{a_{i+1}} \cdots \xrightarrow{a_j} q_j$ , also  $q_i \xrightarrow{y} q_j$
3.  $q_j \xrightarrow{a_{j+1}} \cdots \xrightarrow{a_m} q_m$ , also  $q_j \xrightarrow{z} q_m$

Damit gilt auch

$$q_0 \xrightarrow{x} \underbrace{q_i \xrightarrow{y} q_i \xrightarrow{y} \cdots \xrightarrow{y} q_i}_k \xrightarrow{z} q_m$$

und das zeigt, dass

$$xy^k z \in \mathcal{L} \quad \text{für alle } k \in \mathbb{N}$$

Aus  $i < j$  und  $|y| = j - i$  folgt  $|y| \geq 1$ .

Wegen  $|xy| = j$  und  $j \in \{0, \dots, n\}$  folgt auch

$$|xy| \leq n.$$

□

## Anwendung des Pumping Lemma

**Aufgabe:** Sei  $\Sigma = \{a, b\}$ . Zeigen Sie, dass die Sprache

$$\mathcal{L} := \{w \in \Sigma^* \mid \text{nr}(w, a) = \text{nr}(w, b)\}$$

nicht regulär ist.

Ein Wort in  $\mathcal{L}$  enthält genauso viele Buchstaben a wie b.

**Beweis** indirekt.

**Annahme:** Sei  $\mathcal{L}$  reguläre Sprache der Ordnung  $n$ .  
Betrachte das Wort  $w := a^n b^n$ . Es gilt

$$a^n b^n \in \mathcal{L} \quad \text{und} \quad |a^n b^n| = 2 * n.$$

Also gibt es Worte  $x$ ,  $y$  und  $z$  mit:

$$1. \quad w = a^n b^n = xyz$$

$$2. \quad |xy| \leq n$$

$$3. \quad |y| \geq 1$$

$$4. \quad xy^k z \in \mathcal{L}$$

Aus  $|xy| \leq n$  und  $xyz = a^n b^n$  folgt

$$xy \prec a^n$$

und daraus folgt

$$\text{nr}(x, b) = 0 \quad \text{und} \quad \text{nr}(y, b) = 0.$$

## Fortsetzung der Aufgabe

Aus  $xyz \in \mathcal{L}$  folgt wegen

$$\text{nr}(xyz, a) = \text{nr}(x, a) + \text{nr}(y, a) + \text{nr}(z, a) \text{ und}$$

$$\text{nr}(xyz, b) = \text{nr}(x, b) + \text{nr}(y, b) + \text{nr}(z, b) = \text{nr}(z, b)$$

die Gleichung

$$\text{nr}(x, a) + \text{nr}(y, a) + \text{nr}(z, a) = \text{nr}(z, b) \quad (\text{A})$$

Ebenso folgt aus  $xy^2z \in \mathcal{L}$  und

$$\text{nr}(xy^2z, a) = \text{nr}(x, a) + 2 * \text{nr}(y, a) + \text{nr}(z, a) \text{ und}$$

$$\text{nr}(xy^2z, b) = \text{nr}(x, b) + 2 * \text{nr}(y, b) + \text{nr}(z, b) = \text{nr}(z, b)$$

die Gleichung

$$\text{nr}(x, a) + 2 * \text{nr}(y, a) + \text{nr}(z, a) = \text{nr}(z, b) \quad (\text{B})$$

Ziehen wir (A) von (B) ab, so erhalten wir

$$\text{nr}(y, a) = 0$$

Weil auch

$$\text{nr}(y, b) = 0$$

gilt, folgt

$$y = \varepsilon$$

Das steht im Widerspruch zu

$$|y| \geq 1 .$$

□

**Erkenntnis:** Endliche Automaten können nicht zählen!

## Reguläre Ausdrücke in der Praxis: sed

Zeichen mit Sonderbedeutung bei sed:

1. `"."`: paßt auf jeden Buchstaben außer Zeilenumbruch.  
`"a.c"` paßt auf `"abc"`, `"azc"`, `"a1c"`, ...
2. `"\*"`: Quantor  
beliebig viele Wiederholungen.  
`"fo\*"` paßt auf `"f"`, `"fo"`, `"foo"`, `"fooo"`, etc.
3. `"\+"`: Quantor  
mindestens eine Wiederholung.  
`"fo\+"` paßt auf `"fo"`, `"foo"`, `"fooo"`, ...
4. `"\?"`: Quantor  
keine oder eine Wiederholung.  
`"fo\?"` paßt auf `"f"` und `"fo"`.
5. `"\|"`: Auswahl von zwei Möglichkeiten  
`"eins\|zwei"` paßt sowohl auf `"eins"` als auch auf `"zwei"`
6. `"^"` paßt auf leeren String am Zeilen–Anfang  
`"^int"` paßt auf `"int"` am Zeilen–Anfang.
7. `"$"` paßt auf leeren String am Zeilen–Ende.  
`"^$"` paßt auf Leerzeile.

## Reguläre Ausdrücke

### 8. “[” und “]” begrenzen *Zeichen–Mengen*

“[ad]” paßt auf “a” und “d”.

Spezifikation von *Intervallen* durch “–”

(a) “[a–z]” paßt auf jeden Klein–Buchstaben.

(b) “[a–z\*.]” paßt auf jeden Klein–Buchstaben und auf “\*” und “.”

**Merke:** “\*”, “+”, “?”, “.” verlieren Sonderbedeutung in Zeichen–Mengen.

(c) “^”: Komplement einer Zeichen–Mengen.

“[^0–9]” paßt auf alle Zeichen, die keine Ziffern sind.

### 9. “\ (“ und “\)” : Gruppierung

“(ja)\+” paßt auf

“ja”, “jaja”, “jajaja”, ...

“(J\|j\)a” paßt auf

“Ja” und “ja”

### Anwendung: sed

sed *s/regex/replacement/g*

*replacement* kann “\1”, “\2”, “\3” ... “\3” enthalten.

“\n” steht für *n*-te Klammer

sed 's/\emph{\([^}]\*\)}/<em>\1</em>/g'

ersetzt “\emph{*string*” durch “<em>*string*</em>”.



## Ein komplexeres Beispiel sed

**Gegeben:** Datei der Form

```
8 5
9 3
10 3
11 9
12 3
...
```

**Idee:** Datei spezifiziert Primzahlen

$n \quad m$  spezifiziert  $2^n - m$

Automatisches Ausrechnen mit sed und bc

```
cat prim.txt | \
  sed 's/\([0-9]\+\) \([0-9]\+\)/2^\1 - \2/g' | \
  bc
```

Andere Kommandos, die mit regulären Ausdrücken arbeiten

1. `grep regexp file1 ... filen`
2. *XEmacs*:
  - (a) `isearch-forward-regexp`
  - (b) `query-replace-regexp`
3. `find [directory] -regexp regexp-pattern [action]`
4. Grundlage von *Perl*, *Tcl*, *Python*