

## 1.6 Aakkostot, merkkijonot ja kielet

Automaattiteoria. ~ diskreetin signaalinkäsittelyn perusmallit ja -menetelmät

(~ diskreettien I/O-kuvausten yleinen teoria)



Automaatin käsite on *matematiikan abstraktio*. Yleisellä tasolla suunniteltu automaatti voidaan toteuttaa eri tavoin: esim. sähköpiirinä, mekaanisena laitteena tai (tavallisimmin) tietokoneohjelmana.

## Peruskäsitteitä ja merkintöjä

**Aakkosto** (engl. alphabet, vocabulary): äärellinen, epätyhjä joukko alkioita, joita sanotaan symboliteita. Esim.:

- ▶ binääriraakkoisto  $\{0, 1\}$ ;
- ▶ latinalainen aakkosto  $\{A, B, \dots, Z\}$ .

**Merkkijono** (engl. string): äärellinen järjestetty jono jonkin aakkoston merkkejä. Esim.:

- ▶ “01001”, “0000”: binääriraakkoiston merkkijonoja;
- ▶ “TKTP”, “XYZZY”: latinalaisen aakkoston merkkijonoja.
- ▶ tyhjä merkkijono (engl. empty string).

Tyhjässä merkkijonossa ei ole yhtään merkkiä; siitä voidaan merkitä  $\epsilon$ :llä.

Merkkijonon  $x$  pituus  $|x|$  on sen merkkien määrä. Esim.:  $|01001| = 5$ ,  $|\text{TKTP}| = 4$ ,  $|\epsilon| = 0$ .

Tällä kurssilla keskitytään pääosin automaatteihin, joiden:

1. syötteet ovat äärellisiä, diskreettejä *merkkijonoja*
2. tulokset ovat muotoa “hyväksy”/“hyllä” (~ “syöte OK”/”syöte ei kelppaa”)

Yleistynkseen:

- ▶ äärettömät syötejionot ( $\rightarrow$  “reaktiiviset” järistelmät, Büchi-automaatit)
- ▶ funktioautomaatti ( $\rightarrow$  Moore- ja Mealy-tiakoneet, Turingin funktiokoneet)

Merkkijonojen välinen perusoperaatio on *katenaatio* eli jonojen peräkkäin kirjoittaminen. Katernaation operaatiomerkkina käytetään joskus selkeyden lisäämiseksi symbolia  $\wedge$ . Esim.

- ▶  $KALA \wedge KUKKO = KALAKUKKO$ ;
- ▶ jos  $x = 00$  ja  $y = 11$ , niin  $xy = 0011$  ja  $yx = 1100$ ;
- ▶ kaikilla  $x$  on  $x\epsilon = \epsilon x = x$ ;
- ▶ kaikilla  $x, y, z$  on  $(xy)z = x(yz)$ ;
- ▶ kaikilla  $x, y$  on  $|xy| = |x| + |y|$ .

## Automaatit ja formaalit kielet

Olkoon  $M$  automaatti, jonka syötteet ovat jokin aakkoston  $\Sigma$  merkkijonoja, ja tulos on yksinkertaisesti muotoa ‘syöte hyväksytään’/‘syöte hylätään’. (Merk. lyhyesti 1/0.)

Merkitään  $M$ :n syötteellä  $x$  antamaa tulosta  $M(x)$ :llä ja  $M$ :n hyväksymien syötteiden joukkoaan  $A_M$ :llä, so.

$$A_M = \{x \in \Sigma^* \mid M(x) = 1\}.$$

Sanotaan, että automaatti  $M$  tunnistaa (engl. recognizes) kielentäytyvyyden  $A_M \subseteq \Sigma^*$ .

Automaattiteorian (yksi) idea: *automaatin  $M$  rakenne heijastuu kielentäytyvyydestä*.

Kääntääen: olkoon annettuna jokin toivottu I/O-kuvaus  $f : \Sigma^* \rightarrow \{0, 1\}$ . Tarkastelemalla kielstä

$$A_f = \{x \in \Sigma^* \mid f(x) = 1\}$$

saadaan vihjetä siitä, millainen automaatti tarvitaan kuvauksen  $f$  toteuttamiseen.

Harri Haanpää 2003-2004



**Vakiintuneita merkintöjä**

Em. matematisille käsitteille käytetyt merkinnät ovat periaatteessa vapaasti valittavissa, mutta esityksen ymmärrettävyyden parantamiseksi on tapana pitäätyä tietyissä käytännöissä. Seuraavat merkintätavat ovat vakiintuneet:

**Aakkostot:**  $\Sigma, \Gamma, \dots$  (isoja kreikkalaisia kirjaimia). Esim. binääriaakkosto  $\Sigma = \{0, 1\}$ .

**Aakkoston kokonaisuus (tai yleisemmin joukon mahtavuus):**  $|\Sigma|$ .

**Alkeismerkit:**  $a, b, c, \dots$  (pieniä alkupään latinalaisia kirjaimia). Esim.: Olkoon  $\Sigma = \{\alpha_1, \dots, \alpha_n\}$  aakkosto; tällöin  $|\Sigma| = n$ .

