

4.4 Kielioppien jäsenysongelma

Ratkaistava tehtävä:

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

- Esim. Kuuluuko virke ”jänis joka pelkäsi arkaa peikkoa metsästi suurta sutta” kieleen L_{rel} ?
- Ratkaisumenetelmä: *jäsennysalgoritmi*
- Useita vaihtoehtoisia menetelmiä, erityisesti kun G on jotain rajoitettua (käytännössä esiintyvää) muotoa.

- Tarkastellaan ensin jäsentämiseen liittyviä peruskäsitteitä

4.4.1 Johdot

- Olk. merkkijono $\gamma \in V^*$ kieliopin $G = (V, \Sigma, P, S)$ lausejohdos
- γ :n *johdoksi* G :ssä kutsutaan lähtösymbolista S merkkijonoon γ johtavaa suorien johtojen jonoa

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

- Johdon *pituus* on siihen kuuluvien suorien johtojen määrä (edellä n)
- Esim. lauseen $a + a$ johtoja kieliopissa G_{expr} :

$$\begin{array}{l} \text{(i)} \quad E \Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \\ \quad \quad \Rightarrow a + T \Rightarrow a + F \Rightarrow a + a \\ \text{(ii)} \quad E \Rightarrow E + T \Rightarrow E + F \Rightarrow T + F \\ \quad \quad \Rightarrow F + F \Rightarrow F + a \Rightarrow a + a \\ \text{(iii)} \quad E \Rightarrow E + T \Rightarrow E + F \Rightarrow E + a \\ \quad \quad \Rightarrow T + a \Rightarrow F + a \Rightarrow a + a \end{array}$$

- Johto $\gamma \Rightarrow^* \gamma'$ on

1. *vasen johto*, merk.

$$\gamma \xRightarrow[\text{lm}]{}^* \gamma',$$

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

2. *oikea johto*, merk.

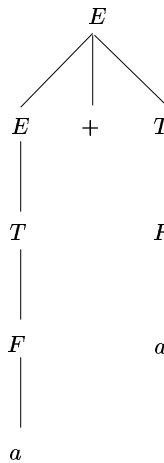
$$\gamma \xRightarrow[\text{rm}]{}^* \gamma',$$

jos — ” — oikeanpuoleiseen välikkeeseen (edellä (iii))

- Suoria vasempia ja oikeita johtoaskelia merkitään $\gamma \xRightarrow[\text{lm}]{} \gamma'$ ja $\gamma \xRightarrow[\text{rm}]{} \gamma'$

4.4.2 Jäsennyspuut

- Suurin osa vaihtoehtoisten johtotapojen eroista muodostuu vain välikkeiden laventamisesta eri järjestyksessä (esim. em. johdot (i) – (iii) ovat kaikki “pohjimmiltaan” samanlaisia)
- Lausejohdoksen *jäsennyspuu* (syntaksipuu, johtopuu) (engl. parse tree, syntax tree, derivation tree) = esitystapa, jossa nämä epäoleelliset erot on abstrahoitu pois
- Jäsennyspuu kertoo ainoastaan, *miten* välikkeet on lavennettu, ei *missä järjestyksessä* lavennukset on tehty
- Esim. kaikkia kolmea em. johtoa vastaa sama jäsennyspuu:



Kuva 4.6: Lauseen $a + a$ jäsennyspuu kieliopissa G_{expr} .

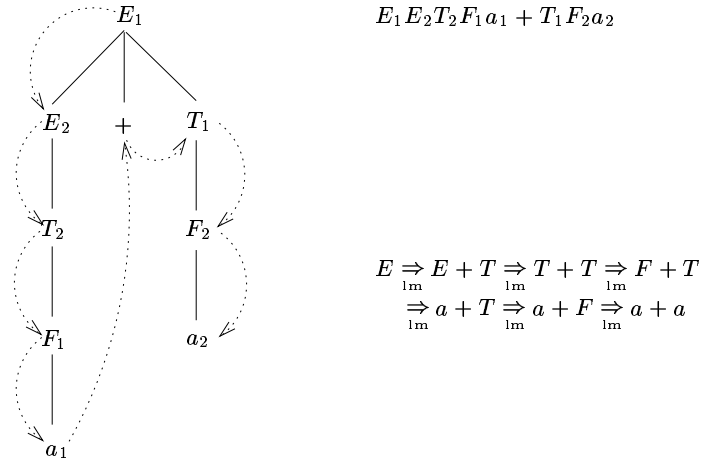
Määritelmä: Olkoon $G = (V, \Sigma, P, S)$ kontekstiton kielioppi. Kieliopin G mukainen *jäsennyspuu* on järjestetty puu, jolla on seuraavat ominaisuudet:

(i) puun solmut on nimetty joukon $V \cup \{\epsilon\}$ 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 τ *tuotos* = merkkijono, joka saadaan liittämällä yhteen sen lehtisolmujen nimet esijärjestyksessä (“vasemmalta oikealle”)



Kuva 4.7: Vasemman johdon muodostaminen jäsennyspuusta.

