

**Lemma (Pumppauslemma).** Olkoon  $A$  säännöllinen kieli. Tällöin on olemassa  $n \geq 1$  siten, että kaikki  $A$ :n merkkijonot  $x$ , joiden pituus  $|x| \geq n$  ovat ilmaistavissa muodossa  $x = uvw$ , missä  $|uv| \leq n$ ,  $|v| \geq 1$  ja merkkijonot muotoa  $uv^i w$  kuuluvat kieleen  $A$  kaikilla  $i \geq 0$ . Kompaktimmin, painottaen pumppauslemman asettamia vaatimuksia, voimme kirjoittaa seuraavasti:

$\forall$  säännöllisille kielille  $A$

$\exists n \geq 1$  s.e.

$\forall x \in A : |x| \geq n$

$\exists$  osinjako  $x = uvw$ , missä  $|uv| \leq n$  ja  $|v| \geq 1$

$\forall i \geq 0 \ uv^i w \in A$ .

Pumppauslemmaa voidaan käyttää hyväksi, kun halutaan osoittaa kieli  $L$  *ei-säännölliseksi*. Tehdään ensin vasta oletus, eli oletetaan  $L$  säännölliseksi kieleksi. Tavoitteena on päästä ristiriitaan tämän oletuksen kanssa seuraten pumppauslemman asettamia vaatimuksia säännöllisille kielille.

Pumppauslemmaa käytettäessä täytyy aina muistaa, että sillä voi osoittaa vain kielen epä-säännöllisyyden, ja sitä *ei* voi käyttää toiseen suuntaan. Esimerkiksi kieli

$$I = \{c^i a^n b^n \mid i > 0 \wedge n \geq 0\} \cup L(a^* b^*)$$

ei ole säännöllinen, mutta kaikki siihen kuuluvat sanat (tyhjää sanaa lukuunottamatta) voidaan osoittaa pumppauslemman ehtojen mukaisesti. Näin ollen kieltä  $I$  ei voida suoraan todistaa epä-säännölliseksi, vaan todistuksessa täytyy käyttää apuna säännöllisten kielten sulkeumaominaisuuksia 5. tehtävän vastauksessa esitettävään tapaan.

Jos halutaan osoittaa kieli säännölliseksi, voidaan muodostaa sen hyväksyvä äärellinen automaatti, sillä pätee: Kieli  $L$  on säännöllinen  $\Leftrightarrow$  on olemassa äärellinen automaatti  $M$ , joka hyväksyy kielen  $L$  (merkitään  $L(M) = L$ ).

4. **Tehtävä:** Modernissa WWW-sivujen kuvaamiseen käytetyssä XML-kielessä on sivujen suunnittelijan mahdollista laatia omia ns. dokumenttityypimäärittäjiä (engl. *Document Type Definition*, lyh. DTD), jotka ovat oleellisesti sivulla esitettävän tekstin tai muun datan rakennetta kuvaavia yhteydettömiä kielioppeja. Tutustu tämän XML/DTD-kuvauskielen notaatioon (esim. WWW-sivulta <http://www.rpbouret.com/xml/xml.dtd.htm>), ja laadi seuraavaa XML/DTD-kuvausta vastaava yhteydetön kielioppi:

```
<!DOCTYPE Book [  
  <!ELEMENT Book (Title, Chapter+)>  
  <!ATTLIST Book Author CDATA #REQUIRED>  
  <!ELEMENT Title (#PCDATA)>  
  <!ELEMENT Chapter (#PCDATA)>  
  <!ATTLIST Chapter id ID #REQUIRED>  
>
```

**Vastaus:** Yllä oleva DTD-kuvaus määrittelee rakenteen kirjalle. Määrittelyssä esiintyy kahdenlaisia asioita: *osia* (element) ja *attribuutteja* (attlist). Perusajatuksena on, että kirja itsessään koostuu osista, ja attribuutit puolestaan liittyvät kirjan osiin ylimääräistä tietoa.

Yleisesti ottaen attribuutteja ei voida esittää puhtailla kieliopilla, vaan niitä varten tarvitaan *attribuuttikieliopit* (ks. opetusmoniste s. 60–68). Näin ollen DTD-kuvauksesta mallinnetaan ensin vain osat, eli käytännössä ainoastaan rivit, jotka alkavat merkinnällä `!ELEMENT`.

