

Kuva 4.10: Kuva kontekstittomien kielten alaluokkien välisistä suhteista.

4.7 CYK-jäsennysalgoritmi

- Rekursiivisesti etenevä jäsentäminen on selkeä ja tehokas jäsenysmenetelmä LL(1)-kieliopille: n merkin mittaisen syötemerkkijonon käsittely sujuu ajassa $O(n)$
- Entäpä mielivaltaisen kontekstittoman kieliopin tuottamien merkkijonojen tunnistaminen?
- Ratkaisu 1: muunnetaan kielioppi ensin em. Greibachin normaalimuotoon
 - Tämänmuotoisella kieliopilla millä tahansa n merkin mittaisella lauseella on $n:n$ askeleen johto
 - Merkkijonon x , $|x| = n$, kuulumisen tuotettuun kieleen voidaan testata käymällä läpi kaikki n askelen mittaiset johdot ja katsomalla, tuottaako jokin niistä $x:n$
 - Tehotonta! (jokaisessa epätriviaalissa kieliopissa on n askelen mittaisia johtoja eksponentiaalinen määrä)

- Ratkaisu 2: *Cocke–Younger–Kasami*- eli *CYK-algoritmi*
 - Toimii ajassa $O(n^3)$, missä n on tutkittavan merkkijonon pituus
 - Kielopin on oltava *Chomskyn normaalimuodossa*
 - Mikä tahansa kontekstiton kielioppi voidaan kuitenkin muuntaa Chomskyn normaalimuotoon, joten sopii kaikille kontekstittomille kielille!

Chomskyn normaalimuoto

Määritelmä: Kontekstiton kielioppi $G = (V, \Sigma, P, S)$ on *Chomskyn normaalimuodossa*, jos

- Välikkeistä enintään S on tyhjentyvä (engl. nullable), eli $S \xRightarrow[G]{*} \epsilon$
- Muut produktiot ovat muotoa $A \rightarrow BC$ tai $A \rightarrow a$, missä A, B ja C ovat välikkeitä ja a on päätemerkki
- Lisäksi vaaditaan yksinkertaisuuden vuoksi, että lähtösymboli S ei esiinny minkään produktio-oikealla puolella

Esim. seuraava kielioppi on Chomskyn normaalimuodossa:

$$\begin{aligned} S &\rightarrow AB|\epsilon \\ A &\rightarrow BA|a \\ B &\rightarrow b \end{aligned}$$

Cocke–Younger–Kasami-algoritmi

- Perustuu dynaamiseen ohjelmointiin (välitulosten taulukointiin)
- Ongelma: Olk. $G = (V, \Sigma, P, S)$ Chomskyn normaalimuodossa (ellei ole, muutetaan ensin Chomskyn normaalimuotoon). Onko $x \in L(G)$ eli $S \xRightarrow[G]{*} x$?

Algoritmi:

1. Jos $x = \epsilon$, niin $x \in L(G) \Leftrightarrow S \rightarrow \epsilon$ on G :n produktio

2. Muuten merk. $x = a_1 \dots a_n$ ja tarkastellaan x :n osajonot tuottavien välikkeiden joukkoja

$$N_{i,j} = \{A \in V - \Sigma \mid A \xrightarrow[G]{*} a_i \dots a_j\}, \quad 1 \leq i \leq j \leq n$$

3. $x \in L(G) \Leftrightarrow S \in N_{1,n}$

Joukkojen $N_{i,j}$ laskeminen

- muodostetaan kolmiomainen taulukko:

a_1	a_2	a_3	a_4	a_5
$N_{1,1}$				
$N_{1,2}$	$N_{2,2}$			
$N_{1,3}$	$N_{2,3}$	$N_{3,3}$		
$N_{1,4}$	$N_{2,4}$	$N_{3,4}$	$N_{4,4}$	
$N_{1,5}$	$N_{2,5}$	$N_{3,5}$	$N_{4,5}$	$N_{5,5}$

- x-akseli vastaa merkkijonon positiota
- taulukon paikka $N_{i,j} \sim$ joukko välikesymboleita A , joille $A \xrightarrow[G]{*} a_i a_{i+1} \dots a_j$
- jokainen diagonaali vastaa tietyn pituisia x :n osajonoja
 - 1. diagonaali \sim välikkeet, jotka tuottavat yhden merkin osajonot
 - 2. diagonaali \sim välikkeet, jotka tuottavat kahden merkin osajonot jne.

- (i) Lasketaan ensin for $i = 1$ to n :

$$N_{i,i} := \{A \in V - \Sigma \mid A \rightarrow a_i \text{ on } G\text{:n produktio}\};$$

- (ii) Lasketaan sitten induktiivisesti

for $k = 1$ to $n - 1$
 for $i = 1$ to $n - k$:

$$N_{i,i+k} := \{A \in V - \Sigma \mid \text{jollakin } j, i \leq j < i + k, \text{ on välikkeet} \\ B \in N_{i,j} \text{ ja } C \in N_{j+1,i+k} \text{ s.e.} \\ A \rightarrow BC \text{ on } G\text{:n produktio}\}$$

- Huom! Tunnetaan jo kaikki lyhyemmät osajonot ja kaikki x :n kelvolliset prefiksit ja suffiksit
- Johdon $A \xRightarrow{G}^* a_i a_{i+1} \dots a_j$ täytyy alkaa askeleella $A \Rightarrow BC$, missä B :stä voidaan johtaa $a_i \dots a_j$:n prefiksi, olk. $B \xRightarrow{G}^* a_i a_{i+1} \dots a_l$, ja C :stä voidaan johtaa loput:
 $C \xRightarrow{G}^* a_{l+1} a_{l+2} \dots a_j$
- Esim. $N_{3,7}$:ta varten tarkistellaan kaikkia pareja
 $(N_{3,3}, N_{4,7}), (N_{3,4}, N_{5,7}), (N_{3,5}, N_{6,7}), (N_{3,6}, N_{7,7})$

Esim. Sovelletaan CYK:iä Chomskyn normaalimuotoiseen kielioppiin G :

$$\begin{array}{l|l} S & \rightarrow AB \quad | \quad BC \\ A & \rightarrow BA \quad | \quad a \\ B & \rightarrow CC \quad | \quad b \\ C & \rightarrow AB \quad | \quad a \end{array}$$

syötemerkkijonolla $x = baaba$:

$i \rightarrow$				
1 : b	2 : a	3 : a	4 : b	5 : a
B				
S, A	A, C			
\emptyset	B	A, C		
\emptyset	B	S, C	B	
S, A, C	S, A, C	B	S, A	A, C

Lähtösymboli S kuuluu joukkoon $N_{1,5}$, joten x kuuluu kieliopin tuottamaan kieleen $L(G)$.

Tutkitaan sitten merkkijonoa $ababab$:

$i \rightarrow$					
1 : a	2 : b	3 : a	4 : b	5 : a	6 : b
A, C					
S, C	B				
\emptyset	A, S	A, C			
B	S	S, C	B		
A, S	\emptyset	\emptyset	A, S	A, C	
S	\emptyset	B	S	S, C	B

Siis $ababab \in L(G)$.

Jäsennystä voi joskus helpottaa valitsemalla pääteakkoston sopivasti. Pääteakkosten ei tarvitse olla yksittäisiä kirjaimia, vaan ne voivat olla myös kielen ”atomisia” sanoja (osasia, joista pidemmät sanat koostuvat). Jos kieliopissa on sääntö $A \rightarrow w$, missä w koostuu vain päätemerkeistä, voidaan w valita uuden aakkoston päätesymboliksi.

Esim. 2. Tarkastellaan seuraavaa Kissastanian kielen mukaista kielioppia G_{miu} :

$$\begin{array}{l} S \rightarrow miuSmau|miu|mau|U \\ U \rightarrow Uu|u \end{array}$$

Valitaan pääteakkostoksi $\Sigma = \{miu, mau, u\}$. Chomskyn normaalimuotoon muunnettuna saadaan kielioppi

$$S \rightarrow IW|UU|miu|mau$$

$$W \rightarrow SA$$

$$U \rightarrow Uu|u$$

$$I \rightarrow miu$$

$$A \rightarrow mau$$

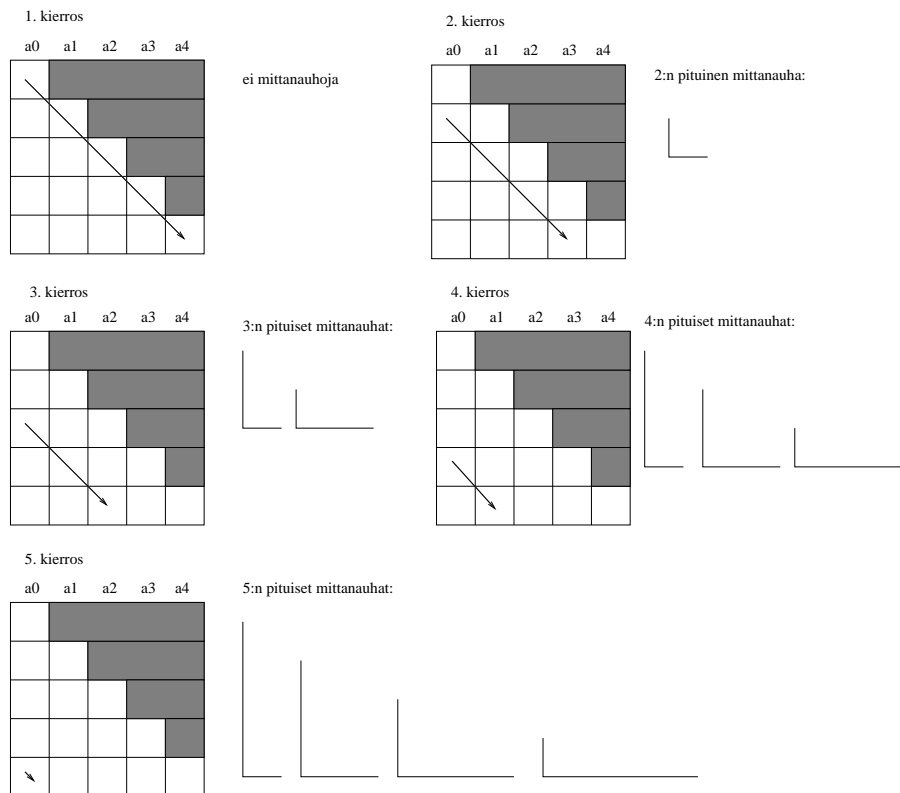
Tutkitaan CYK-algoritmilla, kuuluuko merkkijono $miumiuuumaumau$ kieleen:

$i \rightarrow$					
1 : miu	2 : miu	3 : u	4 : u	5 : mau	6 : mau
S, I					
\emptyset	S, I				
\emptyset	\emptyset	U			
\emptyset	\emptyset	S	U		
\emptyset	S	W	\emptyset	S, A	
S	W	\emptyset	\emptyset	W	S, A

Siis $miumiuuumaumau \in L(G_{miu})$.

Apukeino CYK:in simuloimiseen

CYK-algoritmissa edetään siis diagonaaleittain, kunnes tullaan vasempaan alakulmaan. Mikäli merkkijonon pituus on n , täytyy käydä läpi n diagonaalia. 1. diagonaalin i :teen sarakkeeseen tulevat vain ne välitteet A , joilla on sääntö $A \rightarrow a_i$. Seuraavilla diagonaaleilla ruutu täytetään edellisten diagonaalien sisällön perusteella. Voit kuvitella mielessäsi (tai askarrella piirtoheitinkalvosta!) ”mittanauhat”, joiden avulla muistat, mitä edeltävistä ruuduista täytyy tutkia ruudun täyttämiseksi. 2. diagonaalilla mitan pituus on 2 ja k :nnellä diagonaalilla mitan pituus on k . Huomaa, että kahden pituisen mitan voi asettaa vain yhdellä tapaa taulukkoon, mutta kolmenpituisella mitalla on jo kaksi vaihtoehtoista tapaa. k :nnellä diagonaalilla tarvittavia k :n pituisia ”mittanauhoja” on siis jo $k - 1$ kpl.



Kuva 4.11: CYK:in simulointi. Taulukkoa käydään läpi diagonaaleittain, ja diagonaalin joka ruudussa sovitetaan kaikkia annetuista ”mittanauhoista”. Mittanauhan päiden kohdalla olevat ruudut tutkitaan, ja mikäli niissä oleville välitteille B ja C löytyy sääntö $A \rightarrow BC$, sijoitetaan A senhetkiseen ruutuun.

4.8 Muunnos Chomskyn normaalimuotoon

- Idea:
 1. Poistetaan lähtösymboli produktioiden oikealta puolelta
 2. Poistetaan ϵ -produktiot
 3. Poistetaan yksikköproduktiot $A \rightarrow B$
 4. Poistetaan produktiot $A \rightarrow X_1 X_2 \dots X_k$, $k > 2$
- Produktioiden oikealla puolella olevien lähtösymbolien poistaminen
 - Jos lähtösymboli S on jonkin produktion oikealla puolella lisätään uusi lähtösymboli S' ja sille produktio $S' \rightarrow S$

ϵ -produktioiden poistaminen

Olkoon $G = (V, \Sigma, P, S)$ kontekstiton kielioppi. Välike $A \in V - \Sigma$ on *tyhjentyvä* (engl. nullable), jos $A \xRightarrow{G}^* \epsilon$

Lemma:

Mistä tahansa kontekstittomasta kieliopista G voidaan muodostaa ekvivalentti kielioppi G' , jossa enintään lähtösymboli on tyhjentyvä.
(Huom! Lähtösymbolin tyhjentymistä ei voida välttää, jos $\epsilon \in L(G)$.)
Todistus: Olkoon $G = (V, \Sigma, P, S)$.

1. Selvitetään ensin G :n tyhjentyvät välitteet (NULL-joukko) seuraavasti:
 - (i) asetetaan aluksi

$$\text{NULL} := \{A \in V - \Sigma \mid A \rightarrow \epsilon \text{ on } G\text{:n produktio}\};$$

- (ii) toistetaan sitten seuraavaa NULL-joukon laajennusoperaatiota, kunnes joukko ei enää kasva:

$$\begin{aligned} \text{NULL} &:= \text{NULL} \cup \\ &\quad \{A \in V - \Sigma \mid A \rightarrow B_1 \dots B_k \text{ on } G\text{:n produktio,} \\ &\quad B_i \in \text{NULL kaikilla } i = 1, \dots, k\}. \end{aligned}$$

2. Tämän jälkeen korvataan kukin G :n produktio $A \rightarrow X_1 \dots X_k$ kaikkien sellaisten produktioiden joukolla, jotka ovat muotoa

$$A \rightarrow \alpha_1 \dots \alpha_k, \quad \text{missä } \alpha_i = \begin{cases} X_i, & \text{jos } X_i \notin \text{NULL}; \\ X_i \text{ tai } \epsilon, & \text{jos } X_i \in \text{NULL}. \end{cases}$$

3. Lopuksi poistetaan kaikki muotoa $A \rightarrow \epsilon$ olevat produktiot. Jos poistettavana on myös produktio $S \rightarrow \epsilon$, otetaan muodostettavaan kielioppiin G' uusi lähtösymboli S' ja sille produktiot $S' \rightarrow S$ ja $S' \rightarrow \epsilon$. \square

Esim. Poistetaan ϵ -produktiot seuraavasta kieliopista:

$$\begin{array}{l} S \rightarrow A \mid B \\ A \rightarrow aBa \mid \epsilon \\ B \rightarrow bAb \mid \epsilon \end{array} \quad \Rightarrow \quad (\text{NULL} = \{A, B, S\})$$

$$\begin{array}{l} S \rightarrow A \mid B \mid \epsilon \\ A \rightarrow aBa \mid aa \mid \epsilon \\ B \rightarrow bAb \mid bb \mid \epsilon \end{array} \quad \Rightarrow$$

$$\begin{array}{l} S' \rightarrow S \mid \epsilon \\ S \rightarrow A \mid B \\ A \rightarrow aBa \mid aa \\ B \rightarrow bAb \mid bb \end{array}$$

Yksikköproduktioiden poistaminen

- Produktio muotoa $A \rightarrow B$, missä A ja B ovat välitteitä, on *yksikköproduktio* (engl. unit production)

Lemma:

Mistä tahansa kontekstittomasta kieliopista G voidaan muodostaa ekvivalentti kielioppi G' , jossa ei ole yksikköproduktioita.

Todistus: Olkoon $G = (V, \Sigma, P, S)$.

1. Selvitetään ensin G :n kunkin välitteen "yksikköseuraajat" seuraavasti:

- (i) asetetaan aluksi kullekin $A \in V - \Sigma$:

$$F(A) := \{B \in V - \Sigma \mid A \rightarrow B \text{ on } G\text{:n produktio}\};$$

- (ii) toistetaan sitten seuraavia F -joukkojen laajennusoperaatiota, kunnes joukot eivät enää kasva:

$$F(A) := F(A) \cup \bigcup \{F(B) \mid A \rightarrow B \text{ on } G\text{:n produktio}\}.$$

2. Tämän jälkeen poistetaan G :stä kaikki yksikköproduktiot ja lisätään niiden sijaan kaikki mahdolliset produktiot muotoa $A \rightarrow \omega$, missä $B \rightarrow \omega$ on G :n ei-yksikköproduktio jollakin $B \in F(A)$. \square

Esim. Poistetaan yksikköproduktiot edellä saadusta kieliopista

$$\begin{aligned} S' &\rightarrow S \mid \epsilon \\ S &\rightarrow A \mid B \\ A &\rightarrow aBa \mid aa \\ B &\rightarrow bAb \mid bb. \end{aligned}$$

Välikkeiden yksikköseuraaajat ovat: $F(S') = \{S, A, B\}$, $F(S) = \{A, B\}$, $F(A) = F(B) = \emptyset$. Korvaamalla yksikköproduktiot edellä esitetyllä tavalla saadaan kielioppi

$$\begin{aligned} S' &\rightarrow aBa \mid aa \mid bAb \mid bb \mid \epsilon \\ S &\rightarrow aBa \mid aa \mid bAb \mid bb \\ A &\rightarrow aBa \mid aa \\ B &\rightarrow bAb \mid bb \end{aligned}$$

Produktioiden $A \rightarrow X_1 \dots X_k$, $k > 2$ poisto

- halutaan poistaa oikealta puolelta merkkijonot, joissa on enemmän kuin kaksi välikesymbolia, ja samalla muutetaan joukossa olevat päätesymbolit sopiviksi välikkeiksi (siis esim. $A \rightarrow BbC$, $A \rightarrow abcd$ ovat tällaisia "laittomia" sääntöjä)
- Lisätään kielioppiin kutakin päätemerkkiä a varten uusi välike C_a ja sille produktio $C_a \rightarrow a$
- Korvataan kussakin muotoa $A \rightarrow X_1 \dots X_k$, $k > 2$, olevassa produktiossa ensin kaikki päätemerkit em. uusilla välikkeillä, ja sitten koko produktio produktiojoukolla

$$\begin{aligned} A &\rightarrow X_1 A_1 \\ A_1 &\rightarrow X_2 A_2 \\ &\vdots \\ A_{k-2} &\rightarrow X_{k-1} X_k, \end{aligned}$$

missä A_1, \dots, A_{k-2} ovat jälleen uusia välikkeitä.

(Tarkkaan ottaen uusi produktiojoukko on siis oikeastaan

$$\begin{aligned} A &\rightarrow X'_1 A_1 \\ A_1 &\rightarrow X'_2 A_2 \\ &\vdots \\ A_{k-2} &\rightarrow X'_{k-1} X'_k, \end{aligned}$$

missä

$$X'_i = \begin{cases} X_i, & \text{jos } X_i \in V - \Sigma; \\ C_a, & \text{jos } X_i = a \in \Sigma \end{cases}$$

Lause: Mistä tahansa kontekstittomasta kieliopista G voidaan muodostaa ekvivalentti Chomskyn normaalimuotoinen kielioppi G' .

Todistus: Olkoon $G = (V, \Sigma, P, S)$. Poistetaan ensin G :stä lähtösymboli produktioiden oikealta puolelta, ϵ -produktiot ja yksikköproduktiot em. lemموjen konstruktiolla. Tämän jälkeen kaikki G :n produktiot ovat muotoa $A \rightarrow a$ tai $A \rightarrow X_1 \dots X_k$, $k \geq 2$ (tai $S \rightarrow \epsilon$). Viimeksi mainitut poistetaan ed. konstruktion mukaisesti. \square

Huom! Turhat välitteet ja produktiot voidaan aina poistaa.

Esim. Muunnetaan Chomskyn normaalimuotoon kielioppi

$$\begin{aligned} S &\rightarrow aBCd \mid bbb \\ B &\rightarrow b \\ C &\rightarrow c \end{aligned}$$

Tuloksena saadaan kielioppi

$$\begin{aligned} S &\rightarrow C_a S_1^1 \\ S_1^1 &\rightarrow B S_2^1 \\ S_2^1 &\rightarrow C C_d \\ S &\rightarrow C_b S_1^2 \\ S_1^2 &\rightarrow C_b C_b \\ B &\rightarrow b \\ C &\rightarrow c \\ C_a &\rightarrow a \\ C_b &\rightarrow b \\ (C_c &\rightarrow c) \\ C_d &\rightarrow d \end{aligned}$$

4.9 *Attribuuttikieliopit

4.9.1 Perusidea

Attribuuttikieliopit ovat kätevä tapa ilmaista, mitä toimintoja jäsennyksen yhteydessä suoritetaan. Ne siis lisäävät puhtaasti syntaktiseen kielioppiin semantiikan (merkityksen). Attribuuttikieliopit ovat tarpeellisia esimerkiksi kääntäjissä, joissa lähdekoodin jäsenitys (syntaktisen oikeellisuuden tarkistus) ei riitä, vaan koodi pitäisi myös kääntää suoritettavaksi ohjelmaksi. Ajan puutteen takia emme kuitenkaan käsittele attribuuttikielioppeja kurssillamme, mutta seuraavssa on esitly niiden idea lyhyesti.

- liitetään kontekstittomiin kielioppeihin kielen semantiikan kuvausta
- esim. lauseke $3 + 2 \times 1$ noudattaa kieliopin G_{expr} syntaksia (ja kuuluu siis G_{expr} :in tuottamaan kieleen) – mutta mitä lauseke merkitsee?
- lausekkeen arvon evaluoiminen edellyttää sopivien attribuuttien ja niiden evaluointisääntöjen liittämistä kielioppiin
- Idea:
 - kieliopin mukaisen jäsennyksipuun solmu, jonka nimi on symboli $X \sim$ tietue tyyppiä X .
 - tietue tyyppiin X kuuluvat kentät $\sim X$:n *attribuutit*, merk. $X.s$, $X.t$ jne.
 - kussakin X -tyyppisessä jäsennyksipuun solmussa X :n attribuuteista eri *ilmentymät*
 - produktioihin $A \rightarrow X_1 \dots X_k$ liitetään attribuuttien *evaluointisääntöjä*, jotka ilmaisevat miten annetun jäsennyksipuun solmun attribuutti-ilmentymien arvot määräytyvät sen vanhempi- ja lapsisolmujen attribuutti-ilmentymien arvoista.
- Huom! Säännöt voivat olla periaatteessa minkälaisia funktioita tahansa, kunhan niiden argumentteina esiintyy vain paikallisesti saatavissa olevaa tietoa.
- ts. produktioon $A \rightarrow X_1 \dots X_k$ liitettävissä säännöissä saa mainita vain symbolien A, X_1, \dots, X_k attribuutteja

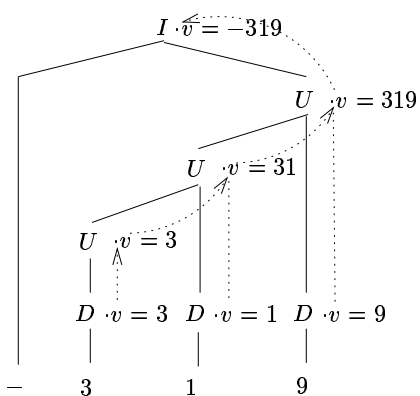
Esim. 1 Liitetään etumerkillisiä kokonaislukuja tuottavaan kontekstittomaan kielioppiin attribuutit ja niiden evaluointisäännöt kieliopin tuottamien lukujen arvojen määrittämiseen.

- kuhunkin jäsenyspuun X -tyyppiseen välikesolmuun liitetään attribuutti-ilmentymä $X.v$, jonka arvoksi tulee X :stä tuotetun numerojonon lukuarvo
- juurisolmun v -ilmentymän arvoksi tulee koko puun tuotoksena olevan numerojonon lukuarvo

<i>Produktiot:</i>	<i>Evaluointisäännöt:</i>
$I \rightarrow +U$	$I.v := U.v$
$I \rightarrow -U$	$I.v := -U.v$
$I \rightarrow U$	$I.v := U.v$
$U \rightarrow D$	$U.v := D.v$
$U \rightarrow UD$	$U_1.v := 10 * U_2.v + D.v$
$D \rightarrow 0$	$D.v := 0$
$D \rightarrow 1$	$D.v := 1$
\vdots	
$D \rightarrow 9$	$D.v := 9$

- produktion evaluointisäännössä saman välikkeen eri esiintymät voidaan erottaa alaindeksillä
- yllä U_1 vastaa ensimmäistä produktiossa $U \rightarrow UD$ esiintyvää U :ta ja U_2 toista

Evaluointisääntöjen mukainen *attributoitu jäsenyspuu* kieliopin tuottamalle lauseelle “-319” (katkoviivoilla on esitetty attribuutti-ilmentymien väliset evaluointiriippuvuudet):



Kuva 4.12: Attributoitu jäsenyspuu.

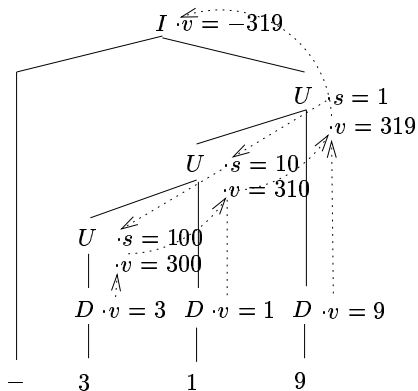
- Attribuuttikieliopin attribuutti t on *synteettinen*, jos sen kuhunkin produktion $A \rightarrow X_1 \dots X_k$ liittyvä evaluointisääntö on muotoa $A.t := f(A, X_1, \dots, X_k)$.
- Tällöin jäsenyspuussa kunkin solmun mahdollisen t -ilmentymän arvo riippuu vain solmun omien ja sen jälkeläisten attribuutti-ilmentymien arvoista
- Muunlaiset attribuutit ovat *periytyviä*
- yllä esim. attribuutti v on synteettinen
- pyritään käyttämään pääasiassa synteettisiä attribuutteja, koska ne voidaan evaluoida helposti yhdellä jäsenyspuun lehdistä juureen suuntautuvalla läpikäynnillä.
- perittyjä attribuutteja voidaan käyttää, kunhan attribuutti-ilmentymien riippuvuusverkkoihin ei tule syklejä

Esimerkiksi edellä olevaan, etumerkillisiä kokonaislukuja tuottavaan kielioppiin voitaisiin liittää lukujen arvot määrittävä semantiikka myös seuraavasti, periytyvää "positiokerroin"-attribuuttia s ja synteettistä "arvo"-attribuuttia v käyttäen:

<i>Produktiot:</i>	<i>Evaluointisäännöt:</i>	
$I \rightarrow +U$	$U.s := 1,$	$I.v := U.v$
$I \rightarrow -U$	$U.s := 1,$	$I.v := -U.v$
$I \rightarrow U$	$U.s := 1,$	$I.v := U.v$
$U \rightarrow D$		$U.v := (D.v) * (U.s)$
$U \rightarrow UD$	$U_2.s := 10 * (U_1.s),$	$U_1.v := U_2.v + (D.v) * (U_1.s)$
$D \rightarrow 0$		$D.v := 0$
$D \rightarrow 1$		$D.v := 1$
\vdots		
$D \rightarrow 9$		$D.v := 9$

Kuvassa 4.13 on esitetty tämän semantiikan mukainen attributoitu jäsenyspuu lauseelle "-319".

Attribuutti-ilmentymien arvot voidaan usein laskea suoraan jäsenysrutiineissa, tarvitsematta muodostaa jäsenyspuuta eksplisiittisesti.



Kuva 4.13: Positiokerrointa käyttäen attributoitu jäsenyspuu.

4.9.2 Aritmeettinen laskin

Esim. 2 Laaditaan aritmeettinen laskin, jolla voi suorittaa etumerkillisten kokonaislukujen yhteen-, kerto- ja vähennyslaskuoperaatioita. Selvyyden vuoksi vaaditaan, että etumerkillisten lukujen ympärillä on aina sulut. Kielioppi on siis seuraava:

$$\begin{aligned}
 E &\rightarrow T + E \mid T - E \mid T \\
 T &\rightarrow F \mid F * T \\
 F &\rightarrow DN \mid D(+DN) \mid (-DN) \mid (E) \\
 N &\rightarrow DN \mid D \\
 D &\rightarrow 0 \mid 1 \mid \dots \mid 9
 \end{aligned}$$

Esitetään lausekkeen arvon evaluointi attribuuttikielioppina:

<i>Produktiot:</i>	<i>Evaluointisäännöt:</i>
$E \rightarrow T + E$	$E_1.v := T.v + E_2.v$
$E \rightarrow T - E$	$E_1.v := T.v - E_2.v$
$E \rightarrow T$	$E.v := T.v$
$T \rightarrow F * T$	$T_1.v := F.v \times T_2.v$
$T \rightarrow F$	$T.v := F.v$
$F \rightarrow DN$	$F.v := 10 \times D.v + N.v$
$F \rightarrow D$	$F.v := D.v$
$F \rightarrow (+DN)$	$F.v := 10 \times D.v + N.v$
$F \rightarrow (-DN)$	$F.v := -1 \times (10 \times D.v + N.v)$
$F \rightarrow (E)$	$F.v := E.v$
$N \rightarrow DN$	$N_1.v := 10 \times D.v + N_2.v$
$N \rightarrow D$	$N.v := D.v$
$D \rightarrow 0$	$D.v := 0$
...	...
$D \rightarrow 9$	$D.v := 9$

Ohjelmassa ei lukuja tarvitse lukea merkki kerrallaan, vaan voimme käyttää jotain sopivaa apurutiinia. Välikkeitä N ja D sisältävissä säännöissä voidaan siis oikaista.

Aritmeettinen kokonaislukulaskin pseudokoodina:

```

/*Apurutiinit: */
char getnext; /* lue seuraava merkki */
int isdigit(char c); /* onko c numeromerkki? */
int readnumber; /* lue koko luku ja palauta sen arvo */
void ERROR; /* virheenkäsittely */
/* Seuraava merkki: */
char next;

int function E {
int op1, op2;
op1 = T;
if(next == '+' || next == '-') { /* rule E → T + E */
    next = getnext;
    op2 = E;
    return op1 + op2;
}
}

```



```

else
  if (next == ' -') { /* rule  $E \rightarrow T - E$  */
    next = getnext;
    op2 = E;
    return op1 - op2;
  }
  else return op1; /* rule  $E \rightarrow T$  */
}

```

```

int function T {
int op1, op2;
op1 = F;
if (next == '*') { /* rule  $T \rightarrow F * T$  */
  next = getnext;
  op2 = T;
  return op1 * op2;
}
else return op1; /* rule  $T \rightarrow F$  */
}

```

```

int function F {
int op
if (isdigit(next)) { /* rule  $F \rightarrow DN$  */
  op = readnumber;
  next = getnext;
  return op;
}
else
  if (next == '(') {
    next = getnext;
    if (next == '+') { /* rule  $F \rightarrow (+DN)$  */
      op = readnumber;
      if (next != ')') ERROR;
      next = getnext;
      return op;
    }
    else { /* rule  $F \rightarrow (-DN)$  */
      if (next == '-') {

```

```
        next = getnext;
        op = readnumber;
        if (next! =')') ERROR;
        next = getnext;
        return -1 * op;
    }
    else{ /* rule  $F \rightarrow (E)$  */
        next = getnext;
        op = E;
        if (next! =')') ERROR;
        next = getnext;
        return op;
    }
}
}
else ERROR;
}
```

4.10 Kontekstittomien kielten ominaisuuksia

4.10.1 Kontekstittomien kielten sulkeumaominaisuudet

Kontekstittomille kielille pätee joitain samantapaisia sulkeumaominaisuuksia kuin säännöllisille kielille.

Lause Olkoon L_1 ja L_2 kontekstittomia kieliä. Tällöin myös

1. $L_1 \cup L_2$ (kielten yhdiste)
2. $L_1 L_2$ (kielten katenaatio)
3. $(L_1)^*$ (kielen sulkeuma)
4. $(L_1)^R$ (kielen käänteiskieli)

ovat kontekstittomia.

Huom! Kontekstittomat kielet eivät kuitenkaan ole suljettuja leikkauksen ja komplementin suhteen! Jos L_1 ja L_2 ovat kontekstittomia, ei $L_1 \cap L_2$ silti ole välttämättä kontekstiton! Samoin $\overline{L_1} = \Sigma^* \setminus L_1$ (kielen komplementti) ei välttämättä ole kontekstiton.

Esim. kielet $L_1 = \{a^n b^n c^k \mid n, k = 0, 1, \dots\}$ ja $L_2 = \{a^k b^n c^n \mid k, n = 0, 1, \dots\}$ ovat kontekstittomia. Kuitenkaan kieli $L_1 \cap L_2 = \{a^n b^n c^n \mid n = 0, 1, \dots\}$ ei ole kontekstiton!

Toisaalta, jos L on kontekstiton kieli ja R on säännöllinen kieli, niin $L \cap R$ on kontekstiton. (Ks. esim. Hopcroft, Motwani, Ullman, teoreema 7.27.)

4.10.2 Ratkeavat ja ratkeamattomat ominaisuudet

Edellä olemme jo oppineet joitain kontekstittomien kielten ominaisuuksia, joita voidaan ratkaista tietokoneella. Tärkein tällainen ominaisuus koskee jäsentämistä: mille tahansa merkkijonolle x ja annetulle kontekstittomalle kieliopille G voimme ratkaista, kuuluuko x kieleen $L(G)$. Ts. ongelma $x \in L(G)$ on algoritmisesti ratkeava. Mikä jäsenysmenetelmä on kulloinkin paras, riippuu annetusta kielestä. Jos kielioppi voidaan esittää LL(1)-muodossa, tarjoaa tämä tehokkaan (ajassa $O(n)$ toimivan) ja helpon jäsenysmenetelmän. Ohjelmointikielten kääntäjät puolestaan käyttävät LR(k)-jäsentäjiä, kuten Unixin yacc- ja bison-ohjelmien tuottamia jäsentäjiä. Mielivaltaisen kieliopin taas voimme aina muuntaa Chomskyn normaalimuotoon ja jäsentää CYK-algoritmillä ajassa $O(n^3)$.

Eräät tärkeistä kontekstittomien kielten ominaisuuksista ovat kuitenkin ratkeamattomia ongelmia. Tällaisia ongelmia ovat:

- Kahden kieliopin ekvivalenssi ts. $L(G_1) = L(G_2)$?
- Annetun kieliopin moniselitteisyys $Ambiguous(G)$?
- Onko annetun mv. kieliopin kuvaama kieli kontekstiton $Context-free(L(G))$?

Viimeiseen mäistä ongelmista liittyy kontekstittomien kielten pumppauslemma, jolla ihminen voi osoittaa joitain kieliä ei-kontekstittomiksi samaan tapaan kuin säännöllisten kielten pumppauslemmalla. Tämäkään pumppauslemma ei kuitenkaan anna kontekstittomuuden välttämättömiä vaan ainoastaan riittävät ehdot, eikä pumppauslemmaa voida soveltaa algoritmisesti. Ongelma pysyy siis yleisessä tapauksessa ratkeamattomana.

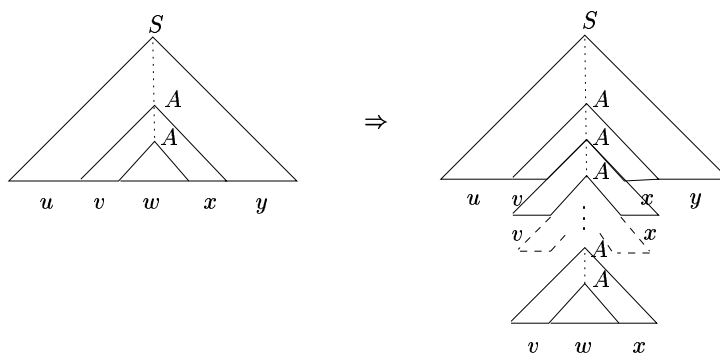
4.11 *Kontekstittomien kielten rajoituksista

- Kontekstittomille kielille voidaan todistaa samantapainen pumppauslemma kuin säännöllisillekin kielille
- Erona on vain se, että nyt merkkijonoa (osat $uvwxy$) on pumpattava samanaikaisesti kahdesta paikasta (osista v ja x)
- *Lemma* [Kontekstittomien kielten pumppauslemma eli $uvwxy$ -lemma]: Olk. L kontekstiton kieli. Tällöin on olemassa sellainen $n \geq 1$, että mikä tahansa $z \in L$, $|z| \geq n$, voidaan jakaa osiin $z = uvwxy$ siten, että
 - (i) $|vx| \geq 1$,
 - (ii) $|vwx| \leq n$,
 - (iii) $uv^iwx^iy \in L$ kaikilla $i = 0, 1, 2, \dots$

Huom: Taas vain äärettömät kielet ovat kiinnostavia

Todistus:

- Olk. $G = (V, \Sigma, P, S)$ Chomskyn normaalimuotoinen kielioppi L :lle. Tällöin missä tahansa G :n jäsenyyspuussa, jonka korkeus (= pisimmän juuresta lehteen kulkevan polun pituus) on h , on enintään 2^h lehteä. Ts. mikä tahansa merkkijonon $z \in L$ jokaisessa jäsenyyspuussa on polku, jonka pituus on vähintään $\log_2 |z|$



Kuva 4.14: Kontekstittoman kielen merkkijonon pumppaus.

