

Luku 4

Kontekstittomat kielet ja pinoautomaatit



- Kuten edellä nähtiin, esim. tasapainoisten sulkulausekkeiden muodostama kieli

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

if-else-parien muodostama kieli

$$L_{\text{if-else}} = \{\text{if}^k \text{else}^l \mid l \leq k\}$$

eivät ole säännöllisiä (Pumppauslemma)

- Mutta näitä kieliä voi kuvata *kontekstittomilla kieliopilla*
- Kuten edellä säännöllisiä kieliä tunnistava kone oli äärellinen automaatti, kontekstitonta kieltä tunnistava kone on *pinoautomaatti*

4.1 Kontekstittomat kieliopit ja kielet

- Esim. sulkulausekekielellä yksinkertainen rekursiivinen kuvaus: Merk. $S =$ “mielivaltainen tasapainoinen sulkumerkkijono”. Tällöin S on tasapainoinen sulkumerkkijono, jos
 - $S = \epsilon$ tai
 - S on muotoa (S') , missä S' on tasapainoinen sulkumerkkijono
- Toinen kuvaustapa: seuraavat *muunnossäännöt tuottavat* täsmälleen kielen L_{match} merkkijonot symbolista S :
 - $S \rightarrow \epsilon$,
 - $S \rightarrow (S)$
- Esim. merkkijonon $((()))$ tuottaminen:

$$S \Rightarrow (S) \Rightarrow ((S)) \Rightarrow (((S))) \Rightarrow (((\epsilon))) = ((()))$$

- Kontekstiton kielioppi on muunnossysteemi, jossa kuvattavat merkkijonot tuotetaan korvaamalla erityisiä muuttuja- t. *välikesymboleita* annettujen sääntöjen mukaan yksi kerrallaan, symbolia ympäröivän merkkijonon rakenteesta riippumatta
- Merkinnällä

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

lyhennetään

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

- Esim. Yksinkertainen kielioppi tietyille aritmeettisille lausekkeille:

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

Esim. lausekkeen $(a + a) * a$ tuottaminen:

$$\begin{array}{lclclcl}
 \underline{E} & \Rightarrow & \underline{T} & \Rightarrow & \underline{T} * F & \Rightarrow & \underline{F} * F \\
 & \Rightarrow & (\underline{E}) * F & \Rightarrow & (\underline{E} + T) * F & \Rightarrow & (\underline{T} + T) * F \\
 & \Rightarrow & (\underline{F} + T) * F & \Rightarrow & (a + \underline{T}) * F & \Rightarrow & (a + \underline{F}) * F \\
 & \Rightarrow & (a + a) * \underline{F} & \Rightarrow & (a + a) * a & &
 \end{array}$$

Määritelmä: Kontekstiton kielioppi (engl. context-free grammar) on nelikko

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

missä

- V on kieliopin aakkosto;
- $\Sigma \subseteq V$ on kieliopin *päätemerkkien* joukko; sen komplementti $N = V \setminus \Sigma$ 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 $A \rightarrow \omega$

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

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

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

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

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

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

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

- Erikoistapaus: $n = 0$, $\gamma \xRightarrow{G}^* \gamma$ millä tahansa $\gamma \in V^*$

- Merkkijono $\gamma \in V^*$ on kieliopin G lausejohdos, jos on $S \xRightarrow[G]{*} \gamma$
- G :n lause on pelkästään päätemerkeistä koostuva G :n lausejohdos $x \in \Sigma^*$
- Kieliopin G tuottama t. kuvaama kieli koostuu G :n lauseista:

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

Määritelmä: Formaali kieli $L \subseteq \Sigma^*$ on *kontekstitton*, jos se voidaan tuottaa jollakin kontekstittomalla kieliopilla

- Esim. Tasapainoisten sulkujonojen muodostaman kielen $L_{\text{match}} = \{(^k)^k \mid k \geq 0\}$ tuottaa kielioppi

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

- Esim. Aiemmin tarkisteltujen yksinkertaisten aritmeettisten lausekkeiden muodostaman kielen L_{expr} tuottaa kielioppi

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

missä

$$\begin{aligned} V &= \{E, T, F, a, +, *, (,)\}, \\ \Sigma &= \{a, +, *, (,)\}, \\ P &= \{E \rightarrow T, E \rightarrow E + T, T \rightarrow F, T \rightarrow T * F, \\ &F \rightarrow a, F \rightarrow (E)\}. \end{aligned}$$

