

## 2.8 Säännöllisten kielten rajoituksista

Kardinaliteettisyydestä 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”.



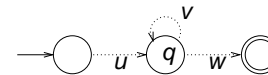
**Lemma 2.6 (Pumppauslemma)** Olkoon  $A$  säännöllinen kieli.

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 alkuosa 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  $i = 0, 1, 2, \dots$  □



**Esimerkki.** Tarkastellaan em. sulkulausekekieltä (merk. ‘(’ =  $a$ , ‘)’ =  $b$ ):

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

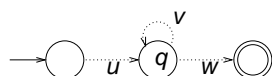
Oletetaan, että  $L$  olisi säännöllinen. Tällöin pitäisi pumppauslemman mukaan olla jokin  $n \geq 1$ , jota pitempiä  $L$ :n merkkijonoja voidaan pumpata. Valitaan  $x = a^n b^n$ , jolloin  $|x| = 2n > n$ . Lemman mukaan  $x$  voidaan jakaa pumpattavaksi osiin  $x = uvw$ ,  $|uv| \leq n$ ,  $|v| \geq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

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

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 uudelleenkirjoitusaskeleessa korvataan yksi erityinen muuttujat. *välikesymboli* jollakin siihen liitetyllä 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).

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

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

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

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

Yhteydettömillä kieliopeilla 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:

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

**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 *välimerkkien* 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

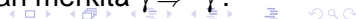
$$\gamma \xrightarrow{G}^* \gamma'$$

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  $\gamma \Rightarrow^* \gamma'$ .



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

Pelkästään päätemerkeistä 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.



Esimerkiksi tasapainoisten sulkujonojen muodostaman kielen

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

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

Yksinkertaisten aritmeettisten lausekkeiden muodostaman

kielen  $L_{\text{expr}}$  tuottaa kielioppi

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

missä

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

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

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



Toinen kielioppi 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 kielioppi  $G'_{\text{expr}}$  näyttää yksinkertaisemmalta kuin kielioppi  $G_{\text{expr}}$ , sen ongelmana on ns. rakenteellinen moniselitteisyys, mikä on monesti ei-toivottu ominaisuus.



**Vakiintuneita merkintätapoja**

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

Päätemerkkejä: kirjaimet  $a, b, c, \dots, s, t$ ;

numerot  $0, 1, \dots, 9$ ;

erikoismerkit; lihavoidut tai alleviivatut varatut sanat (**if, for, end, ...**).

Mielivaltaisia merkkejä (kun välikkeitä ja päätteitä ei erotella):

$X, Y, Z$ .

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

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

**Eräitä konstruktioita**

Olkoon  $L(T)$  välikkeestä  $T$  johdettavissa olevien päättejonon joukko. Olkoon annettu produktiokokoelma  $P$ , 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
$A \rightarrow AB \mid \varepsilon$ (vasen rekursio) tai $A \rightarrow BA \mid \varepsilon$ (oikea rekursio)	Kleenen sulkeuma $L(A) = L(B)^*$



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

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

Kielioppi 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 välikeyholeit edellisten merkintäsopimusten mukaan tai siitä, että ne esiintyvät sääntöjen vasempina puolina; muut esiintyvät merkit ovat päätemerkkejä.

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



Välikeyholeiden *keskeisupotus* on yhteydettömille kieliopille ominainen konstruktio, joka tekee usein (muttei aina) kielestä epä säännöllisen: lisäämällä produktio

$A \rightarrow BAC \mid \varepsilon$  saadaan

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



### 3.3 KIELIOPIEN JÄSENNYSONGELMA

Ratkaistava tehtävä:

“Annettu yhteydetön kielioppi  $G$  ja merkkijono  $x$ . Onko  $x \in L(G)$ ?”

Ratkaisumenetelmä = *jäsennysalgoritmi*.

Useita vaihtoehtoisia menetelmiä, erityisesti kun  $G$  on jotain rajoitettua (käytännössä esiintyvää) muotoa.

Johto  $\gamma \Rightarrow^* \gamma'$  on *vasen johto*, merkitään

$$\gamma \underset{\text{lm}}{\Rightarrow^*} \gamma',$$

jos kussakin johtoaskelella on produktiota sovellettu merkkijonon vasemmanpuoleisimpaan välikkeeseen (edellä johto (i)).

Vastaavasti määritellään *oikea johto* (edellä (iii)), jota merkitään

$$\gamma \underset{\text{rm}}{\Rightarrow^*} \gamma'$$

Suuria vasempia ja oikeita johtoaskelia merkitään  $\gamma \underset{\text{lm}}{\Rightarrow} \gamma'$  ja

$$\gamma \underset{\text{rm}}{\Rightarrow} \gamma'.$$

### Johdot ja jäsenyspuut

Olkon  $\gamma \in V^*$  kieliopin  $G = (V, \Sigma, P, S)$  lausejohdos.

Lähtösymbolista  $S$  merkkijonoon  $\gamma$  johtavaa suorien johtojen jonoa

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \dots \Rightarrow \gamma_n = \gamma$$

sanotaan  $\gamma$ :n *johdoksi*  $G$ :ssä.

Johdon *pituus* on siihen kuuluvien suorien johtojen määrä (edellä  $n$ ).

Esimerkki: lauseen  $a + a$  johtoja kieliopissa  $G_{\text{expr}}$ :

$$\begin{aligned} \text{(i)} \quad E &\Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \\ &\Rightarrow a + T \Rightarrow a + F \Rightarrow a + a \\ \text{(ii)} \quad E &\Rightarrow E + T \Rightarrow E + F \Rightarrow T + F \\ &\Rightarrow F + F \Rightarrow F + a \Rightarrow a + a \\ \text{(iii)} \quad E &\Rightarrow E + T \Rightarrow E + F \Rightarrow E + a \\ &\Rightarrow T + a \Rightarrow F + a \Rightarrow a + a. \end{aligned}$$

Olkon  $G = (V, \Sigma, P, S)$  yhteydetön kielioppi.

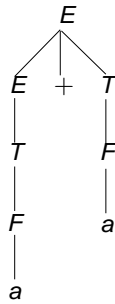
Kieliopin  $G$  mukainen *jäsennyspuu* on järjestetty puu, jolla on seuraavat ominaisuudet:

(i) puun solmut on nimetty joukon  $V \cup \{\varepsilon\}$  alkioilla siten, että sisäsolmujen nimet ovat välikkeitä (so. joukosta  $N = V - \Sigma$ ) ja juurisolmun nimenä on lähtösymboli  $S$ ;

(ii) jos  $A$  on puun jonkin sisäsolmun nimi, ja  $X_1, \dots, X_k$  ovat sen jälkeläisten nimet järjestyksessä, niin  $A \rightarrow X_1 \dots X_k$  on  $G$ :n produktio.

Jäsennyspuun  $\tau$  *tuotos* on merkkijono, joka saadaan liittämällä yhteen sen lehtisolmujen nimet esijärjestyksessä (“vasemmalta oikealle”).

**Esimerkki.** Lauseen  $a + a$  jäsennympuu kieliopissa  $G_{\text{expr}}$ :



Lauseen johto:

$$\begin{aligned} E &\Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \\ &\Rightarrow a + T \Rightarrow a + F \Rightarrow a + a \end{aligned}$$



Johtoa

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \dots \Rightarrow \gamma_n = \gamma$$

vastaavan jäsennympuun muodostaminen:

(i) puun juuren nimeksi tulee  $S$ ; jos  $n = 0$ , niin puussa ei ole muita solmuja; muuten

(ii) jos ensimmäisessä johtoaskelella on sovellettu produktiota  $S \rightarrow X_1 X_2 \dots X_k$ , niin juurelle tulee  $k$  jälkeläissolmua, joiden nimet vasemmalta oikealle ovat

$$X_1, X_2, \dots, X_k;$$

(iii) jos seuraavassa askelella on sovellettu produktiota  $X_i \rightarrow Y_1 Y_2 \dots Y_l$ , niin juuren  $i$ :nnelle jälkeläissolmulle tulee  $l$  jälkeläistä, joiden nimet vasemmalta oikealle ovat  $Y_1, Y_2, \dots, Y_l$ ; ja niin edelleen.

Konstruktioista huomataan, että jos  $\tau$  on jotakin johtoa  $S \Rightarrow^* \gamma$  vastaava jäsennympuu, niin  $\tau$ :n tuotos on  $\gamma$ .



Olkoon  $\tau$  kieliopin  $G$  mukainen jäsennympuu, jonka tuotos on päätimerkkijono  $x$ .

Tällöin  $\tau$ :sta saadaan vasen johto  $x$ :lle käymällä puun solmut läpi esijärjestyksessä ("ylhäältä alas, vasemmalta oikealle") ja laventamalla vastaan tulevat välitteet järjestyksessä puun osoittamalla tavalla.