- Olk. $k = |V - \Sigma|$ kieliopin G välikkeiden määrä. Asetetaan $n = 2^{k+1}$. Tarkastellaan jotakin $z \in L$, $|z| \geq n$, ja sen jotakin jäsenyspuuta
- \Rightarrow puussa on polku, jonka pituus on $\geq k + 1$. Tarkastellaan tämän polun “alinta” $k + 1$:n pituista osaa. Tällä polulla, jossa on $k + 1$ välikettä vastaavaa solmua ja välikkeitä on k kappaletta, on siis jonkin välikkeen toistuttava. Olkoon A joku polun toistuva välike (ks. kuva 4.14).
- Merkkijono z voidaan nyt osittaa $z = uvwxy$, missä w on A :n alimasta ilmentymästä tuotettu osajono ja vwx seuraavaksi ylemmästä A :n ilmentymästä tuotettu osajono; osajonot saadaan johdosta

$$S \Rightarrow^* uAy \Rightarrow^* uvAxy \Rightarrow^* uvwxy$$

- Koska $S \Rightarrow^* uAy$, $A \Rightarrow^* vAx$ ja $A \Rightarrow^* w$, osajonoja v ja x voidaan “pumppata” w :n ympärillä:

$$S \Rightarrow^* uAy \Rightarrow^* uvAxy \Rightarrow^* uv^2Ax^2y \Rightarrow^* \dots \Rightarrow^* uv^iAx^i y \Rightarrow^* uv^iwx^i y$$

- $\Rightarrow uv^iwx^i y \in L$ kaikilla $i = 0, 1, 2, \dots$
- Koska kielioppi G on Chomskyn normaalimuodossa ja $A \Rightarrow^* vAx$, on oltava $|vx| \geq 1$
- Ylemmästä A :n ilmentymästä alkavan jäsenyspuun polun pituus on enintään $k + 1 \Rightarrow$ vastaavan alipuun tuotokselle $|vwx| \leq 2^{k+1} = n$. \square
- Esim. Väite: kieli $L = \{a^k b^k c^k \mid k \geq 0\}$ ei ole kontekstiton. Tod. Vastaväite: L on kontekstiton. Valitaan parametri n lemmän mukaisesti ja tarkastellaan merkkijonoa $z = a^n b^n c^n \in L$. Tällöin z voidaan jakaa pumpattavaksi osiin

$$z = uvwxy, \quad |vx| \geq 1, \quad |vwx| \leq n.$$

Tällöin merkkijono vx ei voi sisältää sekä a :ta, b :tä että c :tä. Merkkijonossa $uv^0wx^0y = uwy$ on siten ylijäämä jotakin merkkiä muihin merkkeihin nähden eli $uwy \notin L$. Ristiriita. \square

4.12 Kontekstittomien kielten sovelluksia

Jäsentäjät

Ohjelmointikieli voidaan kuvata kontekstittomana kielioppina, jota voidaan käsitellä jäsentäjällä (*parser*). Jäsentäjä on kääntäjän osa, joka tutkii ohjelman rakenteen ja esittää sen jäsennyksipuuna. Esim. UNIX tarjoaa komennon `yacc`, joka generoi jäsentäjän annetusta kontekstittomasta kieliopista `yacc`:ille annetaan syötteenä kontekstiton kielioppi (`yacc`-notaatiossa säännön nuoli on korvattu kaksoispisteellä), jonka jokaiseen sääntöön voidaan liittää C-koodina määritelty toiminto. Toiminto suoritetaan aina, kun luodaan vastaava jäsennyksipuun solmu. Tyypillisesti toiminto vain luo jäsennyksipuun solmun, mutta joissain sovelluksissa jäsennyksipuuta ei eksplisiittisesti luoda, vaan toiminto tekee jotain muuta, esim. sisällyttää palan objektikoodia kyseiseen paikkaan.

Esimerkiksi seuraava kielioppi:

$$E \rightarrow I|E + E|E * E|(E)$$

$$I \rightarrow a|b|Ia|Ib|I0|I1$$

annettaisiin `yacc`:ille muodossa:

```

Exp  : Id          {...}
      | Exp' +' Exp {...}
      | Exp' *' Exp {...}
      | '('Exp')'  {...}
      ;
Id    : 'a'         {...}
      | 'b'         {...}
      | Id 'a'      {...}
      | Id 'b'      {...}
      | Id '0'      {...}
      | Id '1'      {...}
      ;

```

missä {...} sisältää jonkin toiminnon C-koodina.

Dokumentin rakenteen kuvaus

- ns. ”mark-up”-kielet kuten HTML ja XML

- kielen merkkijonot dokumentteja, joissa tiettyjä kielen semantiikkaa kuvaavia merkkejä (tags) — esim. `` ja `` ~ järjestetty lista
- HTML: laillisen HTML-dokumentin kuvaus ja dokumentin prosessoinnin ohjaus
- XML: kuvaa tekstin semantiikan — esim.
`<ADDR>Perhospolu 17 C 3, 33000 Kukkamäki</ADDR>` kertoo, että osajono on osoite.