Toinen kielioppi kielen L_{expr} tuottamiseen on

$$G'_{\text{expr}} = (V, \Sigma, P, E),$$

missä

$$\begin{aligned} V &= \{E, a, +, *, (,)\}, \\ \Sigma &= \{a, +, *, (,)\}, \\ P &= \{E \rightarrow E + E, E \rightarrow E * E, E \rightarrow a, E \rightarrow (E)\} \end{aligned}$$

- Esim. (\sim Orponen teht.44) Tarkastellaan suomen kielen virkettä, joka koostuu yksinkertaisesta pääauseesta + 0 tai useammasta sisäkkäisestä relatiivilauseesta:

$$L_{rel} = \{subj(\text{joka } pred \text{ attr } obj)^* pred \text{ attr } obj\}$$

Tällaisia virkkeitä voidaan tuottaa esim. seuraavilla kontekstittoman kieliopin G_{rel} säännöillä:

(Esitetään yksinkertaisuuden vuoksi säännöt merkkijonoilla)

VIRKE \rightarrow SUBJ SL PRED ATTR OBJ

SL \rightarrow joka PRED ATTR OBJ SL | ϵ

SUBJ \rightarrow poika | tyttö | jänis | susi | peikko

PRED \rightarrow pelkäsi | metsästi

ATTR \rightarrow suurta | pientä | vihaista | hirmuista | arkaa

OBJ \rightarrow poikaa | tyttöä | jänistä | sutta | peikkoa

Kielen kuuluvat mm. seuraavat virkkeet:

- poika joka metsästi sutta joka pelkäsi peikkoa joka pelkäsi suurta tyttöä pelkäsi hirmuista jänistä
- tyttö joka metsästi arkaa poikaa pelkäsi vihaista jänistä

Vakiintuneita merkintätapoja

- Välikesymboleita: 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 välitteitä ja päätteitä ei erotella): X, Y, Z
- Päätemerkkijonoja: u, v, w, x, y, z
- Sekamerkkijonoja: $\alpha, \beta, \gamma, \dots, \omega$

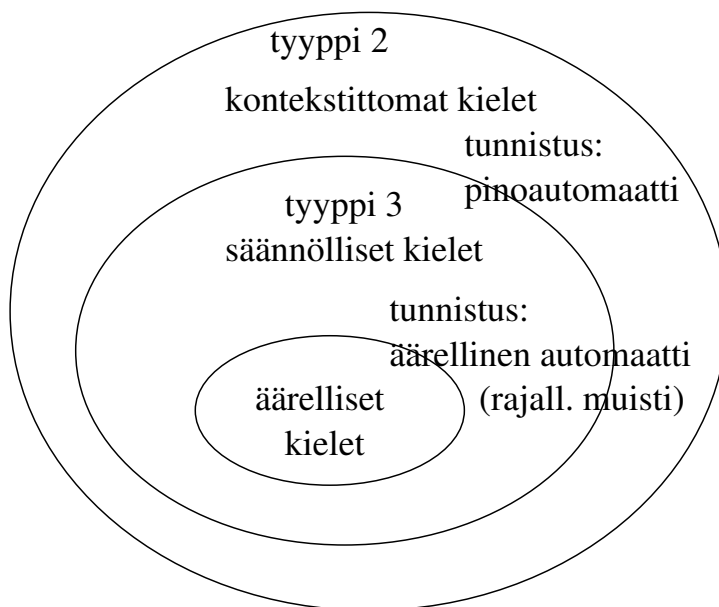
- Kielioppi esitetään usein pelkkänä sääntöjoukkona:

$$\begin{array}{l} A_1 \rightarrow \omega_{11} \quad | \quad \dots \quad | \quad \omega_{1k_1} \\ A_2 \rightarrow \omega_{21} \quad | \quad \dots \quad | \quad \omega_{2k_2} \\ \vdots \\ A_m \rightarrow \omega_{m1} \quad | \quad \dots \quad | \quad \omega_{mk_m} \end{array}$$

- Tällöin päätellään välikesymbolit edellisten merkintäsopimusten mukaan tai siitä, että ne esiintyvät sääntöjen vasempina puolina; muut esiintyvät merkit ovat päätemerkkejä
- *Lähtösymboli* on tällöin *ensimmäisen säännön vasempana puolena* esiintyvä välike; tässä siis A_1