**Merkkijonot:**  $u, v, w, x, y, \dots$  (pieniä loppupään latinalaisia kirjaimia).

Aakkoston  $\Sigma$  kaikkien merkkijonojen joukko:  $\Sigma^*$ . Esim.:

$$\{a, b\}^* = \{\varepsilon, a, b, aa, ab, aaa, aab, \dots\}.$$

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Aakkoston  $\Sigma$  kaikkien merkkijonojen joukko merkitään  $\Sigma$ :lla.

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .

Mielivaltaista merkkijonojoukko  $A \subseteq \Sigma^*$  sanotaan aakkoston  $\Sigma$  (*formaaliksi kieleksi* (engl. formal language)).

Esimerkiksi jos  $\Sigma = \{0, 1\}$ , niin  $\Sigma^* = \{\varepsilon, 0, 1, 00, 01, 10, \dots\}$ .



## Merkkijonoinduktio

Automaattiteoriassa tehdään usein konstruktioita ”induktioilla merkkijonon pituuden suhteen.” Tämä tarkoittaa, että määritellään ensin toiminto tyhjän merkkijonon  $\varepsilon$  ( tai joskus yksittäisen aakkosmerkin) tapauksessa. Sitten oletetaan, että toiminto on määritelty kaikilla annetuilla pituisilla merkkijonoilla  $u$  ja esitetään, miten se täällöin määritellään yhtä merkkiä pitemmillä merkkijonoilla  $w = ua$ .

**Esimerkki.** Olkoon  $\Sigma$  mielivaltainen aakkosto. Merkkijonon  $w \in \Sigma^*$  *käänteisjono* (engl. reversal)  $w^R$  määritellään induktiivisesti säännöillä:

1.  $\varepsilon^R = \varepsilon$ ;
2. jos  $w = ua$ ,  $u \in \Sigma^*$ ,  $a \in \Sigma$ , niin  $w^R = a^R u^R$ .

Induktivista (“rekursiivista”) määritelmää voidaan tieteenkin käyttää laskujen perustana; esim.:

$$\begin{aligned} (011)^R &= 1^{\frown}(01)^R &= 1^{\frown}(1^{\frown}0^R) \\ &= 11^{\frown}(0^R\varepsilon^R) &= 110^{\frown}\varepsilon^R \\ &= 110^{\frown}\varepsilon &= 110. \end{aligned}$$

Tärkeämpää on kuitenkin konstruktioiden ominaisuuksien todistaminen määritelmää noudattelevalla induktiolla. Esimerkki:

1.  $\varepsilon^R = \varepsilon$ ;

2. jos  $w = ua$ ,  $u \in \Sigma^*$ ,  $a \in \Sigma$ , niin  $w^R = a^R u^R$ .

**Väite.** Olkoon  $\Sigma$  aakkosto. Kaikilla  $x, y \in \Sigma^*$  on voimassa  $(xy)^R = y^R x^R$ .

**Todistus.** Induktio merkkijonon  $y$  pituuden suhteen.

1. Perustapaus  $y = \varepsilon$ :  $(x\varepsilon)^R = x^R = \varepsilon^R x^R$ .
2. Induktioaskel: Olkoon  $y$  muotoa  $y = ua$ ,  $u \in \Sigma^*$ ,  $a \in \Sigma$ . Oletetaan, että väite on voimassa merkkijonoilla  $x, u$ . Tällöin on:

$$\begin{aligned} (xy)^R &= (xua)^R &[R:n määritelmä] \\ &= a^R(xu)^R &[\text{induktio-oleitus}] \\ &= a^R(u^R x^R) &[\text{in liitännäisyys}] \\ &= (a^R u^R)x^R &[R:n määritelmä] \\ &= (ua)^R x^R & \\ &= y^R x^R. \square \end{aligned}$$

$$X = \{x_0, x_1, x_2, \dots, x_{n-1}\},$$

jos  $X$  on  $n$ -alkoinen äärellinen joukko ja

$$X = \{x_0, x_1, x_2, \dots\},$$

jos  $X$  on numeroituvasti ääretön.

Numeroituvan joukon kaikki osajoukot ovat myös numeroituvia (HT), mutta ylinumeroituvilla joukoilla on sekä numeroituvia että ylinumeroitavia osajoukkoja. Siten ylinumeroituvat joukot ovat jossain mielessä “isompia” kuin numeroituvat.

**Lause 1.11** Minkä tahansa aakkoston  $\Sigma$  merkkijonojen joukko  $\Sigma^*$  on numerotuvasti ääretön.

**Todistus.** Muodostetaan bijektilo  $f : \mathbb{N} \rightarrow \Sigma^*$  seuraavasti.

Olkoon  $\Sigma = \{a_1, a_2, \dots, a_n\}$ . Kiinnitetään  $\Sigma$ :n merkeille jokin "aakkosjärjestys", olkoon se  $a_1 < a_2 < \dots < a_n$ .

Joukon  $\Sigma^*$  merkkijonot voidaan nyt luetella valitun aakkosjärjestyksen suhteen kanonisessa t. *leksikografisessa järjestyksessä* (engl. canonical t. *lexicographic order*) seuravasti:

- ▶ ensin luetellaan 0:n mittaiset merkkijonot ( $= \varepsilon$ ), sitten 1:n mittaiset ( $= a_1, a_2, \dots, a_n$ ), sitten 2:n mittaiset jne.;
- ▶ kunkin pituusryhmän sisällä merkkijonot luetellaan aakkosjärjestyksessä.

**Bijektilo  $f$  on siis:**

$$\begin{array}{ll} 0 & \mapsto \varepsilon \\ 1 & \mapsto a_1 \\ 2 & \mapsto a_2 \\ \vdots & \vdots \\ n & \mapsto a_n \\ n+1 & \mapsto a_1 a_1 \\ n+2 & \mapsto a_1 a_2 \\ \vdots & \vdots \\ 2n & \mapsto a_1 a_n \end{array}$$

□

**Lause 1.12** Minkä tahansa aakkoston  $\Sigma$  kaikkien formaalien kielten perhe on ylinumeroitava.

**Todistus** (ns. Cantorin diagonaaliargumentti). Merkitään aakkoston  $\Sigma$  kaikkien formaalien kielten perhettä  $\mathcal{P}(\Sigma^*) = \mathcal{A}$ . Tehdään vastaoleitus: oletetaan, että olisi olemassa kaikki  $\Sigma$ :n formaalit kielet kattava numeroointi:

$$\mathcal{A} = \{A_0, A_1, A_2, \dots\}.$$

Olkoot  $\Sigma^*$ :n merkkijonot kanonisessa järjestyksessä  $x_0, x_1, x_2, \dots$  Määritellään em. numeroointeja käyttäen formaali kielet  $\tilde{A}$ :

$$\tilde{A} = \{x_i \in \Sigma^* \mid x_i \notin A_i\}.$$

Koska  $\tilde{A} \in \mathcal{A}$  ja  $A$ :n numeroointi oletettiin kattavaksi, pitäisi olla  $\tilde{A} = A_k$  jollakin  $k \in \mathbb{N}$ . Mutta tällöin  $\tilde{A}$ :n määritelmän mukaan

$$x_k \in \tilde{A} \Leftrightarrow x_k \notin A_k = \tilde{A}.$$

Saatu ristiintää osoittaa, että vastaoletus on väärä. □

$\tilde{A}$	$A_0$	$A_1$	$A_2$	$A_3$	$\dots$
	1	0	0	1	$\dots$
$x_0$	0	1	0	0	$\dots$
$x_1$	0	0	1	0	$\dots$
$x_2$	1	1	1	1	$\dots$
$x_3$	0	0	0	0	$\dots$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

Kuvallisen todistuksen idea voidaan esittää seuraavasti.

Muodostetaan kielten  $A_0, A_1, A_2, \dots$  ja merkkijonojen  $x_0, x_1, x_2, \dots$ , "insidenssimatriisi", jonka rivin  $j$  on arvo 1 jos  $x_j \in A_j$  ja muuten 0. Tällöin kielii  $\tilde{A}$  polikkeaa kustakin kielestä  $A_k$  matriisin "diagonaallilla":

### 1.8 \*Ekskursio: Turingin pysähdytysongelma

Lauseiden 1.11 ja 1.12 mukaan on siis olemassa formaaleja kielitä (I/O-kuvauksia), joita ei voida toteuttaa esim. C-ohjelmilla. Entä jokin konkreettinen esimerkki tälläisestä?

Tunnetuin esimerkki on ns. *Turingin pysähdytysongelma*. (Alan Turing, 1936). C-ohjelma käyttääen tulos voidaan muotoilla seuraavasti:

**Väite.** Ei ole olemassa C-funktioita  $halt(p, x)$ , joka saa syötteenään mielivaltaisen C-funktion tekstin  $p$  ja tälle tarkoitetun syötteen  $x$  ja tuottaa tuloksen 1, jos  $p$ :n suoritus pysähyy syötteellä  $x$ , ja 0 jos  $p$ :n suoritus  $x$ :llä jää ikuiseen silmukkaan.

**Todistus.** Oletetaan väitteen vastaisesti, että tälläinen funktio  $halt$  voitaisiin laatia. Muodostetaan tästä käyttämän toinen funktio  $confuse$ .

Merkitään funktio  $confuse$  ohjelmateksttiä  $c$ :llä ja tarkastellaan funktion  $confuse$  laskentaa tällä omalla kuvauksellaan.

```
void confuse(char *p) {
    int halt (char *p, char *x) {
        confuse (c) pysähyy
        ... /* Funktion halt runko. */ * /
        ⇔
        halt (c, c) == 1
    }
    if (halt (p, p) == 1) while (1);
    ⇔
    confuse (c) ei pysähdy.
}
```

Ristiinräystä seuraa, että oletettua pysähymistestausfunktiota  $halt$  ei voi olla olemassa. □

Samansuksia ns. *rätteamattomia ongelmia* on itse asiassa *paljon*. Asiaan palataan kurssin loppupuolella.