

Limbaje Formale și Tehnici de Compilare

Drăgan Mircea

April 19, 2006

P R E F A Ț Ă

Materialul prezentat constituie notele de curs ținute studenților din primii doi ani la Secția Informatică a Facultății de Matematică și Informatică de la Universitatea de Vest Timișoara.

În primul capitol, *Introducere* se trec în revistă noțiunile fundamentale din teoria limbajelor formale și se prezintă problema generală a compilării.

Capitolul al doilea, *Limbaje regulate și analiza lexicală*, prezintă chestiunile teoretice fundamentale privind teoria automatelor și echivalența cu limbajele de tipul trei, proprietăți speciale și mecanisme echivalente cu automatele finite. Sunt prezentate de asemenea și principiile analizei lexicale, finalizate cu realizarea unui analizor lexical.

Capitolul al treilea, *Limbaje independente de context* este dedicat studiului teoretic al mecanismelor de generare și de recunoaștere a limbajelor de tipul doi. Sunt prezentate formele normale ale gramaticilor independente de context și câteva proprietăți speciale.

Capitolul patru, *Analiza sintactică* este dedicat prezentării principiilor generale de analiză sintactică și a algoritmilor specializați de analiză. Pentru fiecare tip de algoritm analizat s-au considerat exemple practice de aplicare.

Capitolul al cincilea, *Sinteza programelor* tratează formele intermediare uzuale pentru traducerea programelor și generarea codului obiect pornind de la formatul intermediar. Pentru cazul expresiilor aritmetice sunt prezentați și algoritmi direcți de generarea a formatului intermediar, împreună cu proceduri standard de optimizare a codului generat.

Ultimul capitol, *Mașina Turing* prezintă succint mecanismul formal ce definește modelul de calculabilitate.

La sfârșitul fiecărui capitol sunt date câteva exerciții (nerezolvate), aplicații directe la chestiunile teoretice prezentate. Acestea pot fi parcurse în cadrul orelor de seminar și laborator.

În expunerea algoritmilor s-a folosit un limbaj mai puțin rigid, fără prea multe reguli stricte. În general s-a urmărit descrierea cât mai simplă a structurilor care apar în text.

Cuprins

1	Introducere	5
1.1	Limbae formale	5
1.2	Gramatici generative de tip Chomsky	8
1.3	Ierarhia Chomsky	12
1.4	Traducerea programelor	17
1.5	Probleme propuse	20
2	Limbae Regulate și Analiza Lexicală	23
2.1	Automate finite și limbae regulate	23
2.2	Proprietăți speciale ale limbajelor regulate	29
2.3	Expresii regulate și sisteme tranzitionale	34
2.4	Analiza lexicală	37
2.5	Probleme propuse	48
3	Limbae Independente de Context	51
3.1	Arbori de derivare	51
3.2	Ambiguitate în familia \mathcal{L}_2	54
3.3	Forme normale pentru gramatici de tipul 2	56
3.4	Lema Bar-Hillel	62
3.5	Automate push-down (APD)	65
3.6	Automate push-down deterministe	73
3.7	Probleme propuse	76
4	Analiza Sintactică	79
4.1	Algoritmi TOP-DOWN	80
4.1.1	Algoritmul general de analiză top-down	80
4.1.2	Analiza top-down fără reveniri	81
4.1.3	Programarea unui analizor sintactic. Studiu de caz	84
4.2	Algoritmi BOTTOM-UP	90
4.2.1	Gramatici cu precedență simplă	90
4.2.2	Relații de precedență	92
4.2.3	Proprietăți ale gramaticilor cu precedență simplă	93

4.2.4	Determinarea relațiilor de precedență pentru gramatici cu precedență simplă	95
4.2.5	Studiu de caz	95
4.2.6	Gramatici operatoriale	97
4.2.7	Gramatici operatoriale	100
4.2.8	Determinarea relațiilor de precedență pentru gramatici operatoriale	103
4.2.9	Studiu de caz	104
5	Sinteza Programelor	107
5.1	Forme interne ale programelor	107
5.1.1	Tabelele compilatorului	107
5.1.2	Cvadruple și triplete	109
5.2	Generarea formatului intermediar	114
5.2.1	Generarea cvadruplelor	115
5.2.2	Generarea tripletelor	117
5.2.3	Generarea șirului polonez	118
5.3	Generarea formatului intermediar pentru instrucții clasice	121
5.3.1	Instrucțiunea de atribuire	121
5.3.2	Instrucțiunea If	121
5.3.3	Variabile indexate	121
6	Masina Turing	123
6.1	Limbaje de tipul zero	123
6.2	Mașina Turing	126

Capitolul 1

Introducere

1.1 Limbaje formale

Noțiunea generală de limbaj. Se numește **alfabet** sau vocabular orice mulțime finită și nevidă. Elementele unui alfabet V le vom numi **simboluri** (sau litere, caractere, variabile).

Definiție 1.1 *Un cuvânt peste un alfabet V este o secvență $p = a_1 a_2 \dots a_n, a_i \in V, i = 1, \dots, n$.*

Numărul n , deci numărul simbolurilor cuvântului p , se numește *lungimea cuvântului* și va fi notat cu $|p|$ sau $l(p)$. Vom considera și *cuvântul vid* λ sau e , care nu conține nici un simbol; evident $|\lambda| = 0$.

Noțiunea de *cuvânt* este fundamentală în teoria limbajelor formale sau în alte domenii ale informaticii; termeni sinonimi utilizați în literatură sunt *propoziție*, *frază* sau *șir*. Să observăm că nu există o similitudine foarte bună între noțiunile de "alfabet", "cuvânt", etc. din teoria limbajelor formale și noțiunile corespunzătoare din lingvistică.

Mulțimea tuturor cuvintelor peste un alfabet V o notăm cu V^+ . Această mulțime împreună cu cuvântul vid va fi notată cu V^* . În general, vom utiliza litere mari de la sfârșitul alfabetului pentru notarea diverselor alfabete, U, V, W , etc.; litere de la începutul alfabetului (mari sau mici) pentru notarea simbolurilor, $A, B, C, \dots, a, b, c, \dots, i, j, \dots$ (uneori cifre $0, 1, 2, \dots$); pentru notarea cuvintelor vom utiliza litere mici de la sfârșitul alfabetului, $p, q, r, s, t, u, v, w, x, y, z$, etc. (această convenție de notare nu va fi absolută).

Fie $p = a_1 \dots a_n, q = b_1 \dots b_m$. Definim pe mulțimea V^* operația de *concatenare* sau *juxtapunere* prin $pq = a_1 \dots a_n b_1 \dots b_m$.

Se poate verifica ușor că această operație este asociativă. Prin urmare, mulțimea V^+ înzestrată cu această operație este un semigrup (semigrupul liber peste V). Mulțimea V^* cu aceiași operație este un semigrup cu unitate, deci un monoid, unitatea fiind cuvântul vid (monoidul liber peste V). Structura are și

proprietatea de simplificare la stânga și la dreapta, adică:

$$ua = ub \Rightarrow a = b, \text{ respectiv, } au = bu \Rightarrow a = b, \forall a, b, u \in V^*$$

Fie din nou $p, q \in V^*$. Vom spune că q este un subcuvânt sau un *infix* (propriu) al lui p dacă $p = uqv$, $u, v \in V^*$ ($u, v \in V^+$); q este *prefix* (propriu) al lui p dacă $p = qv$, $v \in V^*$ ($v \in V^+$); q este *suffix* (propriu) al lui p dacă $p = uq$, $u \in V^*$ ($u \in V^+$).

Definiție 1.2 *Un limbaj L peste alfabetul V este o parte a mulțimii tuturor cuvintelor peste V , deci $L \subseteq V^*$.*

Să observăm că V^* (sau V^+) este întotdeauna o mulțime infinită (evident numărabilă); în această accepțiune generală, un limbaj poate să fie o mulțime finită sau infinită, uneori chiar vidă.

Exemplu. Fie $V = \{0, 1\}$. Avem

$$V^+ = \{0, 1, 00, 01, 10, 000, \dots\},$$

$$V^* = \{\lambda, 0, 1, 00, 01, 10, 000, \dots\}.$$

Limbaje peste alfabetul V sunt de exemplu mulțimile

$$L_1 = \{\lambda, 00, 11\},$$

$$L_2 = \{1, 11, 111, \dots\} = \{1^n | n \geq 1\}.$$

Observație. Notăția a^n , unde a este un simbol al unui alfabet, înseamnă cuvântul constituit din n simboluri a , adică $\underbrace{aa \dots a}_n$. În particular $a^0 = \lambda$.

Operații cu limbaje Limbajele fiind mulțimi se pot efectua cu limbaje operațiile obișnuite cu mulțimi: *reuniune*, *intersecție*, *diferență*, *complementariere* (*față de V^**). Există și operații specifice limbajelor.

În general, o operație de n -aritate oarecare (cel mai adesea binară sau unară) pe mulțimea V^* definește o operație corespunzătoare pe mulțimea limbajelor. Astfel, dacă

$$\alpha : V^* \longrightarrow \mathcal{P}(V^*) \text{ și } \beta : V^* \times V^* \longrightarrow \mathcal{P}(V^*)$$

sunt două operații pe V^* (unară și respectiv binară) și L_1, L_2 sunt două limbaje peste V , putem defini limbajele $\alpha(L_1)$ respectiv $\beta(L_1, L_2)$ prin

$$\alpha(L_1) = \bigcup_{x \in L_1} \alpha(x), \quad \beta(L_1, L_2) = \bigcup_{x \in L_1, y \in L_2} \beta(x, y).$$

Exemple:

1. *Produsul* (concatenarea) a două limbaje definit prin

$$L_1 L_2 = \{pq \mid p \in L_1, q \in L_2\}.$$

Dacă $L_1 = L_2 = L$ vom nota $LL = L^2$. Prin recurență, se definește L^n astfel

$$L^0 = \{\lambda\}, L^k = L^{k-1}L, k \geq 1.$$

2. *Închiderea (Kleene)* a unui limbaj L este

$$L^* = \bigcup_{k=0}^{\infty} L^k.$$

3. Limbajul $Sub(L)$. Fie V^* și $Sub(x)$ mulțimea tuturor subcuvintelor lui x (evident Sub este o operație unară pe V^*). Dacă L este un limbaj peste V , putem defini limbajul

$$Sub(L) = \bigcup_{x \in L} \{Sub(x)\}.$$

adică limbajul constituit din toate subcuvintele tuturor cuvintelor lui L . Semnificații analoage vor avea și limbajele $Pref(L)$ și $Suf(L)$.

4. Limbajul $Mi(L)$. Fie $x = a_1 \dots a_n$ un cuvânt peste alfabetul V . Cuvântul $Mi(x) = a_n \dots a_1$ se numește *răsturnatul* sau *oglinditul* lui x (Mi este prescurtarea cuvântului englez *mirror*). Se mai notează $Mi(x) = \tilde{x}$. Avem atunci și răsturnatul unui limbaj

$$Mi(L) = \bigcup_{x \in L} \{Mi(x)\}.$$

5. Operația de *substituție*. Fie U și V două alfabeturi și fie aplicația $s : V \longrightarrow \mathcal{P}(U^*)$. Extindem (prelungim) această aplicație la V^* prin

$$s(\lambda) = \{\lambda\}, s(xy) = s(x)s(y), \forall x, y \in V^*.$$

O astfel de prelungire se numește canonică; ea păstrează operația de concatenare, în sensul că dacă $p = xy$, atunci $s(p) = s(x)s(y)$ este concatenarea limbajelor $s(x), s(y)$. Operația de substituție a limbajelor este dată de

$$s(L) = \bigcup_{x \in L} s(x).$$

Să observăm că această operație transformă un limbaj peste un alfabet V într-un limbaj peste un alfabet U și că păstrează operația de concatenare. Dacă $card(a) < \infty, \forall a \in V$, vom spune că substituția este finită, iar dacă $card(a) = 1, \forall a \in V$ vom spune că s este un homomorfism.

Operațiile reuniune, produs și închidere Kleene se mai numesc **operații regulate** asupra limbajelor.

1.2 Gramatici generative de tip Chomsky

Un limbaj peste un alfabet poate să fie o mulțime finită sau infinită. Dacă este o mulțime finită, el poate fi definit prin scrierea efectivă a cuvintelor limbajului. În cazul în care este o mulțime infinită, el poate fi definit în anumite cazuri punând în evidență structura cuvintelor lui. De exemplu

$$L_2 = \{01, 0011, 000111, \dots\} = \{0^n 1^n | n \geq 1\}.$$

Există două procedee mai generale pentru definirea limbajelor:

1. Procedee *generative*, care permit generarea tuturor cuvintelor limbajului. Există mai multe tipuri de mecanisme de generare a limbajelor, între care gramaticile Chomsky, sisteme Lindenmayer, etc.
2. Procedee *analitice*, care determina dacă un cuvânt dat aparține sau nu limbajului. Sunt așa-numitele *automate*, *automate finite*, *automate push-down*, etc.

Un rol deosebit în teoria limbajelor formale îl au gramaticile Chomsky.

Noțiunea de gramatică Fie V_N și V_T două alfabeturi disjuncte, $V_N \cap V_T = \emptyset$ numite respectiv *alfabetul simbolurilor neterminale* (V_N) și *alfabetul simbolurilor terminale* (V_T). Notăm $V_G = V_N \cup V_T$ *alfabetul general* și $P \subset V_G^* V_N V_G^* \times V_G^*$ *alfabetul regulilor*.

Mulțimea P va fi deci formată din perechi de forma (u, v) , unde $u = u' A u''$, $u', u'' \in V_G^*$, $A \in V_N$ iar $v \in V_G^*$, deci u și v sunt cuvinte peste V_G , cu observația că u trebuie să conțină cel puțin un simbol neterminal. Vom spune că o astfel de pereche este o **regulă** (producție, regulă de generare, regulă de rescriere) și o vom nota $u \rightarrow v$ (vom spune: u se transformă în v). Apartenența unei reguli la P o vom nota în mod obișnuit $(u \rightarrow v) \in P$, sau mai simplu, $u \rightarrow v \in P$ (nu va exista confuzia cu faptul că v este un element al lui P).

Definiție 1.3 O gramatică este un sistem $G = (V_N, V_T, X_0, P)$, unde V_N este alfabetul simbolurilor neterminale, V_T este alfabetul simbolurilor terminale, $X_0 \in V_N$ și se numește simbol de start al gramaticii, iar P este mulțimea de reguli.

Observație. Simbolurile alfabetului V_N le vom nota în general cu litere mari A, B, C, \dots, X, Y, Z (mai puțin U, V, W) iar cele ale alfabetului V_T cu litere mici de la început a, b, c, \dots sau cu cifre $0, 1, 2, \dots$.

Fie G o gramatică și $p, q \in V_G^*$. Vom spune că p se **derivează direct** în q și vom scrie $p \Rightarrow_G q$ (sau mai simplu $p \Rightarrow q$) dacă există cuvintele $r, s, u, v \in V_G^*$

astfel încât $p = rus, q = rvs$ iar $u \rightarrow v \in P$. Vom spune că p' se **derivează** în p'' (fără specificația direct) dacă există p_1, p_2, \dots, p_n $n \geq 1$ astfel încât

$$p' = p_1 \Rightarrow_G p_2 \Rightarrow_G \dots \Rightarrow_G p_n = p''.$$

Vom scrie $p' \xRightarrow[G]{+} p''$ (sau $p' \xRightarrow{+} p''$ când nu există nici o confuzie) dacă $n > 1$ și $p' \xRightarrow[G]{*} p''$ (sau $p' \xRightarrow{*} p''$) dacă $n \geq 1$. Șirul de mai sus va fi numit **derivare** iar numărul de derivări directe din șir îl vom numi *lungimea derivării*; se mai spune că p' *derivă* în p'' .

Să observăm că transformările astfel definite $\Rightarrow, \xRightarrow{+}, \xRightarrow{*}$ sunt relații pe V_G^* . Este clar că relația $\xRightarrow{+}$ este închiderea tranzitivă a relației \Rightarrow , iar relația $\xRightarrow{*}$ este închiderea tranzitivă și reflexivă a relației de transformare directă.

Definiție 1.4 . *Limbaajul generat de gramatica G este prin definiție mulțimea*

$$L(G) = \{p \in V_T^*, X_0 \xRightarrow[G]{*} p\}.$$

Observație. Dacă $p \in V_G^*$ și $X_0 \xRightarrow[G]{*} p$ se spune că p este o *formă propozițională* în gramatica G .

Exemple:

1. Fie $G = (V_N, V_T, X_0, P)$, unde $V_N = \{A\}$, $V_T = \{0, 1\}$, $X_0 = A$ (evident) și $P = \{A \rightarrow 0A1, A \rightarrow 01\}$. O derivare în această gramatică este, de exemplu

$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000111 = 0^31^3.$$

Este evident că $L(G) = \{0^n1^n | n \geq 1\}$.

Observație. În cazul în care mai multe reguli au aceeași parte stângă, le vom scrie compact astfel $u \rightarrow v_1|v_2|\dots|v_n$, simbolul $|$ având sensul de sau; în cazul nostru, $A \rightarrow 0A1|01$.

2. $G = (V_N, V_T, X_0, P)$, unde
 $V_N = \{\langle \text{propoziție} \rangle, \langle \text{subiect} \rangle, \langle \text{atribut} \rangle, \langle \text{predicat} \rangle, \langle \text{complement} \rangle,$
 $\langle \text{substantiv} \rangle, \langle \text{adjectiv} \rangle, \langle \text{verb} \rangle, \langle \text{articol} \rangle\},$
 $V_T = \{o, \text{orice}, \text{matrice}, \text{funcție}, \text{derivabilă}, \text{continuă}, \text{este}\},$
 $X_0 = \langle \text{propoziție} \rangle,$
 $P = \{\langle \text{propoziție} \rangle \rightarrow \langle \text{subiect} \rangle \langle \text{atribut} \rangle \langle \text{predicat} \rangle \langle \text{complement} \rangle,$
 $\langle \text{subiect} \rangle \rightarrow \langle \text{articol} \rangle \langle \text{substantiv} \rangle,$
 $\langle \text{atribut} \rangle \rightarrow \langle \text{adjectiv} \rangle,$
 $\langle \text{predicat} \rangle \rightarrow \langle \text{verb} \rangle,$
 $\langle \text{complement} \rangle \rightarrow \langle \text{adjectiv} \rangle,$
 $\langle \text{articol} \rangle \rightarrow o|\text{orice},$
 $\langle \text{substantiv} \rangle \rightarrow \text{matrice}|funcție,$
 $\langle \text{adjectiv} \rangle \rightarrow \text{derivabilă}|continuă,$
 $\langle \text{verb} \rangle \rightarrow \text{este}.$

Observație. În acest exemplu, " $\langle \text{propoziție} \rangle$ ", " $\langle \text{subiect} \rangle$ ", etc., reprezintă fiecare câte un simbol neterminal; de asemenea, "o", "orice", "matrice",

etc., reprezintă simboluri terminale. Se poate ușor observa că această gramatică generează propoziții simple de forma subiect-atribut-predicat-complement care exprimă judecăți asupra conceptelor de "matrice" și "funcție". De exemplu, se poate forma propoziția: "orice funcție derivabilă este continuă", care este din punct de vedere semantic corectă, precum și propoziția "orice funcție continuă este derivabilă", care, după cum se știe, este falsă. Evident, se pot genera și numeroase propoziții care nu au sens. Ceea ce ne interesează în acest moment este aspectul formal, deci din punct de vedere sintactic toate aceste propoziții sunt corecte; semantic, unele propoziții pot să fie incorecte sau chiar să nu aibă sens.

Să mai observăm că o gramatică Chomsky este în măsură să constituie un model matematic pentru sintaxa unei limbi, fără să intereseze aspectele semantice. Este ceea ce a încercat să facă Naom Chomsky pentru limba engleză în lucrările sale din anii '50.

3. $G = (V_N, V_T, X_0, P)$, unde
- $$V_N = \{ \langle \text{program} \rangle, \langle \text{instrucție} \rangle, \langle \text{atribuire} \rangle, \langle \text{if} \rangle, \langle \text{expresie} \rangle, \langle \text{termen} \rangle, \langle \text{factor} \rangle, \langle \text{variabilă} \rangle, \langle \text{index} \rangle \},$$
- $$V_T = \{ \text{begin, end, if, then, stop, t, i, +, *, (,), =, ,, ; } \},$$
- $$X_0 = \langle \text{program} \rangle$$
- $$P = \{ \langle \text{program} \rangle \rightarrow \text{begin } \langle \text{linie} \rangle \text{ end} \\ \langle \text{linie} \rangle \rightarrow \langle \text{linie} \rangle ; \langle \text{instrucție} \rangle \mid \langle \text{instrucție} \rangle \\ \langle \text{instrucție} \rangle \rightarrow \langle \text{atribuire} \rangle \mid \langle \text{if} \rangle \mid \text{stop} \\ \langle \text{atribuire} \rangle \rightarrow \langle \text{variabilă} \rangle = \langle \text{expresie} \rangle \\ \langle \text{if} \rangle \rightarrow \text{if} (\langle \text{expresie} \rangle) \text{ then } \langle \text{atribuire} \rangle \\ \langle \text{expresie} \rangle \rightarrow \langle \text{expresie} \rangle + \langle \text{termen} \rangle \mid \langle \text{termen} \rangle \\ \langle \text{termen} \rangle \rightarrow \langle \text{termen} \rangle * \langle \text{factor} \rangle \mid \langle \text{factor} \rangle \\ \langle \text{factor} \rangle \rightarrow (\langle \text{expresie} \rangle) \mid \langle \text{variabilă} \rangle \\ \langle \text{variabilă} \rangle \rightarrow \text{t} (\langle \text{index} \rangle) \mid \text{i} \\ \langle \text{index} \rangle \rightarrow \langle \text{index} \rangle, \langle \text{expresie} \rangle \mid \langle \text{expresie} \rangle$$

Gramatica din acest exemplu definește un limbaj de programare simplu cu trei tipuri de instrucții: atribuiri, if-then, stop. Expresiile aritmetice au numai operatorii + și *; variabilele pot fi simple sau indexate (tablouri), iar i ține loc de identificator sau constantă. Menționăm că definirea în acest mod a unui limbaj, inclusiv utilizarea croșetelor pentru desemnarea simbolurilor neterminale, poartă adesea denumirea de **notație Backus Naur**; în acest mod s-a definit limbajul ALGOL60.

Tipuri de gramatici După forma regulilor de generare, gramaticile Chomsky se împart în mai multe tipuri; clasificarea obișnuită este următoarea:

- *Gramatici de tipul 0*; sunt gramatici fără restricții asupra regulilor;
- *Gramatici de tipul 1 (dependente de context)*; sunt gramatici care au reguli de forma

$$uAv \rightarrow upv, \quad u, p, v \in V_G^*, \quad p \neq \lambda, \quad A \in V_N$$

sau $A \rightarrow \lambda$ și în acest caz A nu apare în dreapta vreunei reguli.

Observație. Evident, regulile de forma a doua au sens numai dacă A este simbolul de start.

- *Gramaticile de tipul 2 (independente de context);* sunt gramatici care au reguli de forma

$$A \rightarrow p, \quad A \in V_N, \quad p \in V_G^*.$$

- *Gramaticile de tipul 3 (regulate);* sunt gramatici care au reguli de forma

$$\left\{ \begin{array}{l} A \rightarrow Bp \\ C \rightarrow q \end{array} \right. \quad \text{sau} \quad \left\{ \begin{array}{l} A \rightarrow pB \\ C \rightarrow q \end{array} \right.$$

cu $A, B, C \in V_N$ și $p, q \in V_T^*$.

Vom nota cu \mathcal{L}_j , $j = 0, 1, 2, 3$ familiile de limbaje generate de gramaticile de tipurile $j = 0, 1, 2, 3$; vom avea astfel limbaje de tipul 0, limbaje de tipul 1 (sau *dependente de context*), limbaje de tipul 2 (sau *independente de context*) și limbaje de tipul 3 (sau *regulate*). Să observăm că este importantă structura cuvintelor unui limbaj și nu modul în care sunt notate simbolurile terminale. De exemplu, limbajele

$$L'_2 = \{0^n 1^n | n \geq 1\}, \quad L''_2 = \{a^n b^n | n \geq 1\}$$

sunt în mod practic identice. Putem utiliza o unică notație pentru alfabetul simbolurilor terminale, de exemplu, $V_T = \{i_1, \dots, i_n\}$. Clasificarea de mai sus este fundamentală în teoria limbajelor formale, ea a fost introdusă de Naom Chomsky în 1958 și prin tradiție noile clase de limbaje sunt raportate la această clasificare. O altă clasificare este următoarea

- *Gramatici de tipul 0; fără restricții;*
- *Gramatici monotone (de tipul 1):*

$$u \rightarrow v, \quad |u| \leq |v|, \quad u, v \in V_G^*;$$

- *Gramatici dependente de context:*

$$uAv \rightarrow upv, \quad u, p, v \in V_G^*, \quad p \neq \lambda, \quad A \in V_N;$$

- *Gramatici independente de context (de tipul 2):*

$$A \rightarrow p, \quad A \in V_N, \quad p \in V_G^*;$$

- *Gramatici liniare:*

$$A \rightarrow uBv, \quad A \in V_N, \quad B \in V_N \cup \{\lambda\} \quad u, v \in V_T^*;$$

- *Gramatici (stâng) drept liniare:*

$$A \rightarrow uB (A \rightarrow Bv), \quad A \in V_N, B \in V_N \cup \{\lambda\} \quad u, v \in V_T^*;$$

- *Gramatici regulate (de tipul 3); gramatici stâng liniare sau gramatici drept liniare.*

Gramaticile monotone ca și cele dependente de context nu pot avea reguli cu partea dreaptă vidă. Se introduce următoarea convenție de completare: *într-o gramatică monotonă sau dependentă de context se admite o regula de forma $A \rightarrow \lambda$ cu condiția ca A să nu apară în partea dreaptă a vreunei reguli.* După cum vom vedea, existența sau inexistența regulilor de forma $A \rightarrow \lambda$, reguli numite de ștergere, poate modifica esențial puterea generativă a unei gramatici. O gramatică în care nu avem astfel de reguli o vom numi *gramatică λ -liberă*; de asemenea, un limbaj care nu conține cuvântul vid, îl vom numi *limbaj λ -liber*. Să mai observăm că existența regulilor de ștergere într-o gramatică nu implică în mod necesar existența cuvântului vid în limbajul generat.

Două gramatici care generează același limbaj se numesc **echivalente**.

Gramaticile monotone și gramaticile dependente de context sunt echivalente; de asemenea, gramaticile drept și stâng liniare sunt echivalente, justificându-se astfel clasa gramaticilor regulate.

1.3 Ierarhia Chomsky

Lemele care urmează vor avea o utilizare frecventă în cele ce urmează.

Lema 1.1 (*Lema de localizare a gramaticilor independente de context*) Fie G o gramatică independentă de context și fie derivarea

$$x_1 \dots x_m \xRightarrow{*} p, \quad \text{unde } x_j \in V_G, j = \overline{1, m}, p \in V_G^*.$$

Atunci există $p_1, p_2, \dots, p_m \in V_G^*$ astfel încât

$$p = p_1 \dots p_m \quad \text{și} \quad x_j \xRightarrow{*} p_j, \quad j = \overline{1, m}.$$

Demonstrație. Procedăm prin inducție asupra lungimii derivării l .

Dacă $l = 0$ atunci $p = x_1 \dots x_m$ și luăm $p_j = x_j$.

Presupunem că proprietatea este adevărată pentru derivări de lungime l și fie o derivare de lungime $l + 1, x_1 \dots x_m \xRightarrow{*} p$. Punem în evidență ultima derivare

directă $x_1 \dots x_m \xRightarrow{*} q \Rightarrow p$. Conform ipotezei inductive, $q = q_1 \dots q_m$ și $x_j \xRightarrow{*} q_j$, $j = 1, \dots, m$.

Fie apoi $A \rightarrow u$ regula care se aplică în derivarea directă $q \Rightarrow p$ și să presupunem că A intră în subcuvântul q_k , deci $q_k = q'_k A q''_k$. Vom lua

$$p_j = \begin{cases} q_j & , j \neq k \\ q'_k u q''_k & , j = k \end{cases}$$

Este evident că $x_j \xRightarrow{*} p_j$, $j \neq k$, iar pentru $j = k$ avem

$$x_k \xRightarrow{*} q_k = q'_k A q''_k \Rightarrow q'_k u q''_k = p_k. \square$$

Vom pune în evidență în continuare o proprietate asupra structurii regulilor gramaticilor Chomsky. Partea dreaptă a unei reguli, pentru toate tipurile de gramatici, este un cuvânt format din terminale sau neterminale. Este convenabil de multe ori ca partea dreaptă a regulilor să conțină un singur tip de simboluri, terminale sau neterminale. Acest lucru este posibil fără modificarea tipului gramaticii.

Lema 1.2 (*Lema $A \rightarrow i$*) Fie $G = (V_N, V_T, S, P)$ o gramatică de tipul 2. Există o gramatică G' echivalentă, de același tip, cu proprietatea că dacă o regulă are în partea dreaptă un terminal, atunci ea este de forma $A \rightarrow i$, $A \in V_N, i \in V_T$.

Demonstrație. Luăm gramatica G' de forma $G' = (V'_N, V'_T, S, P')$ unde V'_N și P' se construiesc astfel: $V_N \subseteq V'_N$ și includem în P' toate regulile din P care convin. Fie acum o regulă din P care nu convine (punem în evidență apariția terminalelor din partea dreaptă):

$$u \rightarrow v_1 i_1 v_2 i_2 \dots i_n v_{n+1}, \quad i_k \in V_T, \quad v_k \in V_N^*.$$

Vom introduce în P' următoarele reguli:

$$u \rightarrow v_1 X_{i_1} v_2 X_{i_2} \dots X_{i_n} v_{n+1}, \quad X_{i_k} \rightarrow i_k, \quad k = \overline{1, n},$$

unde X_{i_k} sunt neterminale noi pe care le adăugăm la V'_N . Este evident că G' păstrează tipul lui G și că $L(G') = L(G) \square$.

Observație. Construcția din lema se poate extinde ușor la gramatici de tipul 0, înlocuind și terminalele din partea stângă a regulilor (u poate fi un șir arbitrar ce conține minim un neterminale). Gramatica astfel obținută este echivalentă cu cea inițială și păstrează tipul.

Ierarhia Chomsky. Este evident că $\mathcal{L}_3 \subset \mathcal{L}_2$ și că $\mathcal{L}_1 \subset \mathcal{L}_0$, deoarece orice regulă a unei gramatici de tipul 3 respectă prescripțiile unei gramatici de tipul 2; analog pentru familiile \mathcal{L}_1 și \mathcal{L}_0 . Aparent, o regulă de forma $A \rightarrow p$ (de tipul 2) este un caz particular a unei reguli de forma $uAv \rightarrow upv$ (de tipul 1), pentru

$u = v = \lambda$; totuși, realitatea nu este aceasta, deoarece la tipul 2 de gramatici sunt permise reguli de ștergere, pe când la tipul 1 se impune condiția $p \neq \lambda$. Vom arăta că avem $\mathcal{L}_2 \subset \mathcal{L}_1$.

Șirul de incluziuni

$$\mathcal{L}_3 \subset \mathcal{L}_2 \subset \mathcal{L}_1 \subset \mathcal{L}_0$$

poartă denumirea de **ierarhia Chomsky** (vom arăta pe parcursul acestui curs că incluziunile sunt stricte). Această ierarhie caracterizează puterea generativă a celor patru tipuri de gramatici, această putere fiind crescătoare de la 3 la 0. Orice alte mecanisme generative se raportează la această ierarhie fundamentală. Vom demonstra mai întâi următoarea leamnă.

Lema 1.3 (Lema eliminării regulilor de ștergere). *Orice limbaj independent de context λ -liber poate fi generat de o gramatică de tipul 2 fără reguli de ștergere.*

Demonstrație. Fie $G = (V_N, V_T, X_0, P)$ o gramatică independentă de context și $L(G)$ limbajul generat. Prin ipoteză $\lambda \notin L(G)$. Definim prin recurență șirul de mulțimi U_k astfel:

$$U_1 = \{X | X \in V_N, X \rightarrow \lambda \in P\}$$

$$U_{k+1} = \{X | X \in V_N, X \rightarrow p \in P, p \in U_k^*\} \cup U_k, \quad k \geq 1.$$

Să observăm că șirul de mulțimi definite este crescător în raport cu relația de incluziune. Cum toate aceste mulțimi sunt incluse în V_N și V_N este finită, rezultă că există o mulțime maximală U_f astfel încât

$$U_1 \subseteq U_2 \subseteq \dots \subseteq U_f = U_{f+1} = \dots$$

Are loc de asemenea și implicația $X \in U_f \Leftrightarrow X \Rightarrow^* \lambda$.

Vom ilustra această implicație cu un exemplu. Să presupunem că $U_f = U_3$ și fie $X \in U_3$. Atunci în mod necesar trebuie să avem $X \rightarrow p \in P, p \in U_2^*$; de exemplu $p = X_1 X_2$ și $X_1, X_2 \in U_2$. În mod analog $X_1 \rightarrow Y_1 Y_2 Y_3, X_2 \rightarrow Z_1 Z_2$ și $Y_1, Y_2, Y_3, Z_1, Z_2 \in U_1$, prin urmare $Y_1 \rightarrow \lambda, Y_2 \rightarrow \lambda, Y_3 \rightarrow \lambda, Z_1 \rightarrow \lambda, Z_2 \rightarrow \lambda$. Putem scrie derivarea

$$X \Rightarrow X_1 X_2 \Rightarrow^* Y_1 Y_2 Y_3 Z_1 Z_2 \Rightarrow^* \lambda.$$

Definim acum următoarea gramatică independentă de context fără reguli de ștergere $G' = (V_N, V_T, X_0, P')$ unde V_N, V_T, X_0 sunt ca în gramatica dată, iar P' se construiește pornind de la P astfel. Fie $X \rightarrow p \in P, p \neq \lambda$. Includem atunci în P' această regulă precum și toate regulile de forma $X \rightarrow p_j$, unde p_j se obține din p lăsând la o parte, în toate modurile posibile, simbolurile din U_f (se exceptă cazul $p_j = \lambda$). De exemplu dacă $X \rightarrow ABC \in P$ și $A, B \in U_f$, vom induce în P' regulile

$$X \rightarrow ABC, X \rightarrow BC, X \rightarrow AC, X \rightarrow C.$$

Să observăm că în acest fel mulțimea P a fost pe de o parte micșorată (au fost excluse regulile de alegere), iar pe de altă parte îmbogățită cu eventualele noi reguli. Să mai observăm că G' este independentă de context și că nu conține reguli de ștergere.

Vom arăta că $L(G) = L(G')$.

Mai întâi, să arătăm că $L(G) \subseteq L(G')$.

Fie $p \in L(G)$, deci $X_0 \xRightarrow[G]{*} p$; vom arăta că $X_0 \xRightarrow[G']{*} p$. Vom arăta o implicație

ceva mai generală, $X \xRightarrow[G]{*} p$ implică $X \xRightarrow[G']{*} p$, pentru orice $X \in V_N$ (relația cerută

se obține pentru $X = X_0$). Procedăm prin inducție asupra lungimii derivării l . Dacă $l = 1$ avem

$$X \xRightarrow[G]{*} p \text{ deci } X \rightarrow p \in P. \text{ Dar } p \neq \lambda \text{ deci } X \rightarrow p \in P'$$

adică $X \xRightarrow[G']{*} p$.

Presupunem că afirmația este adevărată pentru $l = n$ și luăm o derivare cu lungimea $l = n + 1$. Punem în evidență prima derivare directă

$$X \xRightarrow[G]{*} X_1 \dots X_m \xRightarrow[G]{*} p$$

Conform lemei de localizare avem $p = p_1 \dots p_m$ și $X_j \xRightarrow[G]{*} p_j$, $j = 1, \dots, m$.

Unele din cuvintele p_j pot să fie vide; pentru precizarea ideilor să presupunem că $p_2, p_3, p_5 = \lambda$. Atunci pentru derivările $X_j \xRightarrow[G]{*} p_j$, $j \neq 2, 3, 5$, care au lungimea

de cel mult n , conform ipotezei inductive avem $X_j \xRightarrow[G']{*} p_j$.

Pe de altă parte, pentru $j = 2, 3, 5$ avem $X_2 \xRightarrow[G]{*} \lambda$, $X_3 \xRightarrow[G]{*} \lambda$, $X_5 \xRightarrow[G]{*} \lambda$, deci $X_2, X_3, X_5 \in U_f$. Rezultă că

$$X \rightarrow X_1 X_4 X_6 \dots X_m \in P'$$

așa încât putem scrie

$$X \xRightarrow[G']{*} X_1 X_4 X_6 \dots X_m \xRightarrow[G']{*} p_1 p_4 p_6 \dots p_m = p.$$

Deci $X \xRightarrow[G']{*} p$ și luând $X = X_0$ obținem $p \in L(G')$. Prin urmare $L(G) \subseteq L(G')$.

Să arătăm acum incluziunea inversă, $L(G') \subseteq L(G)$.

Fie $p \in L(G)$, $X_0 \xRightarrow[G']{*} p$. Punem în evidență o derivare directă oarecare

$$X_0 \xRightarrow[G']{*} u \xRightarrow[G']{*} v \xRightarrow[G']{*} p.$$

Dacă în derivarea directă $u \xRightarrow{G'} v$ se aplică o regulă care există și în G , atunci evident pasul respectiv poate fi făcut și în G . Să presupunem că se aplică o regulă nou introdusă de exemplu $X \rightarrow BC$, deci pasul respectiv va avea forma

$$u = u'Xu'' \xRightarrow{G'} u'BCu'' = v$$

Regula $X \rightarrow BC \in P'$ a provenit dintr-o regulă din P , lăsând la o parte simboluri din U_f , în cazul nostru din $X \rightarrow ABC$, lăsându-l la o parte pe $A \in U_f$. Deoarece $A \xRightarrow{G} \lambda$, avem

$$u = u'Xu'' \xRightarrow{G} u'ABCu'' \xRightarrow{G} u'ABCu'' = v.$$

Prin urmare orice pas al derivării considerate se poate obține și în gramatica G , deci $X_0 \xRightarrow{G} p$ și $p \in L(G)$, adică $L(G') \subseteq L(G)$.

Teorema 1.1 $\mathcal{L}_2 \subseteq \mathcal{L}_1$.

Demonstrație. Fie $L \in \mathcal{L}_2$ un limbaj independent de context și G o gramatică de tipul 2 care îl generează, adică $L = L(G)$.

Presupunem că $\lambda \in L$. Construim gramatica G' ca în lema precedentă; orice cuvânt $p \neq \lambda$ din $L(G)$ se poate obține în G' și invers, deci $L(G') = L(G) - \{\lambda\}$. Considerăm atunci o gramatică $G'' = (V_N \cup \{X_0''\}, V_T, X_0'', P' \cup \{X_0'' \rightarrow \lambda, X_0'' \rightarrow X_0\})$. Evident $L(G'') = L(G)$. Toate regulile lui G'' respectă tipul 1 (cu $u = v = \lambda$) și conține o singură regulă de ștergere $X_0'' \rightarrow \lambda$ iar X_0'' nu apare în partea dreaptă a vreunei reguli. Deci G'' este de tipul 1 și orice limbaj independent de context este inclus în \mathcal{L}_1 , adică $\mathcal{L}_2 \subseteq \mathcal{L}_1$. \square

Fiind dată o operație binară notată cu "•" pe o familie de limbaje L , vom spune că familia L este închisă la operația "•" dacă $L_1, L_2 \in L$ implică $L_1 \bullet L_2 \in L$. Definiția noțiunii de închidere pentru operații unare sau cu aritate oarecare este analoagă. Relativ la proprietățile de închidere a familiilor din clasificarea Chomsky menționăm următorul rezultat.

Teorema 1.2 Familiile $\mathcal{L}_j, j = 0, 1, 2, 3$ sunt închise la operațiile regulate.

Demonstrație. Fie $G_k = (V_{N_k}, V_{T_k}, S_k, P_k), k = 1, 2$ două gramatici de același tip $j, j = 0, 1, 2, 3$. Putem presupune că $V_{N_1} \cap V_{N_2} = \emptyset$. Trebuie să arătăm că limbajele $L(G_1) \cup L(G_2)$, $L(G_1)L(G_2)$, $L(G_1)^*$, sunt de același tip j . În acest scop vom construi gramatici de tipul j care să genereze le genereze. Vom indica fără demonstrație modul de construcție al gramaticilor (pentru detalii vezi [?]):

Reuniune

$j = 0, 1, 2, 3$:

$$G = V_{N_1} \cup V_{N_2} \cup \{S\}, V_{T_1} \cup V_{T_2}, P_1 \cup P_2 \cup \{S \rightarrow S_1 | S_2\}.$$

Produs

$j = 0, 1, 2$:

$$G = V_{N_1} \cup V_{N_2} \cup \{S\}, V_{T_1} \cup V_{T_2}, P_1 \cup P_2 \cup \{S \rightarrow S_1 S_2\}.$$

$j = 3$:

$$G = V_{N_1} \cup V_{N_2} \cup \{S\}, V_{T_1} \cup V_{T_2}, S_1, P'_1 \cup P_2,$$

unde P'_1 se obține din P_1 prin înlocuirea regulilor de forma $A \rightarrow p$ cu $A \rightarrow pS_2$.

Închidere Kleene

$j = 0, 1$:

$$G^* = (V_N \cup \{S^*, X\}, V_T, S^*, P \cup \{S^* \rightarrow \lambda | S | X S, X i \rightarrow S i | X S i, \forall i \in V_T\})$$

$j = 2$:

$$G^* = (V_N \cup \{S^*\}, V_T, S^*, P \cup P^{prime} \cup \{S^* \rightarrow S^* S | \lambda\})$$

$j = 3$:

$$G^* = (V_N \cup \{S^*\}, V_T, S^*, P \cup P^{prime} \cup \{S^* \rightarrow S | \lambda\})$$

1.4 Traducerea programelor

Un **limbaj de programare** este un limbaj care are drept scop descrierea unor procese de prelucrare a anumitor date și a structurii acestora (în unele cazuri descrierea structurii datelor este preponderentă), prelucrare care se realizează în general cu ajutorul unui sistem de calcul.

Există în prezent un număr mare de limbaje de programare de *nivel înalt* sau evaluate care se caracterizează printr-o anumită naturalitate, în sensul că descrierea procesului de prelucrare cu ajutorul limbajului este apropiată de descrierea naturală a procesului respectiv precum și prin independența unor astfel de limbaje față de sistemul de calcul. Dintre limbajele de acest tip cu o anumită răspândire în momentul de față menționăm limbajele PASCAL, FORTRAN, C, JAVA, etc.

O altă clasă importantă de limbaje, sunt *limbajele de asamblare*, sau de nivel inferior, ale căror caracteristici depind de sistemul de calcul considerat. În general, fiecare sistem de calcul (sau tip de sistem de calcul), are propriul său limbaj de asamblare; de exemplu, limbajul de asamblare ale sistemelor de calcul echipate cu procesoare de tip Intel Z-80 este denumit în mod curent ASSEMBLER. Instrucțiunile unui limbaj de asamblare corespund cu operațiile simple ale

sistemului de calcul iar stocarea datelor în memorie este realizată direct de utilizator la nivelul locațiilor elementare ale memoriei. Există de asemenea anumite *pseudo-instrucții* sau *directive* referitoare la alocarea memoriei, generarea datelor, segmentarea programelor, etc., precum și macroinstrucții care permit generarea unor secvențe tipice de program sau accesul la bibliotecile de subprograme.

În afară de limbajele evolute și de limbajele de asamblare, există numeroase *limbaje specializate*, numite uneori și de *comandă*, care se referă la o anumită clasă de aplicații. Menționăm de exemplu limbajul pentru prelucrarea listelor LISP, limbajele utilizate în cadrul softwarelui matematic (Mathematica, Maple, MATCAD, etc.) și multe altele. În general însă astfel de limbaje nu sunt considerate limbaje de programare propriu-zise.

Un program redactat într-un limbaj de programare poartă denumirea de **program sursă**. Fiecare sistem de calcul, în funcție de particularitățile sale, posedă un anumit limbaj propriu, numit *cod mașină*, acesta fiind singurul limbaj înțeles de procesorul sistemului. Un astfel de limbaj depinde de structura instrucțiilor procesorului, de setul de instrucții, de posibilitățile de adresare, etc. Un program redactat în limbajul cod mașină al sistemului de calcul îl numim **program obiect**.

Procesul de transformare al unui program sursă în program obiect se numește **compilare** sau *translatare*, uneori chiar *traducere*. De obicei termenul de compilare este utilizat numai în cazul limbajelor evolute, în cazul limbajelor de asamblare fiind utilizat termenul de *asamblare*.

Compilarea (asamblarea) este efectuată de un program al sistemului numit **compilator** (*asamblor*). De multe ori compilatoarele nu produc direct program obiect, ci un text intermediar apropiat de programul obiect, care în urma unor prelucrări ulterioare devine program obiect. De exemplu, compilatoarele sistemelor de operare DOS produc un text numit obiect (fișiere cu extensia obj) care în urma unui proces numit editare de legături și a încărcării în memorie devine program obiect propriu-zis, numit program executabil (fișiere cu extensia exe). Există mai multe rațiuni pentru o astfel de tratare, între care posibilitatea cuplării mai multor module de program realizate separat sau provenite din limbaje sursă diferite, posibilitatea creării unor biblioteci de programe în formatul intermediar și utilizarea lor în alte programe, etc.

Un program sursă poate de asemenea să fie executat de către sistemul de calcul direct, fără transformarea lui prealabilă în program obiect. În acest caz, programul sursă este prelucrat de un program al sistemului numit *interpretor*; acesta încarcă succesiv instrucțiile programului sursă, le analizează din punct de vedere sintactic și semantic și după caz, le execută sau efectuează anumite operații auxiliare.

Procesul de compilare este un proces relativ complex și comportă operații care au un anumit caracter de autonomie. Din aceste motive procesul de compilare este de obicei descompus în mai multe subprocesse sau faze, fiecare fază fiind o operație coerentă, cu caracteristici bine definite. În principiu aceste faze sunt

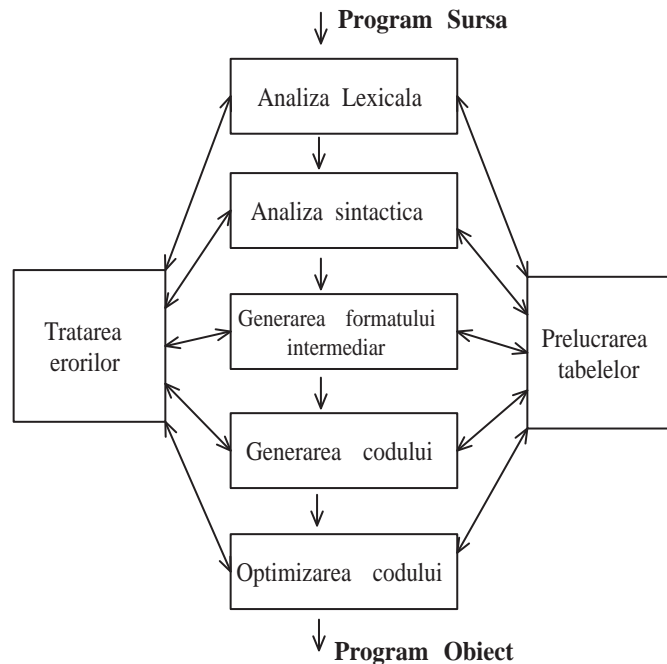


Figura 1.1: Fazele compilării

parcuse secvențial (pot să existe și anumite reveniri) iar programul sursă este transformat succesiv în formate intermediare. Se poate considera că fiecare fază primește de la faza precedentă un fișier cu programul prelucrat într-un mod corespunzător fazei respective, îl prelucrează și furnizează un fișier de ieșire, iarăși într-un format bine precizat, fișier utilizat în faza următoare. Există cinci faze de compilare principale: analiza lexicală, analiza sintactică, generarea formatului intermediar, generarea codului, optimizarea codului și două faze auxiliare, tratarea erorilor și tratarea tabelor (vezi figura 1.1).

Analiza lexicală are ca obiectiv principal determinarea unităților lexicale ale unui program, furnizarea codurilor acestora și detectarea eventualelor erori lexicale. Pe lângă aceste operații de bază, la analiza lexicală se mai pot efectua anumite operații auxiliare precum: eliminarea blank-urilor (dacă limbajul permite utilizarea fără restricții a acestora), ignorarea comentariilor, diferite conversiuni ale unor date (care se pot efectua la această fază), completarea tabelor compilatorului.

Analiza sintactică determină unitățile sintactice ale programului (secvențe de text pentru care se poate genera format intermediar) și verifică programul din punct de vedere sintactic. Este faza centrală a unui compilator, deseori toate celelalte faze sunt rutine apelate de analizorul sintactic pe măsură ce este posibilă efectuarea unei părți din faza respectivă. Tot la analiza sintactică se definitivează

tabelele de informații și se realizează prelucrarea erorilor sintactice.

Faza de *generare a formatului intermediar* se referă la transformarea programului într-o formă numită intermediară pornind de la care se poate, printr-o procedură relativ simplă, să se obțină programul obiect. Structura acestei faze depinde de formatul intermediar ales de către proiectant și de modalitatea de implementare; uzual se folosesc *cvadrupele*, *triplele* sau *șirurile poloneze*.

Generarea codului este o fază în care se realizează codul obiect corespunzător programului. Practic în această fază se obține programul în limbaj de asamblare corespunzător programului sursă redactat într-un limbaj evoluat. În mod obișnuit generatorul de cod este constituit din rutinele generatoare de cod corespunzătoare diverselor unități ale formatului intermediar.

În faza de *optimizare* a codului se realizează o anumită îmbunătățire a codului obiect astfel încât programul rezultat să fie cât mai performant (în privința consumului de timp și memorie). Cele mai tipice exemple sunt eliminarea încărcărilor și a memorărilor redundante sau optimizarea ciclurilor.

Fazele de *Prelucrare a tabelor* și de *Tratare a erorilor* vor fi atinse numai parțial în această curs. Facem însă mențiunea că prelucrarea tabelor poate să aibă o influență importantă asupra performanțelor unui compilator, în mod special din punctul de vedere al timpului de execuție (compilare) și că tratarea erorilor are o anumită implicație în eliminarea erorilor sintactice.

1.5 Probleme propuse

1. Găsiți limbajul generat de gramatica $G = (V_N, V_T, S, P)$, precizând tipul gramaticii (cf. clasificării Chomsky):

- (a) $V_N = \{S\}$; $V_T = \{0, 1, 2\}$;
 $P = \{S \rightarrow 0S0|1S1|2S2|0\}$.
- (b) $V_N = \{S\}$; $V_T = \{a, b\}$;
 $P = \{S \rightarrow aSb|ab\}$.
- (c) $V_N = \{S, A, B\}$; $V_T = \{a, b, c\}$;
 $P = \{S \rightarrow abc|aAbc, Ab \rightarrow bA, Ac \rightarrow Bbcc, bB \rightarrow Bb, aB \rightarrow aaA|aa\}$.
- (d) $V_N = \{S, A, C\}$; $V_T = \{+, -, 0, 1, \dots, 9\}$;
 $P = \{S \rightarrow A|+A|-A, A \rightarrow AC|C, C \rightarrow 0|1|\dots|9\}$.
- (e) $V_N = \{S, A, B\}$; $V_T = \{0, 1\}$;
 $P = \{S \rightarrow 0S|1A, A \rightarrow 0B|1S|0, B \rightarrow 0A|1B|1\}$.
- (f) $V_N = \{A\}$, $V_T = \{a, b\}$, $S = A$, $P = \{A \rightarrow aA|b\}$;
- (g) $V_N = \{x_0\}$, $V_T = \{A, B, \dots, Z\}$, $S = x_0$, $P = \{x_0 \rightarrow x_1D, x_1 \rightarrow x_2N, x_2 \rightarrow E\}$;

- (h) $V_N = \{A\}$, $V_T = \{0, 1, 2\}$, $S = A$, $P = \{A \rightarrow 0A0|1A1|2A2|\lambda\}$;
- (i) $V_N = \{S, A\}$, $V_T = \{0, 1, \dots, 9, \cdot\}$,
 $P = \{S \rightarrow A.A, A \rightarrow 0A|1A| \dots |9A|0|1| \dots |9\}$;
- (j) $V_N = \{S\}$, $V_T = \{PCR, PDAR, UDMR\}$,
 $P = \{S \rightarrow PCR|PDAR|UDMR\}$;
- (k) $V_N = \{A, B, C\}$, $V_T = \{0, 1\}$, $S = A$, $P = \{A \rightarrow 0A|1B|1,$
 $B \rightarrow 0C|1A, C \rightarrow 0B|1C|0\}$;
- (l) $V_N = \{S, A, B, C\}$, $V_T = \{0, 1, \dots, 9, +, -\}$,
 $P = \{S \rightarrow +A|-A|A, A \rightarrow 0A|1A| \dots |9A|0|1| \dots |9\}$;
- (m) $V_N = \{S\}$, $V_T = \{(\cdot)\}$, $P = \{S \rightarrow S(S)S|\lambda\}$;
- (n) $V_N = \{E, T, F\}$, $V_T = \{(\cdot), i, +, *\}$, $S = E$,
 $P = \{E \rightarrow E + T|T, T \rightarrow T * F|F, F \rightarrow (E)|i\}$;
- (o) $V_N = \{S, A, B\}$, $V_T = \{a, b, c\}$, $P = \{S \rightarrow abc|aAbc,$
 $Ab \rightarrow bA, Ac \rightarrow Bbcc, bB \rightarrow Bb, aB \rightarrow aaA|aa\}$;
- (p) $V_N = \{S, A, B, C, D, E\}$, $V_T = \{a\}$, $P = \{S \rightarrow ACaB, Ca \rightarrow aaC, CB \rightarrow$
 $DB|E, aD \rightarrow Da, AD \rightarrow AC, aE \rightarrow Ea, AE \rightarrow \lambda\}$;
- (q) $V_N = \{S, A, B, C, D, E\}$, $V_T = \{a, b\}$, $P = \{S \rightarrow ABC,$
 $AB \rightarrow aAD|bAE, DC \rightarrow BaC, EC \rightarrow BbC, Da \rightarrow aD, Db \rightarrow bD, Ea \rightarrow$
 $aE, Eb \rightarrow bE, AB \rightarrow \lambda, C \rightarrow \lambda, aB \rightarrow Ba, bB \rightarrow Bb\}$;

2. Precizați care dintre gramaticile precedente sunt echivalente.

3. Găsiți gramatici pentru generarea următoarelor limbaje:

- (a) $L = \{\lambda\}$;
- (b) $L = \emptyset$;
- (c) $L = \{0^n | n \in \mathbb{N}\}$;
- (d) $L = \{a, ab, aab, ba, baa\}$.
- (e) $L = \{a^n cb^n | n \geq 1\}$.
- (f) $L = \{w \in \{a, b, c\}^* | w \text{ conține } cc \text{ și se termină cu } a\}$.
- (g) $L = \{BEGIN|END|IF|WHILE|UNTIL\}$.
- (h) $L = \{w \in \{0, 1\}^* | \text{reprezentarea binară pe 8 biți a unui întreg}\}$.
- (i) $L = \{a^n b^{n+3} a^{n+1} | n \geq 0\}$.
- (j) $L = \{w \in \{a, b, c\}^* | w \text{ începe cu } a \text{ și are maxim 4 litere}\}$.
- (k) $L = \{a^i b^j c^k | i, j, k > 0\}$.
- (l) $L = \{w \in \{0, 1\}^* | w \text{ multiplu de 8}\}$.
- (m) $L = \{\text{constante reale în scrierea obișnuită cu punct zecimal } p_i.p_z\}$.

- (n) $L = \{a^i b^j a^i b^j\}$;
 - (o) $L = \{awbbw' | w, w' \in \{0, 1\}^*\}$;
 - (p) $L = \{w \in \{0, 1\}^* | w \text{ conține maxim 2 de } 0 \}$;
 - (q) $L = \{wa\tilde{w} | w \in \{0, 1\}^*\}$;
 - (r) $L = \{w | w \text{ octet ce reprezintă un număr par } \}$;
 - (s) $L = \{A, B, C, \dots, Z\}$;
4. Construiți o gramatică ce conține reguli de ștergere, dar generează un limbaj λ -liber.
 5. Folosind teorema să se construiască o gramatică independentă de context, fără reguli de ștergere care generează limbajul de la punctul precedent.

Capitolul 2

Limbaje Regulate și Analiza Lexicală

2.1 Automate finite și limbaje regulate

Automate finite. Automatele finite sunt mecanisme pentru recunoașterea limbajelor de tipul 3 (regulate). Un automat finit (AF) se compune dintr-o *bandă de intrare* și un *dispozitiv de comandă*.

Pe banda de intrare sunt înregistrate simboluri ale unui *alfabet de intrare*, constituind pe bandă un cuvânt p . La fiecare pas de funcționare banda se deplasează cu o poziție spre stânga.

Dispozitivul de comandă posedă un *dispozitiv de citire* de pe bandă; dispozitivul se află permanent într-o anumită *stare internă*, element al unei *mulțimi finite de stări*. Schema unui automat finit este redată în figura 2.1.

Automatul finit funcționează în pași discreți. Un pas de funcționare constă din: dispozitivul de comandă citește de pe banda de intrare simbolul aflat în dreptul dispozitivului de citire; în funcție de starea internă și de simbolul citit, automatul trece într-o nouă stare și mută banda cu o poziție spre stânga. Automatul își încetează funcționarea după ce s-a citit ultimul simbol înregistrat pe bandă; în acest moment el se va afla într-o anumită stare, care, după cum vom vedea, va juca un rol important în recunoașterea cuvintelor.

Din punct de vedere matematic, un automat finit este un sistem

$$AF = (\Sigma, I, f, s_0, \Sigma_f),$$

unde

Σ este alfabetul (mulțimea) de stări;

I este alfabetul de intrare;

$f : \Sigma \times I \longrightarrow \mathcal{P}(\Sigma)$ este funcția de evoluție;

$s_0 \in \Sigma$ este starea inițială;

$\Sigma_f \subseteq \Sigma$ este mulțimea de stări finale.

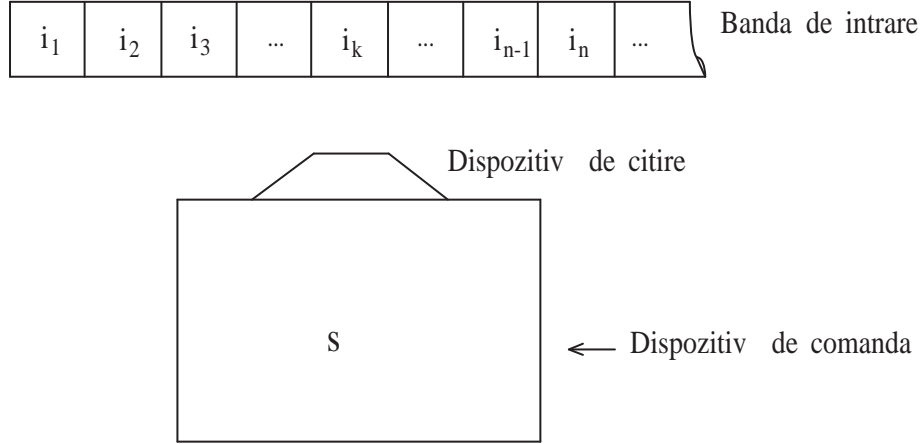


Figura 2.1: Reprezentarea schematică a unui automat finit

Dacă pentru orice $(s, i) \in \Sigma \times I$, avem $|f(s, i)| \leq 1$ automatul se numește *determinist*; în caz contrar se numește *nedeterminist*.

Observații:

1. Funcționarea unui automat finit se poate bloca în situația în care el se află în starea s , citește de pe bandă simbolul i și $f(s, i) = \emptyset$; evident că în acest caz funcționarea în continuare nu mai este posibilă.
2. În cazul unui automat determinist vom scrie $f(s, i) = s'$ (în loc de $f(s, i) \in \{s'\}$).
3. Definiția dată este specifică teoriei limbajelor formale. O altă definiție (mai generală), întâlnită în teoria automatelor este următoarea: un automat finit este un sistem $AF = (\Sigma, I, O, f, g, s_0, F)$ unde, Σ, I, f, s_0, F au semnificația de mai sus, O este *alfabetul de ieșire* iar $g : \Sigma \times I \rightarrow \mathcal{P}(O)$ este *funcție de ieșire*. Funcționalitatea unui astfel de automat finit este analoagă cu cea descrisă mai sus, cu deosebirea că la fiecare pas automatul furnizează o ieșire $o \in g(s, i)$.

Prin *diagrama de stări* a unui automat finit înțelegem un graf orientat care are nodurile etichetate cu stările $s \in \Sigma$ iar arcele se construiesc astfel: nodurile s, s' se unesc cu un arc orientat de la s la s' dacă există $i \in I$ astfel încât $s' \in f(s, i)$; arcul respectiv va fi notat cu i .

Exemplu $AF = (\Sigma, I, f, s_0, \Sigma_f)$ unde $\Sigma = \{s_0, s_1, s_2\}$, $I = \{i_1, i_2\}$, $\Sigma_f = \{s_2\}$, iar f este dată de tabelul

	s_0	s_1	s_2
i_1	$\{s_1\}$	$\{s_2\}$	$\{s_0\}$
i_2	\emptyset	$\{s_0, s_1\}$	$\{s_0, s_1\}$

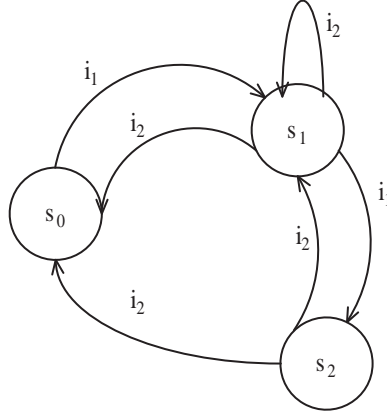


Figura 2.2: Diagrama de stări

Diagrama de stări corespunzătoare este prezentată în figura 2.2.

Funcția de evoluție Funcția f se prelungește de la $\Sigma \times I$ la $\mathcal{P}(\Sigma) \times I^*$ deci definim $f' : \mathcal{P}(\Sigma) \times I^* \rightarrow \mathcal{P}(\Sigma)$, astfel:

- (a) $f'(s, \lambda) = \{s\}, \forall s \in \Sigma,$
- (b) $f'(\emptyset, i) = \emptyset, \forall i \in I,$
- (c) $f'(Z, i) = \bigcup_{s \in Z} f(s, i), \quad Z \in \mathcal{P}(\Sigma), Z \neq \emptyset,$
- (d) $f'(Z, pi) = f'(f'(Z, p), i).$

Prin abuz de limbaj, vom folosi pentru noua funcție tot notația f . Să observăm că relațiile de mai sus constituie o definiție prin recurență, corectă; fiind dat f , putem defini mai întâi toate valorile $f(Z, i)$ (un număr finit, deoarece I este o mulțime finită), apoi $f(Z, p)$ pentru $|p| = 2$, în continuare $f(Z, p)$ pentru $|p| = 3$, etc.

Proprietățile funcției f :

1. Dacă Z_k este o familie de părți a lui Σ , avem

$$f\left(\bigcup_k Z_k, i\right) = \bigcup_k f(Z_k, i)$$

Demonstrație. Utilizând (c) putem scrie

$$\bigcup_k f(Z_k, i) = \bigcup_k \left(\bigcup_{s \in Z_k} f(s, i) \right) = \bigcup_{s \in \bigcup_k Z_k} f(s, i) = f\left(\bigcup_k Z_k, i\right). \square$$

2. $f(Z, p) = \bigcup_{s \in Z} (f(s, p)), \quad p \in I^*.$

Demonstrație. Prin inducție asupra lungimii lui p .

Pentru $|p| = 1$ avem $f(Z, i) = \bigcup_{s \in Z} f(s, i)$, adică (c).

Presupunem că relația este adevărată pentru $|p| = m$ și considerăm un p astfel

încât $|p| = m + 1$, deci $p = p'i$, $|p'| = m$. Putem scrie

$$\begin{aligned} f(Z, p) &= f(Z, p'i) = f(f(Z, p'), i) = f(\bigcup_{s \in Z} f(s, p'), i) = \square \\ &= \bigcup_{s \in Z} f(f(s, p'), i) = \bigcup_{s \in Z} f(s, p'i) = \bigcup_{s \in Z} f(s, p). \end{aligned}$$

3. $f(s, pq) = f(f(s, p), q)$, $p, q \in I^*$.

Demonstrație. Inducție asupra lungimii lui q .

Dacă $|q| = 1$, atunci $q = i$ și relația se reduce la (d).

Presupunem că proprietatea este adevărată pentru r și considerăm $|q| = r + 1$.

Deci $q = q'i$. Avem

$$\begin{aligned} f(s, pq) &= f(s, pq'i) = f(f(s, pq'), i) = f(f(f(s, p), q'), i) = \square \\ &= f(f(s, p), q'i) = f(f(s, p), q). \end{aligned}$$

Limbafe regulate.

Definiție 2.1 *Limbaful recunoscut de automatul finit $AF = (\Sigma, I, f, s_0, \Sigma_f)$ este*

$$L(AF) = \{p | p \in I^*, f(s_0, p) \cap \Sigma_f \neq \emptyset\}$$

Deci $p \in L(AF)$ dacă automatul aflându-se în starea inițială s_0 , după $|p|$ pași de funcționare poate să ajungă într-o stare finală.

În cazul unui automat finit determinist limbaful recunoscut poate fi definit în modul următor. Pentru fiecare $s \in \Sigma$ definim funcția $f_s : I^* \rightarrow \mathcal{P}(\Sigma)$ prin $f_s(p) = f(s, p)$.

Atunci

$$f(s_0, p) \cap \Sigma_f \neq \emptyset \Leftrightarrow f(s_0, p) = f_{s_0}(p) \in \Sigma_f,$$

deci

$$L(AF) = \{p | f(s_0, p) \cap \Sigma_f \neq \emptyset\} = f_{s_0}^{-1}(\Sigma_f)$$

Limbafele recunoscute de automate finite le vom numi *limbafe regulate*; familia acestor limbafe o vom nota cu \mathcal{R} . Evident, familia limbajelor recunoscute de automate finite deterministe, \mathcal{R}_d , este o parte a lui \mathcal{R} , $\mathcal{R}_d \subseteq \mathcal{R}$. Vom arăta că cele două familii coincid.

Teorema 2.1 $\mathcal{R}_d = \mathcal{R}$.

Demonstrație. Fie $AF = (\Sigma, I, f, s_0, \Sigma_f)$ un automat finit (în general nedeterminist). Construim următorul automat finit determinist $AF' = (\Sigma', I, f', \{s_0\}, \Sigma'_f)$ unde $\Sigma' = \mathcal{P}(\Sigma)$, $f' = f$ (prelungirea la $\Sigma \times I^*$), $\Sigma'_f = \{Z | Z \in \mathcal{P}(\Sigma), Z \cap \Sigma_f \neq \emptyset\}$.

Evident, automatul AF' este determinist.

Fie $p \in L(AF)$. Atunci $f(s_0, p) \cap \Sigma_f \neq \emptyset$ și $f(s_0, p) \in \Sigma'_f$. Pe de altă parte, conform cu proprietatea 2 a funcției de evoluție, avem

$$f'(\{s_0\}, p) = f(s_0, p)$$

și deci $f'(s_0, p) \in \Sigma'_f$, adică $p \in L(AF')$ și $L(AF) \subseteq L(AF')$.

Pe o cale analoagă se arată că $L(AF') \subseteq L(AF)$. \square

Observație. Faptul că un cuvânt este recunoscut de un automat finit se poate verifica prin calcul direct sau pe diagrama de stări.

Exemplu. Considerăm automatul din exemplul anterior (figura 2.2) și fie $p = i_1 i_2 i_1$. Prin calcul direct:

$$\begin{aligned} f(s_0, i_1 i_2 i_1) &= f(f(f(s_0, i_1), i_2), i_1) = f(f(\{s_1\}, i_2), i_1) = \\ &= f(\{s_0, s_1\}, i_1) = f(\{s_0\}, i_1) \cup f(\{s_1\}, i_1) = \{s_1\} \cup \{s_2\}. \end{aligned}$$

Astfel că $f(s_0, i_1 i_2 i_1) \cap \Sigma_f = \{s_2\} \neq \emptyset$ și $p \in L(AF)$.

Pe diagrama de stări există traiectoriile:

$$\begin{aligned} s_0 &\xrightarrow{i_1} s_1 \xrightarrow{i_2} s_0 \xrightarrow{i_1} s_1; \\ s_0 &\xrightarrow{i_1} s_1 \xrightarrow{i_2} s_1 \xrightarrow{i_1} s_2; \end{aligned}$$

A doua traiectorie ne duce într-o stare finală, deci $p \in L(AF)$.

Limbaje de tipul trei și limbaje regulate. Vom arăta în cele ce urmează că familia limbajelor de tipul 3 coincide cu familia limbajelor regulate. În prealabil vom pune în evidență o formă specială a limbajelor de tipul 3, pe care convenim să o numim *formă normală*.

Definiție 2.2 Vom spune că o gramatică de tipul 3 este în forma normală dacă are reguli de generare de forma

$$\begin{cases} A \rightarrow iB, \\ C \rightarrow j, \end{cases} \quad \text{unde } A, B, C \in V_N, i, j \in V_T$$

sau regula de completare $S \rightarrow \lambda$ și în acest caz S nu apare în dreapta vreunei reguli.

Lema 2.1 Orice gramatică de tipul 3 admite o formă normală.

Demonstrație. Dacă $G = (V_N, V_T, S, P) \in \mathcal{G}_3$ este gramatica dată, eliminăm în primul rând λ -regulile (ca în lema eliminării regulilor de ștergere) apoi construim gramatica $G' = (V'_N, V_T, S, P')$, unde V'_N și P' se definesc astfel: introducem în V'_N toate simbolurile din V_N iar în P' toate regulile din P care convin; fie acum în P o regulă de forma

$$A \rightarrow pB, \quad p = i_1 \dots i_n$$

Vom introduce în P' regulile:

$$\begin{aligned} A &\rightarrow i_1 Z_1, \\ Z_1 &\rightarrow i_2 Z_2, \\ &\dots, \\ Z_{n-1} &\rightarrow i_n B, \end{aligned}$$

iar simbolurile Z_1, \dots, Z_{n-1} le includem în V'_N .

În cazul unei reguli de forma $A \rightarrow p$ cu $|p| > 1$ procedăm analog, exceptând ultima regulă nou introdusă care va avea forma $Z_{n-1} \rightarrow i_n$. Să mai facem observația că simbolurile Z_1, \dots, Z_{n-1} le luăm distincte pentru fiecare caz.

Se poate arăta ușor că $L(G) = L(G')$. \square

Teorema 2.2 *Familia limbajelor de tipul 3 coincide cu familia limbajelor regulate.*

Demonstrație. Partea I: $E \in \mathcal{L}_3 \Rightarrow E \in \mathcal{R}$.

Fie E un limbaj de tipul 3 și $G = (V_N, V_T, S, P)$ gramatica care îl generează; putem presupune că G este în formă normală. Construim automatul finit $AF = (\Sigma, I, f, s_0, \Sigma_f)$ unde $\Sigma = V_N \cup \{X\}$ (X simbol nou), $I = V_T$, $s_0 = S$ și

$$\Sigma_f = \begin{cases} \{X, S\} & , \text{pentru } \lambda \in E \\ \{X\} & , \text{pentru } \lambda \notin E \end{cases}$$

Funcția de evoluție este definită de:

$$\begin{aligned} &\text{dacă } A \rightarrow iB \in P \text{ luăm } B \in f(A, i), \\ &\text{dacă } C \rightarrow j \in P \text{ luăm } X \in f(C, j), \\ &\text{în rest } \emptyset \end{aligned}$$

Observație. Automatul astfel definit este în general nedeterminist. De exemplu, dacă $A \rightarrow 0B|0$ atunci $f(A, 0) = \{B, X\}$.

Fie $p \in L(G)$, $p = i_1 \dots i_n$, deci $S \xRightarrow{*} p$. Detaliat, această derivare va avea forma

$$(1) \quad S \Rightarrow i_1 A_1 \Rightarrow i_1 i_2 A_2 \Rightarrow i_1 i_2 \dots i_{n-1} A_{n-1} \Rightarrow i_1 i_2 \dots i_n.$$

S-au aplicat regulile:

$$(2) \quad \begin{aligned} &S \rightarrow i_1 A_1, \\ &A_1 \rightarrow i_2 A_2, \\ &\dots \\ &A_{n-1} \rightarrow i_n. \end{aligned}$$

În automat avem corespunzător:

$$(3) : \quad \begin{aligned} &A_1 \in f(S, i_1), \\ &A_2 \in f(A_1, i_2), \\ &\dots \\ &X \in f(A_{n-1}, i_n) \end{aligned}$$

Putem scrie traiectoria

$$(4) \quad S \xrightarrow{i_1} A_1 \xrightarrow{i_2} A_2 \xrightarrow{i_3} \dots \xrightarrow{i_n} X \in \Sigma_f$$

Deci $p \in L(AF)$.

Dacă $p = \lambda$, atunci $S \Rightarrow \lambda$ și $S \rightarrow \lambda \in P$. Dar atunci $\lambda \in L(AF)$, căci automatul este în starea S și rămîne în această stare după "citirea" lui λ ; cum însă în acest caz $S \in \Sigma_f$ rezultă că și în acest caz $p \in L(AF)$. În consecință $L(G) \subseteq L(AF)$.

Fie acum $p = i_1 \dots i_n \in L(AF)$; atunci avem traiectoria (4), relațiile (3), regulile de generare (2) și putem scrie derivarea (1), adică $p \in L(G)$ și $L(AF) \subseteq L(G)$. \square

Partea II $E \in \mathcal{R} \Rightarrow E \in \mathcal{L}_3$.

Vom indica numai modul de construcție a gramaticii. Fie $AF = (\Sigma, I, f, s_0, \Sigma_f)$ automatul finit care recunoaște limbajul E , pe care îl presupunem determinist. Construim gramatica $G = (\Sigma, I, s_0, P)$ unde mulțimea P este definită astfel

$$f(A, i) = B \text{ generează regula } A \rightarrow iB \in P, \\ \text{în plus dacă } B \in \Sigma_f \text{ se generează și regula } A \rightarrow i \in P.$$

Putem arăta că $L(G) = L(AF)$. \square

2.2 Proprietăți speciale ale limbajelor regulate

Caracterizarea algebrică a limbajelor regulate. Limbajele regulate, fiind părți din I^* , se pot caracteriza algebric, independent de mecanismele de generare (gramaticile de tipul 3) sau de cele de recunoaștere (automatele finite).

Teorema 2.3 *Fie $E \subset I^*$ un limbaj. Următoarele afirmații sunt echivalente.*

- (a) $E \in \mathcal{R}$;
- (b) E este o reuniune de clase de echivalențe a unei congruențe de rang finit;
- (c) Următoarea congruență

$$\mu = \{(p, q) \mid \chi_E(r_1 p r_2) = \chi_E(r_1 q r_2), \forall r_1, r_2 \in I^*\},$$

unde χ_E este funcția caracteristică a lui E , este de rang finit.

Demonstrație: Vom arăta următoarele implicații: (a) \Rightarrow (b), (b) \Rightarrow (c), (c) \Rightarrow (a).

(a) \Rightarrow (b).

Fie $AF = (\Sigma, I, f, s_0, \Sigma_f)$ automatul finit care recunoaște limbajul E . Definim pe I^* relația

$$\xi = \{(p, q) \mid f(s, p) = f(s, q), \forall s \in \Sigma\}.$$

Se poate vedea cu ușurință că ξ este o relație de echivalență (reflexivă, simetrică, tranzitivă). În plus, dacă $r \in I^*$ și $(p, q) \in \xi$, atunci $(pr, qr) \in \xi$ și $(rp, rq) \in \xi$. De exemplu, prima apartenență se deduce astfel

$$f(s, pr) = f(f(s, p), r) = f(f(s, q), r) = f(s, qr), \text{ etc.}$$

Prin urmare ξ este o relație de congruență.

Să arătăm că această congruență este de rang finit, adică mulțimea cât I^*/ξ este finită.

Fie $\alpha : \Sigma \longrightarrow \Sigma$ o aplicație oarecare și fie mulțimea

$$I^*(\alpha) = \{p | p \in I^*, f(s, p) = \alpha(s), \forall s \in \Sigma\}.$$

Să observăm că dacă α este funcția identică atunci $\lambda \in I^*(\alpha)$. Deci nu toate $I^*(\alpha)$ sunt vide; în acest caz $I^*(\alpha)$ este o clasă de echivalență. Într-adevăr, fie $p \in I^*(\alpha) \subseteq I^*$ fixat și fie C_p clasa de echivalență a lui p . Arătăm că $C_p = I^*(\alpha)$. Dacă $q \in I^*(\alpha)$, atunci $f(s, q) = \alpha(s), \forall s \in \Sigma$, ceea ce înseamnă că $f(s, q) = f(s, p), \forall s \in \Sigma$, și deci $(p, q) \in \xi$ adică $q \in C_p$ și $I^*(\alpha) \subseteq C_p$.

Invers dacă $q \in C_p$ atunci $f(s, q) = f(s, p) = \alpha(s), \forall s \in \Sigma$ și $q \in I^*(\alpha)$, adică $C_p \subseteq I^*(\alpha)$. Aceasta înseamnă că $I^*(\alpha) = C_p$, adică $I^*(\alpha)$ este o clasă de echivalență.

Între mulțimea cât I^*/ξ și mulțimea funcțiilor definite pe Σ cu valori în Σ putem stabili următoarea corespondență biunivocă: *unei funcții $\alpha : \Sigma \longrightarrow \Sigma$ îi corespunde clasa de echivalență $I^*(\alpha)$* . Invers, fiind dată o clasă de echivalență C , luăm $p \in C$ (oarecare) și atașăm lui C funcția $\alpha(s) = f(s, p), \forall s \in \Sigma$. Ținând cont că dacă $q \in C$ atunci $f(s, p) = f(s, q), \forall s \in \Sigma$, rezultă că funcția α nu depinde de elementul p ales.

Dar mulțimea funcțiilor definite pe Σ cu valori în Σ este finită, deci I^*/ξ este finită, adică congruența ξ este de rang finit.

Fie acum $p \in L(AF)$ și q astfel ca $(p, q) \in \xi$. Avem

$$f(s_0, q) = f(s_0, p) \in \Sigma_f;$$

adică $q \in L(AF)$. Aceasta înseamnă că odată cu elementul p , $L(AF)$ conține clasa de echivalență a lui p . De aici rezultă că $L(AF)$ este constituit dintr-un anumit număr de clase de echivalență a lui ξ . \square

(b) \Rightarrow (c)

Fie ξ o congruență de rang finit și E o reuniune de clase de echivalență. Fie apoi $(p, q) \in \xi$; aceasta înseamnă că $r_1 p r_2 \in E \Leftrightarrow r_1 q r_2 \in E$, deci

$$\chi_E(r_1 p r_2) = \chi_E(r_1 q r_2), \forall r_1, r_2 \in I^*.$$

Prin urmare, $(p, q) \in \mu$. Orice clasă de echivalență din I^*/ξ este inclusă într-o clasă de echivalență din I^*/μ , așa încât $\text{card}(I^*/\mu) < \text{card}(I^*/\xi)$, adică congruența μ este de rang finit. \square

(c) \Rightarrow (a)

Presupunem că μ este o congruență de rang finit; considerăm automatul finit $AF = (I^*/\mu, I, f, C_\lambda, \Sigma_f)$, unde funcția de evoluție f și mulțimea de stări finale sunt

$$f(C_p, i) = C_{pi}, \quad \Sigma_f = \{C_p | p \in E\}.$$

Vom arăta că $E = L(AF)$. În primul rând să observăm că $f(C_p, q) = C_{pq}$ (se poate arăta prin inducție asupra lui $|q|$). Avem

$$p \in E \Leftrightarrow C_p \in \Sigma_f \Leftrightarrow f(C_\lambda, p) = C_{\lambda p} = C_p \in \Sigma_f \Leftrightarrow p \in L(AF)$$

adică $E = L(AF)$ și E este un limbaj regulat. \square

Corolar 2.4 Familia R este închisă la operația de răsturnare.

Demonstrație. Fie $E \in \mathcal{R}$ și \tilde{E} răsturnatul lui E . Conform teoremei de caracterizare, $\text{card}(I^*/\mu^E) < \infty$. Avem

$$(p, q) \in \mu^E \Leftrightarrow \chi_E(r_1 p r_2) = \chi_E(r_1 q r_2) \Leftrightarrow$$

și

$$\chi_{\tilde{E}}(\widetilde{r_1 p r_2}) = \chi_{\tilde{E}}(\widetilde{r_1 q r_2}) \Leftrightarrow (\tilde{p}, \tilde{q}) \in \mu^{\tilde{E}}$$

Cu alte cuvinte dacă C este o clasă de echivalență a lui μ^E atunci \tilde{C} (răsturnatul lui C) este o clasă de echivalență a lui $\mu^{\tilde{E}}$. Aceasta înseamnă că $\text{card}(I^*/\mu^{\tilde{E}}) = \text{card}(I^*/\mu^E) < \infty$ și conform aceleiași teoreme de caracterizare $\tilde{E} \in \mathcal{R}$. \square

Observație. Am văzut că limbajele de tipul 3 pot fi definite de gramatici cu reguli de două categorii: drept liniare sau stâng liniare. Este evident că limbajele stâng liniare sunt răsturnatele limbajelor drept liniare. Cum familia limbajelor regulate (drept liniare) este închisă la operația de răsturnare, rezultă că cele două familii de limbaje coincid.

Închiderea familiei \mathcal{L}_3 la operațiile Pref și complementariere.

Operațiile Pref și complementariere se definesc în modul următor

$$\text{Pref}(E) = \{p \mid \exists r \in I^*, pr \in E\}, \quad C(E) = I^* \setminus E.$$

Teorema 2.5 Familia \mathcal{L}_3 este închisă la operațiile Pref și complementariere.

Demonstrație. Fie $AF = (\Sigma, I, f, s_0, \Sigma_f)$ automatul finit care recunoaște limbajul E . Putem presupune că AF este determinist.

Limbajul $\text{Pref}(E)$. Construim automatul finit $AF' = (\Sigma, I, f, s_0, \Sigma'_f)$ unde

$$\Sigma'_f = \{s \in \Sigma \mid s = f(s_0, q), q \in \text{Pref}(E)\}.$$

Este evident că $\text{Pref}(E) \subseteq L(AF')$, căci dacă $q \in \text{Pref}(E)$ atunci $s = f(s_0, q) \in \Sigma'_f$ conform definiției lui Σ'_f .

Să arătăm acum că $L(AF') \subseteq \text{Pref}(E)$. Fie $r \in L(AF')$, atunci $f(s_0, q) \in \Sigma'_f$, deci există $q \in \text{Pref}(E)$ astfel încât $f(s_0, r) = f(s_0, q)$. Cum q este prefixul unui cuvânt din E , există $w \in I^*$ astfel încât $qw \in E$, adică $f(s_0, q) \in \Sigma_f$. Dar

$$f(s_0, rw) = f(f(s_0, r), w) = f(f(s_0, q), w) = f(s_0, qw) \in \Sigma_f,$$

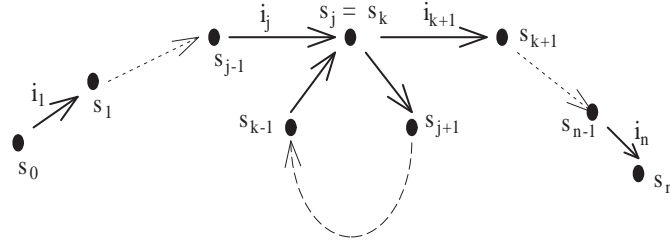


Figura 2.3: Traectoria automatului finit

deci $rw \in E$ și $r \in Pref(E)$. Aceasta înseamnă că $L(AF') \subseteq Pref(E)$ și prin urmare $Pref(E) = L(AF')$, adică $Pref(E)$ este limbaj regulat. \square

Limbajul $C(E)$. Avem

$$C(E) = \{p \in I^* | p \notin E\} = \{p \in I^* | f(s_0, p) \notin \Sigma_f\} = \{p \in I^*, f(s_0, p) \in C(\Sigma_f)\}.$$

Prin urmare $C(E) = L(AF_c)$ unde $AF_c = (\Sigma, I, f, s_0, C(\Sigma_f))$, adică $C(E)$ este un limbaj regulat. \square

Lema de pompare pentru limbaje regulate Sub această denumire (sau lema uvw) este cunoscută o proprietate a limbajelor regulate (ca și a altor familii de limbaje) care ne permite să descompunem cuvintele suficient de lungi ale limbajului în forma uvw și să multiplicăm subcuvântul v de un număr arbitrar de ori, obținând cuvinte care aparțin de asemenea limbajului. Cu alte cuvinte, putem să ”pompăm” în cuvântul dat o anumită parte a sa. Astfel de leme se utilizează deseori pentru a rezolva probleme de neapartenență, adică pentru a arăta că un anumit limbaj nu aparține unei familii date de limbaje.

Lema 2.2 Fie E un limbaj regulat și $AF = (\Sigma, I, f, s_0, \Sigma_f)$ automatul finit care îl recunoaște. Dacă $p \in E$ și $|p| \geq \text{card}(\Sigma)$ atunci p se descompune în forma $p = uvw$, $v \neq \lambda$ și $uv^m w \in E$, $\forall m \in \mathbb{N}$.

Demonstrație. Fie $p = i_1 \dots i_n$, $n \geq \text{card}(\Sigma)$; fie s_0, s_1, \dots, s_n stările parcurse de automat la citirea cuvântului p . Atunci, $s_j = f(s_{j-1}, i_j)$, $j = \overline{1, n}$, $s_n \in \Sigma_f$. Există în mod necesar două stări egale, $s_j = s_k$, $j < k$. Traectoria va avea o buclă (vezi figura 2.3).

Descompunem cuvântul p în forma $p = uvw$ unde $u = i_1 \dots i_j$, $v = i_{j+1} \dots i_k$, $w = i_{k+1} \dots i_n$. Este clar că $v \neq \lambda$, căci $j < k$. Pe de altă parte, putem parcurge traectoria făcând de mai multe ori bucla, adică

$$f(s_0, uv^m w) = s_n \in \Sigma_f.$$

Prin urmare $uv^m w \in E$. \square

Consecință 2.3 *Incluziunea $\mathcal{L}_3 \subseteq \mathcal{L}_2$ este strictă.*

În adevăr, fie limbajul $L_2 = \{0^n 1^n | n \geq 1\}$. Știm că acest limbaj este de tipul 2 și că poate fi generat de gramatica $G = (\{A\}, \{0, 1\}, A, \{A \rightarrow 0A1 | 01\})$. Să arătăm că L_2 nu este de tipul 3.

Să presupunem că L_2 este de tipul 3 și fie $AF = (\Sigma, I, f, s_0, \Sigma_f)$ automatul finit care îl recunoaște. Cum L_2 conține cuvinte oricât de lungi, fie $p \in L_2$ astfel încât $p \geq \text{card}(\Sigma)$. Conform lemei de pompare, p se descompune în forma $p = uvw$, $v \neq \lambda$ și $uv^m w \in E$. Putem avea una din situațiile:

$$\begin{aligned} (1) \quad p &= \underbrace{0 \dots 0}_u \underbrace{0 \dots 0}_v \underbrace{01 \dots 1}_w, \\ (2) \quad p &= \underbrace{0 \dots 01 \dots 1}_u \underbrace{1 \dots 1}_v \underbrace{1 \dots 1}_w, \\ (3) \quad p &= \underbrace{0 \dots 0}_u \underbrace{0 \dots 1}_v \underbrace{1 \dots 1}_w. \end{aligned}$$

Primele două cazuri nu pot avea loc deoarece multiplicându-l pe v , numărul de simboluri 0 și 1 nu s-ar păstra egal. În al treilea caz, luând de exemplu, $m = 2$ obținem

$$p_2 = 0 \dots 00 \dots 10 \dots 11 \dots 1 \in L_2$$

ceea ce din nou nu este posibil, întrucât se contrazice structura cuvintelor lui L_2 . Prin urmare L_2 nu este de tipul 3. \square

Observație. Este interesant de observat că limbajele simple de forma lui L_2 sunt semnificative pentru clasele din clasificarea Chomsky. Astfel

$$\begin{aligned} L_1 &= \{a^n | n \geq 1\}, \quad L_1 \in \mathcal{L}_3; \\ L_2 &= \{a^n b^n | n \geq 1\}, \quad L_2 \in \mathcal{L}_2, \quad L_2 \notin \mathcal{L}_3; \\ L_3 &= \{a^n b^n c^n | n \geq 1\}, \quad L_3 \in \mathcal{L}_1(?), \quad L_3 \notin \mathcal{L}_2; \end{aligned}$$

Ne-am putea aștepta ca limbajul L_3 , un exemplu analog lui L_2 , să fie de tip 1, adică $L_3 \in \mathcal{L}_1$, $L_3 \notin \mathcal{L}_2$. În adevăr, se poate arăta că $L_3 \notin \mathcal{L}_2$, dar după cunoștința autorului, aparența $L_3 \in \mathcal{L}_1$ este o problemă deschisă.

Consecință 2.4 *Fie E un limbaj regulat și $AF = (\Sigma, I, f, s_0, \Sigma_f)$ automatul finit care îl recunoaște. Atunci E este infinit dacă și numai dacă există $p \in E$ astfel încât $|p| \geq \text{card}(\Sigma)$.*

Dacă limbajul este infinit, este clar că există $p \in E$ cu $|p| \geq \text{card}(\Sigma)$. Invers, dacă există $p \in E$ cu $|p| \geq \text{card}(\Sigma)$ atunci $p = uvw$, $v \neq \lambda$ și $uv^m w \in E, \forall m \in \mathbb{N}$, deci limbajul este infinit. \square

2.3 Expresii regulate și sisteme tranzitionale

Expresii regulate Fie V un alfabet. Expresiile regulate sunt cuvinte peste alfabetul $V \cup \{\bullet, *, |\} \cup \{(\cdot)\}$ la care se adaugă simbolul \emptyset . Simbolurile $|$ -sau, \bullet -produs, $*$ -î închidere, le vom considera operatori. Expresiile regulate se definesc astfel:

- (1) λ și \emptyset sunt expresii regulate;
- (2) pentru orice $a \in V$, cuvântul a este o expresie regulată;
- (3) dacă R și S sunt expresii regulate, atunci $(R|S)$ ($R \bullet S$) și (R^*) sunt expresii regulate;
- (4) orice expresie regulată se obține prin aplicarea de un număr finit de ori a regulilor (1)-(3).

Parantezele sunt utilizate pentru a pune în evidență ordinea de aplicare a operatorilor. Pentru simplificarea scrierii vom considera că operatorul $*$ are ponderea cea mai mare, apoi operatorul \bullet și $|$ ponderea cea mai mică. Astfel parantezele redundante pot fi eliminate. De exemplu, prin $R|S^*$ vom înțelege $(R|(S^*))$.

Observație. Expresiile regulate se pot defini cu ajutorul gramaticii $G = (\{E\}, V \cup \{|\cdot, \bullet, *, (\cdot)\}, E, P)$ unde

$$P = \{E \rightarrow (E|E) \mid (E \bullet E) \mid (E^*) \mid a \mid \lambda \mid \emptyset\}.$$

Unei expresii regulate îi putem asocia un anumit limbaj peste V ; vom spune că expresia regulată reprezintă (desemnează, notează) acel limbaj. Modul în care asociem un limbaj unei expresii regulate este

- (1) λ reprezintă limbajul $\{\lambda\}$;
- (1') \emptyset reprezintă limbajul \emptyset ;
- (2) a reprezintă limbajul $\{a\}$;
- (3) dacă R și S sunt expresii regulate ce reprezintă limbajele L_R respectiv L_S atunci
 - (i) $R|S$ reprezintă limbajul $L_R \cup L_S$;
 - (ii) $R \bullet S$ reprezintă limbajul $L_R L_S$;
 - (iii) R^* reprezintă limbajul $(L_R)^*$.

Fie R, S, P trei expresii regulate și L_R, L_S, L_P limbajele reprezentate. Avem :

$$L_{(R|S)|P} = (L_R \cup L_S) \cup L_P = L_R \cup (L_S \cup L_P) = L_{R|(S|P)},$$

întrucât operația de reuniune este asociativă. Vom scrie $(R|S)|P = R|(S|P)$. În mod analog se pot obține și alte proprietăți. De exemplu:

$$\begin{aligned} S|R &= S|R, \\ R|R &= R, \\ (R \bullet S) \bullet P &= R \bullet (S \bullet P), \\ R \bullet (S|P) &= (R \bullet S)|(R \bullet P), \text{ etc.} \end{aligned}$$

În cazul în care nu există pericol de confuzie, vom nota cu L (fără indice) limbajul reprezentat de o anumită expresie regulată.

Exemple.

1. $R = a^*$; $L = \bigcup_{j=0}^{\infty} \{a^j\} = \{\lambda, a, a^2, \dots\}$.
2. $R = aa^*$; $L = a \bullet \{\lambda, a, a^2, \dots\} = \{a, a^2, a^3, \dots\}$.
3. $R = (a|b)^*$; $L = (L_a \cup L_b)^* = (\{a\} \cup \{b\})^* = \{a, b\}^*$;
 $\{a, b\}^* = \{\lambda\} \cup \{a, b\}^1 \cup \{a, b\}^2 \cup \dots = \{\lambda, a, b, aa, ab, ba, bb, \dots\}$, adică
 $(a|b)^*$ reprezintă mulțimea tuturor cuvintelor peste alfabetul $\{a, b\}$.
4. $R = a|ba^*$; $L = \{a\} \cup \{b\} \bullet \{\lambda, a, a^2, \dots\} = \{a, b, ba, ba^2, \dots\}$.

Limbajele reprezentate de expresii regulate constituie o anumită familie de limbae; o vom nota cu \mathcal{L}_{lr} . Apare următoarea problemă: *care este poziția acestei familii în ierarhia Chomsky?* Vom arăta că \mathcal{L}_{lr} coincide cu familia limbajelor regulate.

Sisteme tranziționale

Definiție 2.3 *Un sistem tranzițional este un sistem de forma*

$$ST = (\Sigma, I, f, \Sigma_0, \Sigma_f, \delta)$$

unde:

- Σ este o mulțime (finită) de stări;
- I este alfabetul de intrare;
- $f : \Sigma \times I \longrightarrow \mathcal{P}(\Sigma)$ este funcția de tranziție;
- $\Sigma_0 \subseteq \Sigma$ este mulțimea de stări inițiale;
- $\Sigma_f \subseteq \Sigma$ este mulțimea de stări finale;
- $\delta \subset \Sigma \times \Sigma$ este relația de tranziție.

Exemplu. $\Sigma = \{s_0, s_1, s_2\}$, $I = \{0, 1\}$, $\Sigma_0 = \{s_0, s_1\}$, $\Sigma_f = \{s_2\}$ iar funcția și relația de tranziție sunt date de:

f	s_0	s_1	s_2
0	$\{s_1\}$	$\{s_2\}$	$\{s_0\}$
1	\emptyset	$\{s_0, s_1\}$	$\{s_0, s_2\}$

$\delta = \{(s_0, s_1), (s_2, s_1)\}$.

Ca și în cazul unui automat finit putem construi diagrama de stări completată cu relația δ (arcele punctate). În cazul exemplului nostru diagrama de stări este prezenta în figura 2.4

Observație: δ fiind o relație, are sens δ^* (închiderea tranzitivă și reflexivă).

Fie $i \in I \cup \{\lambda\}$; vom spune că sistemul tranzițional evoluează direct din starea s' în starea s'' dacă:

(1) $i = \lambda$ și $(s', s'') \in \delta^*$. Pe diagrama de stări vom avea o traiectorie punctată de la starea s' la starea s'' ;

$$s' \longrightarrow O \longrightarrow O \longrightarrow \dots \longrightarrow O \longrightarrow s''.$$

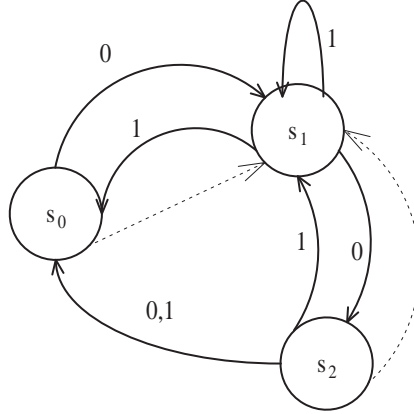


Figura 2.4: Diagrama de stări a sistemului tranzițional

(2) $i \neq \lambda$ și există $s_1, s_2 \in \Sigma$ astfel încât $(s', s_1) \in \delta^*$, $s_2 \in f(s_1, i)$ și $(s_2, s'') \in \delta^*$. Pe diagrama de stări, putem ajunge din s' în s'' pe o traiectorie punctată, apoi un pas pe un arc plin și din nou pe o traiectorie punctată.

$$s' \longrightarrow O \longrightarrow \dots \longrightarrow s_1 \xrightarrow{i} s_2 \longrightarrow O \longrightarrow \dots \longrightarrow s''.$$

Vom scrie $s' \vdash^i s''$.

Fie acum $p = i_1 \dots i_n$. Vom spune că sistemul evoluează din starea s' în starea s'' dacă există s_0, s_1, \dots, s_n astfel încât

$$s' = s_0 \vdash^{i_1} s_1 \vdash^{i_2} \dots \vdash^{i_n} s_n = s''.$$

Vom scrie $s' \vdash^p s''$.

Definiție 2.4 *Limbaajul recunoscut de un sistem tranzițional ST este*

$$L(ST) = \{p \mid p \in I^*, \exists s_0 \in \Sigma_0, s_0 \vdash^p s, s \in \Sigma_f\}.$$

Vom nota cu \mathcal{L}_{ST} familia limbajelor recunoscute de sisteme tranziționale. Este evident că orice automat finit este un sistem tranzițional particular în care $\text{card}(\Sigma_0) = 1$ iar $\delta = \emptyset$ (nu există arce punctate). Prin urmare $\mathcal{R} \subseteq \mathcal{L}_{ST}$.

Teorema 2.6 $\mathcal{R} = \mathcal{L}_{ST}$.

Demonstrație. Evident, trebuie să arătăm incluziunea $\mathcal{L}_{ST} \subseteq \mathcal{R}$. Fie $ST = (\Sigma, I, f, \Sigma_0, \Sigma_f, \delta)$ un sistem tranzițional. Construim automatul finit $AF = (\mathcal{P}(\Sigma), I, f', \Sigma_0, \Sigma'_f)$ unde

$$f'(Z, i) = \{s \mid \exists s' \in Z, s' \vdash^i s\},$$

$$\Sigma'_f = \{Z \mid Z \cap \Sigma_f \neq \emptyset\}.$$

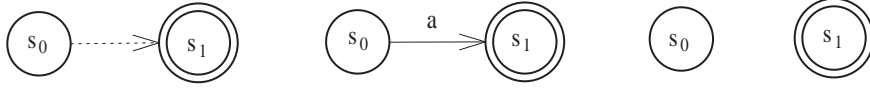


Figura 2.5: Sistemele tranziționale ce recunosc limbajele λ , a , \emptyset .

Fie $p = i_1 \dots i_n \in L(ST)$ și fie următoarea evoluție a sistemului tranzițional

$$s_0 \xrightarrow{i_1} s_1 \xrightarrow{i_2} \dots \xrightarrow{i_n} s_n \in \Sigma_f.$$

Putem construi o traiectorie a lui AF de forma

$$\Sigma_0 \xrightarrow{i_1} Z_1 \xrightarrow{i_2} \dots \xrightarrow{i_n} Z_n,$$

unde $Z_1 = f'(\Sigma_0, i_1)$, $Z_k = f'(Z_{k-1}, i_k)$, $k = 2, \dots, n$. Să observăm că $s_0 \in \Sigma_0$ și că dacă $s_{k-1} \in Z_{k-1}$, atunci conform definiției funcției f' , avem $s_k \in f'(Z_{k-1}, i_k) = Z_k$. Astfel, $s_k \in Z_k$, $k = 1, \dots, n$; pentru $k = n$ avem $s_n \in Z_n$ și cum $s_n \in \Sigma_f$ rezultă că $Z_n \cap \Sigma_f \neq \emptyset$, adică $Z_n \in \Sigma'_f$. Deci automatul ajunge într-o stare finală, $p \in L(AF)$ și $L(ST) \subseteq L(AF)$.

Incluziunea inversă se arată în mod analog. \square

Construcția sistemelor tranziționale pentru expresii regulate. Fiind dată o expresie regulată putem întotdeauna construi un sistem tranzițional care recunoaște limbajul reprezentat de expresia respectivă.

Construcția se face cu ajutorul diagramelor de stări.

Sistemele tranziționale (diagramele de stări) corespunzătoare expresiilor regulate λ , a și \emptyset sunt prezentate în figura 2.5.

Dacă R și S sunt expresii regulate și notăm cu ST_R și ST_S sistemele tranziționale corespunzătoare, atunci sistemele tranziționale pentru $R|S$, $R \bullet S$ și R^* sunt redat în figura 2.6.

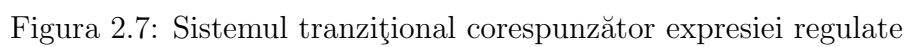
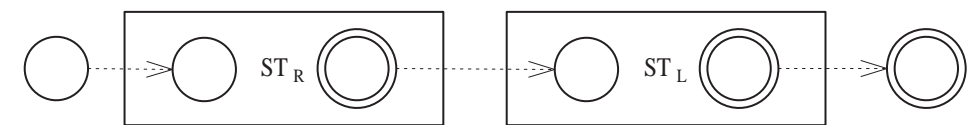
În acest mod putem construi succesiv (recurent) un sistem tranzițional corespunzător unei expresii regulate.

Exemplu. Sistemul tranzițional corespunzător expresiei $R = a|b \bullet a^*$ este redat în figura 2.7.

Consecință Dându-se o expresie regulată, putem construi sistemul tranzițional care recunoaște limbajul reprezentat de expresia respectivă. Cum orice limbaj recunoscut de un sistem tranzițional este regulat, rezultă că limbajele reprezentate de expresii regulate sunt limbaje regulate.

2.4 Analiza lexicală

Procesul de analiză lexicală este o fază a procesului de compilare în care se determină unitățile lexicale (cuvintele, atomii) ale unui program sursă, se furnizează



codurile interne ale acestora și se detectează eventualele erori lexicale. Analizorul lexical mai poate efectua o serie de operații auxiliare precum: eliminarea blank-urilor, ignorarea comentariilor, diferite conversiuni ale unor date, completarea tabelelor compilatorului, gestiunea liniilor textului sursă.

Unități lexicale O unitate lexicală (Lexic = vocabular; totalitatea cuvintelor unei limbi) este o secvență din textul sursă care are o anumită unitate logică. Definiția riguroasă a unităților lexicale ale unui limbaj particular se dă la definirea limbajului de programare respectiv. De obicei, în majoritatea limbajelor de programare, unitățile lexicale sunt: *cuvinte cheie*, *identificatori*, *constante*, *operatori*, *delimitatori*. Din punctul de vedere al analizei lexicale și al modului de prelucrare, unitățile lexicale pot fi de două categorii:

- Unități lexicale simple, sunt unități lexicale care nu comportă atribute suplimentare, de exemplu, cuvintele cheie, operatorii;
- Unități lexicale compuse (atributive), sunt unități lexicale care comportă anumite atribute suplimentare, de exemplu, identificatorii și constantele. Atributele sunt informații specifice, de exemplu, tipul identificatorului sau al constantei (întreg, real, etc.). Este necesară specificarea tipului unui identificator sau al unei constante din startul programului deoarece structura programului obiect sau reprezentarea internă a constantelor depinde de acest tip.

Reprezentarea internă a unităților lexicale se face în funcție de categoria lor. Cele simple se reprezintă printr-un cod specific (număr întreg). Unitățile lexicale compuse se reprezintă prin cod și informații (de natură semantică) asupra sa. De obicei compilatoarele utilizează tabele pentru stocarea atributelor (tabel de constante, tabel de variabile, tabel de etichete, etc.). În acest caz unitatea lexicală se reprezintă intern printr-un cod urmat de o referință într-un tabel. Informația conținută de tabel ar putea fi pentru constante: cod, tip, valoare, iar pentru identificatori: cod, tip, indicator de inițializare.

Este posibil ca o clasă de unități lexicale simple să se reprezinte printr-un cod unic și un atribut pentru distingerea în cadrul clasei. De exemplu, operatorii aritmetici cu aceeași prioritate au o tratare similară din punct de vedere al analizei sintactice. Lista unităților lexicale și definiția riguroasă a acestora se dă la proiectarea compilatorului.

Un exemplu de unități lexicale și coduri asociate ar putea fi cele din figura 2.8.

Analizorul primește textul sursă și produce șirul de unități lexicale în codificarea internă. De exemplu secvența de text sursă următoare:

```
{if (unu < 2) return 0;  
  a=33;  
}
```


Unitate lexicală	COD	ATRIBUT	Exemplu
if	<u>if</u> = 1	-	if, If, IF
else	<u>else</u> = 2	-	else ElSe
identificator	<u>ID</u> = 3	referință	Nelu v tabel
constantă întreagă	<u>NUM</u> = 4	referință	4 -3 233
constantă reală	<u>FLOAT</u> = 5	referință	4.32 -3.233
+	<u>op</u> = 6	1	
-	<u>op</u> = 6	2	
×	<u>op</u> = 6	3	
/	<u>op</u> = 6	4	
<	<u>opr</u> = 7	1	
>	<u>opr</u> = 7	2	
<=	<u>opr</u> = 7	3	
>=	<u>opr</u> = 7	4	
(<u>LPAR</u> = 8	-	
)	<u>RPAR</u> = 9	-	
{	<u>LBRACE</u> = 10	-	
}	<u>RBRACE</u> = 11	-	

Figura 2.8: Coduri asociate unităților lexicale

va produce sirul de unități lexicale

LBRACE If LPAR [ID,22] [opr,1] [NUM, 40], RPAR RETURN [NUM, 42] SEMI [ID,24] opAssign [NUM,44] SEMI RBRACE.

În lista codurilor interne găsite atributul identificatorului este adresa relativa din tabela de identificatori, analog atributele constantelor numerice sunt adrese relative în tabelele de constante.

Majoritatea limbajelor evaluate conțin și secvențe de text care nu sunt unități lexicale, dar au acțiuni specifice asociate. De exemplu:

```
comentarii          /* text */

directive de preprocesare  #include<stdio.h>
                           #define MAX 5.6
```

Înainte de analiza lexicală (sau ca subrutină) se preprocesează textul și abia apoi se introduce rezultatul în analizorul lexical.

Specificarea unităților lexicale Definirea riguroasă a unităților lexicale se face de către proiectantul limbajului. O posibilitate de descriere este limbajul natural. De exemplu, pentru **C** și **JAVA**:

”Un identificator este o secvență de litere și cifre: primul caracter trebuie să fie literă. Linia de subliniere contează ca literă. Literele mari și mici sunt diferite. Dacă șirul de intrare a fost împărțit în unități lexicale până la un caracter dat, noua unitate lexicală se consideră astfel încât să includă cel mai lung șir de caractere ce poate constitui o unitate lexicală. Spațiile, taburile, newline și comentariile sunt ignorate cu excepția cazului când servesc la separarea unităților lexicale. Sunt necesare spații albe pentru separarea identificatorilor, cuvintelor cheie și a constantelor.”

Orice limbaj rezonabil poate fi folosit pentru implementarea unui analizor lexical.

Unitățile lexicale se pot specifica cu ajutorul limbajelor regulate, deci folosind gramatici de tipul 3 sau expresii regulate ce notează limbajele. Ambele specificații conduc la construirea de automate finite echivalente, care se pot ușor programa. În cele ce urmează, vom folosi ambele variante de specificare pentru un set uzual de unități lexicale întâlnit la majoritatea limbajelor evaluate.

Considerăm gramatica regulată ce generează identificatori, constante întregi, cuvinte cheie, operatori relaționali și aritmetici.

$$G : \left\{ \begin{array}{l} < ul > \rightarrow < id > \mid < num > \mid < cc > \mid < op > \mid < opr > \\ < id > \rightarrow l < id1 > \mid l, < id1 > \rightarrow l < id1 > \mid c < id1 > \mid l \mid c \\ < num > \rightarrow c < num > \mid c \\ < cc > \rightarrow if \mid do \mid else \mid for \\ < op > \rightarrow + \mid - \mid * \mid / \\ < opr > \rightarrow < \mid < = \mid > \mid > = \end{array} \right. ,$$

unde l -litera, c -cifra.

Pornind de la această gramatică se poate construi una echivalentă în formă normală, apoi se extrage funcția de evoluție a automatului finit determinist echivalent ce recunoaște unitățile lexice.

Descrieri echivalente ale unităților lexice cu ajutorul expresiilor regulate sunt

`cuvinte cheie = if | do | else | for`

`identificatori = (a|b|c|...z)(a|b|c|...z|0|1|...|9)*`

`Numar = (0|1|...|9)(0|1|...|9)*`

`Operatori aritmetici = + | - | * | /`

`Operatori relationali = < | <= | > | >=`

Limbajul generat de gramatica precedentă se obține prin suma expresiilor regulate. Menționăm că există programe specializate (Lex, Flex, JavaCC) ce generează un analizor lexical (în C sau Java) pornind de la expresiile regulate. Sintaxa folosită în scrierea expresiilor regulate este dependentă de programul folosit.

Programarea unui analizor lexical Realizarea efectivă a unui analizor revine la simularea funcționării unui automat finit. O variantă de programare este atașarea unei secvențe de program la fiecare stare a automatului. Dacă starea nu este stare finală atunci secvența citește următorul caracter din textul sursă și găsește următorul arc din diagrama de stări. Depinzând de rezultatul căutării se transferă controlul altei stări sau se returnează eșec (posibilă eroare lexicală). Dacă starea este finală atunci se apelează secvența de returnare a codului unității lexice și eventuala instalare a unității lexice în tabelele compilatorului.

Pentru simplificarea implementării se caută următoarea unitate lexicală prin încercarea succesivă a diagramelor corespunzătoare fiecărei unități lexice (într-o ordine prestabilită). Eroarea lexicală se semnalează doar atunci când toate încercările se încheie cu eșec.

De obicei textul sursă conține și secvențe ce se pot descrie cu ajutorul expresiilor regulate, dar care nu sunt unități lexice (de exemplu comentariile). Aceste secvențe nu vor genera cod lexical, dar au asociate diverse acțiuni specifice. Pentru a evita apariția unor caractere necunoscute în textul sursă se consideră și limbajul ce constă din toate simbolurile ASCII. Astfel, indiferent de unde începe analiza textului, programul de analiză lexicală găsește o potrivire cu o descriere. Spunem că specificația este **completă**.

Există posibilitatea ca mai multe secvențe cu aceeași origine să corespundă la diferite descrieri ale unităților lexice. Se consideră unitate lexicală cel mai lung șir ce se potrivește unei descrieri (**longest match rule**). Dacă sunt două reguli

care se potrivesc la același șir de lungime maximă atunci se consideră o prioritate asupra descrierilor (**rule priority**).

De exemplu, în textul următor,

```
i  if  if8
```

unitățile lexicale delimitate vor fi **i**–identificator, **if**–cuvânt cheie, **if8**–identificator. Regula de prioritate se aplică pentru potrivirea lui **if** cu identificator și cuvânt cheie, iar criteriul de lungime maximă pentru **if8**.

Pentru depistarea celei mai lungi potriviri, din punct de vedere al programării analizorului, este suficient să prevedem un pointer suplimentar pe caracterul ce corespunde ultimei stări finale atinse pe parcursul citirii.

Studiu de caz. Se consideră problema realizării unui analizor lexical ce delimitează într-un text sursă cuvinte din limbajul ce conține identificatori, cuvinte cheie (pentru simplificare folosim doar cuvântul cheie **if**), constante numerice (întregi fără semn). De asemenea se face salt peste spațiile albe și se ignoră comentariile. Presupunem că un comentariu începe cu două caractere slash și se termină cu newline. Orice caracter ce nu se potrivește descrierii este semnalat ca și caracter ilegal în text.

Etapa I. Descriem secvențele cu ajutorul expresiilor regulate

IF = "if"

ID = (a|b|c|...|z)(a|b|c|...|z|0|1|...|9)*

NUM = (0|1|...|9)(0|1|...|9)* = (0|1|...|9)+

WS = (\n|\t|" ")+

COMMENT = "//"(a|b|c|...|z|0|1|...|9|" ")*\n

ALL = a|b|...|z|0|...|9|\t| ... toate caracterele ASCII

Etapa II. Corespunzător expresiilor avem următoarele automate finite deterministe echivalente, prezentate în figura 2.9 (stările au fost notate prin numere întregi):

Etapa III. Se construiește sistemul tranzițional (vezi figura 2.10) ce recunoaște limbajul reuniune, adăugând o nouă stare inițială (notată cu 1), pe care o conectăm prin arce punctate (ce corespund λ -tranzițiilor). Construcția provine din *legarea* sistemelor tranziționale *în paralel*. S-au renumerotat stările sistemului tranzițional, asignând nume simbolice stărilor finale.

Etapa III. Construcția automatului finit determinist general ce recunoaște reuniunea limbajelor. Pentru aceasta sistemul tranzițional se transformă cu teorema de echivalență în automat finit determinist (practic se trece de la stări ale sistemului tranzițional la submulțimi de stări, ce devin stările automatului finit).

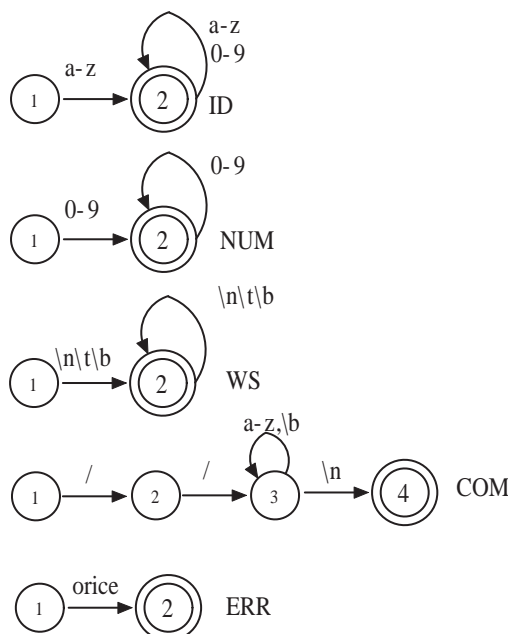


Figura 2.9: Automatele finite corespunzătoare expresiilor regulate

În cazul particular al automatului nostru, diagrama de stări este dată în figura 2.11.

Etapa IV. Programarea analizorului lexical.

Automatul finit obținut are stările finale asociate cu clasele de cuvinte recunoscute. Se asociază acțiuni stărilor finale, corespunzător definițiilor unităților lexicale (de exemplu pentru constante numerice se generează reprezentarea internă, se memorează în tabelul de constante și se returnează codul unității lexicale NUM). Pentru programarea analizorului se folosesc trei variabile de tip pointer în textul sursă: **FirstSymbol**, **CurrentSymbol**, **LastFinalSymbol**, ce rețin indicii caracterului de început al secvenței, indicele caracterului ce urmează la citire, indicele ultimului caracter ce corespunde atingerii unei stări finale. Cei trei pointeri au fost reprezentați prin semnele grafice $|$, \perp respectiv \top . De asemenea considerăm o variabilă *State*, ce reține starea curentă a automatului.

În tabelul 2.12 este indicată evoluția analizorului lexical (inclusiv acțiunile asociate) pentru cazul analizei următorului text sursă.

```
if if8%// ha\n
```

Pentru simplificarea codificării, stările automatului finit determinist au fost redenumite prin numere întregi începând cu starea inițială 1, starea $\{3, 6, 16\} = 2$, $\{4, 6\} = 3$, $\{6, 16\} = 4$, ș.a.m.d. de la stânga la dreapta și de sus în jos. De obicei

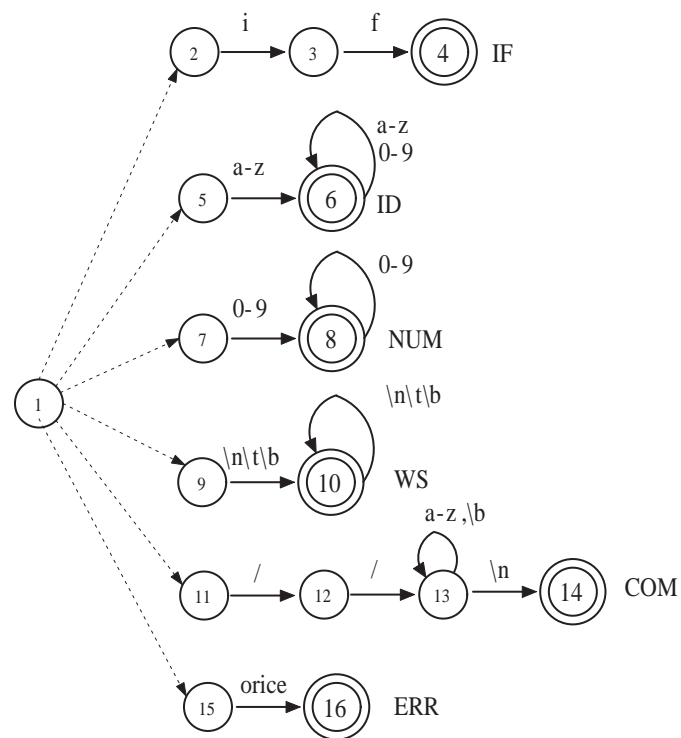


Figura 2.10: Sistemul tranzițional ce recunoaște reuniunea limbajelor

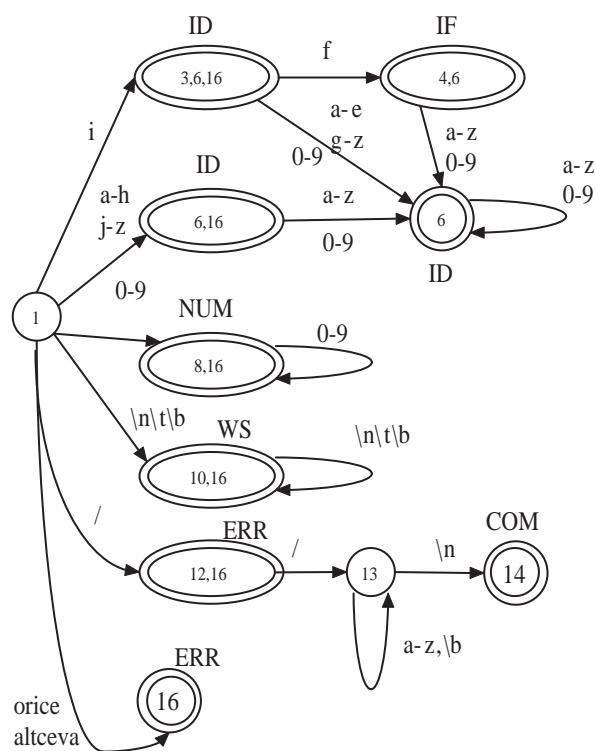


Figura 2.11: Automatul finit determinist ce recunoaște reuniunea limbajelor

Last Final	Current State	Current Input	Accept Action
0	1	$\begin{array}{l} \top i f \quad i f 8 \% // \quad h a \backslash n \\ i \top f \quad i f 8 \% // \quad h a \backslash n \\ i f \top i f 8 \% // \quad h a \backslash n \\ i f \top \quad i f 8 \% // \quad h a \backslash n \end{array}$	return $cc = < if >$ <i>resume</i>
2	2		
3	3		
	0		
0	1	$\begin{array}{l} i f \top i f 8 \% // \quad h a \backslash n \\ i f \top i f 8 \% // \quad h a \backslash n \\ i f \top i f 8 \% // \quad h a \backslash n \end{array}$	<i>resume</i>
7	7		
	0		
0	1	$\begin{array}{l} i f \top i f 8 \% // \quad h a \backslash n \\ i f i \top f 8 \% // \quad h a \backslash n \\ i f i f \top 8 \% // \quad h a \backslash n \\ i f i f 8 \top \% // \quad h a \backslash n \\ i f i f 8 \top \% \perp // \quad h a \backslash n \end{array}$	return $id = < if8 >$ <i>resume</i>
2	2		
3	3		
5	5		
	0		
0	1	$\begin{array}{l} i f \quad i f 8 \top \% // \quad h a \backslash n \\ i f \quad i f 8 \% \top // \quad h a \backslash n \\ i f \quad i f 8 \% \top / \perp \quad h a \backslash n \end{array}$	print ("illegal character: %"); <i>resume</i>
11	11		
	0		
0	1	$\begin{array}{l} i f \quad i f 8 \% \top // \quad h a \backslash n \\ i f \quad i f 8 \% \top / \perp \quad h a \backslash n \\ i f \quad i f 8 \% \top // \perp \quad h a \backslash n \end{array}$	
0	8		
0	9		
...	

Figura 2.12: Evoluția analizorului lexical pentru textul $if\ if8\ \%//\ ha\ \backslash n$


```

int edges[] [] = { /* ... 0 1 2 ... e f g h i j ... */
/* state 0 */      {0,0, ...,0,0,0, ...,0,0,0,0,0,0, ... },
/* state 1 */      {0,0, ...,6,6,6, ...,4,4,4,4,2,4, ... },
/* state 2 */      {0,0, ...,5,5,5, ...,5,3,5,5,5,5, ... },
etc
}

```

Figura 2.13: Reprezentarea funcției de evoluție a automatului finit

funcția de evoluție asociată automatului finit determinist se memorează sub forma unui tablou bidimensional de întregi, ca în figura 2.13. Starea 0 este asociată cu blocarea automatului. Ajungerea în această stare echivalează cu găsirea ultimei unități lexicale, între pointerii $|$ și \top . Se execută acțiunea asociată stării finale și se reia căutarea (*resume*) următoarei unități lexicale începând cu caracterul imediat următor pointerului \top .

Observație: Cele mai costisitoare operațiuni (ca timp) din analiza lexicală sunt ignorarea comentariilor și tratarea erorilor lexicale. Primele generatoare automate de analizare lexicale și sintactice au apărut în anii '70 și au fost incluse în sistemul de operare Unix.

2.5 Probleme propuse

1. Construiți automate finite pentru recunoașterea limbajelor:
 - (a) $L = \{PSDR, PNL, PUNR\}$;
 - (b) $L = \{w \mid \text{șiruri de 0 și 1 terminate cu 1}\}$;
 - (c) $L = \{w \mid w \text{ identificator PASCAL}\}$;
 - (d) $L = \{w \mid w \text{ constantă întreagă cu semn în PASCAL}\}$;
 - (e) $L = \{w \in \{0,1\}^* \mid w \text{ multiplu de 3}\}$;
 - (f) $L = \{a^i b^j \mid i, j > 0\}$;
 - (g) $L = \emptyset$.
2. Construiți automate finite echivalente cu gramaticile de tipul trei de la problema 1 capitolul 1.

3. Construiți automate finite deterministe echivalente cu cele nedeterminate obținute la problema precedentă.
4. Găsiți gramatici regulate echivalente cu automatele de la problema 1.
5. Folosind lema de pompă pentru limbaje regulate dovediți că următoarele limbaje nu sunt regulate:
 - (a) $L = \{0^{i^2} | i \geq 1\}$;
 - (b) $L = \{0^{2^n} | n \geq 1\}$;
 - (c) $L = \{0^n | n \text{ este număr prim } \}$;
 - (d) $L = \{0^m 1^n 0^{m+n} | m \geq 1, n \geq 1\}$;
6. Specificați limbajele denotate de următoarele expresii regulate:
 - (a) $(11|0)^*(00|1)^*$;
 - (b) $(1|01|001)^*(\lambda|0|00)$;
 - (c) $10|(0|11)0^*1$;
 - (d) $((0|1)(0|1))^*$;
 - (e) $01^*|1$;
 - (f) $((11)^*|101)^*$.
7. Construiți sisteme tranziționale ce recunosc limbajele specificate la problema precedentă. Pentru fiecare sistem tranzițional construiți un automat finit determinist echivalent.

Capitolul 3

Limbaje Independente de Context

3.1 Arbori de derivare

Caracterizarea familiei \mathcal{L}_2 cu arbori de derivare. Una din caracteristicile de bază ale limbajelor independente de context este aceea că o derivare într-un astfel de limbaj poate fi reprezentată de un arbore, numit în acest context *arbore de derivare*. Această reprezentare este importantă în mod special pentru faptul că permite o imagine intuitivă simplă a unei derivări și deci posibilitatea de a lucra ușor cu limbaje de tipul 2.

Vom prezenta în primul rând câteva noțiuni elementare de teoria grafurilor, cu scopul de a preciza notațiile și terminologia.

Un *graf orientat* \mathcal{G} este o pereche $\mathcal{G} = (V, \Gamma)$ unde V este o mulțime finită iar Γ o aplicație $\Gamma : V \longrightarrow \mathcal{P}(V)$. Mulțimea V se numește mulțimea vârfurilor (nodurilor) grafului iar dacă $v_2 \in \Gamma(v_1)$, perechea (v_1, v_2) este un arc în graf; v_1 este *originea* iar v_2 este *extremitatea* arcului. Un *drum* de la vârful v' la vârful v'' în graful \mathcal{G} este o mulțime de arce $(v_1, v_2)(v_2, v_3) \dots (v_{n-1}, v_n)$ cu $v' = v_1$ și $v'' = v_n$. Numărul $n - 1$ este *lungimea drumului*. Un drum pentru care $v_1 = v_n$ se numește *circuit*. Un circuit de lungime 1 poartă numele de *bucă*.

Definiție 3.1 *Un arbore este un graf fără circuite, cu $\text{card}(V) \geq 2$ și care satisface următoarele două condiții :*

1. $\exists v_0 \in V$ astfel încât $v_0 \notin \Gamma(v), \forall v \in V$; v_0 se numește **rădăcina arborelui**;
2. $\forall v \in V \setminus \{v_0\}, \exists ! w$ cu $v \in \Gamma(w)$; altfel spus orice vârf diferit de v_0 este extremitatea unui singur arc.

Exemplu. $V = \{v_0, v_1, v_2, v_3, v_4\}$ iar funcția Γ este dată de:

x	v_0	v_1	v_2	v_3	v_4
$\Gamma(x)$	$\{v_1, v_2\}$	\emptyset	$\{v_3, v_2\}$	\emptyset	\emptyset

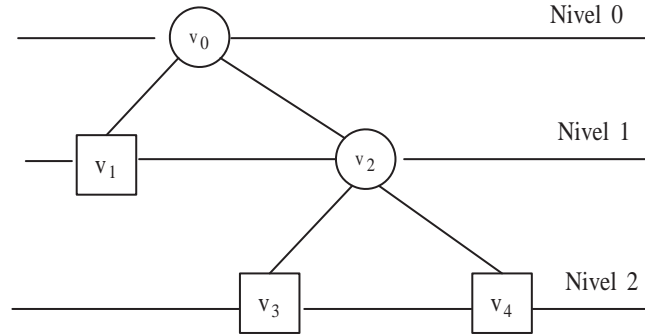


Figura 3.1: Reprezentarea grafică a arborelui $\mathcal{G} = (V, \Gamma)$.

Reprezentarea în plan a acestui arbore este dată în figura 3.1

Nodurile v pentru care $\Gamma(v) = \emptyset$ se numesc noduri *terminale* (finale); celelalte se numesc *interne*. Mulțimea nodurilor terminale constituie *frontiera* arborelui. În general vom nota un arbore cu litere mari, specificînd ca indici rădăcina și frontiera; de exemplu $\mathcal{A}_{v_0, v_1v_2v_3}$. Un arbore comportă mai multe *ramuri*; în exemplu avem următoarele ramuri : v_0v_1 , $v_0v_2v_3$, $v_0v_2v_4$.

Fie $G = (V_N, V_T, S, P)$ o gramatică de tipul 2.

Definiție 3.2 *Un arbore de derivare în gramatica G este un arbore cu următoarele două proprietăți .*

- (1) *Nodurile sunt etichetate cu elementele din V_G ;*
- (2) *Dacă un nod v are descendenți direcți v_1, \dots, v_n atunci $v \rightarrow v_1v_2 \dots v_n \in P$.*

Exemplu. $G = (\{A, B\}, \{a, b\}, A, P)$ unde

$$P = \{A \rightarrow aBA|Aa|a, B \rightarrow AbB|ba|abb\}.$$

Arborele $\mathcal{A}_{A, aabbaa}$ reprezentat în figura 3.2 (Varianta 1) este un arbore de derivare (poate fi desenat coborînd frontiera pe nivelul ultim , Varianta 2):

Teorema 3.1 *Fie G o gramatică de tipul 2. Atunci $X \xRightarrow{*} p$ dacă și numai dacă există un arbore $\mathcal{A}_{X, p}$.*

Demonstrație. $X \xRightarrow{*} p$ implică $\exists \mathcal{A}_{X, p}$.

Procedăm prin inducție asupra lungimii derivării l .

Dacă $l = 1$, $X \Rightarrow p = i_1 \dots i_n$ și $X \rightarrow i_1 \dots i_n \in P$. Arborele din figura 3.3 corespunde cerințelor teoremei.

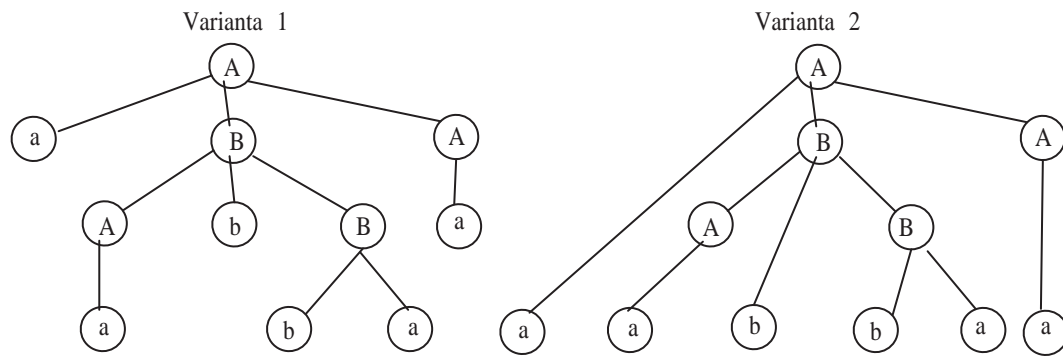


Figura 3.2: Variante de reprezentare a arborelui $\mathcal{A}_{A, aabbbaa}$.

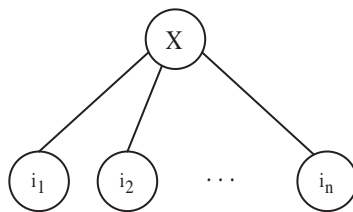


Figura 3.3: Arbore corespunzător unei derivări directe.

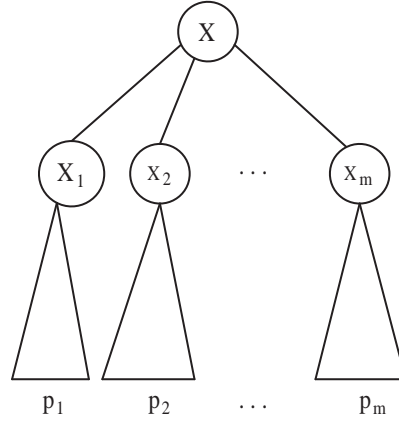


Figura 3.4: Construcția arborelui $\mathcal{A}_{X,p_1...p_m}$.

Presupunem că proprietatea este adevărată pentru derivări de lungime l și considerăm o derivare de lungime $l+1$, $X \xRightarrow{*} p$. Punem în evidență prima derivare directă

$$X \Rightarrow X_1 \dots X_n \xRightarrow{*} p$$

Conform lemei de localizare, $p = p_1 \dots p_n$ și $X_j \xRightarrow{*} p_j$, $j = \overline{1, n}$. Putem face următoarea construcție: conform ipotezei inductive, fiecărei derivări $X_j \xRightarrow{*} p_j$ îi corespunde câte un arbore \mathcal{A}_{X_j,p_j} ; unim apoi toate nodurile X_j în nodul X plasat la nivelul zero. Obținem astfel un arbore $\mathcal{A}_{X,p_1...p_m} = \mathcal{A}_{X,p}$ (vezi figura 3.4) care corespunde cerințelor teoremei.

Pentru implicația $\exists \mathcal{A}_{X,p} \Rightarrow X \xRightarrow{*} p$, se parcurge o cale inversă, făcând o inducție asupra numărului de nivele. De exemplu, dacă acest număr este 2, arborele de derivare trebuie să arate ca în 3.3 și deci avem $X \rightarrow i_1 i_2 \dots i_n = p \in P$ și $X \xRightarrow{*} p$, etc. \square

3.2 Ambiguitate în familia \mathcal{L}_2 .

Fie G o gramatică de tipul 2. O derivare $S = u_0 \Rightarrow u_1 \Rightarrow \dots \Rightarrow u_n$ în care la fiecare derivare directă se înlocuiește simbolul neterminal cel mai din stânga (dreapta) se numește *derivare extrem stângă (dreaptă)*. Să observăm că în particular într-o gramatică de tipul 3 orice derivare este extrem dreaptă (scrierea drept liniară).

Definiție 3.3 O gramatică G de tipul 2 în care există un cuvânt $p \in L(G)$ care se poate obține cu două derivări extrem stângi (drepte) distincte, se numește *ambiguă*. În caz contrar este *neambiguă*.

Exemplu. Gramatica $A \rightarrow aBA|Aa|a$, $B \rightarrow AbB|ba|abb$ este ambiguă. Într-adevăr, avem

$$A \Rightarrow aBA \Rightarrow aBa \Rightarrow aAbBa \Rightarrow aAbbaa \Rightarrow aabbaa;$$

$$A \Rightarrow Aa \Rightarrow aBAa \Rightarrow aBaa \Rightarrow aabbaa.$$

Definiție 3.4 Un limbaj este ambigu dacă toate gramaticile care îl generează sunt ambigu. În caz contrar (adică dacă există o gramatică neambiguă care să îl genereze) limbajul este neambigu.

Dacă G este ambigu și $p \in L(G)$ este un cuvânt care se poate obține cu două derivări extrem stângi distincte, atunci există arborii $\mathcal{A}_{S,p}$ și $\mathcal{A}'_{S,p}$, diferiți, dar care au aceeași rădăcină și frontieră.

Teorema 3.2 Dacă L_1 și L_2 sunt limbaie disjuncte neambigue, atunci $L_1 \cup L_2$ este neambigu.

Demonstrație. Fie $G_k = (V_{N_k}, V_{T_k}, S_k, P_k)$, $k = 1, 2$ două gramatici de tipul 2 neambigue și fie $G = (V_{N_1} \cup V_{N_2}, V_{T_1} \cup V_{T_2}, S, P_1 \cup P_2 \cup \{S \rightarrow S_1|S_2\})$ gramatica ce generează limbajul $L(G_1) \cup L(G_2)$.

Să presupunem că $L(G)$ este ambigu. Atunci există $p \in L(G)$ care se poate obține cu două derivări extreme stângi diferite. Să presupunem că $p \in L(G_1)$, $p \notin L(G_2)$. Atunci obținem două derivări distincte în gramatica G

$$(1) S \Rightarrow_G S_1 \xRightarrow{*}_G p, \text{ deci } S_1 \xRightarrow{*}_{G_1} p;$$

$$(2) S \Rightarrow_G S_1 \xRightarrow{*}_G p, \text{ deci } S_1 \xRightarrow{*}_{G_2} p,$$

deci și două derivări extrem stângi în gramatica G_1 . Aceasta ar însemna că G_1 este ambigu. Contradicție cu ipoteza! \square

Teorema 3.3 Limbajele de tipul 3 sunt neambigue.

Demonstrație. Fie L un limbaj de tipul 3 și G gramatica care îl generează; fie apoi AF automatul finit care recunoaște limbajul L și AFD automatul finit echivalent determinist. Construim gramatica G' astfel încât $L(G') = L(AFD)$. Reamintim că regulile lui G' se construiesc astfel $f(A, a) = B \Rightarrow A \rightarrow aB$, $f(A, a) \in \Sigma_f \Rightarrow A \rightarrow a$.

Să presupunem acum că L este ambigu; atunci orice gramatică care îl generează, inclusiv G' , este ambigu. Aceasta înseamnă că există un $p \in L(G')$ astfel încât

$$S \Rightarrow i_1 A_1 \Rightarrow i_1 i_2 A_2 \Rightarrow \dots \Rightarrow i_1 \dots i_{n-1} A_{n-1} \left\{ \begin{array}{l} \Rightarrow i_1 \dots i_n A'_n \Rightarrow \\ \Rightarrow i_1 \dots i_n A''_n \Rightarrow \end{array} \right\} \xRightarrow{*} p.$$

Deci există regulile $A_{n-1} \rightarrow i_n A'_n$ și $A_{n-1} \rightarrow i_n A''_n$, adică în automatul AFD avem

$$f(A_{n-1}, i_n) = A'_n, \quad f(A_{n-1}, i_n) = A''_n,$$

ceea ce contrazice faptul că AFD este determinist. \square

3.3 Forme normale pentru gramatici de tipul 2

Forma normală Chomsky.

Definiție 3.5 *O gramatică în forma normală Chomsky este o gramatică cu reguli de forma*

$$\begin{aligned} A &\rightarrow BC, \\ D &\rightarrow i, \end{aligned}$$

unde $A, B, C, D \in V_N$ și $i \in V_T$. Se acceptă și regula de completare $S \rightarrow \lambda$ cu condiția ca S să nu apară în dreapta vreunei reguli.

Lema 3.1 (lema substituției). *Fie G o gramatică de tipul 2 și $X \rightarrow uYv$ precum și $Y \rightarrow p_1 \dots p_n$ toate regulile din G care au Y în stânga. Atunci G este echivalentă cu o gramatică G' în care am făcut "substituțiile"; adică facem următoarea înlocuire*

$X \rightarrow uYv$ se înlocuiește cu $X \rightarrow up_1v | \dots | up_nv$
(Regulile $Y \rightarrow p_1 | \dots | p_n$ le vom păstra neschimbate).

Demonstrație. Fie $p \in L(G)$ și $S \xRightarrow{*} p$. Punem în evidență doi pași consecutivi oarecare:

$$G : S \xRightarrow{*} r \Rightarrow s \Rightarrow t \xRightarrow{*} p.$$

Dacă în $r \Rightarrow s$ se utilizează regula $X \rightarrow uYv$ atunci în mod necesar în pasul următor se utilizează una din regulile $Y \rightarrow p_1 | \dots | p_n$, să presupunem $Y \rightarrow p_j$ (evident, este posibil ca această regulă să nu se aplice în pasul imediat următor, dar ea poate fi "adusă" în această poziție). Prin urmare

$$(A) \quad G : r = r'Xr'' \Rightarrow r'uYvr'' \Rightarrow r'up_jvr'' = t.$$

Acești doi pași se pot obține și în G' (într-un singur pas):

$$(B) \quad G' : r = r'Xr'' \Rightarrow r'up_jr'' = t.$$

Deci $S \xRightarrow[G']{*} p, p \in L(G')$ și $L(G) \subseteq L(G')$.

Invers, dacă $p \in L(G')$ și $S \xRightarrow[G']{*} p$, atunci dacă la un pas se utilizează o regulă nouă introdusă (pasul (B)), transformarea respectivă se poate obține și în G cu doi pași (pașii (A)); deci $p \in L(G)$ și $L(G') \subseteq L(G)$. \square

Corolar 3.4 *Orice gramatică de tipul 2 este echivalentă cu o gramatică de același tip în care mulțimea de reguli nu conține redenumiri. (O redenumire este o regulă de forma $A \rightarrow B$, $A, B \in V_N$).*

Intr-adevăr, dacă $A \rightarrow B \in P$ este o redenumire și $B \rightarrow p_1 | \dots | p_n$ sunt toate regulile care au B în stânga, efectuăm substituțiile, deci înlocuim regula $A \rightarrow B$ cu $A \rightarrow p_1 | \dots | p_n$. În cazul în care printre acestea apare o nouă redenumire, repetăm procedeul. \square

Exemplu. Gramatica G_E care generează expresii aritmetice $E \rightarrow E+T | T, T \rightarrow T * F | F, F \rightarrow (E) | i$ se poate pune sub următoarea formă (fără redenumiri) :

$$\begin{aligned} E &\rightarrow E + T | T * F | (E) | i \\ T &\rightarrow T * F | (E) | i \\ F &\rightarrow (E) | i \end{aligned}$$

Teorema 3.5 (teorema lui Chomsky de existență a formei normale). *Orice gramatică independentă de context este echivalentă cu o gramatică în forma normală Chomsky.* *ema2*

Demonstrație. Putem porni de la o gramatică G care nu are redenumire și ale cărei reguli cu terminale au forma $A \rightarrow i, A \in V_N, i \in V_T$. De asemenea presupunem că G nu are reguli de ștergere.

Rezultă că regulile lui G au una din formele:

- (1) $A \rightarrow BC$,
- (2) $D \rightarrow i$,
- (3) $X \rightarrow X_1 \dots X_n, n > 2$.

Construim o gramatică $G' = (V'_N, V_T, S, P')$ unde $V_N \subseteq V'_N$ și P' conține toate regulile din P de forma (1) și (2). Fiecare regulă de forma (3) o înlocuim cu:

$$\begin{aligned} X &\rightarrow X_1 Z_1, \\ Z_1 &\rightarrow X_2 Z_2, \\ &\dots \\ Z_{n-2} &\rightarrow X_{n-1} X_n \end{aligned}$$

și includem neterminalele Z_1, \dots, Z_{n-2} (altele pentru fiecare regulă de forma (3) în V'_N).

Se poate relativ ușor arăta că $L(G) = L(G')$. De exemplu, dacă $u \Rightarrow v$ (direct) în G și de aplică o regulă de forma (1) sau (2), atunci evident derivarea respectivă se poate obține și în G' ; în cazul în care se aplică o regulă de forma (3), avem

$$G : u = u' X u'' \Rightarrow u' X_1 \dots X_n u'' = v.$$

Această derivare se poate obține și în G' în mai mulți pași și anume

$$G' : u = u' X u'' \Rightarrow u' X_1 Z_1 u'' \Rightarrow u' X_1 X_2 Z_2 u'' \Rightarrow \dots \Rightarrow u' X_1 \dots X_n u'' = v. \square$$

Observație. O gramatică ce are reguli de forma $A \rightarrow BC, A \rightarrow B, A \rightarrow a$ unde $A, B, C \in V_N$ și $a \in V_T$ spunem că este în *forma 2-canonică*. Este evident că orice gramatică de tip 2 este echivalentă cu o gramatică în formă 2-canonică.

Gramatici recursive

Definiție 3.6 Un simbol neterminal X al unei gramatici de tipul 2 este recursiv dacă există o regulă de forma $X \rightarrow uXv$, $u, v \in V_G^*$.

Dacă $u = \lambda$ ($v = \lambda$) simbolul X este stâng (drept) recursiv. O gramatică ce are cel puțin un simbol recursiv se numește recursivă. De exemplu, gramatica G_E care generează expresiile aritmetice are două simboluri stâng recursive, E și T .

Existența simbolurilor stâng recursive poate provoca dificultăți în aplicarea algoritmilor de analiză top-down. Într-adevăr, într-o astfel de gramatică, încercarea de a construi arborele de derivare corespunzător unui cuvânt p prin aplicarea întotdeauna a primei reguli pentru simbolul cel mai din stânga, poate să conducă la un ciclu infinit (de exemplu în G_E s-ar obține $E \Rightarrow E + T \Rightarrow E + T + T \Rightarrow \dots$).

Teorema 3.6 Orice limbaj de tipul 2 poate să fie generat de o gramatică fără recursie stângă.

Demonstrație. Fie $G = (V_N, V_T, S, P)$ o gramatică de tipul 2; presupunem că G are un singur simbol recursiv X și fie

$$(A) \quad X \rightarrow u_1|u_2|\dots|u_n|Xv_1|\dots|Xv_m$$

toate regulile care au X în stânga. Construim gramatica $G' = (V'_N, V_T, S, P')$, unde $V_N \subset V'_N$, $P \subset P'$ cu excepția regulilor (A); acestea se înlocuiesc cu

$$\begin{aligned} X &\rightarrow u_1|u_2|\dots|u_n|u_1Y|u_2Y|\dots|u_nY, \\ Y &\rightarrow v_1|\dots|v_m|v_1Y|\dots|v_mY \end{aligned}$$

G' este de tipul 2 și nu are simboluri stâng recursive; se vede însă că Y este un simbol drept recursiv.

Fie $p \in L(G)$, $S \xRightarrow[G]{*} p$. Dacă în această derivare nu intervine simbolul recursiv, atunci evident că $S \xRightarrow[G']{*} p$. Să presupunem că X intervine la un anumit pas: $S \Rightarrow u \Rightarrow p$, unde $u = u'Xu''$. Putem aplica, începând de la u spre dreapta, în primul rând regulile pentru X și să urmărim numai subarborele respectiv, deci

$$G: \quad X \xRightarrow[G]{*} Xv_{j_1} \xRightarrow[G]{*} Xv_{j_2}v_{j_1} \xRightarrow[G]{*} \dots \xRightarrow[G]{*} Xv_{j_s}\dots v_{j_1} \xRightarrow[G]{*} u_jv_{j_s}\dots v_{j_1}.$$

Aceeași formă propozițională o putem obține și în gramatica G' astfel

$$G': \quad X \xRightarrow[G']{*} u_jY \xRightarrow[G']{*} u_jv_{j_s}Y \xRightarrow[G']{*} \dots \xRightarrow[G']{*} u_jv_{j_s}\dots v_{j_1}.$$

Prin urmare avem $S \xRightarrow[G']{*} u \xRightarrow[G']{*} p$, adică $p \in L(G')$ și $L(G) \subseteq L(G')$. Analog,

$L(G') \subseteq L(G)$. \square

Forma normală Greibach.

Definiție 3.7 *O gramatică în forma normală Greibach este o gramatică cu reguli de forma*

$$A \rightarrow ip, \text{ unde } A \in V_N, i \in V_T, p \in V_N^*.$$

Se acceptă și regula de completare $S \rightarrow \lambda$ cu condiția ca S să nu apară în dreapta vreunei reguli.

Teorema 3.7 *(Teorema de existență a formei normale Greibach). Orice gramatică de tipul 2 este echivalentă cu o gramatică în forma normală Greibach.*

Demonstrație. Fie G o gramatică de tipul 2 în forma normală Chomsky și fie $V_N = \{S = X_1, X_2, \dots, X_n\}$. Vom construi o gramatică echivalentă care să satisfacă cerințele din forma normală Greibach în mai multe etape.

Etapa I. Vom modifica regulile de generare astfel încât toate regulile care nu sunt de forma $X \rightarrow i$ să satisfacă condiția $X_j \rightarrow X_k p$, $j < k$, $p \in V_N^*$. Acest lucru îl facem cu un algoritm pe care îl prezentăm într-un limbaj nestandard de publicare (tip PASCAL):

```

      j := 1;
e1:   begin
      Se elimină recursiile stângi; neterminalele
      noi le notăm cu  $Y_1, Y_2, \dots$ 
      end
      if  $j = n$  then STOP;
       $j := j + 1$ ;
       $l := 1$ ;
e2:   begin
      Fie  $X_j \rightarrow X_l p$ ,  $p \in V_N^*$  și  $X_l \rightarrow p_1 \dots p_m$ 
      toate regulile care au  $X_l$  în stânga; se efectuează toate substituțiile.
      end
       $l := l + 1$ ;
      if  $l < j - 1$  then goto e2
      goto e1

```

Să observăm că pentru $j = 1$ și după eliminarea recursiilor stângi condiția cerută este evident îndeplinită; în plus, dacă au fost recursii, vom avea reguli cu partea stângă neterminale noi Y_1, Y_2, \dots . Mai departe, luăm toate regulile care au în stânga X_2 ($j := j + 1 = 2$) și efectuăm substituțiile; acestea se vor transforma în $X_2 \rightarrow X_k p$ cu $k \geq 2$ și după o nouă eliminare a recursiilor stângi vom avea $k > 2$ plus reguli cu partea stângă neterminale noi. În felul acesta toate regulile care au în stânga X_1 și X_2 satisfacă condiția cerută; în continuare $j := j + 1 = 3$, etc.

Etapa II. Avem acum trei categorii de reguli:

- (1) $X_j \rightarrow i$;
- (2) $X_j \rightarrow X_k p$, $j < k$, $p \in V_N^*$;
- (3) $Y \rightarrow iq$, $q \in V_N^*$, $i = 1, \dots, m$.

Aranjăm toate neterminalele într-un șir unic, la început Y_1, \dots, Y_m apoi X_1, \dots, X_n și le redenumim, de exemplu cu X_1, \dots, X_{m+n} :

$$\begin{array}{cccccccc} Y_1, & Y_2, & \dots, & Y_m, & X_1, & X_2, & \dots, & X_n \\ X_1, & X_2, & \dots, & X_m, & X_{m+1}, & X_{m+2}, & \dots, & X_{m+n} \end{array}$$

Vom nota $n + m = N$. În felul acesta regulile gramaticii vor avea numai formele (1) și (2).

Etapa III. Toate regulile care au X_N în stânga vor avea forma (1). Fie $X_{n-1} \rightarrow X_N p_1 | \dots | X_N p_n$ toate regulile care au X_{N-1} în stânga și care nu sunt de forma (1). Efecuăm substituțiile lui X_N ; în acest fel regulile care au X_N și X_{N-1} în stânga satisfac cerințele din forma normală Greibach. În continuare, considerăm toate regulile care au X_{N-2} în stânga și efectuăm substituțiile, etc. \square

Forma normală operator

Una din formele importante pentru gramatici independente de context, utilizată în analiza sintactică prin metoda precedenței, este forma operator a acestor gramatici.

Definiție 3.8 *O gramatică independentă de context $G = (V_N, V_T, S, P)$ se spune că este în forma normală operator dacă oricare ar fi producția $A \rightarrow \beta \in P$, în β nu apar două neterminale (variabile) consecutive, adică*

$$P \subseteq V_N \times [(V_N \cup V_T)^* \setminus (V_N \cup V_T)^* V^2 (V_N \cup V_T)^*].$$

Teorema 3.8 *Orice gramatică independentă de context este echivalentă cu o gramatică în forma normală operator.*

Demonstrație. Fie $G = (V_N, V_T, S, P)$ o gramatică de tipul 2 și $L(G)$ limbajul generat. Fără a restrânge generalitatea presupunem că $\lambda \notin L(G)$ și G este în forma 2-canonică (regulile sunt de forma $A \rightarrow BC$, $A \rightarrow B$, $A \rightarrow a$ vezi teorema ??). Definim o gramatică echivalentă $G' = (V'_N, V_T, S, P')$ astfel: $V'_N = \{S\} \cup (V_N \times V_T)$, iar $P' = P_1 \cup P_2 \cup P_3 \cup P_4$ unde

- i) $P_1 = \{S \rightarrow (S, a)a \mid a \in V_T\}$;
- ii) $P_2 = \{(A, a) \rightarrow \lambda \mid A \in V_N, a \in V_T, A \rightarrow a \in P\}$;
- iii) $P_3 = \{(A, a) \rightarrow (B, a) \mid A, B \in V_N, a \in V_T, A \rightarrow B \in P\}$;
- iv) $P_4 = \{(A, a) \rightarrow (B, b)b(C, a) \mid A, B, C \in V_N, a, b \in V_T, A \rightarrow BC \in P\}$.

Să observăm că G' este în forma normală operator. Pentru a demonstra că $L(G) = L(G')$ vom defini mai întâi o funcție $\phi : P_2 \cup P_3 \cup P_4 \rightarrow P$ astfel:

- $\phi((A, a) \rightarrow \lambda) = A \rightarrow a$ pentru $(A, a) \rightarrow \lambda \in P_2$;
- $\phi((A, a) \rightarrow (B, a)) = A \rightarrow B$ pentru $(A, a) \rightarrow (B, a) \in P_3$;
- $\phi((A, a) \rightarrow (B, b)b(C, a)) = A \rightarrow BC$ pentru $(A, a) \rightarrow (B, b)b(C, a) \in P_4$.

Funcția ϕ se extinde în mod natural la $\phi' : (P_2 \cup P_3 \cup P_4)^* \longrightarrow P^*$. Vom arăta că în gramatica G , $\forall w \in V_T^*$, $\forall a \in V_T$ are loc derivarea extrem dreaptă $A \xRightarrow[G]{*} wa$ folosind producțiile $\pi_1, \pi_2, \dots, \pi_n$ dacă și numai dacă există în P' producțiile $\pi'_1, \pi'_2, \dots, \pi'_n$ astfel ca $\phi(\pi'_i) = \pi_i$, $1 \leq i \leq n$ și în G' are loc derivarea extrem dreaptă $(A, a) \xRightarrow[G']{*} w$ folosind producțiile $\pi'_1, \pi'_2, \dots, \pi'_n$.

Să demonstrăm afirmația prin inducție după n , lungimea derivării.

Dacă $n = 1$ atunci $w = \lambda$, $A \rightarrow a \in P$ și în P' există producția $(A, a) \rightarrow \lambda$, deci $(A, a) \Rightarrow \lambda$ și $\phi((A, a) \rightarrow \lambda) = A \rightarrow a$.

Invers, dacă $(A, a) \Rightarrow w$ în G' atunci $w = \lambda$ (după forma producțiilor din G') și are loc proprietatea enunțată.

Să presupunem afirmația adevărată pentru derivări de lungime cel mult $n - 1$ și să o demonstrăm pentru derivări de lungime $n > 1$. Fie așadar $A \xRightarrow[G]{*} wa$

o derivare de lungime n în gramatica G și punem în evidență prima derivare directă. Distingem două cazuri:

I. Prima producție utilizată în derivare este $A \rightarrow B$. Atunci,

$$A \Rightarrow_G B \xRightarrow[G]{*} wa$$

și conform ipotezei inductive avem în G' o derivare $(B, a) \xRightarrow[G']{*} w$ (de lungime $n - 1$)

cu producții satisfăcând condițiile arătate. Dar cum $A \rightarrow B \in P$, în P' avem producția $(A, a) \rightarrow (B, a)$ deci $(A, a) \xRightarrow[G']{*} w$ în gramatica G' .

II. Prima producție este de forma $A \rightarrow BC$. Atunci

$$A \Rightarrow_G BC \xRightarrow[G]{*} wa.$$

În acest caz $wa = ubva$ (conform lemei de localizare), astfel că $B \xRightarrow[G]{*} ub$ și

$C \xRightarrow[G]{*} va$. După ipoteza inductivă, vom avea în G' derivările:

$$(B, b) \xRightarrow[G']{*} u, \quad (C, a) \xRightarrow[G']{*} v.$$

Cum $A \rightarrow BC \in P$ vom avea în P' producția $(A, a) \rightarrow (B, b)b(C, a)$ și în G' putem scrie derivarea extrem dreaptă

$$(A, a) \Rightarrow (B, b)b(C, a) \xRightarrow[G']{*} (B, b)bv \xRightarrow[G']{*} ubv = w$$

și producțiile care s-au aplicat îndeplinesc condițiile din enunț.

În mod analog se demonstrează reciproca.

Din această afirmație, luând în particular $A = S$, obținem:

$$S \xRightarrow[G]{*} wa \Leftrightarrow (S, a) \xRightarrow[G']{*} w, \quad \forall w \in V_T^*, \forall a \in V_T.$$

Cum în G' există și producția $S \rightarrow (S, a)a$, am găsit: $wa \in L(G)$ dacă și numai $wa \in L(G')$, deci cele două gramatici sunt echivalente.

Pentru a încheia demonstrația trebuie să considerăm și cazul $\lambda \in L(G)$. Aplicăm construcția de mai sus unei gramatici ce generează $L(G) \setminus \{\lambda\}$ și obținem $G' = (V'_N, V_T, S, P')$ gramatica operator corespunzătoare. Considerăm acum gramatica $G_1 = (V_{N_1}, V_T, S_1, P_1)$ unde $V_1 = V' \cup \{S_1\}$, $P_1 = P' \cup \{S_1 \rightarrow \lambda, S_1 \rightarrow S\}$ care este în forma normală operator și $L(G_1) = L(G)$. \square

3.4 Lema Bar–Hillel

Ca și în cazul limbajelor regulate, lema Bar–Hillel pune în evidență următoarea proprietate a unui limbaj independent de context: orice cuvânt al unui astfel de limbaj, suficient de lung, conține un subcuvânt (nevid) pe care multiplicându-l de un număr arbitrar de ori, obținem noi cuvinte care aparțin de asemenea limbajului. Mai poartă denumirea de "lema de pompă" sau "lema $uvwxy$ ".

Ne vom referi în cele ce urmează la gramatici de tipul 2 în formă normală Chomsky. Într-o gramatică de această formă arborii de derivare sunt întotdeauna arbori binari.

Lema 3.2 *Fie G o gramatică în forma normală Chomsky și $X \xRightarrow{*} p, p \in V_G^*$. Dacă cea mai lungă ramură din arborele $\mathcal{A}_{X,p}$ conține m noduri, atunci $|p| \leq 2^{m-1}$.*

Demonstrație. Procedăm prin inducție asupra lui m .

Pentru $m = 2$ arborele are două nivele și în mod evident $|p| \leq 2 = 2^{m-1}$.

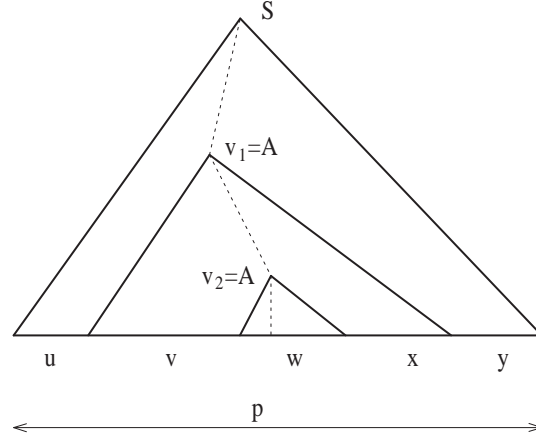
Presupunem că proprietatea este adevărată pentru un m oarecare și considerăm un arbore care are pe cea mai lungă ramură $m + 1$ noduri. În derivarea $X \xRightarrow{*} p$ punem în evidență prima derivare directă

$$X \Rightarrow Y Z \xRightarrow{*} p.$$

Conform lemei de localizare, $p = p_1 p_2$ și $Y \xRightarrow{*} p_1, Z \xRightarrow{*} p_2$. Arborii de derivare \mathcal{A}_{Y,p_1} și \mathcal{A}_{Z,p_2} conțin, fiecare pe cea mai lungă ramură cel mult m noduri; conform ipotezei de inducție, avem $|p_1| \leq 2^{m-1}, |p_2| \leq 2^{m-1}$. Prin urmare

$$|p| = |p_1 p_2| = |p_1| + |p_2| \leq 2^{m-1} + 2^{m-1} = 2^m. \square$$

Observație. Dacă $X \xRightarrow{*} p$ și $|p| > 2^{m-1}$ atunci există în arborele $\mathcal{A}_{X,p}$ cel puțin o ramură cu $m + 1$ noduri.

Figura 3.5: Descompunerea cuvântului p .

Lema 3.3 Fie G o gramatică de tipul 2 și fie $X \Rightarrow X_1 \dots X_m \xRightarrow{*} p$. Atunci, conform lemei de localizare, $p = p_1 \dots p_m$ și $X_j \xRightarrow{*} p_j, j = 1, \dots, m$. Fie $\mathcal{A}_{Y,q} \subseteq \mathcal{A}_{X,p}$. atunci q este subcuvânt într-unul din cuvintele p_j .

Demonstrație. Neterminalul Y aparține unuia din subarborii \mathcal{A}_{X_j,p_j} , fie acesta \mathcal{A}_{X_k,p_k} ; atunci $\mathcal{A}_{Y,q} \subseteq \mathcal{A}_{X_k,p_k}$ și $q \in \text{Sub}(p_k)$. \square

Lema 3.4 (Lema Bar–Hillel) Fie E un limbaj independent de context. Atunci există un număr natural k astfel încât, dacă $p \in E$ și $|p| > k$ atunci p se poate descompune în forma $p = uvwxy$ cu următoarele proprietăți:

1. $vx \neq \lambda$;
2. $|vwx| \geq k$;
3. $uv^jwx^jy \in E, \forall j \in \mathbb{N}$,

Demonstrație. Fie $n = \text{card}(V_N)$. Luăm $k = 2^n$. Cum $|p| > k = 2^n$, conform observației de la prima lea, există în arborele $\mathcal{A}_{S,p}$ cel puțin o ramură cu $n + 2$ noduri; pe această ramură ultimul nod este terminal, deci există $n + 1$ noduri etichetate cu același simbol A . Descompunem cuvântul p ca în figura 3.5, $p = uvwxy$.

(1) Considerăm subarborile $\mathcal{A}_{A,vwx}$ căruia îi corespunde derivarea $A \xRightarrow{*} vwx$. Punem în evidență primul pas

$$A \Rightarrow BC \xRightarrow{*} vwx,$$

și vwx se descompune în $vwx = p_B p_C, B \xRightarrow{*} p_B, C \xRightarrow{*} p_C$ și $p_B, p_C \neq \lambda$. Cum $\mathcal{A}_{A,w} \subseteq \mathcal{A}_{A,vwx}$, rezultă că w este subcuvânt în p_B sau în p_C . Să presupunem că $w \in \text{Sub}(p_C)$; atunci $p_B \in \text{Sub}(v)$ și cum $p_B \neq \lambda$ rezultă $v \neq \lambda$.

(2) Putem alege nodurile v_1 și v_2 astfel încât pe ramura punctată începând de la v_1 în jos până la frontieră să avem exact $n + 1$ noduri. Rezultă conform lemei precedente, $|vwx| \geq 2^n = k$.

(3) putem scrie derivările $A \xRightarrow{*} w, A \xRightarrow{*} vAx$. Deci

$$S \xRightarrow{*} uAy \xRightarrow{*} uvAxy \xRightarrow{*} uv^2Ax^2y \xRightarrow{*} \dots \xRightarrow{*} uv^jwx^jy. \square$$

Problema închiderii familiei \mathcal{L}_2 la intersecție

Teorema 3.9 *Familia \mathcal{L}_2 nu este închisă la intersecție.*

Demonstrație. Considerăm limbajul $L_3 = \{a^n b^n c^n | n \geq 1\}$. Să rătăm că $L_3 \notin \mathcal{L}_2$. Într-adevăr, să presupunem că $L_3 \in \mathcal{L}_2$ și fie k constanta din lema Bar-Hillel. Luăm $n > k/3$ și fie $p = a^n b^n c^n$; avem $|p| > k$, deci conform acestei leme p se descompune în forma $p = uvwxy$ cu $vx \neq \lambda$ și $uv^jwx^jy \in L_3, j \in \mathbb{N}$.

Vom analiza posibilitățile de constituire a subcuvintelor v și x . Să presupunem că în v (sau x) intră două din simbolurile a, b, c ; de exemplu $v = aabbb$. atunci considerăm cuvântul $p_2 = uv^2wx^2y = uaabbaabbwx^2y$ care nu are structura cuvintelor lui L_3 ("a" nu poate să urmeze după "b") dar conform lemei Bar-Hillel aparține lui L_3 . Deci, în v (sau x) nu pot intra două (sau trei) din simbolurile a, b, c . Să presupunem că intră un singur simbol; de exemplu $v \in \{a\}^*$ și $x \in \{b\}^*$. Atunci multiplicând în p subcuvintele v și x , puterile simbolurilor "a" și "b" se măresc iar "c" rămâne pe loc, contrazicându-se din nou structura cuvintelor din L_3 . În concluzie L_3 nu este independent de context.

Fie acum limbajele

$$L'_3 = \{a^m b^n c^n | m, n \geq 1\}$$

$$L''_3 = \{a^n b^n c^m | m, n \geq 1\}.$$

Se poate vedea ușor că aceste limbaje sunt de tipul 2; gramaticile care le generează sunt $S \rightarrow aS|aA, A \rightarrow bAc|bc$ și respectiv $S \rightarrow Sc|Ac, A \rightarrow aAb|ab$.

Avem $L'_3 \cap L''_3 = L_3$, deci intersecția a două limbaje de tipul 2 este un limbaj care nu aparține lui \mathcal{L}_2 . \square

Corolar 3.10 *Familia \mathcal{L}_2 nu este închisă la complementariere.*

Într-adevăr, să presupunem că \mathcal{L}_2 este închisă la complementariere și fie $E_1, E_2 \in \mathcal{L}_2$. Cum familia \mathcal{L}_2 este închisă la reuniune, ar rezulta $C(E_1) \cup C(E_2) \in \mathcal{L}_2$. Dar (de Morgan) $C(E_1) \cup C(E_2) = C(E_1 \cap E_2) \in \mathcal{L}_2$.

Complementul limbajului $C(E_1 \cap E_2)$ este $E_1 \cap E_2$ și conform presupunerii ar trebui ca $E_1 \cap E_2 \in \mathcal{L}_2$, oricare ar fi E_1, E_2 , ceea ce nu este adevărat. \square

Generalizări ale lemei Bar-Hillel Lema lui Bar-Hillel reprezintă o condiție necesară ca un limbaj să fie independent de context. cu ajutorul ei se poate demonstra relativ ușor că anumite limbaje nu sunt independente de context. Ideea este aceea că prin multiplicarea subcuvintelor v și x , de un număr suficient de ori, se contrazice structura limbajului.

Totuși, nu orice limbaj care nu este de tipul 2 poate fi respins de lema lui Bar-Hillel. De exemplu, această leamnă nu poate respinge limbajul

$$L = \{a^n b^{2^m} | n, m \geq 1\} \cup \{b\}^*,$$

care nu este de tipul 2 (se demonstrează pe altă cale). Într-adevăr, pentru orice $p = a^n b^{2^m}$ luăm

$$u = \lambda, v = a, w = \lambda, x = \lambda, y = a^{n-1} b^{2^m}.$$

Putem itera subcuvintele v și x fără să contrazicem structura cuvintelor limbajului.

O generalizare a lemei Bar-Hillel este

Lema 3.5 (*Lema lui Ogden*) *pentru orice limbaj independent de context L există o constantă k astfel încât orice $p \in L$ pentru care cel puțin k simboluri sunt "marcate", se poate descompune în forma $p = uvwxy$ astfel încât*

1. *sau u, v, w sau w, x, y conțin fiecare câte un simbol marcat;*
2. *vw conține cel mult k simboluri marcate;*
3. *$uv^jwx^jy \in L, \forall j \in \mathbb{N}$*

Cu ajutorul acestei leme se poate arăta că limbajul de mai sus nu este de tipul 2, considerând marcate simbolurile b .

3.5 Automate push-down (APD)

Automatele push-down sunt mecanisme pentru recunoașterea limbajelor independente de context.

Un APD se compune din (vezi figura 3.6):

1. O bandă de intrare care conține simboluri ale unui *alfabet de intrare*; aceste simboluri constituie pe o bandă un anumit cuvânt peste alfabetul de intrare. Banda se mișcă numai spre stânga;
2. O *memorie push-down* (memorie inversă, stivă, pilă, etc) care conține simboluri ale unui alfabet propriu, numit *alfabetul memoriei push-down*. Această memorie funcționează ca o stivă - ultimul introdus, primul extras (Last In, First Out);
3. Un *dispozitiv de comandă* care se află permanent într-o anumită stare internă aparținând unei mulțimi finite de stări. Dispozitivul de comandă posedă un *dispozitiv de citire* de pe banda de intrare și un *dispozitiv de scriere-citire* în memoria push-down.

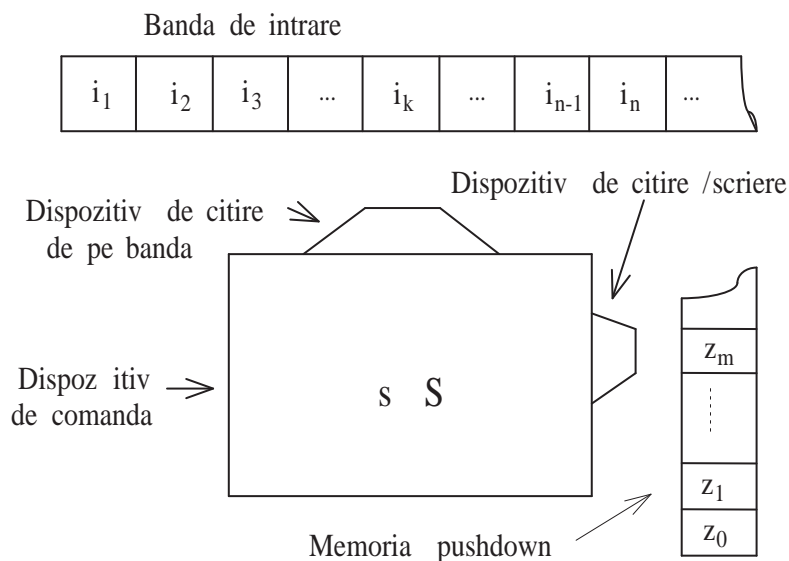


Figura 3.6: Reprezentare schematică a unui automat push-down.

Ca și un automat finit, un automat push-down funcționează în pași discreți; un pas de funcționare compoartă:

1. Dispozitivul de comandă citește simbolul de pe banda de intrare din dreptul dispozitivului de citire și mută banda cu o poziție spre stânga.
2. În funcție de starea internă, de simbolul citit și de simbolul din vârful memoriei push-down dispozitivul de comandă efectuează operațiile:
 - (a) Trece într-o nouă stare;
 - (b) Scrie în memoria push-down un anumit cuvânt peste alfabetul memoriei push-down; în particular, acesta poate să fie cuvântul vid, ceea ce are ca efect ștergerea simbolului din vârful memoriei push-down.

Funcționarea unui APD se termină în general după ce s-a citit ultimul simbol al cuvântului scris pe banda de intrare dar este posibil ca el să efectueze un anumit număr de pași, citind de fiecare dată de pe bandă cuvântul vid λ . De asemenea, este posibil ca în timpul funcționării, deci înainte de a ajunge la ultimul simbol, automatul să se blocheze. De exemplu, automatul ajunge într-o configurație (stare, simbol pe bandă, simbol în vârful memoriei push-down) inadmisibilă sau se golește memoria push-down dar nu s-a epuizat cuvântul de pe bandă, etc.

Matematic, un APD se definește astfel:

Definiție 3.9 *Un automat push-down este un sistem*

$$APD = (\Sigma, I, Z, f, s_0, z_0)$$

unde:

Σ este mulțimea de stări (finită și nevidă);
 I este alfabetul de intrare;
 Z este alfabetul memoriei push-down;
 $f : \Sigma \times (I \cup \{\lambda\}) \times Z \longrightarrow \mathcal{P}(\Sigma \times Z^*)$ este funcția de evoluție;
 $s_0 \in \Sigma$ este starea inițială;
 $z_0 \in Z$ este simbolul inițial al memoriei push-down.

Un APD are în general o funcționare nedeterministă, $\text{card } f(s, i, z) \geq 1$; mulțimea automatelor push-down deterministe formează o clasă specială.

În termenii funcției de evoluție, un pas de evoluție comportă citirea simbolului i de pe bandă, citirea simbolului z din vârful memoriei push-down, apoi, în funcție de starea internă s și de aceste două simboluri, automatul trece într-o nouă stare s' și scrie în memoria push-down un cuvânt $q \in Z^*$ astfel încât $(s', q) \in f(s, i, z)$. În cazul în care $f(s, i, z) = \emptyset$ evoluția este oprită; este situația în care automatul se blochează.

O stare a automatului (sau configurație) este un sistem $\delta = (s, p, q)$ unde $s \in S$ este starea internă, $p \in I^*$ este subcuvântul de pe banda de intrare rămas de citit (inclusiv simbolul din dreptul dispozitivului de citire), iar $q \in Z^*$ este subcuvântul din memoria push-down.

Vom spune că un APD trece direct din starea $\delta_1 = (s_1, p_1, q_1)$ în starea $\delta_2 = (s_2, p_2, q_2)$ și vom scrie $\delta_1 \mapsto \delta_2$ dacă se execută un pas de evoluție; dacă $p_1 = ip'_1, q_1 = zq'_1$ putem avea $(s_2, q) \in f(s_1, i, z)$ și atunci $p_2 = p'_1, q_2 = qq'_1$ sau $(s_2, q) \in f(s_1, \lambda, z)$ și atunci $p_2 = p_1, q_2 = qq'_1$.

Vom spune că automatul evoluează (fără specificația direct) din starea δ' în stare δ'' și vom scrie $\delta' \xrightarrow{*} \delta''$ dacă:

- (1) $\delta' = \delta''$;
- (2) $\exists \delta_1, \dots, \delta_n$ astfel încât $\delta' = \delta_1 \mapsto \delta_2 \mapsto \dots \mapsto \delta_{n-1} \mapsto \delta_n = \delta''$.

Limbaajul recunoscut de un APD se poate defini în două moduri:

- (1) Limbaajul recunoscut de un APD cu golirea memoriei push-down, este, prin definiție

$$L(APD) = \{p | p \in I^* (s_0, p, z_0) \xrightarrow{*} (s, \lambda, \lambda)\}.$$

Aceasta înseamnă că, pornind din starea internă s_0 și având în vârful memoriei push-down simbolul z_0 , cu ajutorul cuvântului p de pe banda de intrare, automatul poate să evolueze astfel încât să golească memoria push-down după citirea cuvântului. Menționăm că golirea memoriei push-down nu trebuie neapărat să coincidă cu citirea ultimului simbol al lui p ; este posibil ca automatul să mai efectueze câțiva pași citind de pe bandă simbolul λ .

- (2) Limbaajul recunoscut de un APD cu stări finale; în definiția automatului se adaugă o submulțime Σ_f a lui Σ numită mulțimea de stări finale. Prin definiție, limbaajul recunoscut de un APD cu stări finale este:

$$L(APD) = \{p | p \in I^* (s_0, p, z_0) \xrightarrow{*} (s, \lambda, q), s \in \Sigma_f, q \in Z^*\}.$$

Prin urmare, este necesar ca după citirea lui p , eventual după încă câțiva pași, APD să ajungă într-o stare finală. Vom vedea că cele două definiții sunt echivalente.

Limbaje recunoscute de automate push-down cu golirea memoriei. Vom arăta că familia limbajelor recunoscute de APD cu stări finale coincide cu familia limbajelor independente de context. În felul acesta, APD constituie mecanisme analitice de definire a limbajelor de tipul 2.

Teorema 3.11 *Un limbaj este independent de context dacă și numai dacă este recunoscut de un automat push-down cu golirea memoriei push-down.*

Demonstrație. Partea I $E \in \mathcal{L}_2 \Rightarrow E = L(APD)$.

Fie $G = (V_N, V_T, S, P)$ o gramatică de tipul 2 în forma normală Greibach care generează limbajul E . Construim un automat pushdown astfel:

$APD = (\{s\}, V_T, V_N, f, s, S)$, funcția de evoluție fiind definită de:

$$A \rightarrow ip \in P \Rightarrow (s, p) \in f(s, i, A),$$

altfel \emptyset .

Fie $p \in L(G)$, $p = i_1 \dots i_n$, $S \xrightarrow[G]{*} p$. Această derivare trebuie să aibă forma (extrem stângă):

$$(A) \quad S \Rightarrow i_1 X_1 u_1 \Rightarrow i_1 i_2 X_2 u_2 u_1 \Rightarrow i_1 i_2 i_3 X_3 u_3 u_2 u_1 \Rightarrow \dots \Rightarrow i_1 \dots i_n,$$

unde $u_1, u_2, u_3, \dots \in V_N^* = Z^*$.

Observație. Aparent, partea $u_s u_{s-1} \dots u_1$ se mărește cu fiecare derivare directă. În realitate, unele din cuvintele u_j sunt vide, și anume atunci când se aplică o regulă de forma $X \rightarrow i$; în particular, în ultimele derivări directe se aplică numai reguli de această formă.

Avem

$$\begin{aligned} S \rightarrow i_1 X_1 u_1 &\Rightarrow (s, X_1 u_1) \in f(s, i_1, S), \\ X_1 \rightarrow i_2 X_2 u_2 &\Rightarrow (s, X_2 u_2) \in f(s, i_2, X_1), \\ X_2 \rightarrow i_3 X_3 u_3 &\Rightarrow (s, X_3 u_3) \in f(s, i_3, X_2), \\ &\dots \end{aligned}$$

Prin urmare automatul poate să aibă următoarea evoluție:

$$\begin{aligned} (s, i_1 i_2 i_3 i_4 \dots i_n, S) &\mapsto (s, i_2 i_3 i_4 \dots i_n, X_1 u_1) \mapsto \\ &\mapsto (s, i_3 i_4 \dots i_n, X_2 u_2 u_1) \mapsto (s, i_4 \dots i_n, X_3 u_3 u_2 u_1) \mapsto \dots \end{aligned}$$

Dacă comparăm această evoluție cu derivarea (A) putem observa că pe banda de intrare avem la fiecare pas partea complementară a cuvântului (față de derivare) iar în memoria push-down se reproduce partea de neterminale din formele propoziționale

ale derivării. Cum în derivare se ajunge la $i_1 \dots i_n$, în evoluție se va ajunge la (s, λ, λ) .

Deci $p \in L(APD)$ și $L(G) \subseteq L(APD)$.

Fie acum $p \in L(APD)$; va trebui să arătăm că $p \in L(G)$, deci că $S \xrightarrow[G]{*} p$.

Vom arăta o implicație ceva mai generală, și anume, pentru orice $u \in V_N^*$, avem

$$(s, p, u) \xrightarrow{*} (s, \lambda, \lambda) \Rightarrow u \xrightarrow[G]{*} p.$$

În particular, dacă $u = S$ obținem implicația dorită.

Procedăm prin inducție asupra lungimii lui p .

Dacă $|p| = 1$, atunci $p = i, u = X$ iar evoluția va avea un singur pas $(s, i, X) \xrightarrow{*} (s, \lambda, \lambda)$, deci $(s, \lambda) \in f(s, i, X)$ și $X \rightarrow i \in P$. Putem scrie $u = X \xRightarrow[G]{*} i = p$.

Presupunem că implicația este adevărată pentru un cuvânt $|p| = l$ și considerăm un p astfel încât $|p| = l + 1$. Fie i și X primele simboluri din p și u , deci $p = ip'$ și $u = Xu'$. În evoluția $(s, p, u) \xrightarrow{*} (s, \lambda, \lambda)$ punem în evidență prima evoluție directă

$$(s, p, u) = (s, ip', Xu') \xrightarrow{*} (s, p', vu') \xrightarrow{*} (s, \lambda, \lambda).$$

Din definiția evoluției directe rezultă că $(s, v) \in f(s, i, X)$ deci $X \rightarrow iv \in P$. Pe de altă parte din ipoteza inductivă rezultă că $vu' \xRightarrow[G]{*} p'$. Avem

$$u = Xu' \xRightarrow[G]{*} ivu' \xRightarrow[G]{*} ip' = p,$$

ceea ce demonstrează implicația.

Prin urmare $p \in L(G)$ și $L(APD) \subseteq L(G)$, de unde $L(G) = L(APD)$. \square

Partea II. $E = L(APD)$ $E \in \mathcal{L}_2$

Fie $APD = (\Sigma, I, Z, f, s_0, z_0)$. Construim G de forma $G = (V_N, V_T, S, P)$ unde $V_N = \{s_0\} \cup \{(s, z, s') | s, s' \in \Sigma, z \in Z\}$, $V_T = I$, $S = s_0$, iar regulile de generare le definim astfel:

- (1) $s_0 \rightarrow (s_0, z_0, s)$, $\forall s \in \Sigma$;
- (2) Dacă $(s_1, z_1 \dots z_m) \in f(s, i, z)$ vom introduce în P reguli de forma

$$(s, z, s') \rightarrow i(s_1, z_1, s_2)(s_2, z_2, s_3) \dots (s_m, z_m, s'),$$

unde $s', s_2, \dots, s_m \in \Sigma$;

- (3) Dacă $(s', \lambda) \in f(s, i, z)$ vom introduce în P reguli de forma

$$(s, z, s') \rightarrow i,$$

unde $s' \in \Sigma$.

Să observăm că gramatica astfel construită este independentă de context, și anume în forma normală Greibach.

Fie $p \in L(APD)$, deci $(s_0, p, z_0) \xrightarrow{*} (s', \lambda, \lambda)$; trebuie să arătăm că $s_0 \xrightarrow[G]{*} p$.

Vom arăta implicația ceva mai generală

$$(s, p, z) \xrightarrow{*} (s', \lambda, \lambda) \Rightarrow (s, z, s') \xrightarrow[G]{*} p.$$

În particular pentru $s = s_0, z = z_0$ rezultă $(s_0, z_0, s') \xrightarrow[G]{*} p$ și putem scrie $s_0 \Rightarrow (s_0, z_0, s') \xrightarrow[G]{*} p$,

adică $p \in L(G)$.

Procedăm prin inducție asupra lungimii evoluției l .

Dacă $l = 1$ atunci $(s, p, z) \xrightarrow{*} (s', \lambda, \lambda)$, deci $p = i$ și $(s', \lambda) \in f(s, i, z)$ și $(s, z, s') \rightarrow i$ este o regulă, adică putem scrie $(s, z, s') \Rightarrow i = p$.

Presupunem că implicația este adevărată pentru evoluții de lungime oarecare l și considerăm o evoluție de lungime $l + 1$; punem în evidență prima evoluție directă

$$(s, p, z) = (s, i_1 p', z) \xrightarrow{*} (s_1, p', z_1 \dots z_m) \xrightarrow{*} (s', \lambda, \lambda).$$

Descompunem cuvântul p' în forma $p' = p_1 \dots p_m$ astfel încât

$$\begin{aligned} (s_1, p_1, z_1) &\xrightarrow{*} (s_2, \lambda, \lambda), \\ (s_2, p_2, z_2) &\xrightarrow{*} (s_3, \lambda, \lambda), \\ &\dots \\ (s_m, p_m, z_m) &\xrightarrow{*} (s', \lambda, \lambda). \end{aligned}$$

Observație. Putem pune în evidență felul în care se definește cuvântul p_1 urmărind evoluția lui APD;

$$\begin{array}{ccc} (s_1, i_1 i_2 \dots i_n, z_1 z_2 \dots z_m) & \xrightarrow{\quad} & (s'_1, i_2 i_3 \dots i_n, q z_2 \dots z_m) \\ (a) & & (b) \\ \dots & \xrightarrow{\quad} & (s_2, i_{j_1} \dots i_n, z_2 \dots z_m) \\ & & (c) \end{array}$$

La primul pas (situația a) automatul este în starea s_1 , pe bandă este i_1 iar în memoria push-down este z_1 . După efectuarea unui pas, automatul trece în starea s'_1 , mută banda cu o poziție spre stânga, extrage pe z_1 și scrie în memoria push-down un cuvânt q (situația b). Se poate observa că z_2 a "coborât"; cum știm că memoria push-down se golește ($p \in L(APD)$), trebuie ca la un moment dat z_2 să ajungă în vârful stivei (situația c). În acest moment partea din p citită va fi p_1 iar starea în care a ajuns automatul o notăm cu s_2 . Este clar că dacă pe bandă am avea scris numai p_1 am avea evoluția $(s_1, p_1, z_1) \xrightarrow{*} (s_2, \lambda, \lambda)$.

Analog p_2, \dots, p_m .

Din definiția derivării directe $(s, i_1 p', z) \mapsto (s_1, p', z_1 \dots z_m)$ avem $(s_1, z_1 \dots z_m) \in f(s, i_1, z)$ iar în P va exista regula

$$(s, z, s') \rightarrow i_1(s_1, z_1, s_2)(s_2, z_2, s_3) \dots (s_m, z_m, s')$$

unde luăm stările s_2, \dots, s_m cele rezultate la descompunerea lui p' . Pe de altă parte, din ipoteza inductivă, avem

$$\begin{aligned} (s_1, z_1, s_2) &\xRightarrow{*} p_1, \\ (s_2, z_2, s_3) &\xRightarrow{*} p_2, \\ &\dots \\ (s_m, z_m, s') &\xRightarrow{*} p_m. \end{aligned}$$

Putem scrie derivarea

$$(s, z, s') \Rightarrow i_1(s_1, z_1, s_2)(s_2, z_2, s_3) \dots (s_m, z_m, s') \xRightarrow{*} i_1 p_1 \dots p_m = i_1 p' = p.$$

După cum am văzut, rezultă mai departe $p \in L(G)$ și $L(APD) \subseteq L(G)$.

Pentru a demonstra incluziunea inversă, vom arăta mai întâi implicația

$$(s, z, s') \xRightarrow{*} p(s_1, z_1 s_2) \dots (s_m, z_m, s') \text{ implică } (s, p, z) \xrightarrow{*} (s_1, \lambda, z_1 \dots z_m).$$

Procedăm prin inducție asupra lungimii derivării l .

Dacă $l = 1$ atunci $p = i$ și se aplică regula

$$(s, z, s') \rightarrow i(s_1, z_1, s_2) \dots (s_m, z_m, s')$$

deci $(s_1, z_1 \dots z_m) \in f(s, i, z)$ și $(s, i, z) \mapsto (s_1, \lambda, z_1 \dots z_m)$.

Presupunem că implicația este adevărată pentru l oarecare și considerăm o derivare de lungime $l + 1$. Fie $p = p'i$ și punem în evidență ultimul pas.

$$\begin{aligned} (s, z, s') &\xRightarrow{*} p'(s'_{j-1}, z'_{j-1}, s'_j)(s_j, z_j, s_{j+1}) \dots (s_m, z_m, s') \\ &\Rightarrow p'i(s_1, z_1, s_2) \dots (s_{j-1}, z_{j-1}, s_j)(s_j, z_j, s_{j+1}) \dots (s_m, z_m, s'), \end{aligned}$$

unde $s'_j = s_j$; la ultimul pas s-a aplicat regula

$$(s'_{j-1}, z'_{j-1}, s_j) \rightarrow i(s_1, z_1, s_2) \dots (s_{j-1}, z_{j-1}, s_j).$$

Rezultă $(s_1, z_1 \dots z_{j-1}) \in f(s'_{j-1}, i, z'_{j-1})$ și putem scrie evoluția $(s'_{j-1}, i, z'_{j-1}) \mapsto (s_1, \lambda, z_1 \dots z_{j-1})$.

Pe de altă parte, conform ipotezei inductive, avem

$$(s, p', z) \xrightarrow{*} (s'_{j-1}, \lambda, z'_{j-1} z_j \dots z_m)$$

Prin urmare

$$(s, p, z) = (s, p'i, z) \xrightarrow{*} (s'_{j-1}, i, z'_{j-1} z_j \dots z_m) \mapsto (s_1, \lambda, z_1 \dots z_m)$$

și implicația este demonstrată.

Fie acum $p \in L(G)$, deci $s_0 \xRightarrow[G]{*} p$. Ținând seama de forma regulilor din G , în această derivare se va aplica prima dată o regulă de forma (1), apoi regula de forma (2) iar la sfârșit reguli de forma (3). La aplicarea regulilor (2) putem rescrie la fiecare pas simbolul neterminal cel mai din stânga, deci să obținem o derivare extrem stângă. Să observăm că în acest caz structura formelor propoziționale intermediare este cea menționată, $p(s_1, z_1, s_2)(s_2, z_2, s_3) \dots (s_m, z_m, s')$.

Prin urmare, derivarea va avea forma

$$s_0 \Rightarrow (s_0, z_0, s') \xRightarrow{*} p(s_1, z_1, s_2) \dots (s_m, z_m, s') \xRightarrow{*} p.$$

Trebuie să avem regulile $(s_j, z_j, s_{j+1}) \rightarrow \lambda, j = 1, \dots, m, s_{m+1} = s'$ și putem scrie

$$(s_0, p, z_0) \xrightarrow{*} (s_1, \lambda, z_1 \dots z_m) \xrightarrow{} (s_2, \lambda, z_2 \dots z_m) \xrightarrow{} \dots \xrightarrow{} (s', \lambda, \lambda)$$

adică $p \in L(APD)$ și $L(G) \subseteq L(APD)$. \square

Automate push-down cu stări finale. Vom nota un automat push-down cu stări finale cu APD_f .

Teorema 3.12 *Un limbaj este recunoscut de un automat push-down dacă și numai dacă este recunoscut de un automat push-down cu stări finale.*

Demonstrație. Partea I $E = l(APD) \Rightarrow E \in L(APD_f)$.

Dacă $APD = (\Sigma, I, Z, f, s_0, z_0)$ construim un automat push-down cu stări finale astfel

$$APD_f = (\Sigma \cup \{s'_0, s_f\}, I, Z \cup \{z'_0\}, f', s'_0, z'_0)$$

unde mulțimea de stări finale este $\{s\}$ iar funcția de evoluție este definită de:

$$\begin{aligned} f'(s, i, z) &= f(s, i, z), s \in \Sigma, i \in I \cup \{\lambda\}, z \in Z; \\ f'(s'_0, \lambda, z'_0) &= (s_0, z_0 z'_0); \\ f(s, \lambda, z'_0) &= (s_f, \lambda), s \in \Sigma; \\ &\text{în rest } \emptyset. \end{aligned}$$

Fie $p \in L(APD)$; atunci $(s_0, p, z_0) \xrightarrow{*} (s, \lambda, \lambda)$. Evident că aceeași evoluție o poate avea și automatul push-down cu stări finale. Putem scrie în APD_f evoluția

$$(APD_f) : (s'_0, p, z'_0) \xrightarrow{*} (s_0, \lambda, z'_0) \xrightarrow{} (s, \lambda, \lambda),$$

deci $p \in L(APD_f)$ și $L(APD) \subseteq L(APD_f)$.

Invers, fie $p \in L(APD_f)$, atunci (în APD_f)

$$(s'_0, p, z'_0) \xrightarrow{} (s, p, z_0 z'_0) \xrightarrow{*} (s_f, \lambda, q).$$

Ultimul pas trebuie să fie de forma $(s, \lambda, z'_0) \xrightarrow{} (s_f, \lambda, \lambda)$ pentru că nu există altă valoare a lui f care să ne ducă într-o stare finală. Deci (în APD_f)

$$(s'_0, p, z_0 z'_0) \xrightarrow{*} (s, \lambda, z'_0) \xrightarrow{} (s_f, \lambda, \lambda)$$

și putem scrie în APD evoluția $(s_0, p, z_0) \xrightarrow{*} (s, \lambda, \lambda)$, adică $p \in L(APD)$ și $L(APD_f) \subseteq L(APD)$. \square

Partea II $E = L(APD_f) \Rightarrow E \in L(APD)$.

Fie $APD_f = (\Sigma, I, Z, f, s_0, z_0, \Sigma_f)$ un automat push-down cu stări finale (mulțimea stărilor finale este Σ_f) și construim un APD astfel

$$APD = (\Sigma \cup \{s'_0, s'\}, I, Z \cup \{z'_0\}, f', s'_0, z'_0)$$

unde

$$f'(s, i, z) = f(s, i, z), s \in \Sigma, i \in I \cup \{\lambda\}, z \in Z;$$

$$f(s'_0, \lambda, z'_0) = (s, z_0 z'_0);$$

$$f(s, \lambda, z) = (s', \lambda), s \in \Sigma_f \cup \{s'\}, z \in Z \cup \{z'_0\};$$

în rest \emptyset .

Fie $p \in L(APD_f)$, atunci $(s_0, p, z_0) \xrightarrow{*} (s, \lambda, q), s \in \Sigma_f$. Este evident că în APD avem evoluția $(s_0, p, z_0) \xrightarrow{*} (s, \lambda, q)$. Putem scrie

$$APD : (s'_0, p, z'_0) \xrightarrow{*} (s_0, p, z_0 z'_0) \xrightarrow{*} (s, \lambda, q z'_0) \xrightarrow{*} (s', \lambda, \lambda),$$

deci $p \in L(APD)$ și $L(APD_f) \subseteq L(APD)$.

Invers, fie $p \in L(APD)$. Avem

$$APD : (s'_0, p, z'_0) \xrightarrow{*} (s_0, p, z_0 z'_0) \xrightarrow{*} (s, \lambda, \lambda).$$

Simbolul z'_0 nu poate fi șters decât cu o regulă de forma $f(s, \lambda, z) = (s', \lambda)$, $s \in \Sigma_f \cup \{s'\}$, deci APD trebuie să ajungă într-o stare $s \in \Sigma_f$, apoi să rămână în s' .

$$APD : (s_0, p, z_0 z'_0) \xrightarrow{*} (s, \lambda, q z'_0), s \in \Sigma_f.$$

Putem scrie

$$APD_f : (s_0, p, z_0) \xrightarrow{*} (s, \lambda, q), s \in \Sigma_f$$

și deci $p \in L(APD_f)$, adică $L(APD) \subseteq L(APD_f)$. \square

3.6 Automate push-down deterministe

Funcționarea unui APD este în general nedeterministă, $\text{card } f(s, i, z) \geq 1$. Pentru ca un APD să aibă o funcționare deterministă nu este suficient să impunem condiția $\text{card } f(s, i, z) = 1$, deoarece dacă pentru un anumit $s \in \Sigma$ și $z \in Z$ avem $f(s, \lambda, z) \neq \emptyset$ și $f(s, i, z) \neq \emptyset$, putem face un pas citind λ sau citind i .

Definiție 3.10 *Un automat push-down este determinist dacă*

- (1) $\text{card } f(s, i, z) \geq 1, s \in \Sigma, i \in I \cup \{\lambda\}, z \in Z$;
- (2) dacă $f(s, \lambda, z) \neq \emptyset$, atunci $f(s, i, z) = \emptyset, \forall i \in I$.

Un limbaj recunoscut de un APD determinist îl vom numi limbaj independent de context (sau de tipul doi) determinist. Familia limbajelor independente de context deterministe este inclusă (strict) în familia limbajelor de tipul 2 (după cum vom vedea).

Un APD se poate bloca în următoarele două situații

1. Automatul ajunge în starea s , în vârful memoriei push-down se află simbolul z , pe banda de intrare urmează i și $f(s, i, z) = f(s, \lambda, z) = \emptyset$;
2. Intră într-un ciclu infinit citind de pe bandă λ ; de exemplu $f(s, \lambda, z) = (s, z)$ și $f(s, i, z) = \emptyset$ pentru o anumită pereche (s, z) .

Definiție 3.11 *Un APD determinist este neblocabil dacă pentru orice cuvânt $p \in I^*$ există o evoluție de forma $(s_0, p, z_0) \xrightarrow{*} (s, \lambda, q)$.*

Într-un APD determinist neblocabil orice cuvânt peste I poate fi citit. Evident, de aici nu rezultă că orice cuvânt este recunoscut de APD.

Lema 3.6 *Un APD determinist cu stări finale este echivalent cu un APD determinist cu stări finale neblocabil (relativ la prima situație de blocare).*

Demonstrație. Fie $APD_f = (\Sigma, I, Z, f, s_0, z_0, \Sigma_f)$. Construim $APD'_f = (\Sigma \cup \{s'_0, s'\}, I, Z \cup \{z'_0\}, f', s'_0, z'_0, \Sigma_f)$ unde:

- (1) $f'(s, i, z) = f(s, i, z)$ dacă $f(s, i, z) \neq \emptyset, s \in \Sigma, i \in I \cup \{\lambda\}, z \in Z$;
- (2) $f'(s, i, z) = (s', z)$ dacă $f(s, i, z) = f(s, \lambda, z) = \emptyset, s \in \Sigma, i \in I, z \in Z$;
- (3) $f'(s', i, z) = (s', z), i \in I, z \in Z$;
- (4) $f'(s', \lambda, z'_0) = (s_0, z_0 z'_0)$.

Avem

$$p \in L(APD_f) \Leftrightarrow (s_0, p, z_0) \xrightarrow{APD_f} (s, \lambda, q), s \in \Sigma_f \Leftrightarrow$$

$$(s'_0, p, z'_0) \xrightarrow{APD'_f} (s_0, p, z_0 z'_0) \xrightarrow{APD'_f} (s, \lambda, q z'_0), s \in \Sigma_f \Leftrightarrow p \in L(APD'_f).$$

Deci $L(APD_f) = L(APD'_f)$. \square

Observație. Într-o situație de blocare (s, ip, zq) și $f(s, i, z) = f(s, \lambda, z) = \emptyset$ putem scrie

$$(s, ip, zq) \xrightarrow{APD'_f} (s', p, zq) \xrightarrow{APD'_f} \dots \xrightarrow{APD'_f} (s', \lambda, zq).$$

Teorema 3.13 *Orice limbaj independent de context determinist este recunoscut de un APD determinist neblocabil.*

Demonstrație. Fiind dat un APD determinist vom construi un APD determinist neblocabil echivalent. Conform lemei anterioare putem presupune că nu are loc prima situație de blocare.

Dacă are loc a doua situație de blocare, putem avea două cazuri:

1. conținutul memoriei push-down se mărește nelimitat;
2. lungimea cuvintelor scrise în memoria push-down nu depășește un anumit număr.

Fie $\text{card}(\Sigma) = n$, $\text{card}(Z) = k$, $l = \max\{|q|, q \in Z^* \mid (s', q) \in f(s, i, z)\}$.

Cazul 1. Există în total un număr nk de perechi de forma (s, z) . Dacă în evoluția lui APD s-ar succede numai configurații cu perechi de forma (s, z) distincte atunci lungimea cuvântului din memoria push-down ar crește la maximum nkl simboluri. Prin urmare, dacă $(s', \lambda, \alpha') \xrightarrow{*} (s'', \lambda, \alpha'')$ și $|\alpha''| - |\alpha'| > nkl$, atunci în această evoluție trebuie să existe două configurații cu aceeași stare s și cu același simbol z în vârful memoriei push-down. Deci

$$(s', \lambda, \alpha') \xrightarrow{*} (s, \lambda, zr) \xrightarrow{*} (s, \lambda, zqr) \xrightarrow{*} (s'', \lambda, \alpha''),$$

de unde urmează că

$$(s', \lambda, \alpha') \xrightarrow{*} (s, \lambda, zr) \xrightarrow{*} (s, \lambda, q^m r), \forall m \in N.$$

Cazul II Lungimea cuvântului scris în memoria push-down nu depășește nkl , căci dacă $|\alpha''| - |\alpha'| > nkl$, ar rezulta că în memoria push-down am avea $zq^m, \forall m \in N$ și lungimea cuvântului nu ar fi finită. \square

Teorema 3.14 *Dacă E este un limbaj independent de context determinist atunci limbajul complementar $I^* \setminus E$ este de asemenea independent de context determinist.*

Demonstrație. Fie $APD_f = (\Sigma, I, Z, f, s_0, z_0, \Sigma_f)$ automatul push-down determinist care recunoaște limbajul E . Construim un automat push-down determinist care va recunoaște limbajul $I^* \setminus E$ în modul următor

$$APD'_f = (\Sigma \times \{\sigma_1, \sigma_2, \sigma_3\}, I, Z, f', s'_0, z_0, \Sigma'_f)$$

unde

$$s'_0 = \begin{cases} (s_0, \sigma_1) & \text{pentru } s_0 \in \Sigma_f, \\ (s_0, \sigma_2) & \text{pentru } s_0 \notin \Sigma_f, \end{cases}$$

mulțimea de stări finale este $\Sigma'_f = \{(s, \sigma_3) \mid s \in \Sigma\}$ iar funcția de evoluție este definită de

- (1) dacă $f(s, i, z) = (s', q)$ atunci
- $$\begin{aligned} f'((s, \sigma_1), i, z) &= ((s', k'), q), \\ f'((s, \sigma_2), i, z) &= ((s, \sigma_3), q), \\ f'((s, \sigma_3), i, z) &= ((s', k'), q); \end{aligned}$$
- (2) dacă $f(s, \lambda, z) = (s', q)$ atunci
- $$\begin{aligned} f'((s, \sigma_1), \lambda, z) &= ((s', k'), q), \\ f'((s, \sigma_3), \lambda, z) &= ((s', k'), q); \end{aligned}$$

unde $k' = 1$ dacă $s' \in \Sigma_f$ și $k' = 2$ dacă $s' \notin \Sigma_f$.

Fie $p \in L(APD'_f)$, atunci $((s_0, \sigma_k), p, z_0) \xrightarrow{*} ((s, \sigma_3), \lambda, q), k = 1, 2$. În mod necesar, ultima configurație trebuie să fie precedată de o configurație de forma $((s', \sigma_2), \lambda, q)$, deci

$$(A) ((s_0, \sigma_k), p, z_0) \xrightarrow{*} ((s', \sigma_2), \lambda, q) \xrightarrow{*} ((s, \sigma_3), \lambda, q)$$

de unde rezultă că $((s_0, p, z_0) \xrightarrow{*} (s', \lambda, q)$ și $s' \notin \Sigma_f$. Deci $p \in I^* \setminus E$ și $L(APD'_f) \subseteq I^* \setminus E$.

Invers, fie $p \in I^* \setminus E$, atunci în APD'_f avem evoluția $((s_0, p, z_0) \xrightarrow{*} (s, \lambda, q), s \notin \Sigma_f$ și putem scrie evoluția (A). Prin urmare $p \in L(APD'_f)$, adică $I^* \setminus E \subseteq L(APD'_f)$. □

Consecință. Putem enunța proprietatea de mai sus astfel: Familia limbajelor independente de context deterministe este închisă la complementariere. Cum familia limbajelor independente de context nu este închisă la complementariere și cum ea coincide cu familia limbajelor recunoscute de automatele push-down nedeterministe putem mai departe obține următorul rezultat:

APD nedeterministe nu sunt echivalente cu APD deterministe.

3.7 Probleme propuse

1. Să se arate că $L(G) = \{(ab)^n a \mid n \geq 0\}$ unde G are producțiile $S \rightarrow SbS, S \rightarrow a$. Să se construiască o gramatică echivalentă cu G care să fie neambiguă.
2. Eliminați redenumirile din gramatica G_E ce generează expresiile aritmetice simple.
3. Aduceți la forma normală Chomsky gramaticile ce au regulile:
 - (a) $S \rightarrow TbT, T \rightarrow TaT \mid ca;$
 - (b) $S \rightarrow aAC, A \rightarrow aB \mid bAB, B \rightarrow b, C \rightarrow c;$
 - (c) $S \rightarrow A.A, A \rightarrow 0A \mid 1A \mid \dots \mid 9A \mid \lambda.$
4. Găsiți forma normală Greibach pentru gramaticile:

- (a) $A_1 \rightarrow A_2 A_3, A_2 \rightarrow A_1 A_2 | 1, A_3 \rightarrow A_1 A_3 | 0$.
 - (b) G_E care generează expresiile aritmetice simple.
5. Construiți un automat push-down pentru recunoașterea limbajului:
- (a) $L = \{w | w \in \{a, b\}^*, \text{ numărul literelor } a \text{ în } w \text{ este egal cu numărul literelor } b \text{ în } w\}$;
 - (b) $L = \{w | w \in \{a, b\}^*, w = \tilde{w}\}$;
 - (c) $L = \{w | w \in \{(,)\}^*, w \text{ este un cuvânt în care fiecare paranteză deschisă are o pereche, paranteză închisă}\}$.
6. Folosind lema de pompare să se arate că următoarele limbaje nu sunt independente de context:
- (a) $L = \{a^i b^j c^k | i < j < k\}$;
 - (b) $L = \{a^i b^j | j = i^2\}$;
 - (c) $L = \{a^n b^n c^n | n \geq 1\}$;
 - (d) $L = \{a^i | i \text{ prim}\}$;
 - (e) $L = \{a^i b^i c^j | j \geq i\}$;

Capitolul 4

Analiza Sintactică

Analiza sintactică este o fază a procesului de compilare care are următoarele două obiective principale:

- Stabilește dacă un cuvânt dat aparține sau nu limbajului, deci dacă cuvântul este *corect* din punct de vedere sintactic. În particular limbajul poate fi definit de o gramatică generativă de tip Chomsky și deci termenul de analiză sintactică trebuie înțeles în sensul teoriei limbajelor formale. Mai menționăm că prin *cuvânt* înțelegem orice structură constituită cu simbolurile acceptate de limbaj, în particular un întreg program, dar de obicei ne vom mărgini la anumite entități, de exemplu, o linie sau un rând.
- Determină derivarea (arborele de derivare) corespunzător cuvântului. Odată cu această operație sunt degajate anumite structuri pentru care se poate genera cod intermediar, structuri pe care le vom numi *unități sintactice*.

Pe lângă aceste obiective principale se mai efectuează și alte operații, de exemplu, analiza și tratarea erorilor, prelucrarea tabelelor, etc. Rezultatul analizei sintactice va fi un fișier care conține derivările (arborii de derivare) corespunzătoare unităților sintactice în care este descompus programul: expresii aritmetice, declarații, etc. Acest fișier este utilizat în faza de generare a formatului intermediar. În mod curent însă, generatoarele de format intermediar sunt niște rutine, apelate de analizorul sintactic, astfel încât formatul intermediar se obține succesiv.

Teoretic, problema analizei sintactice este rezolvată de automatele corespunzătoare diverselor tipuri de limbaje; această cale conduce însă la algoritmi cu complexitate mare (număr mare de stări, funcție de tranziție complexă, etc.). Există algoritmi speciali de analiză sintactică cu eficiență superioară. În continuare ne vom ocupa cu două clase de astfel de algoritmi:

- Algoritmi top-down (de sus în jos);
- Algoritmi bottom-up (de jos în sus).

4.1 Algoritmi TOP-DOWN

4.1.1 Algoritmul general de analiză top-down

Algoritmul general de analiză top-down are un principiu foarte simplu: *se aplică în mod sistematic regulile de generare, începând cu simbolul de start; în cazul unui eșec se revine în sus și se încearcă o altă regulă. Regulile se aplică în ordinea în care sunt scrise în gramatică, fără să existe o anumită ordine preferențială de scriere, întrucât natura algoritmului nu permite nici un fel de ierarhizare a regulilor din punctul de vedere al frecvenței de utilizare.*

Pentru descrierea riguroasă a acestui algoritm am urmat modelul din cartea lui D. Gries, *Compiler construction for digital computers*.

Fie $G = (V_N, V_T, x_0, \mathcal{P})$ o gramatică de tipul 2 și $p \in V_T^*$. Ne interesează următoarele două probleme:

- (a) $p \in L(G)?$,
- (b) Dacă $p \in L(G)$ atunci să se determine derivarea $x_0 \Rightarrow p$.

Considerăm toate regulile care au X_0 în stânga:

$$x_0 \rightarrow x_{11} \dots x_{1n_1} | x_{21} \dots x_{2n_2} | \dots | x_{p1} \dots x_{pn_p}$$

Inițial x_0 devine *activ* și alege prima regulă $x_0 \rightarrow x_{11} \dots x_{1n_1}$. Dacă această alegere este corectă trebuie să avem $x_0 \Rightarrow x_{11} \dots x_{1n_1} \xRightarrow{*} p$ și în conformitate cu *lema de localizare* a gramaticilor de tipul 2, cuvântul p se poate descompune în forma $p = p_1 p_2 \dots p_{n_1}$, unde $x_{1j} \xRightarrow{*} p_j$, $j = 1, \dots, n_1$.

Simbolul x_0 îl activează pe x_{11} și îi cere să găsească derivarea $x_{11} \xRightarrow{*} p_1$; dacă x_{11} reușește, transmite lui x_0 *succes*. Simbolul x_0 îl dezactivează pe x_{11} , îl activează pe x_{12} și îi cere să găsească derivarea $x_{12} \xRightarrow{*} p_2$, etc. Dacă toate simbolurile activate de x_0 transmit succes, construcția este terminată. Să presupunem că x_{1j} transmite *eșec*; atunci x_0 îl dezactivează pe x_{1j} , îl reactivează pe x_{1j-1} căruia îi transmite: *Mi-ai dat o derivare, dar aceasta nu este bună, încearcă alta*. Dacă x_{1j-1} reușește, procesul se continuă spre dreapta; dacă nu, atunci x_0 îl dezactivează pe x_{1j-1} , îl reactivează pe x_{1j-2} căruia îi cere o altă derivare. Procesul se continuă în acest mod fie spre dreapta, fie spre stânga. Dacă se ajunge la x_{11} și acesta nu reușește să găsească o altă derivare, x_0 decide că prima regulă aleasă nu este bună și încearcă cu următoarea regulă, adică $x_0 \rightarrow x_{21} \dots x_{2n_2}$, ș. a. m. d.

Observații

- Fiecare simbol devenit activ, procedează exact ca și părintele său, alege prima regulă, activează primul simbol, etc.
- Nu se cunoaște anticipat descompunerea $p = p_1 p_2 \dots p_{n_1}$. Deci x_{1j} transmite succes dacă reușește să găsească o derivare $x_{1j} \xRightarrow{*} p_j$, unde p_j este un subcuvânt oarecare al lui p , cu singura condiție ca p_{1j} să înceapă din punctul unde s-a terminat p_{1j-1} . De exemplu, dacă $p = i_1 i_2 \dots i_8 \dots$ și $p_1 = i_1 i_2 i_3 i_4$,

$p_2 = i_5 i_6$, atunci x_{13} trebuie să găsească o derivare de forma $x_{13} \xrightarrow{*} i_7 i_8 \dots$. În particular, dacă $x_{1j} \in V_T$ decizia de *succes* sau *eșec* depinde de faptul dacă x_{1j} coincide sau nu cu simbolul din p care urmează (În exemplul de mai sus, dacă x_{13} coincide sau nu cu i_7).

- Când se reactivează un simbol și i se cere o nouă derivare, acesta reactivează ultimul simbol fiu și îi cere același lucru.

Exemplu Considerăm următoarea gramatică G_E care generează expresii aritmetice simple. Operanzii sunt notați simbolic cu a , operatorii sunt $+$ și $*$ iar ordinea naturală a operațiilor este completată de paranteze.

$$\begin{cases} E \rightarrow T + E | T \\ T \rightarrow F * T | F \\ F \rightarrow (E) | a \end{cases}$$

Procesul de analiza sintactică top-down pentru cuvântul $p = a * a$ este prezentat în figura 4.1.

În această figură, revenirile în sus au fost marcate prin încadrarea într-un dreptunghi a subarborului care a condus la un eșec. Dreptunghiurile interioare au semnificația unor eșecuri realizate mai devreme (este vorba de timpul de desfășurare al procesului). Arborele corespunzător cuvântului este cel neîncadrat în vreun dreptunghi (în figură acesta este situat în partea dreaptă).

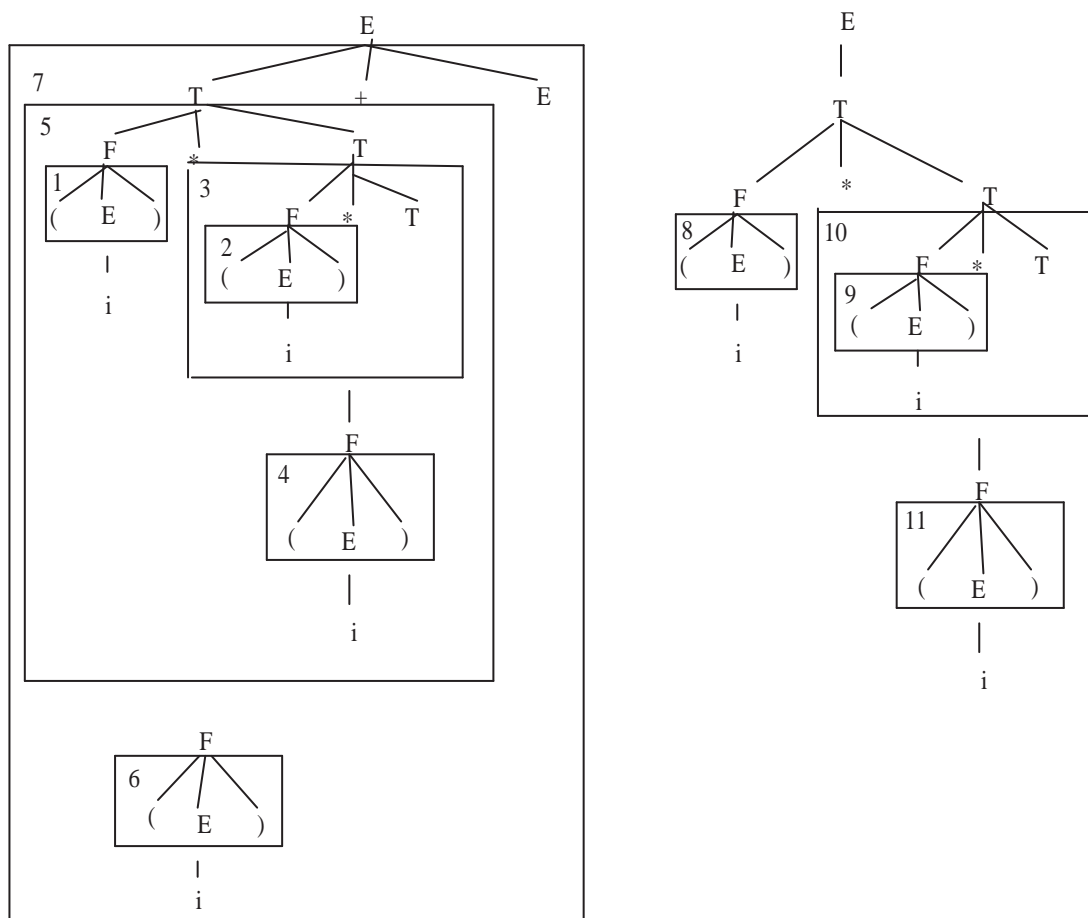
4.1.2 Analiza top-down fără reveniri

În cazul unor gramatici cu o formă specială se poate face o analiză de tip top-down fără reveniri. Principala condiție este ca în cazul mai multor alternative (reguli cu același simbol în stânga), să se poată decide cu precizie ramura corectă. În general o astfel de decizie se poate realiza prin analiza unor simboluri care urmează în cuvântul de analizat.

Exemplu Considerăm următoarea gramatică G care generează secvențe de declarații și instrucțiuni cuprinse între cuvintele cheie **Program** și **EndProgram**. Declarațiile sunt notate simbolic cu d , iar instrucțiunile cu i . Fiecare instrucție sau declarație se încheie cu $;$, exceptând cazul celei ce precede **EndProgram**.

$$G \begin{cases} < program > \rightarrow Program D I EndProgram \\ D \rightarrow d; X \\ X \rightarrow d; X | \lambda \\ I \rightarrow iY \\ Y \rightarrow ; iY | \lambda \end{cases}$$

Să considerăm urmatorul cuvânt de analizat

Figura 4.1: Arborele sintactic pentru $p = i * i$

```

Program
  d;
  d;
  i;
  i;
  i;
EndProgram

```

O derivare extrem stângă pentru acest cuvânt este

$$\begin{aligned}
< program > &\Rightarrow Program\ D\ I\ EndProgram \Rightarrow Program\ d; X\ I\ EndProgram \\
&\Rightarrow Program\ d; d; X\ I\ EndProgram \Rightarrow Program\ d; d; I\ EndProgram \\
&\Rightarrow Program\ d; d; iY\ EndProgram \Rightarrow Program\ d; d; i; iY\ EndProgram \\
&\Rightarrow Program\ d; d; i; i; iY\ EndProgram \Rightarrow Program\ d; d; i; i; i\ EndProgram
\end{aligned}$$

Pentru construcția derivării extrem stângi se procedează astfel: inițial se consideră simbolul de start $< program >$ pentru care se alege singura regulă disponibilă (dacă primul simbol din șirul de analizat nu coincide cu primul terminal al regulii se poate decide imediat *eroare sintactică*). Următorul simbol neterminal pentru care trebuie aleasă o regulă este D iar șirul rămas de analizat (*de generat*) începe cu d , astfel că regula aleasă va fi $D \rightarrow d; X$. Din cuvântul inițial trebuie generată în continuare secvența $d; i; i; i; EndProgram$ ce începe cu d iar neterminalul cel mai din stânga este X , astfel că se alege regula ce începe cu d . În continuare, din cuvântul inițial rămâne de generat secvența $i; i; i; EndProgram$ ce începe cu i iar neterminalul cel mai din stânga este X , astfel că se alege regula ce $X \rightarrow \lambda$, ș.a.m.d.

Dacă se consideră gramatica ce generează expresii aritmetice simple, în forma considerată la algoritmul general top-down (4.2.9), atunci la primul pas al unei derivări extrem stângi pentru cuvântul $(a + a * a) + a$ nu se poate decide regula de ales prin citirea primului terminal (deoarece ar fi necesar să consultăm șirul rezultat până la întâlnirea operatorului $+$ aflat după paranteza închisă).

O gramatică pentru care alegerea regulii de aplicat este **unic determinată** de următorul simbol terminal din șirul de analizat se numește gramatică $LL(1)$ (Left to right parsing, Leftmost derivation, 1 symbol lookahead). În multe cazuri există posibilitatea de a transforma gramatica într-una echivalentă de tip $LL(1)$. Pentru cazul particular al gramaticii pentru generarea expresiilor aritmetice, o gramatică echivalentă este următoarea:

$$\begin{cases} E \rightarrow TE' \\ E' \rightarrow +TE' \mid -TE' \mid \lambda \\ T \rightarrow FT' \\ T' \rightarrow *FT' \mid /FT' \mid \lambda \\ F \rightarrow (E) \mid id \mid num \end{cases}$$

$$\left\{ \begin{array}{l} \langle \text{bloc} \rangle \rightarrow \{ \langle \text{lista} \rangle \} \\ \langle \text{lista} \rangle \rightarrow \langle \text{instr} \rangle L \\ L \rightarrow ; \langle \text{instr} \rangle L \mid \lambda \\ \langle \text{instr} \rangle \rightarrow id = E \mid \text{if}(E) \text{ then } \langle \text{instr} \rangle \mid \text{while}(E) \text{ do } \langle \text{instr} \rangle \mid \{ \langle \text{lista} \rangle \} \\ E \rightarrow TE' \\ E' \rightarrow +TE' \mid -TE' \mid \lambda \\ T \rightarrow FT' \\ T' \rightarrow *FT' \mid /FT' \mid \lambda \\ F \rightarrow (E) \mid id \mid num \end{array} \right.$$

Figura 4.2: Gramatica pentru generarea blocurilor de instrucțiuni.

```
{
  a = 2;
  b = b + 1;
  if ( b-a ) then { x = x-1;
                  a = 3 };
  c = b*72
}
```

Figura 4.3: Program sursă de analizat sintactic

4.1.3 Programarea unui analizor sintactic. Studiu de caz

Programarea unui analizor sintactic top-down fără reveniri se poate face relativ ușor asociind câte o funcție la fiecare neterminal. Analiza unui cuvânt revine la un șir de apeluri corespunzătoare tratării simbolurilor ce apar pe parcursul construcției derivării extrem stângi. Fiecare funcție asociată conține o singură instrucțiune *switch* cu clauze ce corepund regulilor gramaticii. Alegerea regulii se face după următoarea unitate lexicală din textul de analizat.

Să considerăm următoarea gramatică (vezi figura 4.2) ce generează un bloc de instrucțiuni de atribuire, condiționale de tip **if** (expresie aritmetică) **then** instrucțiune sau repetitive de tip **while**. Încușă instrucțiune poate fi o instrucțiune simplă sau bloc de instrucțiuni separate prin delimitatorul ;(SEMICOLON).

Un text sursă ce poate fi analizat de această gramatică este cel din figura 4.3.

Lista de unități lexicale furnizate de analizorul lexical este pentru acest caz

LBRACE id LET num SEMI id LET id PLUS num SEMI if LPAR id MINUS id
RPAR then LBRACE id LET id LET id minus num SEMI id LET num RBRACE
SEMI id LET id ORI num RBRACE.

Un program pentru analiza sintactică top-down fără reveniri corespunzător gramaticii din figura 4.2 este reprezentat parțial în figura 4.4. Lista codurilor de unități lexicale (terminalele gramaticii) conține SEMI LPAR RPAR ș.a.m.d. Presupunem că analizorul lexical funcționează ca funcție ce returnează codul următoarei unități lexicale din textul sursă. Variabila de tip întreg *token* conține codul returnat de analizorul lexical *ALEX()*. S-au mai definit două funcții: *err()* pentru tratarea erorilor de sintaxă depistate și *eat(int tok)* ce consumă din textul sursă unitatea lexicală *tok* pe care ne așteptăm să o găsim în text.

Rularea programului pentru exemplul din figura 4.3 va produce o secvență de apeluri recursive ca în figura 4.5.

Un pic de algebră

Vom da în cele ce urmează o definiție riguroasă a gramaticilor $LL(k)$ (k simboluri citite în avans) vom da în cele ce urmează, împreună cu condițiile necesare și suficiente ca o gramatică să intre în această categorie.

Definiție Fie $G = (V_N, V_T, S, \mathcal{P})$ o gramatică independentă de context. G se zice de tip $LL(k)$ dacă și numai dacă oricare ar fi două derivări extrem stângi

$$\begin{aligned} S \xRightarrow{*} uXv &\Rightarrow u\alpha v \xRightarrow{*} u\gamma \\ S \xRightarrow{*} uXv &\Rightarrow u\beta v \xRightarrow{*} uv \quad \text{unde } X \rightarrow \alpha \in \mathcal{P}, X \rightarrow \beta \in \mathcal{P}, \gamma, \nu \in V_T^+ \end{aligned}$$

din $\gamma^k = \nu^k$ rezultă $\alpha = \beta$ (notația γ^k înseamnă primele k simboluri din șirul γ).

Pentru $k = 1$ se obține cazul gramaticilor din secțiunile precedente. Restricția asupra numărului de simboluri citite în avans restrânge drastic mulțimea limbajelor ce pot fi analizate cu astfel de gramatici. Un inconvenient major pentru $k > 1$ este creșterea exponențială a dimensiunii tabelului de predicție (câte o intrare în tabel pentru fiecare combinație de k simboluri). O variantă mai eficientă de analiză se obține cu gramatici $LR(1)$ (Left to right parsing, Rightmost derivation, 1 symbol lookahead).

Pentru orice $\alpha \in V_G^+$, $X \in V_N$, $X \rightarrow \alpha \in \mathcal{P}$ vom considera următoarele mulțimi:

$$\begin{aligned} Prim(\alpha) &= \{a \in V_T \mid \alpha \xRightarrow{*} a\beta, \beta \in V_G^*\} \\ Urm(\alpha) &= \{a \in V_T \mid S \xRightarrow{*} uXv, a \in Prim(v)\} \\ SD(X, \alpha) &= \{a \in V_T \mid a \in Prim(\alpha) \text{ sau } (\alpha \xRightarrow{*} \lambda \text{ si } a \in Urm(X))\} \\ NULL(G) &= \{X \in V_N \mid X \xRightarrow{*} \lambda\} \end{aligned}$$

Dacă un terminal $a \in Prim(\alpha)$ atunci a poate apărea pe prima poziție într-un șir derivat din α . Mulțimea $SD(X, \alpha)$ poartă denumirea de *mulțimea simbolurilor directe* asociate regulii $X \rightarrow \alpha$, iar $NULL(G)$ conține neterminalele ce se pot șterge (în unul sau mai mulți pași) cu reguli din G .

Teoremă.

$$G \in LL(k) \Leftrightarrow SD(X, \alpha) \cap SD(X, \beta) = \Phi, \forall X \rightarrow \alpha, X \rightarrow \beta \in \mathcal{P}, \beta \neq \alpha$$

```

final int if=1, then=2, LPAR=3, RPAR=4, LBRACE=5, RBRACE=6,
PLUS=7, MINUS=8, ORI=9, DIV=10, SEMI=11, id=12,
while=13, do=14, LET=15;
void err();
int ALEX();
int token = ALEX();
void eat(int tok){
    if (tok==token) token=ALEX else err()
};

void bloc(){
    eat(LBRACE); lista(); eat(RBRACE)
};

void lista(){
    instr();L()
};

void L(){
    switch(token){
        case SEMI: eat(SEMI); instr(); L(); break;
        default: break }
};

void instr(){
    switch(token){
        case if: eat(if); eat(LPAR); E(); eat(RPAR); eat (then); instr(); break;
        case id: eat(id); eat(LET); E(); break;
        case while: eat(while); eat(LPAR); E(); eat(RPAR); eat(do); instr(); break;
        case LBRACE: eat(LBRACE); instr(); eat(RBRACE); break;
        default: printf("syntax error: if, identif, while or left brace expected");
                err()}
};

...
void E(){
    T(); Eprime();
};

void Eprime(){
    switch(token) {
        case PLUS: eat(PLUS);T();Eprime();break;
        case MINUS: eat(MINUS);T();Eprime();break;
        default: break;
    }
};

...

```

Figura 4.4: Cod C corespunzător analizorului sintactic

```

bloc()
  eat(LBRACE);
  lista()
    instr()
      eat(id);
      eat(LET);
      E()
        T()
          F()
            eat(num);
            return_F;
          Tprime()
            return_Tprime;
          return_T;
        Eprime()
          return_Eprime;
      return_E;
    return_instr;
  L()
    eat(SEMI)
    instr()
      .... aici se recunoaste ; b = b + 1
    return_instr;
  L()
    eat(SEMI);
    instr()
      .... aici se recunoaste ; if ( ... }
    return_instr;
  L()
    .... aici se recunoaste ; c = b * 72
    return_L();
  return_L;
return_L;
return_lista;
eat(RBRACE);
return_bloc;

```

Figura 4.5: Execuția analizei lexicale pentru textul sursa

Determinarea mulțimilor definite anterior se poate face ușor. Pentru $NULL(G)$, se poate aplica algoritmul definit la eliminarea regulilor de ștergere pentru o gramatică independentă de context.

```

 $NULL = \{X \in V_N | X \rightarrow \lambda \in \mathcal{P}\}$ 
repeat
  (1)  $AUX = NULL$ ;
  (2)  $NULL = NULL \cup \{X \in V_N | X \rightarrow p, p \in AUX^+\}$ ;
until  $NULL = AUX \square$ 

```

Calculul mulțimilor $Prim(X)$ și $Urm(X)$ se poate face cu algoritmul următor:

Pas 1: { Inițializare }

$Prim(a) = \{a\}; \forall a \in V_T$

$Prim(X) = Urm(X) = \Phi; \forall X \in V_N$

Pas 2: **Repeat**

Pentru fiecare regulă $X \rightarrow Y_1 Y_2 \dots Y_k$ **Do**

For ($i = 1; i \leq k; i++$)

(2.1) **if** ($i = 1$ sau $Y_1 Y_2 \dots Y_i \in NULL^+$)

$Prim(X) = Prim(X) \cup Prim(Y_i)$;

(2.2) **if** ($i = k$ sau $Y_{i+1} \dots Y_k \in NULL^+$)

$Urm(Y_i) = Urm(Y_i) \cup Urm(X)$;

(2.3) **For** ($j = i + 1; j \leq k; j++$)

if ($j = i + 1$ sau $Y_{i+1} \dots Y_{j-1} \in NULL^+$)

$Urm(Y_i) = Urm(Y_i) \cup Prim(Y_j)$;

end for

end for

end do

Until $Prim(X), Urm(X)$ nu se schimbă la o iterație. \square

Calculul mulțimilor $Prim(\alpha)$ revine la definiția recursivă

$$Prim(X\alpha) = \begin{cases} Prim(X) & \text{dacă } \alpha \notin NULL(G) \\ Prim(X) \cup Prim(\alpha) & \text{dacă } \alpha \in NULL(G) \end{cases}$$

Dacă se consideră gramatica ce generează expresii aritmetice 4.1.2, atunci rezultatul aplicării algoritmilor este prezentat în tabelul următor:

	$NULL$	$Prim$	Urm
E	DA	$(, i$	$)$
E'		$+, -$	$)$
T		$(, i$	$+, -,)$
T'		$*, /$	$+, -,)$
F		$(, i$	$*, / ,)$

iar mulțimile simbolurilor directoare sunt

$$\begin{aligned} SD(E, TE') &= \{ (, i \} \\ SD(E', +TE') &= \{ + \} \\ SD(E', \lambda) &= \{ \} \\ SD(T, FT') &= \{ (, i \} \\ SD(T', *FT') &= \{ * \} \\ SD(T', \lambda) &= \{ + \} \\ SD(F, (E)) &= \{ (\} \\ SD(F, id) &= \{ id \} \\ SD(F, num) &= \{ num \} \end{aligned}$$

Funcțiile *C* sau *Java* corespunzătoare neterminalelor se modifică puțin, în sensul că vom avea clauze corespunzătoare simbolurilor terminalelor ce se regăsesc în mulțimile simbolurilor directoare asociate regulilor. Dacă la reconstrucția derivării trebuie aplicată o regulă pentru *X* și următorul terminal din textul de analizat nu face parte din mulțimea simbolurilor directoare asociate vreunei reguli cu *X* în stânga, atunci se apelează modulul de tratare a erorilor sintactice.

Pentru cazul gramaticii anterioare, câteva din secvențele de cod asociate neterminalelor sunt prezentate în continuare.

```
...
void E(){
switch(token) {
    case LPAR:
    case id:
    case num: T(); Eprime(); break;
    default: printf("syntax error: (,identifier or number expected");
              err()
}
}

void Tprime(){
switch(token) {
    case ORI: eat(ORI);F();Tprime();break;
    case DIV: eat(DIV);F();Tprime();break;
    case PLUS:
    case RPAR: break;
    default: printf("syntax error: *,/,+,) expected");
              err()
}
}
```

Tratarea erorilor este dependentă de proiectantul compilatorului și limbaj. Principal, există trei modalități de tratare.

- Se oprește analiza sintactică în punctul unde s-a depistat o eroare, cu transmiterea unui mesaj prietenos către utilizator, de exemplu: *Eroare sintactică: așteptam să urmeze delimitator de instrucțiune. Mai învață sintaxa limbajului! BYE!*
- Se încearcă repararea greșelii de sintaxă inserând în textul sursă un simbol din mulțimea de simboluri directe asociate regulii ce s-a aplicat. Desigur că este suficient să presupunem că am inserat în text, astfel încât analiza poate continua (desigur nu vom uita să anunțăm utilizatorul că nu cunoaște regulile de sintaxă și l-am corectat noi!). Inserarea poate conduce la cicluri infinite, astfel că nu este recomandabilă totdeauna.
- Se încearcă găsirea unui simbol ce se potrivește ignorând toate terminalele textului sursă până se întâlnește un simbol din mulțimea simbolurilor directe. Se transmite același mesaj prietenos către autorul textului sursă, apoi se continuă analiza.

4.2 Algoritmi BOTTOM-UP

4.2.1 Gramatici cu precedență simplă

Fie $G = (V_N, V_T, x_0, \mathcal{P})$ o gramatică de tipul doi și $p \in L(G)$ o formă propozițională, adică un cuvânt peste alfabetul general V_G astfel încât $x_0 \xrightarrow{*} p$. Vom nota cu $\mathcal{A}_{x_0, p}$ arborele de derivare care are rădăcina x_0 și frontiera p .

Definiția 1. Vom spune că f este o frază simplă a lui p dacă f este un subcuvânt al lui p și în plus:

- $\exists x \in V_N$, cu $x \rightarrow f \in \mathcal{P}$;
- $\mathcal{A}_{x, f} \subset \mathcal{A}_{x_0, p}$.

Exemplu. Considerăm gramatica G_E 4.2.9 și forma propozițională $p = T * F + a * (E + T)$. Arborele de derivare este prezentat în figura 4.6. Se poate observa că frazele simple sunt $T * F, a, E + T$.

Observație. Subcuvântul F respectă prima condiție dar nu este frază simplă, întrucât nu este satisfăcută condiția a doua din definiție.

Definiția 2. Fraza simplă cea mai din stânga poartă denumirea de *frază simplă stângă*.

Vom nota fraza simplă stângă corespunzătoare lui p cu f_p în exemplul nostru, $f_p = T * F$. După cum vom vedea în continuare, fraza simplă stângă are un rol important în algoritmi de analiză sintactică bottom-up. În principiu, acești algoritmi prevăd următorii pași (descriși de figura 4.7):

Regulile determinate la pasul (3), aplicate în ordine inversă, ne vor da derivarea corespunzătoare lui p . Să mai observăm că în acest mod arborele de derivare se

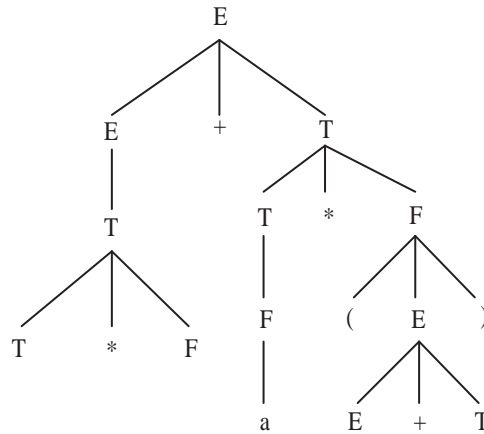


Figura 4.6: Arborele de derivare corespunzător derivării lui $p = T * F + a * (E + T)$

- (1) Inițializare p ;
- (2) Se determină f_p , fraza simplă stângă a lui p ;
- (3) Se determină regula $x \rightarrow f_p$;
- (4) Se **reduce** fraza simplă stângă, adică dacă $p = r f_p s$, se pune $p = r x s$;
GOTO(2).

Figura 4.7: Algoritmul bottom-up

construiește de jos în sus (bottom-up). Problema cea mai dificilă este desigur determinarea frazei simple stângi de la pasul (2). În principiu ar trebui cunoscut anticipat arborele de derivare corespunzător lui p , dar tocmai acesta este scopul analizei sintactice. După cum vom vedea, această dificultate se poate elimina, deci putem determina fraza simplă stângă fără a cunoaște acest arbore. De fapt, fraza simplă stângă este legată de anumite proprietăți ale gramaticii pe care le tratăm în continuare.

4.2.2 Relații de precedență

Definiția 3. Fie $x, y \in V_G$ două simboluri oarecare ale gramaticii.

Vom spune că:

- (1) $x \prec y$ (x precede pe y) dacă $\exists p$ astfel încât $p = rxys$, $x \notin f_p$, $y \in f_p$;
- (2) $x \pm y$ (x este egal cu y) dacă $\exists p$ astfel încât $p = rxys$, $x \in f_p$, $y \in f_p$;
- (2) $x \succ y$ (y succede lui x) dacă $\exists p$ astfel încât $p = rxys$, $x \in f_p$, $y \notin f_p$;

Relațiile $\prec \pm \succ$ se numesc *relații de precedență simple* (atunci când nu există posibilitatea de confuzie se folosesc denumirile *mai mic*, *egal*, *mai mare* pentru relațiile de precedență). Să observăm că aceste relații nu se exclud reciproc, deci că putem avea, de exemplu, $x \prec y$ și în același timp $x \succ y$, întrucât pot exista cazuri de gramatici în care putem găsi două forme propoziționale p_1, p_2 astfel încât să avem simultan cele două relații. Să mai facem de asemenea observația că există mai multe tipuri de astfel de relații, cu definiții asemănătoare, întrucât relațiile trebuie să satisfacă condiții suplimentare care depind și de gramatică; pentru toate aceste tipuri de relații vom utiliza aceleași notații.

Definiția 4. O gramatică de tipul 2 spunem că este cu *precedență simplă* dacă între oricare două simboluri ale gramaticii **există cel mult o relație** de precedență.

Pentru algoritmul de analiză sintactică ascendentă este important ca să nu existe reguli cu părțile drepte identice; altminteri, în cadrul pasului (4) din algoritmul de analiză sintactică nu se poate decide la cine să se facă reducerea frazei simple stângi. De asemenea regulile de ștergere nu sunt admise, pentru acest caz noțiunea de frază simplă neavând sens. Aceste condiții presupunem în mod sistematic că sunt satisfăcute.

Fie acum $V_G = \{x_1, x_2, \dots, x_n\}$ alfabetul general, unde am adoptat o anumită ordine a simbolurilor. Definim matricea de precedență a gramaticii,

$$M = (a_{ij}), \text{ unde } a_{ij} = \begin{cases} \prec & \text{daca } x_i \prec x_j \\ \pm & \text{daca } x_i \pm x_j \\ \succ & \text{daca } x_i \succ x_j \end{cases}$$

Exemplu. Considerăm următoarea gramatică

$$G = (\{A, B\}, \{0, 1\}, A, \{A \rightarrow A0 \mid 1B, B \rightarrow 1\}).$$

Matricea de precedență va fi

	A	B	0	1
A			\prec	
B			\succ	
0			\succ	
1		\pm	\succ	\prec

4.2.3 Proprietăți ale gramaticilor cu precedență simplă

Proprietatea 1. Fie $p = a_1 \dots a_n$ o formă propozițională într-o gramatică cu precedență simplă (a_i sunt simboluri neterminale sau terminale ale gramaticii). Atunci $f_p = a_i \dots a_j$ dacă și numai dacă

$$a_1 \overset{\pm}{\prec} a_2 \dots \overset{\pm}{\prec} a_{i-1} \prec a_i \pm a_{i+1} \dots \pm a_j \succ a_{j+1} (*)$$

Demonstrația acestei proprietăți se efectuează considerând toate structurile de arbori de derivare posibile într-o situație dată. Vom exemplifica această idee pentru cuvântul $p = a_1 \dots a_7$. Dacă presupunem că $f_p = a_3 a_4 a_5$, atunci, conform definiției frazei simple stângi, avem

$$a_2 \prec a_4, a_4 \pm a_5, a_5 \succ a_6$$

Mai trebuie arătat că $a_1 \overset{\pm}{\prec} a_2$. În acest scop vom considera toate posibilitățile de structuri arborescente; o parte din ele sunt prezentate în figura 4.8.

Se poate observa că în ipoteza că $f_p = a_3 a_4 a_5$, singurele situații posibile sunt (a) și (b). Celelalte situații contrazic această ipoteză, de exemplu, în cazurile (c) și (d) fraza simplă stângă ar fi , respectiv, $a_1 a_2$ sau a_2 . Acum, în cazul (a), prin reducerea succesivă a frazei simple stângi, vom obține un arbore în care fraza simplă stângă va fi $a_1 a_2$, adică $a_1 \pm a_2$ iar în cazul (b), fraza simplă stângă va fi a_2 și atunci $a_1 \prec a_2$.

Implicația inversă se bazează pe cea directă; se presupune că $f_p = a_l \dots a_k$ cu $l \neq i$ și $k \neq j$ și se consideră în continuare toate posibilitățile de poziționare a indicilor l, k față de i, j . De exemplu, dacă $f_p = a_2 a_3$ atunci, conform implicației directe, avem

$$a_1 \prec a_2 \pm a_3 \succ a_4$$

ceea ce contrazice relația (*), conform căreia $a_3 \pm a_4$ iar gramatica are precedență simplă.

Observație La delimitarea frazei simple stângi se procedează astfel: parcurgem șirul de simboluri de la stânga la dreapta până se găsește relația \succ (sau marginea dreaptă a formei propoziționale) pentru determinarea ultimului simbol din f_p , apoi se parcurge șirul spre stânga, trecând peste relații \pm , până la găsirea unei

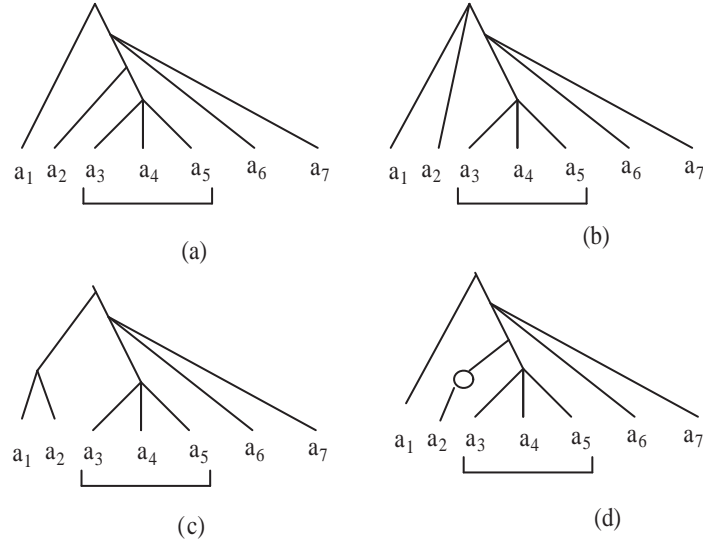
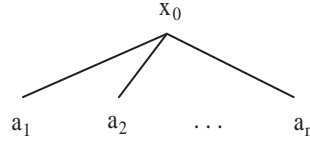


Figura 4.8: Arbori de derivare posibili

Figura 4.9: Arborele de derivare pentru $ni = 1$

relații \prec (sau marginea dreaptă a formei propoziționale) pentru determinarea primului simbol din f_p .

Proprietatea 2. O gramatică cu precedență simplă este neambiguă.

Schiță de demonstrație. Fie p o formă propozițională și $\mathcal{A}_{x_0,p}$ arborele de derivare corespunzător. Vom arăta că acest arbore este unicul arbore de derivare care are rădăcina x_0 și frontiera p . Procedăm prin inducție asupra numărului de noduri interne ni . Dacă $ni = 1$ atunci arborele de derivare va avea forma din figura 4.9 și unicitatea acestuia este evidentă. Presupunem acum că proprietatea este adevărată pentru un ni oarecare și considerăm cazul $ni + 1$. Să presupunem că ar exista doi arbori diferiți cu rădăcina x_0 și frontiera p ; fie aceștia $\mathcal{A}'_{x_0,p}$ și $\mathcal{A}_{x_0,p}$.

Deoarece gramatica este cu precedență simplă, rezultă că fraza simplă stângă, care este aceeași în cei doi arbori, este situată în aceeași poziție, $p = rf_p s$. Efectuăm reducerea frazei simple f_p (conform regulii aplicate $x \rightarrow f_p$) și vom

obține arborii $\mathcal{A}'_{x_0,q}$ și $\mathcal{A}_{x_0,q}$ unde $q = rXs$ iar cei doi arbori trebuie să fie diferiți. Dar numărul de noduri interne este ni , în contradicție cu ipoteza.

4.2.4 Determinarea relațiilor de precedență pentru gramatici cu precedență simplă

Așa cum s-a precizat deja, relațiile de precedență sunt proprietăți intrinseci ale gramaticii, nu depind de contextul în care se află cele două simboluri. Prin urmare, cunoașterea anticipată a acestor relații împreună cu proprietatea 1, rezolvă complet problema pasului 2 de la algoritmul de analiză bottom-up, adică determinarea frazei simple stângi. În acest subparagraf vom prezenta o teoremă de caracterizare pe baza căreia se pot determina relațiile de precedență simple și implicit faptul dacă gramatica este sau nu cu precedență simplă.

Vom defini două relații specifice gramaticilor și care vor fi utilizate în continuare, numite, respectiv, **F** (First) și **L** (Last). Fie $x \in V_N$ și $y \in V_G$ două simboluri ale gramaticii. Vom spune că

$$\begin{aligned} xFy \text{ dacă } \exists x \rightarrow yu \in \mathcal{P}, u \in V_G^* \\ xLy \text{ dacă } \exists x \rightarrow uy \in \mathcal{P}, u \in V_G^* \end{aligned}$$

Închiderile tranzitive, respectiv, tranzitive și reflexive ale acestor relații le vom nota cu F^+ , L^+ și respectiv, F^* , L^* . Există numeroși algoritmi direcți care permit calcularea acestei închideri, cu complexitate redusă.

Teorema 2. Avem

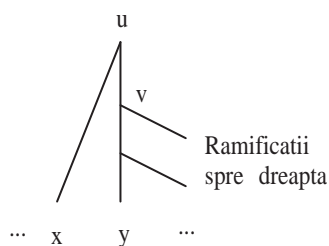
- (1) $x \prec y \Leftrightarrow \exists u \rightarrow \dots xv \dots \in \mathcal{P}, vF^+y$;
- (2) $x \pm y \Leftrightarrow \exists u \rightarrow \dots xy \dots \in \mathcal{P}$;
- (3) $x \succ y \Leftrightarrow \exists u \rightarrow \dots vw \dots \in \mathcal{P}, vL^+x, wF^*y$;

Schiță de demonstrație. Se analizează posibilitățile de arbori de derivare care răspund cerințelor teoremei. De exemplu, pentru cazul (1), trebuie să existe structura de arbore ca în figura 4.10.

Este clar că dacă $x \prec y$, atunci conform definiției frazei simple stângi, trebuie să existe p astfel încât $p = rxys$, $x \notin f_p$, $y \in f_p$, adică să existe structura din 4.10. Invers, dacă există o astfel de structură, atunci fie o formă propozițională care conține u . Efectuăm reduceri succesive până când vom obține o formă propozițională pentru care u aparține frazei simple stângi; la arborele astfel obținut, adăugăm subarborele de mai sus. Vom avea în mod evident $x \notin f_p$, $y \in f_p$.

4.2.5 Studiu de caz

Să considerăm cazul gramaticii pentru generarea expresiilor aritmetice simple G_E , cu regulile obișnuite

Figura 4.10: Arborele de derivare corespunzator cazului $x \prec y$

	E	T	F	$+$	$*$	$($	$)$	a
E				\pm			\pm	
T				\prec	\pm			
F				\prec	\prec			
$+$		\prec	\pm	\prec			\prec	\prec
$*$			\pm			\prec		\prec
$($	\prec	\pm	\prec			\prec		\prec
$)$				\prec	\prec		\prec	
a				\prec	\prec		\prec	

Figura 4.11: Matricea de precedență pentru G_E

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

Relațiile First, Last și închiderile acestora sunt:

$$\begin{array}{lll} E F \{E, T\} & E F^+ \{E, T, F, (, a\} & E F^* \{E, T, F, (, a\} \\ T F \{T, F\} & T F^+ \{T, F, (, a\} & T F^* \{T, F, (, a\} \\ F F \{(, a\} & F F^+ \{(, a\} & F F^* \{F, (, a\} \\ \\ E L \{T\} & E L^+ \{T, F, a,)\} & E F^* \{E, T, F, a,)\} \\ T L \{F\} & T L^+ \{F, a,)\} & T F^* \{T, F, a,)\} \\ F L \{a,)\} & F L^+ \{a,)\} & F F^* \{F, a,)\} \end{array}$$

Conform teoremei de calcul a relațiilor de precedență vom obține matricea de precedență (vezi figura 4.11) asociată simbolurilor gramaticii G_E .

$$G : \begin{cases} A \rightarrow A0 \mid B1 \\ B \rightarrow 1 \end{cases} \quad \begin{array}{c|c|c|c} & S & a & b \\ \hline S & \pm & \prec & \pm \\ \hline a & \pm & \prec & \pm \\ \hline b & \succ & \succ & \succ \end{array}$$

Figura 4.12: Exemplu de gramatică cu precedență simplă

Regula	Forma propozitională
	$a \prec a \prec \underline{a \pm b} \succ aababbbabb$
$S \rightarrow ab$	$a \prec a \pm S \prec a \prec \underline{a \pm b} \succ abbbabb$
$S \rightarrow ab$	$a \prec a \pm S \prec a \pm S \prec \underline{a \pm b} \succ bbabb$
$S \rightarrow ab$	$a \prec a \pm S \prec \underline{a \pm S \pm S \pm b} \succ babb$
$S \rightarrow aSSb$	$a \prec \underline{a \pm S \pm S \pm b} \succ abb$
$S \rightarrow aSSb$	$a \pm S \prec \underline{a \pm b} \succ b$
$S \rightarrow ab$	$a \pm S \pm S \pm b$
$S \rightarrow aSSb$	S

Figura 4.13: Reconstrucția derivării cuvântului $p = aaabaababbbabb$

Deci gramatica G_E nu are proprietatea de precedență simplă, astfel că nu poate fi aplicat direct algoritmul de analiză sintactică bottom up.

Vom exemplifica algoritmul de analiză pentru gramatica din figura 4.12, matricea de precedență asociată este cea alăturată.

Reconstituirea derivării cuvântului $p = aaabaababbbabb$ este prezentată în tabelul 4.13, unde fiecare linie corespunde formei propoziționale, în care fraza simplă stângă este subliniată, precedată de regula identificată.

4.2.6 Gramatici operatoriale

Gramaticile operatoriale sunt gramatici cu o structură specială a regulilor de rescriere, în care anumite simboluri sunt considerate operatori iar celelalte operanzi, operatorii trebuind să separe operanzii. Această structură este preluată de la gramaticile care generează expresiile aritmetice iar principiul de analiză este de asemenea preluat de la algoritmul de evaluare a expresiilor aritmetice. În continuare vom prezenta acest principiu pentru evaluarea expresiilor aritmetice fără paranteze, cazul expresiilor cu paranteze putându-se reduce la cel fără paranteze, de exemplu utilizând tehnica de modificare a ponderilor operatorilor la întâlnirea parantezelor.

În primul rând se atribuie operatorilor anumite ponderi care vor defini ordinea de efectuare a operațiilor. De obicei, se consideră $p(+) = p(-) = 1 < 2 = p(*) =$

$p(/)$, ceea ce înseamnă că mai întâi se fac adunările și scăderile, apoi înmulțirile și împărțirile, etc. Efectuarea unei operații depinde de contextul în care se află operatorul respectiv, și anume, dacă ponderea operatorului precedent este mai mare sau egală cu ponderea operatorului curent, atunci se poate efectua operația definită de operatorul precedent. Din acest motiv, acest tip de gramatici se numesc *cu operator de precedență* (operator precedence grammars).

Putem realiza un algoritm simplu de evaluare utilizând două stive

- **P**- stiva operatorilor;
- **S**- stiva operanzilor.

Vom utiliza indicii i și k pentru a indica nivelele celor două stive. Prin urmare, notația p_i , respectiv s_k au sensul de operatorul, respectiv operandul din stivă aflat pe nivelul i , respectiv k . Pentru simplificare, vom utiliza aceeași notație p_i atât pentru operator cât și pentru ponderea operatorului respectiv.

Algoritmul de evaluare:

PAS 1: *Inițializări:* Se introduce primul operator și primul operand în stivele P și S și se inițializează nivelele celor două stive, $i = k = 2$;

PAS 2: Se introduce următorul operator în stiva P și se actualizează nivelul, $i = i + 1$;

PAS 3: Dacă $p_{i-1} < p_i$ atunci se introduce următorul operand în lista S , se actualizează nivelul stivei, $k = k + 1$, și se revine la pasul 2;

PAS 4: Dacă $p_{i-1} \geq p_i$ atunci:

în stiva P : p_{i-1} se înlocuiește cu p_i ;

în stiva S : s_{k-1} se înlocuiește cu $s_{k-1}p_{i-1}s_k$;

se actualizează nivelele celor două stive, $i = i - 1$; $k = k - 1$ și se revine la pasul 3.

Observație. Pentru uniformitatea algoritmului se introduce operatorul marcaj, notat cu $\#$, cu ponderea cea mai mică; orice expresie se încadrează între două marcaje iar oprirea algoritmului are loc la pasul 4 în cazul în care $p_1 = p_2 = \#$. Schematic, acest algoritm este prezentat în figura 4.14.

Exemplu. Tabelul 4.15 conține starea dinamică a celor două stive pentru expresia $\#a + b * c/d - f\#$. Menționăm că în cadrul pașilor 2 și 4 s-a completat câte un rând nou în tabel.

Observație. Ponderile operatorilor sunt legate de ierarhia lor în gramatică, cu cât un operator este mai jos cu atât ponderea lui este mai mare. Acest fapt este ilustrat în figura 4.16.

Pentru cazul expresiilor aritmetice cu paranteze algoritmul se păstrează, dar este aplicat după preprocesarea expresiei prin modificarea dinamică a ponderilor operatorilor și renunțarea la paranteze. Se parcurge expresia de la stânga la dreapta alocând ponderi operatorilor; la întâlnirea unei paranteze deschise ponderile se incrementează cu 10, iar la întâlnirea unei paranteze închise se decrementează ponderile operatorilor cu 10. În continuare parantezele sunt ignorate.

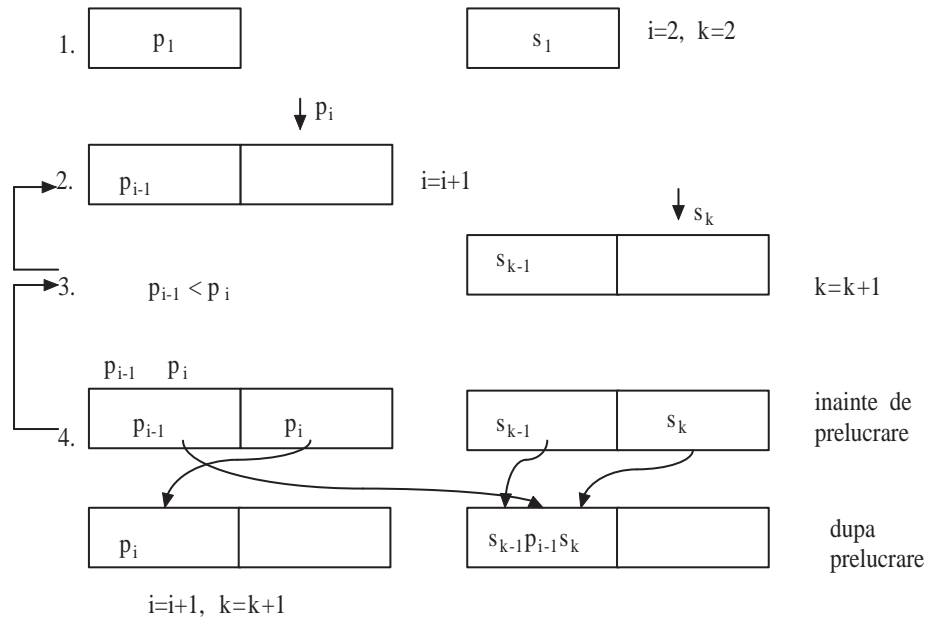


Figura 4.14: Prelucrările efectuate asupra celor două stive

	1	2	3	4	1	2	3	4
1	#0				a			
2	#0	+1			a	b		
3	#0	+1	-1		a	b		
4	#0	-1			$a + b$	c		
5	#0	-1	*2		$a + b$	c	d	
6	#0	-1	*2	/2	$a + b$	c	d	
7	#0	-1	/2		$a + b$	$c * d$	f	
8	#0	-1	/2	#0	$a + b$	$c * d$	f	
9	#0	-1	#0		$a + b$	$c * d / f$		
10	#0	#0			$a + b - c * d / f$			

Figura 4.15: Evaluarea expresiei $\#a + b - c * d / f\#$

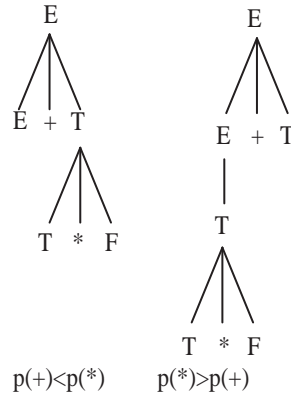


Figura 4.16: Ierarhia operatorilor

De exemplu, expresia

$$a * (b - c * d) - (x + y / (z - t)) * h$$

devine prin preprocesare

$$a *_2 b -_{11} c *_{12} d -_1 x +_{11} y /_{12} z -_{21} t *_2$$

ceea ce va asigura ordinea corectă de evaluare.

4.2.7 Gramatici operatoriale

Fie $G = (V_N, V_T, x_0, \mathcal{P})$ o gramatică de tipul 2. Vom considera că simbolurile neterminale V_N sunt operanzii limbajului iar simbolurile terminale V_T sunt operatori. Să observăm că în cazul gramaticii G_E care generează expresii aritmetice simple o asemenea convenție este justificată, deoarece oricare din neterminalele E, T, F derivează în terminalul a (generic, acest terminal reprezintă operanzii), sau într-o expresie (eventual între paranteze) care are rolul de operand.

Definiția 1. Vom spune că o gramatică este operatorială dacă regulile de rescriere au forma

$$X \rightarrow N_i T_i N_{i+1} T_{i+1} \dots T_j N_{j+1},$$

unde N_k sunt neterminale (*operanzi*), inclusiv cuvântul vid, iar T_k sunt terminale (*operatori*).

Observație. Într-o gramatică operatorială orice formă propozițională are structura:

$$p = N_1 T_1 N_2 T_2 \dots T_n N_{n+1}$$

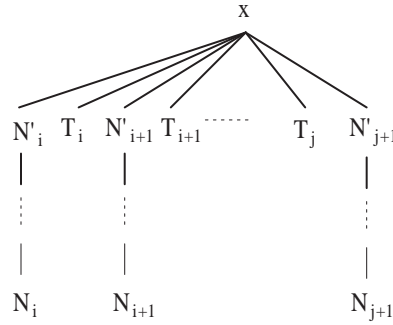


Figura 4.17: Structura subarborelui

Aceasta înseamnă că într-o formă propozițională a unei gramatici operatoriale între oricare doi operanzi există cel puțin un operator; pot însă exista mai mulți operatori alăturați, de exemplu, $\dots((a \dots$

Definiția 2. Fie $p = N_1 T_1 N_2 T_2 \dots T_n N_{n+1}$ o formă propozițională într-o gramatică operatorială. Vom spune că f este o *frază simplă* a lui p dacă $f = N_i T_i N_{i+1} T_{i+1} \dots T_j N_{j+1}$ este un subcuvânt al lui p care conține *cel puțin un simbol terminal* și în plus:

- (1) $\exists x \rightarrow N'_i T_i N'_{i+1} T_{i+1} \dots T_j N'_{j+1} \in \mathcal{P}$, și $N'_k \xrightarrow{*} N_k \forall k = i, \dots, j$.
- (2) $\mathcal{A}_{x,f} \subset \mathcal{A}_{x_0,p}$.

Structura de arbore corespunzătoare unei fraze simple este prezentată în figura 4.17.

Ca și în cazul gramaticilor cu precedență simplă, fraza simplă cea mai din stânga poartă denumirea de *frază simplă stângă*. Principiul de analiză sintactică este identic cu cel de la gramatici cu precedență simplă. Problema principală este și aici determinarea frazei simple stângi. Aceasta operație poate fi făcută utilizând relațiile de precedență dintre simbolurile gramaticii care pentru gramatici operatoriale au o definiție puțin modificată.

Definiția 3. Fie $x, y \in V_T$. Vom spune că:

- (1) $x \overset{\circ}{<} y$ (x precede pe y) dacă $\exists p$ astfel încât $p = rxNys$, $N \in V_N$, $x \notin f_p$, $y \in f_p$;
- (2) $x \overset{\circ}{=} y$ (x este egal cu y) dacă $\exists p$ astfel încât $p = rxNys$, $N \in V_N$, $x \notin f_p$, $y \in f_p$;
- (2) $x \overset{\circ}{>} y$ (y succede lui x) dacă $\exists p$ astfel încât $p = rxy s$, $N \in V_N$, $x \in f_p$, $y \notin f_p$;

Matricea de precedență se definește numai pe mulțimea simbolurilor terminale, ceea ce conduce la o simplificare a algoritmului de analiză.

Definiția 4. O gramatică operatorială se spune că are *precedență simplă* dacă între două simboluri ale gramaticii există cel mult o relație de precedență.

Ca și la celelalte tipuri de gramatici, în vederea efectuării pasului 3 din algoritmul de analiză, vom presupune că nu există reguli cu părțile drepte identice.

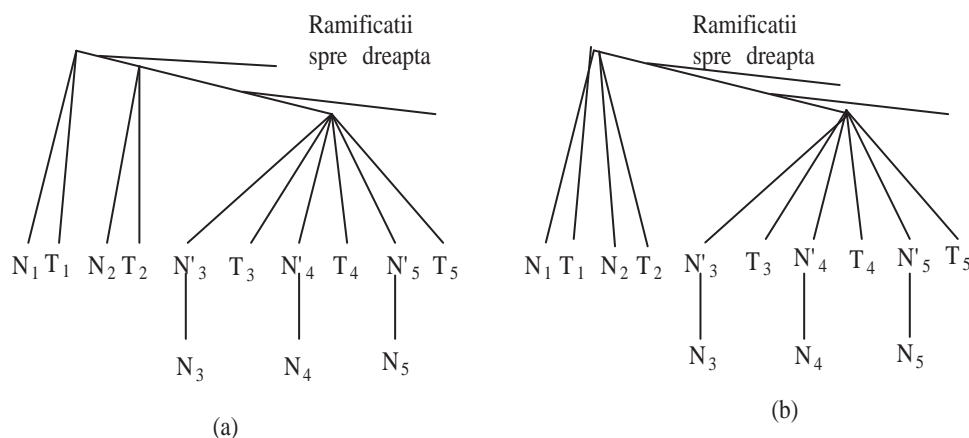


Figura 4.18: Structurile posibile de arbori

De asemenea, sunt valabile proprietățile de caracterizare și de neambiguitate.

Proprietatea 1. Fie $p = N_1T_1N_2T_2 \dots T_nN_{n+1}$ o formă propozițională într-o gramatică cu precedență simplă. Atunci $f_p = N_iT_iN_{i+1}T_{i+1} \dots T_jN_{j+1}$ dacă și numai dacă

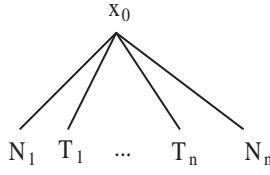
$$T_1 \overset{\circ}{\leq} T_2 \dots \overset{\circ}{\leq} T_{i-1} \overset{\circ}{<} T_i \overset{\circ}{=} T_{i+1} \dots \overset{\circ}{=} T_j \overset{\circ}{>} T_{j+1} \quad (*)$$

Schiță de demonstrație. Presupunem că $p = N_1T_1N_2T_2 \dots T_nN_{n+1}$ și că fraza simplă stângă este $f_p = N_3T_3N_4T_4N_5T_5$. Atunci din definiția frazei simple stângi, rezultă $T_2 \overset{\circ}{<} T_3 \overset{\circ}{=} T_4 \overset{\circ}{=} T_5 \overset{\circ}{>} T_6$. Mai trebuie să arătăm că $T_1 \overset{\circ}{<} T_2$ sau $T_1 \overset{\circ}{=} T_2$. În acest scop analizăm structurile posibile de arbori (figura 4.18).

Efectuăm reduceri succesive până când $T_2 \in f_p$. În acest moment avem $T_1 \overset{\circ}{<} T_2$ pentru cazul (a) și $T_1 \overset{\circ}{=} T_2$ pentru cazul (b). Pentru implicația inversă presupunem că $f_p = N_kT_kN_{k+1}T_{k+1} \dots T_lN_{l+1}$, $l \neq j$, $k \neq i$ și considerăm toate cazurile de poziționare a indicilor k, l în raport cu i, j . De exemplu, dacă $l < j$ atunci, în conformitate cu implicația directă, rezultă $T_l \overset{\circ}{>} T_{l+1}$ iar conform relației (*) din ipoteză, avem $T_l \overset{\circ}{<} T_{l+1}$ sau $T_l \overset{\circ}{=} T_{l+1}$ ceea ce contrazice ipoteza inițială conform căreia gramatica este cu precedență simplă.

Proprietatea 2. O gramatică operatorială cu precedență simplă este neambiguă.

Schiță de demonstrație. Fie p o formă propozițională și $\mathcal{A}_{x_0,p}$ arborele de derivare corespunzător. Vom arăta că acest arbore este unicul arbore de derivare care are rădăcina x_0 și frontiera p . Procedăm prin inducție asupra numărului de noduri interne ni . Dacă $ni = 1$ atunci arborele de derivare va avea forma din figura 4.19 și unicitatea acestuia este evidentă. Presupunem acum că proprietatea

Figura 4.19: Arborele de derivare pentru $ni = 1$

este adevărată pentru un ni oarecare și considerăm cazul $ni + 1$. Să presupunem că ar exista doi arbori diferiți cu rădăcina x_0 și frontiera p ; fie aceștia $\mathcal{A}'_{x_0,p}$ și $\mathcal{A}_{x_0,p}$.

Deoarece gramatica este cu precedență simplă, rezultă că fraza simplă stângă, care este aceeași în cei doi arbori, este situată în aceeași poziție, $p = rf_p s$. Efectuăm reducerea frazei simple f_p (conform regulii aplicate $x \rightarrow f_p$) și vom obține arborii $\mathcal{A}'_{x_0,q}$ și $\mathcal{A}_{x_0,q}$ unde $q = rXs$ iar cei doi arbori trebuie să fie diferiți. Dar numărul de noduri interne este cel mult ni , ceea ce contrazice ipoteza inductivă.

4.2.8 Determinarea relațiilor de precedență pentru gramatici operatoriale

Vom defini mai întâi relațiile $F_{1,2}$ și $L_{1,2}$, analoage cu F și L de la gramaticile cu precedență simplă.

Definiție Fie $x \in V_N$ și $y \in V_G$. Vom spune că

$$\begin{aligned} xF_{1,2}y & \text{daca } \exists x \rightarrow yu \in \mathcal{P}, \text{ sau } x \rightarrow Nyu \in \mathcal{P}, u \in V_G^*, N \in V_N \\ xL_{1,2}y & \text{daca } \exists x \rightarrow uy \in \mathcal{P}, \text{ sau } x \rightarrow uyN \in \mathcal{P}, u \in V_G^*, N \in V_N \end{aligned}$$

Închiderile tranzitive, respectiv, tranzitive și reflexive ale acestor relații le vom nota cu $F_{1,2}^+$, $L_{1,2}^+$ și respectiv, $F_{1,2}^*$, $L_{1,2}^*$. Acum, determinarea relațiilor de precedență operatoriale se poate face cu ajutorul următoarei teoreme.

Teorema (*determinarea relațiilor de precedență*). Avem

- (1) $x \overset{\circ}{<} y \Leftrightarrow \exists u \rightarrow \dots xv \dots \in \mathcal{P}, vF_{1,2}^+y$;
- (2) $x \overset{\circ}{=} y \Leftrightarrow \exists u \rightarrow \dots xNy \dots \in \mathcal{P}$;
- (3) $x \overset{\circ}{>} y \Leftrightarrow \exists u \rightarrow \dots vy \dots \in \mathcal{P}, vL_{1,2}^+x$;

Schiță de demonstrație. Demonstrația se poate face prin analiza structurilor posibile ale arborilor de derivare. Pentru cazul (1), trebuie să existe structura de arbore ca în figura 4.20.

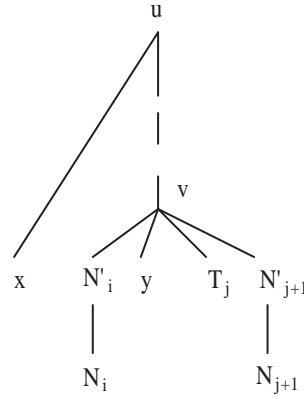


Figura 4.20: Structura posibilă de arbore

Această teoremă ne permite să determinăm matricea de precedență a gramaticii pe baza căreia putem apoi aplica algoritmul de analiză bottom-up. Problema neterminalelor de la capetele frazei simple stângi se poate rezolva analizând părțile drepte ale regulilor de rescriere. Într-adevăr, teorema ne dă numai terminalele care intră în fraza simplă stângă și neterminalele interioare; cele două eventuale neterminale de la capete vor aparține sau nu frazei simple stângi depinzând de forma părții drepte ale regulii respective. Evident, se pot face și ipoteze de neambiguitate suplimentare pentru gramatică, similare cu condiția de a nu exista reguli cu părțile drepte identice.