4.1.1 Säännölliset kielet ja kontekstittomat kieliopit

- Kontekstittomilla kieliopeilla voidaan kuvata joitakin ei-säännöllisiä kieliä (esimerkiksi kielet L_{match} ja L_{expr})
- Osoitetaan, että myös kaikki säännölliset kielet voidaan kuvata kontekstittomilla kieliopeilla
- Kontekstittomat kielet ovat siten säännöllisten kielten aito ylikuokka



4.1.2 Oikealle ja vasemmalle lineaariset kieliopit

- *Määritelmä:* Kontekstiton kielioppi on *oikealle lineaarinen*, jos sen kaikki produktiot ovat muotoa $A \rightarrow \epsilon$ tai $A \rightarrow aB$, ja *vasemmalle lineaarinen*, jos sen kaikki produktiot ovat muotoa $A \rightarrow \epsilon$ tai $A \rightarrow Ba$
- Osoittautuu, että sekä vasemmalle että oikealle lineaarisilla kielioppeilla voidaan tuottaa täsmälleen säännölliset kielet
- \Rightarrow lineaarisia kielioppeja nimitetään myös yhteisesti *säännöllisiksi* kielioppeiksi
- Todistetaan tässä väite vain oikealle lineaarisille kielioppeille:

Lause: Jokainen säännöllinen kieli voidaan tuottaa oikealle lineaarisella kieliopilla. Todistus: Olkoon L aakkoston Σ säännöllinen kieli, ja olkoon $M = (Q, \Sigma, \delta, q_0, F)$ sen tunnistava (deterministinen tai epädeterministinen) äärellinen automaatti. Muodostetaan kielioppi G_M , jolla on $L(G_M) = L(M) = L$.

- G_M :n pääteakkosto = M :n syöteakkosto Σ
- G_M :n välikeakkostoon otetaan yksi välike A_q kutakin M :n tilaa q kohden.
- G_M :n lähtösymboli on A_{q_0}
- G_M :n produktiot vastaavat M :n siirtymiä:
 - (i) kutakin M :n lopputilaa $q \in F$ kohden kielioppiin otetaan produktio $A_q \rightarrow \epsilon$;
 - (ii) kutakin M :n siirtymää $q \xrightarrow{a} q'$ (so. $q' \in \delta(q, a)$) kohden kielioppiin otetaan produktio $A_q \rightarrow aA_{q'}$
- Tarkastetaan konstruktion oikeellisuus:
 - Merk. A_q :sta tuotettavien päätejonojen joukkoa

$$L(A_q) = \{x \in \Sigma^* \mid A_q \xRightarrow{G_M}^* x\}$$

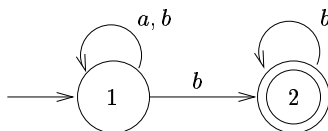
- Induktiolla merkkijonon x pituuden suhteen voidaan osoittaa, että kaikilla q on

$$x \in L(A_q) \Leftrightarrow (q, x) \vdash_M^* (q_f, \epsilon) \text{ jollakin } q_f \in F$$

– Erityisesti on siis

$$\begin{aligned} L(G_M) = L(A_{q_0}) &= \{x \in \Sigma^* \mid (q_0, x) \vdash_M^* (q_f, \epsilon) \\ &\quad \text{jollakin } q_f \in F\} \\ &= L(M) = L. \quad \square \end{aligned}$$

Esim.



Kuva 4.1: Yksinkertainen äärellinen automaatti.

Automaattia vastaava kielioppi on:

$$\begin{aligned} A_1 &\rightarrow aA_1 \mid bA_1 \mid bA_2 \\ A_2 &\rightarrow \epsilon \mid bA_2 \end{aligned}$$

Lause: Jokainen oikealle lineaarisella kieliopilla tuotettava kieli on säännöllinen.

Todistus: Olkoon $G = (V, \Sigma, P, S)$ oikealle lineaarinen kielioppi. Muodostetaan kielen $L(G)$ tunnistava epäterministinen äärellinen automaatti $M_G = (Q, \Sigma, \delta, q_S, F)$ seuraavasti:

- M_G :n tilat vastaavat G :n välitteitä:

$$Q = \{q_A \mid A \in V - \Sigma\}$$

- M_G :n alkutila on lähtösymbolia S vastaava tila q_S
- M_G :n syöteaakkosto on G :n pääteaakkosto Σ
- M_G :n siirtymäfunktio δ jäljittelee G :n produktioita siten, että kutakin produktiota $A \rightarrow aB$ kohden automaatissa on siirtymä $q_A \xrightarrow{a} q_B$ (so. $q_B \in \delta(q_A, a)$)
 M_G :n lopputiloja ovat ne tilat, joita vastaaviin välitteisiin liittyy G :ssä ϵ -produktio:

$$F = \{q_A \in Q \mid A \rightarrow \epsilon \in P\}$$

- Konstruktion oikeellisuus voidaan jälleen tarkastaa induktiolla G :n tuottamien ja M_G :n hyväksymien merkkijonojen pituuden suhteen. \square

Esimerkkejä

Esimerkkejä säännöllisistä lausekkeista ja vastaavista kontekstittomista kieliopeista:

1. Lauseke: a^*
Kielioppi: $S \rightarrow aS|\epsilon$
2. Lauseke: $a^+ = aa^*$
Kielioppi: $S \rightarrow aS|a$
3. Lauseke: $(aa)^*$
Kielioppi: $S \rightarrow aSa|\epsilon$
(merkkijono koostuu vain parillisesta määrästä a :ta)
4. Lauseke: $(b^*ab^*ab^*)^*$
Kielioppi:
 $S \rightarrow BaBaBS|\epsilon$
 $B \rightarrow bB|\epsilon$
(merkkijono sisältää parillisen määrän a :ta, lisäksi saa olla b :tä missä tahansa)
5. Lauseke: $(0 \cup 1 \cup \dots \cup 9)(0 \cup 1 \cup \dots \cup 9)^*$
Kielioppi:
 $S \rightarrow DN$
 $N \rightarrow DN|\epsilon$
 $D \rightarrow 0|1|\dots|9$
(vähintään yhdestä digitistä koostuva numero)

4.2 *Ekskursio: Kasvikieliopit (Lindenmayer Systems)

Esimerkki 1 Olkoon symbolit $\Sigma = \{SIEMEN, SIRKKALEHDET, LEHDET, VARSI, NUPPU, PUNKUKKA, SINKUKKA\}$. Kukkaketo voidaan nyt esittää Σ^* :n merkkijonona. Kukin kasvi kehittyy seuraavien sääntöjen mukaan:

SIEMEN \rightarrow SIRKKALEHDET
 SRKKALEHDET \rightarrow LEHDET | VARSI
 LEHDET \rightarrow VARSI
 VARSI \rightarrow NUPPU | LEHDET
 NUPPU \rightarrow PUNKUKKA | SINKUKKA
 PUNKUKKA \rightarrow SIEMEN | SIEMEN SIEMEN $|\epsilon$
 SINKUKKA \rightarrow SIEMEN | SIEMEN SIEMEN $|\epsilon$

missä ϵ kuvaa kasvin kuolemaa (se katoaa kylvämättä edes siemeniä).

Kyseessä on *kontekstiton kielioppi*, sillä säännön seuraus määräytyy ainoastaan edeltävän symbolin perusteella (jos useita vaihtoehtoisia seurauksia, niin valitaan jokin niistä).

Esimerkki 2 Laajennetaan symbolijoukkoa: $\Sigma_2 = \Sigma \cup \{PAIVA, YO\}$ ja määritellään uudet säännöt:

SIEMEN \rightarrow SIRKKALEHDET
 SRKKALEHDET PAIVA \rightarrow LEHDET | VARSI
 LEHDET PAIVA \rightarrow VARSI
 VARSI PAIVA \rightarrow NUPPU | LEHDET
 NUPPU YO \rightarrow PUNKUKKA | SINKUKKA
 PUNKUKKA \rightarrow SIEMEN | SIEMEN SIEMEN $|\epsilon$
 SINKUKKA \rightarrow SIEMEN | SIEMEN SIEMEN $|\epsilon$
 PAIVA \rightarrow YO
 YO \rightarrow PAIVA

Nyt kukat kasvavat vain päivisin, mutta kukat aukeavat vain öisin. Tämä kielioppi on *kontekstillinen*, sillä myös symbolin konteksti (ympäristötekijät) vaikuttavat säännön seuraukseen.