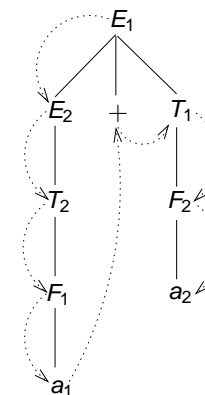
Oikea johto saadaan käymällä puu läpi käänteisessä esijärjestyksessä ("ylhäältä alas, oikealta vasemmalle").

Muodostamalla annetusta vasemmasta johdosta  $S \Rightarrow^* x$  ensin jäsennympuu edellä esitetyllä tavalla, ja sitten jäsennympuusta vasen johto, saadaan takaisin alkuperäinen johto; vastaava tulos pätee myös oikeille johdoille.



**Esimerkki.** Lauseen  $a + a$  vasemman johdon muodostaminen jäsennympuusta.

Jäsennympuu:



Solmut esijärjestyksessä:

$$E_1 E_2 T_2 F_1 a_1 + T_1 F_2 a_2$$

Vasen johto:

$$\begin{aligned} E &\Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \\ &\stackrel{\text{lm}}{\Rightarrow} a + T \stackrel{\text{lm}}{\Rightarrow} a + F \stackrel{\text{lm}}{\Rightarrow} a + a \end{aligned}$$



**Lause 3.3** Olkoon  $G = (V, \Sigma, P, S)$  yhteydetön kielioppi. Tällöin:

(i) jokaisella  $G$ :n lausejohdoksella  $\gamma$  on  $G$ :n mukainen jäsennyspuu  $\tau$ , jonka tuotos on  $\gamma$ ;

(ii) jokaista  $G$ :n mukaista jäsennyspuuta  $\tau$ , jonka tuotos on päättemerkkijono  $x$ , vastaavat yksikäsitteiset vasen ja oikea johto  $S \xRightarrow{lm}^* x$  ja  $S \xRightarrow{rm}^* x$ .

**Seuraus 3.4** Jokaisella  $G$ :n lauseella on vasen ja oikea johto.

Siis: yhteydetön kieliopin tuottamien lauseiden jäsennyspuut, vasemmat ja oikeat johdot vastaavat yksikäsitteisesti toisiaan.

Jäsennysongelman ratkaisuun katsotaan usein kuuluvan pelkän päätösongelman "Onko  $x \in L(G)$ ?" ratkaisemisen lisäksi jonkin näistä jäsennysesityksistä tuottaminen.

Moniselitteisyys on tietojenkäsittelysovelluksissa yleensä ei-toivottu ominaisuus, koska se merkitsee että annetulla lauseella on kaksi vaihtoehtoista "tulkintaa."

Yhteydetön kieli, jonka tuottavat kieliopit ovat kaikki moniselitteisiä, on *luonnostaan moniselitteinen*.

Esimerkiksi kielioppi  $G'_{\text{expr}}$  on moniselitteinen, kieliopit  $G_{\text{expr}}$  ja  $G_{\text{match}}$  yksiselitteisiä. Kieli  $L_{\text{expr}} = L(G'_{\text{expr}})$  ei ole luonnostaan moniselitteinen, koska sillä on myös yksiselitteinen kielioppi  $G_{\text{expr}}$ . Luonnostaan moniselitteinen on esimerkiksi kieli

$$\{a^i b^j c^k \mid i = j \text{ tai } j = k\}.$$

(Todistus sivuutetaan.)

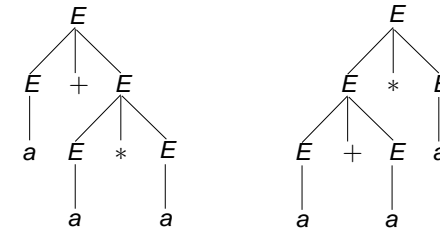
### Kieliopin moniselitteisyys

Lauseella voi olla kieliopissa useita jäsennyksiä.

**Esimerkki.** Tarkastellaan yksinkertaisten aritmeettisten lausekkeiden kielioppia:

$$G'_{\text{expr}} = \{E \rightarrow E + E, E \rightarrow E * E, E \rightarrow a, E \rightarrow (E)\}.$$

Lauseella  $a + a * a$  on tässä kieliopissa kaksi jäsennystä:



Yhteydetön kielioppi  $G$  on *moniselitteinen*, jos jollakin  $G$ :n lauseella  $x$  on kaksi erilaista  $G$ :n mukaista jäsennyspuuta. Muuten kielioppi on *yksiselitteinen*.