Näistä riveistä ensimmäinen:

```
<!ELEMENT Book (Title, Chapter+)>
```

kertoo, että kirja (Book) sisältää otsikon (Title) ja listan lukuja (Chapter). Lukuja täytyy olla vähintään yksi. Seuraava rivi:

```
<!ELEMENT Title (#PCDATA)>
```

puolestaan määrittelee otsikon merkkijoukoksi (#PCDATA). Merkkijoukkoon voi kuulua periaattessa mitä tahansa tietokoneen merkkijärjestelmään kuuluvia symboleita.

Lopuksi, rivi:

```
<!ELEMENT Chapter (#PCDATA)>
```

kertoo luvun olevan taas joukko merkkejä.

Kirjan rakenteen kuvaa siis kielioppi:<sup>1</sup>

$$\begin{aligned} \textit{Book} &\rightarrow \textit{Title Chapters} \\ \textit{Title} &\rightarrow \mathbf{data} \\ \textit{Chapters} &\rightarrow \textit{Chapter Chapters} \mid \textit{Chapter} \\ \textit{Chapter} &\rightarrow \mathbf{data} \end{aligned}$$

XML-kielessä erotetaan dokumentin osat toisistaan käyttäen apuna `<A>` ja `</A>` koodeja. Kun nämä koodit lisätään yllä olevaan kielioppiin, saadaan tulokseksi seuraavanlainen yhteydetön kielioppi:

$$\begin{aligned} \textit{Book} &\rightarrow \langle \mathbf{Book} \textit{Title Chapters} \langle / \mathbf{Book} \rangle \\ \textit{Title} &\rightarrow \langle \mathbf{Title} \mathbf{data} \langle / \mathbf{Title} \rangle \\ \textit{Chapters} &\rightarrow \textit{Chapter Chapters} \mid \textit{Chapter} \\ \textit{Chapter} &\rightarrow \langle \mathbf{Chapter} \mathbf{data} \langle / \mathbf{Chapter} \rangle \end{aligned}$$

Vaikka attribuutteja ei voidakaan täysin esittää yhteydettömällä kieliopilla, niiden syntaksi voidaan kuitenkin kuvata. XML-kielessä osan attribuutit kirjoitetaan osan aloittavan koodin sisään. Lisäämällä nämä ylläolevaan kielioppiin saadaan tulokseksi:

$$\begin{aligned} \textit{Book} &\rightarrow \langle \mathbf{Book} \textit{BookAttributes} \rangle \textit{Title Chapters} \langle / \mathbf{Book} \rangle \\ \textit{Title} &\rightarrow \langle \mathbf{Title} \mathbf{data} \langle / \mathbf{Title} \rangle \\ \textit{Chapters} &\rightarrow \textit{Chapter Chapters} \mid \textit{Chapter} \\ \textit{Chapter} &\rightarrow \langle \mathbf{Chapter} \textit{ChapterAttributes} \rangle \mathbf{data} \langle / \mathbf{Chapter} \rangle \\ \textit{BookAttributes} &\rightarrow \mathbf{author} = \mathbf{data} \\ \textit{ChapterAttributes} &\rightarrow \mathbf{id} = \mathbf{data} \end{aligned}$$

Tässä on huomattava, että yhteydettömällä kieliopilla ei voida määritellä ehtoa, jonka mukaan kaikilla luvuilla on eri tunnus. Tällaisten ehtojen toteutuminen täytyy tarkistaa jollain erillisellä ohjelmakoodilla.

---

<sup>1</sup> *Kursiivilla* kirjoitetut sanat ovat välitteitä, **vahvennetut** päätemerkkejä.

5. **Tehtävä:** Osoita, että kieli  $\{w \in \{a, b\}^* \mid w\text{:ssä on yhtä monta } a\text{:ta ja } b\text{:tä}\}$  ei ole säännöllinen, ja laadi yhteydetön kielioppi sen kuvaamiseen.

**Vastaus:** Kielen  $L = \{w \in \{a, b\}^* \mid w\text{:ssä on yhtä monta } a\text{:ta ja } b\text{:tä}\}$  voisi todistaa ei-säännölliseksi suoraan pumppauslemmalla. Tässä esitetään kuitenkin hieman monimutkaisempi ratkaisu esimerkkinä siitä, miten ”hankalia” kieliä voidaan käsitellä.