- L-järjestelmät \sim yksinkertaistettu abstrakti kielioppi, jonka tuottama kieli kuvaa kasvien (tms. elollisten organismien) kasvua ja kehitystä
- alkujaan botanisti Aristid Lindenmayerin kehittämä matemaattinen malli, jolla voi ennustaa kasvien kasvua
- sittemmin käytetty keinotekkoisten organismien luomiseen ja kaikkeen mahdolliseen!
- L-järjestelmä tuottaa vain kieliopin mukaisen merkkijonon \rightarrow voidaan tulkita graafisena kuvana
- muistuttavat fraktaaleja, mutta eivät välttämättä fraktaaleja
- Yleinen idea: jokaisella aika-askelella merkkijonon symboli lavennetaan naapurisymbolien ja ehtoihin sopivan säännön perusteella (voi pysyä ennallaankin) \leftrightarrow soluautomaatti, jossa kukin hilan (vektorin, 2- tai useampiulotteisen taulukon) solu on äärellinen automaatti
- Leikkaavat Chomskyn hierarkian kaikkia luokkia:
 - voivat noudattaa säännöllisen kielen sääntöjä, ts. ovat muotoa $A \rightarrow aB$ tai $A \rightarrow Ba$ (muunnoksen tuloksena yksi päätesymboli ja yksi väliesymboli) tai $A \rightarrow \epsilon$ ("solun kuolema", symboli katoaa) tai
 - kontekstitonta kielioppia, ts. säännöt muotoa $A \rightarrow V^*$ (ts. muunnetaan vain yhtä symbolia kerrallaan, riippumatta sen naapureista, muunnoksen tuloksena voi tulla mitä tahansa välies- ja päätesymboleja, myös ϵ) tai
 - kontekstillista kielioppia, ts. säännöt muotoa $\alpha A \beta \rightarrow \alpha \omega \beta$ (ts. muunnos vain tietyssä "kontekstissa", kun naapureina α ja β)
 - rajoittamatonta kielioppia, ts. säännöt muotoa $\omega \rightarrow \omega'$ (ts. mikä tahansa merkkijono voi muuntua miksi tahansa merkkijonoksi, merkkijonot voivat sisältää niin pääte- kuin väliesymboleja)

- esim. erään levämuodon kasvun mallinnus:
 - alkuvaiheessa voi olla kahdenlaisia soluja, merk. A ja B ("aksiomat")
 - säännöt:
 - $A \rightarrow AB$
 - $B \rightarrow A$
 - esim. $AB \Rightarrow ABA \Rightarrow ABAAB \Rightarrow ABAABABA$ (3 aika-askelta)
- yksinkertaisimmat L-järjestelmät ns. DOL-järjestelmiä (kontekstiton ja deterministinen L-järjestelmä)
- esim. 2: säännöt:
 - $A \rightarrow CB$
 - $B \rightarrow A$
 - $C \rightarrow DA$
 - $D \rightarrow C$
 - lähtökohta: mikä tahansa aakkoston sallittu "siemen" (lähtösymboli)
 - esim. $A \Rightarrow CB \Rightarrow DAA \Rightarrow CCBCB \Rightarrow DADAADAA$
- esim. 3: Puurakenteen muodostaminen:
 - $A \rightarrow C[B]D$
 - $B \rightarrow A$
 - $C \rightarrow C$
 - $D \rightarrow C(E)A$
 - $E \rightarrow D$
 - tulkitaan hakasulut haarana vasemmalle ja tavalliset sulut haarana oikealle
 - esim. $A \Rightarrow C[B]D \Rightarrow C[A]C(E)A$ (2 aika-askelta)
 - puuna:

- Sovellus (WH): muodostetaan kukkaketoa grafiikkaohjelmalla, joka noudattaa seuraavia sääntöjä:

SIEMEN → *SIRKKALEHDET|VVARSI|OVARSI|KVARSI*

VVARSI → *VLEHTI|VNUPPU* (satunnainen väri)

OVARSI → *OLEHTI|ONUPPU*

VNUPPU → *KUKKA*

ONUPPU → *KUKKA*

KNUPPU → *KUKKA*

KUKKA → *SIRKKALEHDET*

KVARSI → *VVARSI|OVARSI|KNUPPU|SIRKKALEHDET*

|*VKIELONVARSI|OKIELONVARSI|HAVU*

SIRKKALEHTI → *VVARSI|OVARSI|KVARSI|LEHDET*

VKIELOVARSI → *VKIELO*

OKIELOVARSI → *OKIELO*

VKIELO → *VMARJAT*

OKIELO → *OMARJAT*

- Kuvaruutu on jaettu soluihin, joita käydään läpi jossain järjestyksessä ja piirretään ko. paikkaan (tai naapurisoluihin) symbolin mahdollinen seuraaja

Kirjallisuutta

- Prusinkiewicz & Lindenmayer: The Algorithmic Beauty of Plants
- Prusinkiewicz & Hanan: Lindenmayer Systems, Fractals, and Plants
- Rozenberg & Salomaa (toim.): Lindenmayer Systems. Impacts on Theoretical Computer Science, Computer Graphics, and Developmental Biology