4.4.3 Jäsennyspuun muodostaminen

- Johtoa

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

vastaava jäsennyspuu muodostetaan seuraavasti:

- puun juuren nimeksi tulee S ; jos $n = 0$, niin puussa ei ole muita solmuja; muuten
- jos ensimmäisessä johtoaskelessa on sovellettu produktiota $S \rightarrow X_1X_2 \dots X_k$, niin juurelle tulee k jälkeläissolmua, joiden nimet vasemmalta oikealle ovat X_1, X_2, \dots, X_k ;
- jos seuraavassa askelessa on sovellettu produktiota $X_i \rightarrow Y_1Y_2 \dots Y_l$, niin juuren i :nnelle jälkeläissolmulle tulee l jälkeläistä, joiden nimet vasemmalta oikealle ovat Y_1, Y_2, \dots, Y_l ; ja niin edelleen

- Ts. jos τ on jotakin johtoa $S \Rightarrow^* \gamma$ vastaava jäsennyspuu, niin τ :n tuotos on γ
- Johdon muodostaminen: Olk. τ kieliopin G mukainen jäsennyspuu, jonka tuotos on päättemerkkijono x .
 - τ :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”)
- Jos muodostetaan annetusta vasemmasta (oikeasta) johdosta $S \xRightarrow{*}_{lm} x$ ($S \xRightarrow{*}_{rm} x$) ensin jäsennyspanu em. tavalla, ja sitten jäsennyspanuista vasen (oikea) johto, saadaan takaisin alkuperäinen johto

Lause: Olk. $G = (V, \Sigma, P, S)$ kontekstiton kielioppi. Tällöin:

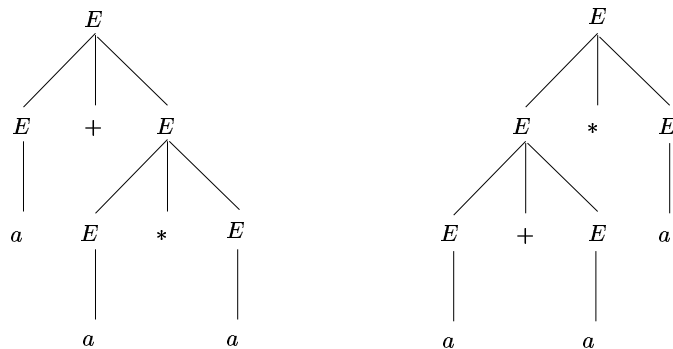
- (i) jokaisella G :n lausejohdoksella γ on G :n mukainen jäsennyspanu τ , jonka tuotos on γ ;
- (ii) jokaista G :n mukaista jäsennyspanua τ , jonka tuotos on päätmerkkijono x , vastaavat yksikäsitteiset vasen ja oikea johto $S \xRightarrow{*}_{lm} x$ ja $S \xRightarrow{*}_{rm} x$

Seuraus: Jokaisella G :n lauseella on vasen ja oikea johto

- Ts. kontekstittoman kieliopin tuottamien lauseiden jäsennyspanut, vasemmat ja oikeat johdot vastaavat yksikäsitteisesti toisiaan
- Jäsennyspangelman ratkaisuksi katsotaan usein päätöspangelman “Onko $x \in L(G)$?” ratkaisemisen lisäksi jonkin näistä jäsennyspanesityksistä tuottaminen x :lle

4.4.4 Kieliopin moniselitteisyys

- Lauseella voi olla kieliopissa useita jäsennyksii (esim. lause $a + a * a$ kieliopissa G'_{expr})



Kuva 4.8: Lauseen $a + a * a$ kaksi erilaista jäsennyttä.

- Kontekstiton kielioppi G on *moniselitteinen*, jos jollakin G :n lauseella x on kaksi erilaista G :n mukaista jäsenyspuuta
- Muuten kielioppi on *yksiselitteinen*
- Kontekstiton kieli, jonka tuottavat kieliopit ovat kaikki moniselitteisiä, on *luonnostaan moniselitteinen*
- Esimerkkejä:
 - Kielioppi G'_{expr} on moniselitteinen
 - Kieliopit G_{expr} ja G_{match} yksiselitteisiä
 - Kieli $L_{\text{expr}} = L(G'_{\text{expr}})$ ei ole luonnostaan moniselitteinen, koska sillä on myös yksiselitteinen kielioppi G_{expr}
 - Kieli $\{a^i b^j c^k \mid i = j \text{ tai } j = k\}$ on luonnostaan moniselitteinen (todistus sivuutetaan)

4.5 Rekursiivisesti etenevä jäsentäminen

4.5.1 LL(1)-kielioppi

- Tarkastellaan seuraavaa kielioppia G :

$$\begin{aligned} E &\rightarrow T + E \mid T - E \mid T \\ T &\rightarrow a \mid (E) \end{aligned}$$

- Muokataan G :stä välkkeen E produktiot ”tekijöimällä” ekvivalentti kielioppi G' :

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +E \mid -E \mid \epsilon \\ T &\rightarrow a \mid (E) \end{aligned}$$

- G' lauseille voidaan helposti muodostaa vasemmat johdot suoraan lähtösymbolista E alkaen: Jäsennyksen joka vaiheessa tavoitteena olevan lauseen *seuraava merkki* määrää yksikäsitteisesti sen, mikä produktio valitaan lavennettavana olevaan välkkeeseen
- Kielioppia, jolla on tämä ominaisuus, sanotaan *LL(1)-kieliopiksi* (”left to right scan, producing left parse with 1 symbol lookahead” – vast LL(k)- ja LR(k)-kieliopit)

4.5.2 LL(1)-kieliopin jäsenysohjelma

- Jäsenysohjelma muodostuu välikkeitä vastaavista rekursiivisista funktioista
- Esim. G' :n mukainen jäsentäjä pseudokoodina:

```
function  $E$   
begin  
 $T$ ;  
 $E'$ ;  
end
```

```
function  $E'$   
begin  
if ( $next == '+'$ ) then  
begin  
    tee jotakin;  
     $next = getnext$ ;  
     $E$ ;  
end  
else  
    begin  
        if ( $next == '-'$ ) then  
            begin  
                tee jotakin;  
                 $next = getnext$ ;  
                 $E$ ;  
            end  
        else tee jotakin;  
    end  
end
```

```
function  $T$   
begin  
if ( $next == 'a'$ ) then  
    begin  
        tee jotakin;  
         $next = getnext$ ;  
        return;  
    end  
else
```

```

begin
  if (next == '(') then
    begin
      tee jotakin;
      next = getnext;
      E;
      if (next ≠ ')')
        ERROR;
      next = getnext;
    end
  else ERROR;
end
end
end

```

PÄÄOHJELMASSA:

```

next = getnext;
E;

```

Esim. syötejonon a-(a+a) jäsennessä, kun käsky ”tee jotakin” tulostaa produktiot:

```

E → TE'
T → a
E' → -E
E → TE'
T → (E)
E → TE'
T → a
E' → +E
E → TE'
T → a
E' → ε
E' → ε.

```

Tulostus vastaa vasenta johtoa:

$$\begin{aligned}
 E &\Rightarrow TE' \Rightarrow aE' \Rightarrow a - E \Rightarrow a - TE' \\
 &\Rightarrow a - (E)E' \Rightarrow a - (TE')E' \\
 &\Rightarrow a - (aE')E' \Rightarrow a - (a + E)E' \\
 &\Rightarrow a - (a + TE')E' \Rightarrow a - (a + aE')E' \\
 &\Rightarrow a - (a + a)E' \Rightarrow a - (a + a)
 \end{aligned}$$

4.5.3 LL(1)-kielioppien yleinen muoto

LL(1)-kieliopissa sallitaan rajoitetussa määrin myös

- Produktioita, joiden oikeat puolet alkavat väliskeellä, sekä
- Tyhjentyviä väliskeitä A , joilla $A \Rightarrow^* \epsilon$

Esim. kielen $a^*b \cup c^*d$ tuottava kielioppi:

$$\begin{array}{l} S \rightarrow Ab \mid Cd \\ A \rightarrow aA \mid \epsilon \\ C \rightarrow cC \mid \epsilon. \end{array}$$

Kielioppi on LL(1), vaikka ensimmäiseksi sovellettavaa produktiota ei voikaan päätellä pelkästään S :n produktioiden perusteella

4.5.4 Apukäsite: päättemerkkijoukot FIRST ja FOLLOW

Määritellään annetun kieliopin $G = (V, \Sigma, P, S)$ väliskeisiin liittyvät päättemerkkijoukot:

- $\text{FIRST}(A) = A$:sta johdettavien päätejonojen 1. merkit + ϵ , jos A tyhjentyvä
- $\text{FOLLOW}(A) =$ ne päättemerkit, jotka voivat seurata A :ta jossain G :n lausejohdoksessa + ϵ , jos A voi sijaita lausejohdoksen lopussa
- Formaalisti:

$$\begin{aligned} \text{FIRST}(A) &= \{a \in \Sigma \mid A \Rightarrow^* ax \text{ jollakin } x \in \Sigma^*\} \\ &\quad \cup \{\epsilon \mid A \Rightarrow^* \epsilon\}; \\ \text{FOLLOW}(A) &= \{a \in \Sigma \mid S \Rightarrow^* \alpha A a \beta \text{ joillakin } \alpha, \beta \in V^*\} \\ &\quad \cup \{\epsilon \mid S \Rightarrow^* \alpha A \text{ jollakin } \alpha \in V^*\}. \end{aligned}$$

FIRST-joukkojen määritelmä laajennetaan mielivaltaisille merkkijonoille, ja edelleen merkkijonojoukkoihin:

$$\begin{aligned} \text{FIRST}(\epsilon) &= \{\epsilon\}; \\ \text{FIRST}(a) &= \{a\} \quad \text{kaikilla } a \in \Sigma; \\ \text{FIRST}(X_1 \dots X_k) &= \begin{cases} \text{FIRST}(X_1) \cup \dots \cup \text{FIRST}(X_i) - \{\epsilon\}, \\ \quad \text{jos } \epsilon \in \text{FIRST}(X_1), \dots, \text{FIRST}(X_{i-1}), \\ \quad \epsilon \notin \text{FIRST}(X_i); \\ \text{FIRST}(X_1) \cup \dots \cup \text{FIRST}(X_k), \\ \quad \text{jos } \epsilon \in \text{FIRST}(X_i) \text{ kaikilla } i = 1, \dots, k; \end{cases} \end{aligned}$$

Määritelmä: Kielioppi on LL(1)-muotoinen, jos sen millä tahansa kahdella samaan välikkeeseen A liittyvällä produktiolla $A \rightarrow \omega_1$ ja $A \rightarrow \omega_2$, $\omega_1 \neq \omega_2$, on voimassa:

$$\begin{aligned} &\text{FIRST}(\{\omega_1\}\text{FOLLOW}(A)) \\ &\cap \text{FIRST}(\{\omega_2\}\text{FOLLOW}(A)) = \emptyset. \end{aligned}$$

LL(1)-kieliopin rekursiivisesti etenevän jäsentäjän muodostaminen yleisesti:

Apurutiinit:

token *getnext*;

Seuraavan päätesymbolin selaus

... — voi olla > 1 kirjainmerkkiä.*

void ERROR(char* *msg*);

Virheiden käsittely; parametri *msg*

... sisältää virheilmoitustekstin.*

Välikettä A vastaava funktio:

function *A* /**A*:n produktiot $A \rightarrow \omega_1 \mid \dots \mid \omega_n$ */

begin

if (*next in* [*a*₁₁, ..., *a*_{1*m*₁}]) /*FIRST({ ω_1 }FOLLOW(*A*)) = {*a*₁₁, ..., *a*_{1*m*₁}}*/

begin /*produktio $A \rightarrow \omega_1$ */

parse(ω_1);

end

else

```

:
if (next in [ $a_{n1}, \dots, a_{nm_n}$ ])          /*FIRST( $\{\omega_n\}$ FOLLOW( $A$ )) =  $\{a_{n1}, \dots, a_{nm_n}\}$ */
begin /*produktio  $A \rightarrow \omega_n$ */
    parse( $\omega_n$ );
end
else
    ERROR("A cannot start with this.")
end

```

Tässä lyhennemerkintä "*parse*(ω_i)" tarkoittaa seuraavalla tavalla muodostettavaa käskyjonoa:

$parse(X_1 \dots X_k) \equiv parse(X_1); \dots ; parse(X_k)$,

missä

$parse(a) \equiv$
if *next* $\neq a$ **then** ERROR('a expected. ');
next := *getnext*; (*a* on päätemerkki)

$parse(B) \equiv B$; (*B* on välike)

4.5.5 Kielioppien muokkaaminen LL(1)-muotoon

- LL(1)-kieliopit ovat varsin suppea kielioppiluokka
- Seuraavilla muunnoksilla voidaan joitakin "melkein" LL(1)-muotoisia kielioppeja muuntaa tähän muotoon:

1. VASEN TEKIJÖINTI

- Kielioppi, jossa on produktiot

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2, \quad \alpha \neq \epsilon, \beta_1 \neq \beta_2$$

ei voi olla LL(1)-muotoinen

- Parannus: otetaan käyttöön uusi välike A' ja korvataan em. produktiot produktioilla:

$$\begin{aligned} A &\rightarrow \alpha A' \\ A' &\rightarrow \beta_1 \mid \beta_2, \end{aligned}$$

missä α on $\alpha\beta_1$:n ja $\alpha\beta_2$:n pisin yhteinen alkuosa

2. VÄLITTÖMÄN VASEMMAN REKURSION POISTAMINEN

- Kielioppi on *vasemmalle rekursiivinen*, jos jollakin välikkeellä A ja merkkijonolla γ on

$$A \Rightarrow^+ A\gamma,$$

missä merkintä $\alpha \Rightarrow^+ \beta$ tarkoittaa, että α :sta voidaan johtaa β johdolla, jonka pituus on vähintään yksi askel

- Vasemmalle rekursiivinen kielioppi ei voi täyttää LL(1)-ehtoa
- *Välitön* vasen rekursio, so. suorat johdot $A \Rightarrow A\gamma$, voidaan välttää korvaamalla produktiot

$$A \rightarrow A\beta \mid \alpha, \quad \beta \neq \epsilon,$$

produktioilla

$$\begin{aligned} A &\rightarrow \alpha A' \\ A' &\rightarrow \beta A' \mid \epsilon \end{aligned}$$

- Muotoa $A \rightarrow A$ olevat produktiot voidaan yksinkertaisesti jättää pois

Jokainen kontekstiton kielioppi voidaan teoriassa muuntaa vasemman rekursion välttävään *Greibachin normaalimuotoon*, missä kaikki produktiot ovat muotoa

$$A \rightarrow aB_1 \dots B_k, \quad k \geq 0,$$

tai $S \rightarrow \epsilon$, missä a on päätimerkki, B_1, \dots, B_k välikkeitä ja S lähtösymboli

4.5.6 *LIITE: FIRST- ja FOLLOW-joukkojen laskeminen

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

Ensin lasketaan FIRST-joukot:

1. Asetetaan aluksi kaikille kieliopin päätteille $a \in \Sigma$:

$$\text{FIRST}(a) := \{a\},$$

ja kaikille välikkeille $A \in V - \Sigma$:

$$\begin{aligned} \text{FIRST}(A) := & \{a \in \Sigma \mid A \rightarrow a\beta \text{ on } G\text{:n produktio}\} \\ & \cup \{\epsilon \mid A \rightarrow \epsilon \text{ on } G\text{:n produktio}\}. \end{aligned}$$

2. Käydään sitten kieliopin produktioita läpi jossakin järjestyksessä ja toistetaan, kunnes FIRST-joukot eivät enää kasva: kullekin produktiolle $A \rightarrow X_1 \dots X_k$ asetetaan:

$$\begin{aligned} \text{FIRST}(A) := & \text{FIRST}(A) \cup \\ & \bigcup \{\text{FIRST}(X_i) \mid 1 \leq i \leq k, \epsilon \in \text{FIRST}(X_j) \text{ kaikilla } j < k\} \\ & \cup \{\epsilon \mid \epsilon \in \text{FIRST}(X_j) \text{ kaikilla } j = 1, \dots, k\}. \end{aligned}$$

FOLLOW-joukot määritetään FIRST-joukkojen avulla seuraavasti:

1. Asetetaan aluksi kaikille välikkeille $B \in V - \Sigma^*$:

$$\begin{aligned} \text{FOLLOW}(B) := & \\ & \bigcup \{\text{FIRST}(\beta) - \{\epsilon\} \mid A \rightarrow \alpha B \beta \text{ on } G\text{:n produktio}\}, \end{aligned}$$

ja lähtösymbolille S lisäksi:

$$\text{FOLLOW}(S) := \text{FOLLOW}(S) \cup \{\epsilon\}.$$

2. Sitten toistetaan, kunnes FOLLOW-joukot eivät enää kasva: kullekin produktiolle $A \rightarrow \alpha B \beta$, missä $\epsilon \in \text{FIRST}(\beta)$, asetetaan:

$$\text{FOLLOW}(B) := \text{FOLLOW}(B) \cup \text{FOLLOW}(A).$$