Määritellään kieli  $L' = L \cap L(a^*b^*)$ . Oletetaan, että  $L$  on säännöllinen. Koska  $L(a^*b^*)$  on säännöllinen ja säännöllisten kielten joukko on suljettu leikkauksen suhteen, täytyy myös  $L'$ :n olla säännöllinen. (Toisinpäin ehto ei päde:  $L'$  voi olla säännöllinen vaikka  $L$  ei olisi, sillä esim.  $A \cap \emptyset = \emptyset$  kaikille kielille  $A$ ).

Huomataan, että  $L' = \{a^k b^k \mid k \geq 0\}$ . Tarkastellaan sanaa  $w = a^n b^n$ , missä  $n$  on pumppauslemmassa esiintyvä parametri. Yritetään osoittaa  $w$  lemman ehtojen mukaisesti. Koska  $|xy| \leq n$ , osituksen täytyy olla muotoa;

$$\begin{aligned}x &= a^{n-i-k} \\y &= a^i \\z &= a^k b^n,\end{aligned}$$

missä  $0 < i \leq n$  ja  $i + k \leq n$ . Nyt  $xz = a^{n-i} b^n$ , joten  $xz \notin L'$ . Näin ollen sanaa  $w$  ei voida pumpata, eikä  $L'$  ole säännöllinen, joten myöskään  $L$  ei ole säännöllinen.

Alla kielen  $L$  kuvaava yhteydetön kielioppi  $G$ :

$$\begin{aligned}G &= (V, \Sigma, P, S), \text{ missä} \\V &= \{S, T, a, b\}, \\ \Sigma &= \{a, b\}, \\ P &= \{ S \rightarrow SS \mid aT \mid Ta \mid \varepsilon, \\ &\quad T \rightarrow ST \mid TS \mid b \}\end{aligned}$$

Tässä kieliopissa välikkeellä  $S$  johdetaan merkkijonot, joissa on molempia aakkosia yhtä monta, ja  $T$ :llä kaikki joissa on  $b$ -merkkejä yksi enemmän kuin  $a$ -merkkejä.

**Esimerkki.** Annetaan johto merkkijonolle  $aababb \in L$ .

$$\begin{aligned}S &\Rightarrow aT \\ &\Rightarrow aST \\ &\Rightarrow aaTT \\ &\Rightarrow aabT \\ &\Rightarrow aabST \\ &\Rightarrow aabaTT \\ &\Rightarrow aababT \\ &\Rightarrow aababb\end{aligned}$$

6. **Tehtävä:** Laadi yhteydetön kielioppi, joka tuottaa kaikki seuraavan esimerkin tapaiset, yksinkertaisista sisäkkäisistä `for`-silmukoista, `begin`- ja `end`-sulkeilla kootuista lauseista ja alkeisoperaatioista `a` rakentuvat ”ohjelmat”:

```
a;
for 3 times do
begin
  for 5 times do a;
  a; a
end
```

Silmukkalaskureiden voit olettaa olevan kokonaislukuja väliltä  $0, \dots, 9$ .

**Vastaus:** Ohjelmointikielten kieliopit määritellään useimmiten siten, että aakkostoksi otetaan kielessä esiintyvät syntaktiset elementit (lekseemit). Tässä tapauksessa niitä ovat numerot, `a` sekä varatut sanat. Ohjelman jäsentäminen jaetaan kahteen osaan:

- (a) Muutetaan ohjelman teksti jonoksi lekseemeitä tilakoneiden avulla.
- (b) Muodostetaan lekseemijonon jäsenyspuu.

Tehtävän kieliopin voi määrittellä monellakin eri tapaa, tässä on yksi mahdollinen tulkinta:

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

$$V = \{C, S, N, \mathbf{begin}, \mathbf{do}, \mathbf{end}, \mathbf{for}, \mathbf{times}, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, ;, a\}$$

$$\Sigma = \{\mathbf{begin}, \mathbf{do}, \mathbf{end}, \mathbf{for}, \mathbf{times}, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, ;, a\}$$

Tässä välikkeen  $S$  tulkintana on "lause" (*statement*),  $C$ :n "yhdistetty lause" (*compound statement*) ja  $N$ :n "numero". Kieliopin säännöt määritellään seuraavasti:

$$P = \{C \rightarrow S \mid S; C$$

$$S \rightarrow a \mid \mathbf{begin} C \mathbf{end} \mid \mathbf{for} N \mathbf{times} \mathbf{do} S$$

$$N \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9\}$$

**Esimerkki.** Tehtävänannossa esiintyneen ohjelman jäsenyspuu:

