

## 2.8 Säännöllisten kielten rajoituksista

Kardinaliteettisyistä on oltava olemassa (paljon) ei-säännöllisiä kieliä: kieliä on ylinumeroituva määrä, säännöllisiä lausekkeita vain numeroituvasti.

Voidaanko löytää konkreettinen, *mielenkiintoinen* esimerkki kielestä, joka ei olisi säännöllinen? Helposti.

Säännöllisten kielten perusrajoitus: äärellisillä automaateilla on vain rajallinen "muisti". Siten ne eivät pysty ratkaisemaan ongelmia, joissa vaaditaan mielivaltaisen suurten lukujen tarkkaa muistamista.

**Esimerkki:** sulkulausekekieli

$$L_{\text{match}} = \{(^k)^k \mid k \geq 0\}.$$

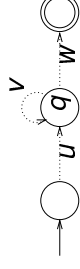
**Formalisointi:** "pumppauslemma".

**Lemma 2.6 (Pumppauslemma)** Olkoon  $A$  säännöllinen kieli. Tällöin on olemassa sellainen  $n \geq 1$ , että mikä tahansa  $x \in A$ ,  $|x| \geq n$ , voidaan jakaa osiin  $x = uvw$  siten, että  $|uv| \leq n$ ,  $|v| \geq 1$ , ja  $uv^i w \in A$  kaikilla  $i = 0, 1, 2, \dots$

**Todistus.** Olkoon  $M$  jokin  $A$ :n tunnistava deterministinen äärellinen automaatti, ja olkoon  $n$   $M$ :n tilojen määrä.

Tarkastellaan  $M$ :n läpikäymiä tiloja syötteellä  $x \in A$ ,  $|x| \geq n$ .

Koska  $M$  jokaisella  $x$ :n merkillä siirtyy tilasta toiseen, sen täytyy kulkea jonkin tilan kautta (ainakin) kaksi kertaa — itse asiassa jo  $x$ :n  $n$ :n ensimmäisen merkin aikana. Olkoon  $q$  ensimmäinen toistettu tila.



Olkoon  $u$   $M$ :n käsittelemä  $x$ :n alkiosa sen tullessa ensimmäisen kerran tilaan  $q$ ,  $v$  se osa  $x$ :stä jonka  $M$  käsittelee ennen ensimmäistä paluutaan  $q$ :hun, ja  $w$  loput  $x$ :stä. Tällöin on  $|uv| \leq n$ ,  $|v| \geq 1$ , ja  $uv^i w \in A$  kaikilla  $i = 0, 1, 2, \dots$   $\square$

**Esimerkki.** Tarkastellaan em. sulkulausekekieltä (merk. '( = a, ' ) = b):

$$L = L_{\text{match}} = \{a^k b^k \mid k \geq 0\}.$$

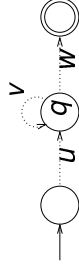
Oletetaan, että  $L$  olisi säännöllinen. Tällöin pitäisi pumppauslemman mukaan olla jokin  $n \geq 1$ , jota pitempiä  $L$ :n merkijonoja voidaan pumpata. Valitaan  $x = a^n b^n$ , jolloin  $|x| = 2n > n$ . Lemman mukaan  $x$  voidaan jakaa pumppattavaksi osiin  $x = uvw$ ,  $|uv| \leq n$ ,  $|v| \geq 1$ ; siis on oltava

$$u = a^i, v = a^j, w = a^{n-(i+j)} b^n, \quad i \leq n-1, j \geq 1.$$

Mutta esimerkiksi "0-kertaisesti" pumppaessa:

$$uv^0 w = a^i a^{n-(i+j)} b^n = a^{n-j} b^n \notin L.$$

Siten  $L$  ei voi olla säännöllinen.



### 3. KIELIOPIT JA MERKKIJONOJEN TUOTTAMINEN

Kielioppi = muunnossysteemi merkkijonojen (kielen "sanojen") tuottamiseen tietystä lähtöjonosta alkaen, osajonoja toistuvasti annettujen sääntöjen mukaan uudelleenkirjoittamalla.

Kielioppi on *yhteydetön*, jos kussakin uudelleenkirjoitusaskellessa korvataan yksi erityinen muuttuja t. *välisymboli* jollakin siihen liitettyä korvausjonolla, ja korvaus voidaan aina tehdä symbolia ympäröivän merkkijonon rakenteesta riippumatta.

Sovelluksia: rakenteisten tekstien kuvaaminen (esim. ohjelmointikielten BNF-syntaksikuvaukset, XML:n DTD/Schema-määrittelyt), yleisemmin rakenteisten "olioiden" kuvaaminen (esim. syntaktinen hahmontunnistus).

Yhteydetömillä kielioppeilla voidaan kuvata (tuottaa) myös ei-säännöllisiä kieliä.

*Esimerkki:* yhteydetön kielioppi kielelle  $L_{\text{match}}$  (lähtösymboli  $S$ ):

- (i)  $S \rightarrow \varepsilon$ ,
- (ii)  $S \rightarrow (S)$ .

Esimerkiksi merkkijonon  $((()))$  tuottaminen:

$$S \Rightarrow (S) \Rightarrow ((S)) \Rightarrow (((S))) \Rightarrow (((\varepsilon))) = (((()))).$$

*Toinen esimerkki:* kielioppi C-tyyppisen ohjelmointikielen aritmeettisille lausekkeille (yksinkertaistettu).

$$\begin{array}{l} E \rightarrow T \mid E + T \\ T \rightarrow F \mid T * F \\ F \rightarrow a \mid (E). \end{array}$$

Esimerkiksi lausekkeen  $(a + a) * a$  tuottaminen:

$$\begin{array}{l} E \Rightarrow \underline{I} \Rightarrow \underline{I} * F \Rightarrow \underline{I} * F \\ \Rightarrow (\underline{E}) * F \Rightarrow (\underline{E} + T) * F \Rightarrow (\underline{I} + T) * F \\ \Rightarrow (\underline{E} + T) * F \Rightarrow (\underline{a} + \underline{I}) * F \Rightarrow (\underline{a} + \underline{E}) * F \\ \Rightarrow (\underline{a} + \underline{a}) * \underline{E} \Rightarrow (\underline{a} + \underline{a}) * a. \end{array}$$

### Määritelmä 3.1 Yhteydetön kielioppi on nelikko

$$G = (V, \Sigma, P, S),$$

missä

- ▶  $V$  on kieliopin aakkosto;
- ▶  $\Sigma \subseteq V$  on kieliopin päätemerkkien joukko; sen komplementti  $N = V - \Sigma$  on kieliopin välikemerkkien t. *-symbolien* joukko;
- ▶  $P \subseteq N \times V^*$  on kieliopin sääntöjen t. *produktioiden* joukko;
- ▶  $S \in N$  on kieliopin lähtösymboli.

Produktiota  $(A, \omega) \in P$  merkitään tavallisesti  $A \rightarrow \omega$ .

Merkkijono  $\gamma \in V^*$  tuottaa t. johtaa suoraan merkkijonon  $\gamma' \in V^*$  kieliopissa  $G$ , merkitään

$$\gamma \xRightarrow{G} \gamma'$$

jos voidaan kirjoittaa  $\gamma = \alpha A \beta$ ,  $\gamma' = \alpha \omega \beta$  ( $\alpha, \beta, \omega \in V^*$ ,  $A \in N$ ), ja kieliopissa  $G$  on produktio  $A \rightarrow \omega$ .

Jos kielioppi  $G$  on yhteydestä selvä, voidaan merkitä  $\gamma \Rightarrow \gamma'$ .

Merkkijono  $\gamma \in V^*$  tuottaa t. johtaa merkkijonon  $\gamma' \in V^*$  kieliopissa  $G$ , merkitään

$$\gamma \xRightarrow{G^*} \gamma'$$

jos on olemassa jono  $V$ :n merkkijonoja  $\gamma_0, \gamma_1, \dots, \gamma_n$  ( $n \geq 0$ ), siten että

$$\gamma = \gamma_0 \xRightarrow{G} \gamma_1 \xRightarrow{G} \dots \xRightarrow{G} \gamma_n = \gamma'$$

Erikoistapauksena  $n = 0$  saadaan  $\gamma \xRightarrow{G^*} \gamma$  millä tahansa  $\gamma \in V^*$ . Jälleen, jos  $G$  on yhteydestä selvä, voidaan merkitä  $\gamma \Rightarrow^* \gamma'$ .

Esimerkiksi tasapainoisten sulkujonojen muodostaman kielen  $L_{\text{match}} = \{(^k)^k \mid k \geq 0\}$  tuottaa kielioppi

$$G_{\text{match}} = (\{S, (\cdot)\}, \{(\cdot)\}, \{S \rightarrow \varepsilon, S \rightarrow (S)\}, S).$$

Yksinkertaisten aritmeettisten lausekkeiden muodostaman kielen  $L_{\text{expr}}$  tuottaa kielioppi

$$G_{\text{expr}} = (V, \Sigma, P, E),$$

missä

$$V = \{E, T, F, a, +, *, (\cdot)\},$$

$$\Sigma = \{a, +, *, (\cdot)\},$$

$$P = \{E \rightarrow T, E \rightarrow E + T, T \rightarrow F, T \rightarrow T * F, F \rightarrow a, F \rightarrow (E)\}.$$

Merkkijono  $\gamma \in V^*$  on kieliopin  $G$  lausejohdos, jos on  $S \xRightarrow{G^*} \gamma$ .

Peikästään päätämerkeistä koostuva  $G$ :n lausejohdos  $x \in \Sigma^*$  on  $G$ :n lause.

Kieliopin  $G$  tuottama t. kuvaama kieli koostuu  $G$ :n lauseista:

$$L(G) = \{x \in \Sigma^* \mid S \xRightarrow{G^*} x\}.$$

Formaali kieli  $L \subseteq \Sigma^*$  on yhteydetön, jos se voidaan tuottaa jollakin yhteydettömällä kieliopilla.

Toinen kieliooppi kielen  $L_{\text{expr}}$  tuottamiseen on

$$G_{\text{expr}} = (V, \Sigma, P, E),$$

missä

$$V = \{E, a, +, *, (, )\},$$

$$\Sigma = \{a, +, *, (, )\},$$

$$P = \{E \rightarrow E + E, E \rightarrow E * E, E \rightarrow a, E \rightarrow (E)\}.$$

*Huom:* Vaikka kieliooppi  $G_{\text{expr}}$  näyttää yksinkertaisemmalta kuin kieliooppi  $G_{\text{expr}}$ , sen ongelmana on ns. rakenteellinen moniselitteisyys, mikä on monesti ei-toivottu ominaisuus.

### Vakiintuneita merkintätapoja

Välikesymboleita:  $A, B, C, \dots, S, T$ .

Päätemerkkejä: kirjaimet  $a, b, c, \dots, s, t$ ;  
numerot  $0, 1, \dots, 9$ ;  
erikoismerkit; lihavoituidut tai alleviivatut varatut sanat (**if**, **for**, **end**,  $\dots$ ).

Mielivaltaisia merkkejä (kun välitteitä ja päätteitä ei erotella):  $X, Y, Z$ .

Päätemerkkijonoja:  $u, v, w, x, y, z$ .

Sekamerkkijonoja:  $\alpha, \beta, \gamma, \dots, \omega$ .

Toinen kieliooppi kielen  $L_{\text{expr}}$  tuottamiseen on

$$G_{\text{expr}} = (V, \Sigma, P, E),$$

missä

$$V = \{E, a, +, *, (, )\},$$

$$\Sigma = \{a, +, *, (, )\},$$

$$P = \{E \rightarrow E + E, E \rightarrow E * E, E \rightarrow a, E \rightarrow (E)\}.$$

*Huom:* Vaikka kieliooppi  $G_{\text{expr}}$  näyttää yksinkertaisemmalta kuin kieliooppi  $G_{\text{expr}}$ , sen ongelmana on ns. rakenteellinen moniselitteisyys, mikä on monesti ei-toivottu ominaisuus.

Produktiot, joilla on yhteinen vasen puoli  $A$ , voidaan kirjoittaa yhteen: joukon

$$A \rightarrow \omega_1, A \rightarrow \omega_2, \dots, A \rightarrow \omega_k$$

sijaan kirjoitetaan

$$A \rightarrow \omega_1 \mid \omega_2 \mid \dots \mid \omega_k.$$

Kieliooppi esitetään usein pelkkänä sääntöjoukkona:

$$A_1 \rightarrow \omega_{11} \mid \dots \mid \omega_{1k_1}$$

$$A_2 \rightarrow \omega_{21} \mid \dots \mid \omega_{2k_2}$$

$\vdots$

$$A_m \rightarrow \omega_{m1} \mid \dots \mid \omega_{mk_m}.$$

Tällöin päätellään välikesymbolit edellisten merkintäsopimusten mukaan tai siinä, että ne esiintyvät sääntöjen vasempina puolina; muut esiintyvät merkit ovat päätemerkkejä.

*Lähtösymboli* on tällöin *ensimmäisen säännön vasempana puolena* esiintyvä välike; tässä siis  $A_1$ .

### Eräitä konstruktioita

Olkoon  $L(T)$  välikkeestä  $T$  johdettavissa olevien päätejonojen joukko. Olkoon meillä produktiokoeelma  $P$ , jossa ei esiinny välikettä  $A$ , ja jolla  $B$ :stä voidaan johtaa  $L(B)$  (ja vastaavasti  $C$ :stä  $L(C)$ ).

Lisäämällä  $P$ :hen jokin seuraavista produktioista saadaan uusia kieliä:

produktio	kieli
$A \rightarrow B \mid C$	yhdiste $L(A) = L(B) \cup L(C)$
$A \rightarrow BC$	katenaatio $L(A) = L(B)L(C)$ , ja
$A \rightarrow AB \mid \epsilon$ (vasen rekursio) tai	Kleenen sulkeuma $L(A) = L(B)^*$
$A \rightarrow BA \mid \epsilon$ (oikea rekursio)	

Välikkeiden keskeisyypotus on yhteydettömille kieliopeille ominainen konstruktio, joka tekee usein (muttei aina) kielestä epääännöllisen: lisäämällä produktio

$A \rightarrow BAC \mid \epsilon$  saadaan

$$L(A) = \bigcup_{i=0}^{\infty} L(B)^i L(C)^i.$$