4.2.9 Studiu de caz

Să considerăm cazul gramaticii pentru generarea expresiilor aritmetice simple G_E , cu regulile obișnuite

$$\begin{cases} E \rightarrow E + T | T \\ T \rightarrow T * F | F \\ F \rightarrow (E) | a \end{cases}$$

Relațiile $F_{1,2}$, $L_{1,2}$ și închiderile acestora sunt:

$$\begin{array}{lll} E F_{1,2} \{E, T, +\} & E F_{1,2}^+ \{E, T, F, +, *, (, a\} & E F_{1,2}^* \{E, T, F, +, *, (, a\} \\ T F_{1,2} \{T, F, *\} & T F_{1,2}^+ \{T, F, *, (, a\} & T F_{1,2}^* \{T, F, *, (, a\} \\ F F_{1,2} \{(, a\} & F F_{1,2}^+ \{(, a\} & F F_{1,2}^* \{F, (, a\} \end{array}$$

$$\begin{array}{lll} E L_{1,2} \{T, +\} & E L_{1,2}^+ \{T, +, F, *, a,)\} & E F_{1,2}^* \{E, T, +, F, *, a,)\} \\ T L_{1,2} \{F, *\} & T L_{1,2}^+ \{F, *, a,)\} & T F_{1,2}^* \{T, F, *, a,)\} \\ F L_{1,2} \{a,)\} & F L_{1,2}^+ \{a,)\} & F F_{1,2}^* \{F, a,)\} \end{array}$$

	+	*	()	<i>a</i>
+	$\overset{\circ}{>}$	$\overset{\circ}{<}$	$\overset{\circ}{<}$	$\overset{\circ}{>}$	$\overset{\circ}{<}$
*	$\overset{\circ}{>}$	$\overset{\circ}{>}$	$\overset{\circ}{<}$	$\overset{\circ}{>}$	$\overset{\circ}{<}$
($\overset{\circ}{<}$	$\overset{\circ}{<}$	$\overset{\circ}{<}$	$\overset{\circ}{=}$	$\overset{\circ}{<}$
)	$\overset{\circ}{>}$	$\overset{\circ}{>}$		$\overset{\circ}{>}$	
<i>a</i>	$\overset{\circ}{>}$	$\overset{\circ}{>}$		$\overset{\circ}{>}$	

Figura 4.21: Matricea de precedență operatorială pentru G_E

<i>Regula</i>	<i>Forma propozitională</i>
$F \rightarrow a$	$\underline{a \overset{\circ}{>}} + \overset{\circ}{<} (a + a) * a * a$
$F \rightarrow a$	$F + \overset{\circ}{<} (\overset{\circ}{<} \underline{a \overset{\circ}{>}} + a) * a * a$
$F \rightarrow a$	$F + \overset{\circ}{<} (\overset{\circ}{<} F + \overset{\circ}{<} \underline{a \overset{\circ}{>}}) * a * a$
$F \rightarrow a$	$F + \overset{\circ}{<} (\overset{\circ}{<} F + F \overset{\circ}{>}) * a * a$
$T \rightarrow T * F, T \rightarrow F$	$F + \overset{\circ}{<} (\underline{T \overset{\circ}{=}}) \overset{\circ}{>} * a * a$
$F \rightarrow (E), E \rightarrow T$	$F + \overset{\circ}{<} F * \overset{\circ}{<} \underline{a \overset{\circ}{>}} * a$
$F \rightarrow a$	$F + \overset{\circ}{<} \underline{F * F} \overset{\circ}{>} * a$
$T \rightarrow T * F, T \rightarrow F$	$F + \overset{\circ}{<} T \overset{\circ}{<} * \overset{\circ}{<} \underline{a}$
$F \rightarrow a$	$F + \overset{\circ}{<} \underline{T * F}$
$T \rightarrow T * F$	$\underline{F + T}$
$E \rightarrow E * T, E \rightarrow T, T \rightarrow F$	\underline{E}

Figura 4.22: Reconstrucția derivării cuvântului $p = a + (a + a) * a * a$

Conform teoremei de calcul a relațiilor de precedență vom obține matricea de precedență (vezi figura 4.21) asociată simbolurilor gramaticii G_E .

Deci gramatica G_E are proprietatea de precedență operatorială simplă, astfel că se poate aplica direct algoritmul de analiză sintactică bottom up.

Reconstituirea derivării cuvântului $p = a + (a + a) * a * a$ este prezentată în tabelul 4.22, unde fiecare linie corespunde formei propoziționale, în care fraza simplă stângă este subliniată, precedată de regula identificată.

Capitolul 5

Sinteza Programelor

5.1 Forme interne ale programelor

5.1.1 Tabelele compilatorului

Compilatoarele prelucrează un număr important de informații, variabile de diverse tipuri, date organizate în anumite structuri standard sau definite de utilizator, etc. De cele mai multe ori aceste informații sunt stocate în tabele și memorate pe suporturile de informații ale sistemului. Organizarea și prelucrarea tabelelor constituie o problemă distinctă, cu influență importantă asupra performanțelor unui compilator, mai ales dacă acestea sunt memorate pe un suport extern. În cele ce urmează vom prezenta pe scurt principalele tabele utilizate în mod obișnuit în cadrul unui compilator și modul de organizare a informațiilor stocate în aceste tabele. Menționăm că structura acestor tabele este definită de proiectantul de compilator și că prin urmare există un anumit grad de subiectivitate în modul de organizare al acestora.

- *Tabelul de variabile.* Conține variabilele utilizate într-un program și informații asupra lor. Menționăm că variabilele conțin datele prelucrate de program, în conformitate cu tipul acestora, conținut care se modifică în mod dinamic pe parcursul evoluției programului. Informațiile din fiecare linie a tabelului pot fi:
 - numele variabilei; așa cum s-a precizat la analiza lexicală, numele unei variabile este de obicei un șir de caractere, litere sau cifre, primul caracter trebuind să fi o literă.
 - tipul variabilei; relativ la tipurile de variabile utilizate în programare trebuie să facem mențiunea că acestea depind în mod esențial de domeniul de aplicabilitate al limbajului. De exemplu, pentru un limbaj de tip matematic (Pascal), tipul unei variabile este caracterizat în general de trei atribute: *structura* (scalar, tablou), *calitatea* (intreg, flotant,

logic), *dimensiunea* (simplă precizie, dublă precizie, etc). În principiu, fiecare variabilă posedă toate aceste trei caracteristici, putând exista și anumite restricții, de exemplu, variabilele întregi pot avea dimensiunea numai simplă precizie sau dublă precizie.

- adresa; adresa de început a zonei de memorie rezervată variabilei;
- Alte informații: sunt informații de tip semantic, utilizate la depanarea programului; de exemplu, dacă variabila a fost sau nu inițializată (adică dacă înainte de prima utilizare i s-a atribuit valoare), dacă a fost utilizată sau nu, etc.

Pentru uniformitatea prelucrării tabelului se poate prevedea o anumită lungime fixă pentru elementele tabelului, chiar dacă spațiul nu este folosit în întregime. De exemplu, în cazul variabilelor de tip tablou, se pot utiliza două sau mai multe linii pentru a stoca valorile maxime ale indicilor (aceste dimensiuni vor fi utilizate atât la rezervarea zonei de memorie cât și la calculul deplasamentului unui element.

Observație asupra funcțiilor. Denumirile funcțiilor sunt deseori utilizate ca variabile care conțin valoarea pe care funcția respectivă o returnează. Este natural ca în acest caz numele funcției să fie considerat variabilă și să fie stocat în acest tabel împreună cu atributele respective; de exemplu, unul din aceste atribute trebuie să fie contorul de amplasare, (CA), adică adresa primei instrucții executabile a rutinei corespunzătoare funcției.

- *Tabelul de etichete.* Conține etichetele utilizate într-un program și adresele primelor instrucții ale secvențelor de program corespunzătoare. În esență, etichetele nu diferă de variabilele obișnuite, dar există anumite particularități; de exemplu, numele unei etichete trebuie să fie constituit dintr-un identificator urmat de un delimitator special (două puncte, etc.) sau numai din caractere numerice (cazul limbajului BASIC) iar valoarea unei etichete este întodeauna o valoare de adresă. De asemenea, în cazul etichetelor trebuie făcută analiza compatibilității dintre etichetele definite și etichetele referite.
- *Tabelul de variabile temporare.* Conține variabile ale compilatorului necesare păstrării unor rezultate intermediare. Având în vedere că rezultatele intermediare, la fel ca și variabilele definite de utilizator, pot avea tipuri diferite, se pot constitui zone separate în care se vor memora rezultate de un anumit tip bine precizat. În acest caz, variabilele temporare pot avea o structură liniară iar adresa unei astfel de variabile se determină cu o formulă de tipul

$$adr = adr_0 + c \times i$$

unde adr_0 este adresa de început a zonei, c este lungimea elementelor iar i numărul de ordine al variabilei temporare. Se poate adopta un sistem de alocare dinamică a variabilelor temporare în scopul economisirii memoriei.

- *Tabelul de constante.* Conține constantele de diverse tipuri utilizate într-un program. În general este un tabel de mai mici dimensiuni, de aceea nu comportă probleme deosebite de organizare și prelucrare.

5.1.2 Cvadrupe și triplete

Cvadrupele și tripletele sunt forme intermediare ale programelor în care fiecare operație elementară împreună cu operanzii ei se stochează într-o secvență de patru sau trei elemente.

Cvadrupe

Ne vom referi în continuare numai la operații binare, pentru celelalte operații putându-se proiecta structuri similare. Considerăm deci o secvență de text de forma $A \text{ op } B$, unde A și B sunt operanzi (variabile, constante, variabile temporare, etc.) iar op este un operator binar. Cvadruplul corespunzător acestei secvențe va fi

$$\overline{op}, \overline{A}, \overline{B}, t_i$$

unde \overline{op} , \overline{A} , \overline{B} sunt codurile unităților lexicale respective iar t_i este o variabilă temporară. Semnificația acestui cvadrupe este următoarea: se execută operația definită de operatorul op între operanzii A și B , în această ordine, iar rezultatul este memorat la t_i . În continuare vom folosi aceeași notație pentru un operand și pentru codul operandului, posibilitatea de confuzie fiind redusă, deci nu vom mai utiliza notații de forma \overline{A} , etc.

Exemplu. Cvadrupele corespunzătoare expresiei aritmetice $a - b + (c + d)/e$ sunt

- (1) $-, a, b, t_1$
- (2) $+, c, d, t_2$
- (3) $/, t_2, e, t_3$
- (4) $+, t_1, t_3, t_4$

În cazul altor operații, chiar de n -aritate diferită de doi, se stabilesc structurile specifice ale cvadrupelelor; pentru operații cu mai mult de doi operanzi, se pot utiliza două sau chiar mai multe cvadrupe. Un exemplu de cvadrupe pentru structurile uzuale ale limbajelor de programare este prezentat în tabelul 5.1.

Observație asupra instrucției **if**. Presupunem că structura instrucției este
if (expresie logica) then instrucțiune_ bloc;

Text Sursă	Cvadruple
$a+b$	$+, a, b, t, \text{ idem } -, *, /$
$a=b$	$=, a, b,$
goto i	$BR, i, ,$
if (a==b) then	BE, a, b, n

Figura 5.1: Exemplu de codificare prin cvadrupe

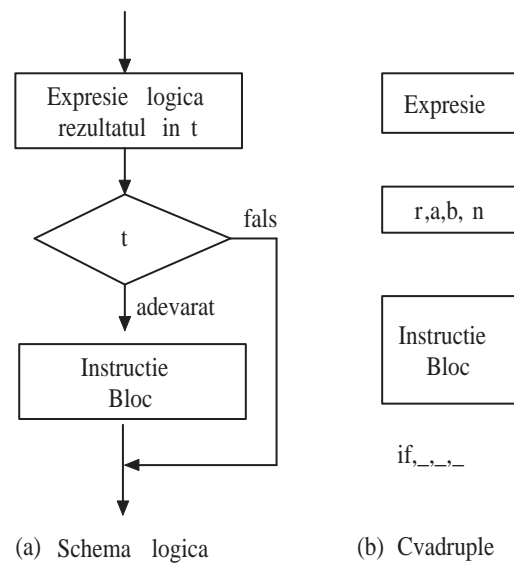


Figura 5.2: Instrucțiunea If logic

Semantica acestei instrucții este prezentată în figura 5.2. Menționăm că expresie logică va avea forma particulară $E_1 r E_2$, unde r este un operator relațional.

În figură (a) prezintă schema logică a instrucției iar (b) cvadrupele corespunzătoare. Partea de cod corespunzătoare expresiei logice va stoca valoarea (logică) a expresiei într-o variabilă temporară t ; această valoare va fi utilizată în instrucția de salt, în conformitate cu tipul de salt prevăzut în programul sursă ($<$, $=$, etc.). În consecință, cvadruplul de salt r va trebui să prevadă această variabilă temporară, plus numărul de ordine n al primului cvadruplu care urmează după cvadrupele corespunzătoare instrucției sau blocului. De exemplu, acest cvadruplu poate avea următoarea formă r, a, b, n unde n este numărul cvadruplului la care trebuie să se facă saltul.

Triplete

Este o altă posibilitate de format intermediar, asemănător cu cvadrupele, singura deosebire fiind aceea că nu se prevede explicit o variabilă temporară pentru stocarea rezultatului. Prin urmare, tripletul corespunzător unei operații binare de forma $A \text{ op } B$ va fi op, A, B . În cazul în care ca operand al unui triplet trebuie specificat rezultatul unui triplet anterior, se va scrie numărul de ordine al tripletului respectiv.

Exemplu. Pentru secvența de program $xx = a - b + (c + d)/e$ șirul de triplete generat va fi

- (1) $-$, a, b
- (2) $+$, c, d
- (3) $/$, (2), e
- (4) $+$, (1), (3)
- (5) $=$, xx , (4)

Notăția poloneză

Notăția poloneză a fost introdusă inițial pentru a se putea evalua expresii aritmetice printr-o singură parcurgere secvențială. Din punctul de vedere al proiectării compilatoarelor, această notație poate fi utilizată ca un format intermediar al programelor, în mod special pentru limbajele de programare algoritmice, Fortran, Pascal. În continuare vom ilustra ideea de evaluare secvențială a expresiilor aritmetice printr-un exemplu.

Considerăm expresia aritmetică $a*b+(c-d*e)*f$. Ideea șirului polonez constă în scrierea operandilor și a operatorilor acestei expresii (fără a utiliza paranteze) într-un șir cu următoarea proprietate: se parcurge șirul secvențial de la stânga la dreapta și fiecare operator întâlnit provoacă efectuarea operației respective între cei doi operanzi din față, în ordinea apariției lor. Pentru exemplul considerat, un

șir care satisface această condiție este

$$ab * cde * -f * +$$

Vom numi acest șir notația poloneză corespunzătoare expresiei, sau *notația poloneză inversă*. Evaluarea se poate face utilizând o stivă cu următoarea prelucrare: se parcurge secvențial șirul, operanzii se introduc în stivă fără nici un control, iar fiecare operator citit provoacă efectuarea operației între operanzii aflați în penultimul și ultimul element al stivei, în această ordine și memorarea rezultatului în penultimul element. Starea dinamică a stivei pentru exemplul considerat este următoarea:

$$\begin{array}{llllll} (1) & a & a*b & a*b & a*b & a*b \\ (2) & b & c & c & c-d*e & (c-d*e)*f \\ (3) & & d & d*e & f & \\ (4) & & e & & & \end{array}$$

Conținutul stivei a fost rescris atunci când un operator provoacă efectuarea unei operații. Se poate observa că șirul considerat de noi (notația poloneză a expresiei) satisface condiția de evaluare secvențială. De fapt, această condiție este esențială și ea poate fi luată chiar ca o definiție a șirului polonez.

Fie V un alfabet și B o mulțime de operatori binari.

Definiția 1. Vom numi șir polonez orice cuvânt peste alfabetul $V \cup B$ obținut recursiv cu ajutorul regulilor:

- (1) dacă $a \in V$ atunci a este șir polonez;
- (2) dacă p și q sunt șiruri poloneze și $b \in B$, atunci pqb este șir polonez.

Să observăm că utilizând aceste două reguli, putem să obținem succesiv toate șirurile poloneze având operanzii V și operatorii binari B . De exemplu, dacă $V = \{a, b, c\}$ și $B = \{+, *\}$, atunci șirurile poloneze vor fi $a, b, c, ab+, ac+, bc+, ab*, \dots, aab+++, aac+++, \dots$

În cele ce urmează vom da o altă definiție, tot constructivă, folosind gramaticile generative Chomsky. Este necesar să precizăm expresiile aritmetice pentru care introducem aceste șiruri.

Definiția 2. Fie gramatica $G_{exp} = (\{E, T, F\}, \{a, +, *\}, E, \mathcal{P})$ unde regulile de rescriere sunt

$$\left\{ \begin{array}{l} E \rightarrow E + T | T \\ T \rightarrow T * F | F \\ F \rightarrow a \end{array} \right.$$

Observație. Terminalul a desemnează operanzii, deci a este de fapt o notație simbolică pentru variabile de forma a, b, c , etc. Este ușor de văzut că această gramatică generează expresii aritmetice fără paranteze, în care singurele operații sunt $+$ și $*$. Utilizarea a trei neterminale E, T, F se face numai pentru a preciza ponderea operațiilor, creșterea făcându-se odată cu adâncimea, deci $p(+) < p(*)$.

Același limbaj al expresiilor aritmetice fără paranteze poate fi generat cu o gramatică cu un singur neterminal, fie acesta A , și având regulile

$$A \rightarrow A + A \mid A * A \mid a$$

Dar în acest caz trebuie introduse alte convenții pentru a preciza ordinea operațiilor; de exemplu, pentru expresia $a + b * c$, în acest ultim caz, nimic nu arată că operația $*$ trebuie făcută înainte de $+$. Trebuie să precizăm că prin considerarea a numai doi operanzi nu se impune o constrângere asupra noțiunii de expresie aritmetică, foarte ușor putem generaliza această gramatică la un număr oarecare de operatori binari, de exemplu

$$\begin{aligned} A_1 &\rightarrow A_1 op_1 A_2 \mid A_2 \\ A_2 &\rightarrow A_2 op_2 A_3 \mid A_3 \\ &\text{etc.} \end{aligned}$$

De asemenea, putem considera operatori de n -aritate oarecare precum și reguli de rescriere în concordanță cu sintaxa pe care dorim să o satisfacem aceștia. Utilizarea unor paranteze pentru a indica o anumită ordine de efectuare a operațiilor, iarăși nu diferă esențial de cazul expresiilor fără paranteze. Putem folosi algoritmul de modificare dinamică a ponderilor descris la gramatici operatoriale: parcurgem expresia de la stânga la dreapta și atașăm fiecărui operator ponderea sa; în momentul în care întâlnim o paranteză deschisă, mărim toate ponderile cu o constantă mai mare decât cea mai mare pondere iar la întâlnirea unei paranteze închise, scădem această constantă. Apoi eliminăm toate parantezele. Este evident că în acest fel operațiile din interiorul parantezelor vor avea prioritatea mai mare, mergându-se spre parantezele din interior.

Definiția 3. Fie $G_p = (V_N, V_T, x_0, \mathcal{P})$, unde $V_N = \{A\}$, $V_T = \{a, +, *\}$, $x_0 = A$ iar regulile sunt

$$A \rightarrow AA + \mid AA * \mid a$$

Orice cuvânt $p \in L(G_p)$ va fi numit șir polonez. Este natural ca aceste șiruri poloneze să fie puse în legătură cu expresiile aritmetice fără paranteze prezentate mai sus. Gramatică G_p nu stabilește o ordine intrinsecă de efectuare a operațiilor, întocmai ca și gramatica cu un singur neterminal de generare a expresiilor. Convenim însă ca această ordine să fie cea *secvențială*.

Dacă $p \in L(G_p)$ atunci convenim să notăm cu $aes(p)$ cuvântul obținut din p cu ajutorul algoritmului de evaluare secvențială.

Următoarea teoremă stabilește legătura între cele două entități.

Teoremă. Pentru orice cuvânt $e \in L(G_{exp})$ există un $p_e \in L(G_p)$ astfel încât $aes(p_e) = e$.

Demonstrație. Procedăm prin inducție asupra numărului n de operatori din e . Dacă $n = 1$ atunci $e = a + b$ sau $e = a * b$ și șirul polonez care satisface condiția din teoremă va fi $p_e = ab +$ sau $p_e = ab *$. Să presupunem că proprietatea este

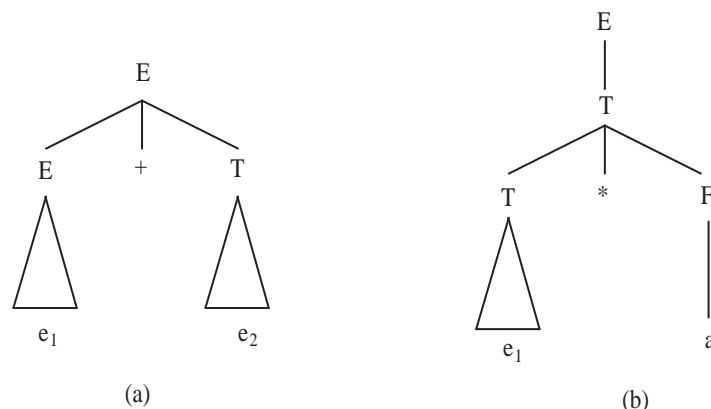


Figura 5.3: Structurile posibile ale arborelui de derivare

adevărată pentru o expresie care conține un număr de operatori mai mic sau egal cu n și considerăm o expresie e cu n operatori. Arborele de derivare al acesteia va avea una din formele prezentate în figura 5.3.

În cazul (b) putem considera $e_2 = a$. În ambele cazuri expresiile e_1 și e_2 au un număr de operanzi inferior lui n și în conformitate cu ipoteza inductivă avem $aes(p_{e_1}) = e_1$ și $aes(p_{e_2}) = e_2$. Luăm $p_e = p_{e_1}p_{e_2} +$ pentru cazul (a) și $p_e = p_{e_1}p_{e_2}*$ pentru cazul (b). Este evident că $aes(p_e) = e$.

5.2 Generarea formatului intermediar

În această secțiune vom analiza problema generării formatului intermediar pentru limbaje algoritmice de tip Fortran, Pascal, Basic. Formele intermediare avute în vedere sunt cvadruplele și notația poloneză, într-un singur caz ne vom referi și la triplete, având în vedere că principal generarea acestora nu diferă esențial de cea a cvadruplelor.

Schema generală a algoritmilor este următoarea.

Fie G gramatica (de tipul 2) care generează limbajul de programare. Prima fază constă în stabilirea unor elemente semantice în concordanță cu structura limbajului; astfel, fiecărui simbol $x \in V_G$ i se atribuie un anumit *sens semantic*, notat $S(x)$, sens care conține o anumită informație depinzând de simbolul respectiv, de contextul în care acesta apare, etc. De asemenea, fiecărei reguli de generare i se atașează o anumită *rutină semantică*. Operațiile obișnuite care se efectuează în aceste rutine sunt: analize specifice, atribuirea unui sens semantic simbolului din stânga, generarea unui cvadruplu sau a unei secvențe de șir polonez, etc. Faza a doua este faza de generare propriu-zisă și ea comportă în principiu următoarele operații: Se efectuează analiza sintactică a tex-

tului sursă cu un algoritm de analiză pe care proiectantul de compilator a decis să îl utilizeze. Indiferent de algoritmul ales, analiza poate fi concepută în două etape, mai întâi construcția arborelui de derivare și apoi reducerea lui succesivă prin eliminarea unor subarbori. În mod natural acești subarbori sunt cei corespunzători frazelor simple stângi așa încât acțiunile care se realizează cu acest prilej sunt cele corespunzătoare rutinelor semantice din prima fază. Menționăm că frontierele subarborilor care se reduc vor constitui, prin definiție, unitățile sintactice ale limbajului. În momentul în care se ajunge la simbolul de start (s-a efectuat întreaga reducere) într-un fișier de ieșire se obține programul în format intermediar, cvadrupele, șir polonez, etc. Se poate observa că analiza sintactică joacă în acest proces un rol central. În cele ce urmează vom presupune că analiza sintactică se efectuează cu un algoritm de tip bottom-up și deci subarborile care se reduce are drept frontieră fraza simplă stângă.

5.2.1 Generarea cvadrupelelor

Gramatica care generează expresiile aritmetice G_E este cunoscută; pentru a ne situa într-un context cât mai apropiat de cazul real, vom completa această gramatică cu operațiile - și / (diferență și împărțire). Prin urmare regulile acestei gramatici vor fi

$$\begin{cases} E \rightarrow E + T | E - T | T \\ T \rightarrow T * F | T / F | F \\ F \rightarrow (E) | a \end{cases}$$

Sensul semantic este prestabilit numai pentru simbolul terminal a , și anume, $S(a)$ este adresa zonei de memorie atașată identificatorului desemnat de a . Pentru celelalte simboluri neterminale, E, T, F sensul semantic se atribuie pe parcursul efectuării analizei. Rutinele semantice corespunzătoare regulilor sunt prezentate în tabelul 5.4. Menționăm că acțiunile prevăzute în aceste rutine sunt stabilite în conformitate cu înțelesul pe care proiectantul dorește să îl atribuie diverselor structuri și nu au o rațiune intrinsecă. Ca o regulă generală însă, putem să facem precizarea că în cazul cvadrupelelor, variabila temporară în care se reține rezultatul intermediar, va fi transmisă ca sens semantic simbolului din stânga.

Exemplu. Considerăm următoarea expresie aritmetică $(a + b) * c - d * (e - f)$. Arborele de derivare corespunzător este prezentat în figura 5.5.

În dreapta fiecărui nod al arborelui (nodurile sunt etichetate cu simboluri ale gramaticii), se află înscris sensul semantic corespunzător variabilei și contextului în care aceasta se află. Procesul de generare al cvadrupelelor se desfășoară astfel: la începutul procesului, după cum se poate vedea pe arborele de generare, fraza simplă stângă este a . Se efectuează reducerea $F \rightarrow a$ și se aplică rutina semantică corespunzătoare acestei reguli (în tabelul rutinelor poziția (8)); singura operație care se face este transmiterea sensului semantic lui F . În continuare, acest sens

<i>Regula</i>	<i>Cvadruplu</i>	<i>Sens</i>
(1) $E_1 \rightarrow E_2 + T$	$+, S(E_2), S(T), t_i$	$S(E_1) = t_i$
(2) $E_1 \rightarrow E_2 - T$	$-, S(E_2), S(T), t_i$	$S(E_1) = t_i$
(3) $E \rightarrow T$		$S(E) = S(T)$
(4) $T_1 \rightarrow T_2 * F$	$*, S(T_2), S(F), t_i$	$S(T_1) = t_i$
(5) $T_1 \rightarrow T_2 / F$	$/, S(T_2), S(F), t_i$	$S(T_1) = t_i$
(6) $T \rightarrow F$		$S(T) = S(F)$
(7) $F \rightarrow (E)$		$S(F) = S(E)$
(8) $F \rightarrow a$		$S(F) = \text{Adr}(a)$

Figura 5.4: Rutine semantice asociate regulilor

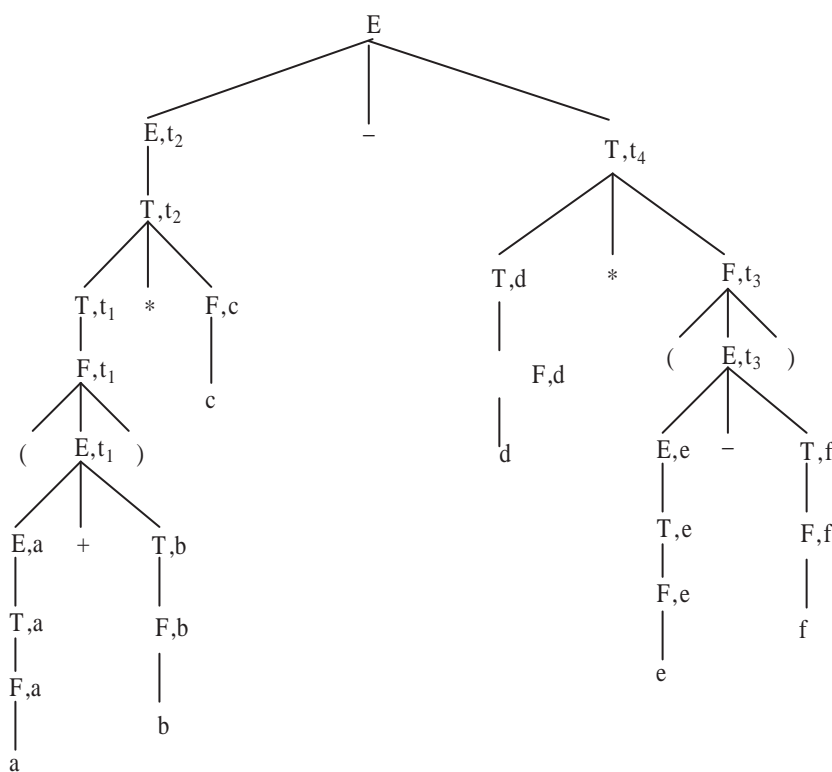


Figura 5.5: Arborele de derivare

<i>Regula</i>	<i>Cvadruplu</i>	<i>Sens</i>
(1) $E_1 \rightarrow E_2 + T$	$+, S(E_2), S(T)$	$S(E_1) = (n)$
(2) $E_1 \rightarrow E_2 - T$	$-, S(E_2), S(T)$	$S(E_1) = (n)$
(3) $E \rightarrow T$		$S(E) = S(T)$
(4) $T_1 \rightarrow T_2 * F$	$*, S(T_2), S(F)$	$S(T_1) = (n)$
(5) $T_1 \rightarrow T_2 / F$	$/, S(T_2), S(F)$	$S(T_1) = (n)$
(6) $T \rightarrow F$		$S(T) = S(F)$
(7) $F \rightarrow (E)$		$S(F) = S(E)$
(8) $F \rightarrow a$		$S(F) = \text{Adr}(a)$

Figura 5.6: Rutine semantice pentru generarea tripletelor

se transmite până la E , apoi fraza simplă stângă este b și se transmite sensul semantic până la T . În acest moment fraza stângă este $E + T$ iar cele două simboluri au sensurile $S(E) = a, S(T) = b$; se efectuează reducerea $E \rightarrow E + T$ și se generează cvadruplul $+, a, b, t_1$ iar sensul semantic care se transmite lui E este t_1 . Procesul continuă în același mod și se obține în final șirul de cvadruple:

$+, a, b, t_1$
 $*, t_1, c, t_2$
 $-, e, f, t_3$
 $*, d, t_3, t_4$
 $-, t_2, t_4, t_5$

5.2.2 Generarea tripletelor

Modificarea ceva mai importantă intervine la rutinele semantice precum și la faptul că după generarea unui triplet, sensul semantic care se transmite va fi numărul de ordine al tripletului generat. Aceste rutine sunt prezentate în tabelul 5.6.

Generarea tripletelor se desfășoară la fel cu cea a cvadruplelor. Pentru expresia considerată la generarea cvadruplelor se obține următorul șir de triplete

(1) $+, a, b$
(2) $*, (1), c$
(3) $-, e, f$
(4) $*, d, (3)$
(5) $-, (2), (4)$

	Regula	Notăția poloneză
(1)	$E \rightarrow E + T$	$p_i = +$
(2)	$E \rightarrow E - T$	$p_i = -$
(3)	$E \rightarrow T$	
(4)	$T \rightarrow T * F$	$p_i = *$
(5)	$T \rightarrow T / F$	$p_i = /$
(6)	$T \rightarrow F$	
(7)	$F \rightarrow (E)$	
(8)	$F \rightarrow a$	$p_i = \text{Adr}(a)$

Figura 5.7: Metoda rutinelor semantice

5.2.3 Generarea șirului polonez

Tehnica rutinelor semantice

După cum s-a precizat, notația poloneză constituie o interesantă și utilă formă intermediară a programelor. În paragraful precedent această notație s-a introdus numai pentru expresii aritmetice fără paranteze, dar această construcție se poate face pentru majoritatea structurilor sintactice ale limbajelor de programare. Desigur că cea mai naturală tehnică de generare a notației poloneze pentru o expresie este procedura generală prezentată la începutul capitolului, bazată pe rutinele semantice atașate regulilor de generare. Spre deosebire de procedura de generare a cvadruplelor și a tripletelor, pentru șirurile poloneze nu mai avem nevoie de sensul semantic al neterminalelor; acestea sunt utilizate în momentul apariției lor și deci nu mai este necesară stocarea unor informații suplimentare. În continuare sunt prezentate rutinele semantice pentru generarea notației poloneze corespunzătoare expresiilor aritmetice cu paranteze (vezi figura 5.7), generate de gramatica G_E . Notația p_i o utilizăm pentru a desemna elementul i din șirul polonez; indicele i se actualizează după fiecare prelucrare a șirului.

Exemplu. Vom aplica rutinele semantice de mai sus pentru generarea șirului polonez corespunzător expresiei aritmetice $a * b + (c - d * e) * f$. Analiza sintactică se efectuează cu un algoritm de tip bottom-up, așa încât trasarea explicită a arborelui de derivare nu mai este necesară. Prelucrarea efectivă se poate organiza într-un tabel ca în figura 5.8; menționăm că fraza simplă stângă este indicată prin sublinierea secvenței corespunzătoare de text.

Tehnica parantezelor

Ideea de bază constă în aceea că se plasează între două paranteze fiecare operator împreună cu cei doi operanzi ai săi; apoi generarea notației poloneze se poate face simplu, folosind o stivă de operatori.

Textul sursă	Regula aplicată	Notăția poloneză
$\underline{a} * b + (c - d * e) * f$	$F \rightarrow a$	a
$F * \underline{b} + (c - d * e) * f$	$F \rightarrow b$	ab
$\underline{F * F} + (c - d * e) * f$	$T \rightarrow T * F, T \rightarrow F$	ab*
$T + (\underline{c} - d * e) * f$	$F \rightarrow c$	ab*c
$T + (F - \underline{d} * e) * f$	$F \rightarrow d$	ab*cd
$T + (F - F * \underline{e}) * f$	$F \rightarrow e$	ab*cde
$T + (F - \underline{F * F}) * f$	$T \rightarrow T * F, T \rightarrow F$	ab*cde*
$T + (\underline{F - T}) * f$	$E \rightarrow E - T, E \rightarrow T, T \rightarrow F$	ab*cde*-
$T + \underline{(E)} * f$	$F \rightarrow (E)$	ab*cde*-
$T + F * \underline{f}$	$F \rightarrow f$	ab*cde*-f
$T + \underline{F * F}$	$T \rightarrow T * F, T \rightarrow F$	ab*cde*-f*
$\underline{T + T}$	$E \rightarrow E + T, E \rightarrow T$	ab*cde*-f*+
E		ab*cde*-f*+

Figura 5.8: Generarea sirului polonez: tehnica analizei sintactice

Plasarea parantezelor. Se atribuie fiecărui operator câte o pondere în conformitate cu convenția de priorități sau cu gramatica care generează expresiile considerate. În cazul unei expresii aritmetice fără paranteze se poate proceda în următorii pași:

- (1) Se atașează fiecărui operator ponderea sa.
- (2) Se parcurge expresia de la stânga la dreapta până când între doi operatori se găsește relația \geq ; de exemplu $\dots a * b + c \dots$ și, cu convenția uzuală, $* \geq +$. Se plasează între paranteze operatorul $*$ împreună cu cei doi operanzi ai săi, adică $\dots (a * b) + c \dots$
- (3) Se șterge ponderea operatorului curent (cel care a fost inclus între paranteze) și se reia procesul de la punctul (2) începând cu operatorul următor (în exemplul nostru de la +).

De exemplu, expresia $a/b * c + d - e$ va avea forma $(((((a/b)*c)+d)-e))$.

În cazul unor expresii aritmetice cu paranteze se folosește un algoritm asemănător în care se modifică pasul (1) și anume: se parcurge expresia de la stânga la dreapta atașând fiecărui operator ponderea sa. La întâlnirea unei paranteze deschise se modifică provizoriu ponderile prin adăugarea unei constante întregi (această constantă trebuie să fie mare decât cea mai mică pondere) și se atribuie în continuare noile ponderi. La întâlnirea unei paranteze închise, această constantă este scăzută, astfel încât operatorii din interiorul parantezelor vor avea ponderi superioare. Apoi toate parantezele se șterg. De exemplu, în expresia

$$(a + b) * (c / (d - e) - f),$$

după efectuarea acestei operații, presupunând că inițial ponderile operatorilor $+$, $-$, $*$, $/$ sunt, respectiv, 1, 2, 3, 4 și constanta de majorare are valoarea 10, operatorilor li se vor atribui ponderile:

a	+	b	*	c	/	d	-	e	-	f
11			3		14		22		12	

Generarea șirului polonez. Pornind de la expresiile aritmetice în care au fost plasate parantezele în conformitate cu prima parte, șirul polonez poate fi construit cu ajutorul algoritmului:

- Se parcurge expresia de la stânga la dreapta, parantezele deschise sunt ignorate și se execută operațiile
 - operanzii se scriu succesiv în șirul polonez;
 - operatorii se scriu într-o stivă.
- La citirea unei paranteze închise, operatorul din vârful stivei se scrie în șirul polonez.

Tehnica stivei de operatori

Este probabil cea mai răspândită metodă de generare a șirului polonez pentru expresii aritmetice. Se utilizează două liste liniare și o stivă. Prima listă conține expresia aritmetică dată iar în cealaltă se obține succesiv șirul polonez corespunzător. Operatorii sunt plasați în stivă ținându-se cont de ponderile lor, apoi sunt extrași și scriși în șir conform cu un algoritm care în esență reproduce acțiunea rutinelor semantice. Mecanismul care realizează acest proces este prezentat în figura 5.9.

Algoritmul este următorul:

- (1) Operanzii se scriu direct în șirul polonez.
- (2) Operatorii și parantezele deschise se scriu în stivă respectându-se următoarele reguli

- Parantezele deschise au pondere 0 și se scriu fără condiții;
- Pentru ceilalți operatori, se compară ponderea operatorului de introdus cu ponderea operatorului din vârful stivei; dacă aceasta este mai mare, introducerea are loc; dacă este mai mică sau egală, se extrage vârful stivei (acesta trece în șirul polonez), se reia comparația cu operatorul din vârful stivei până la prima paranteză deschisă sau până la baza stivei, și apoi se introduce noul operator;

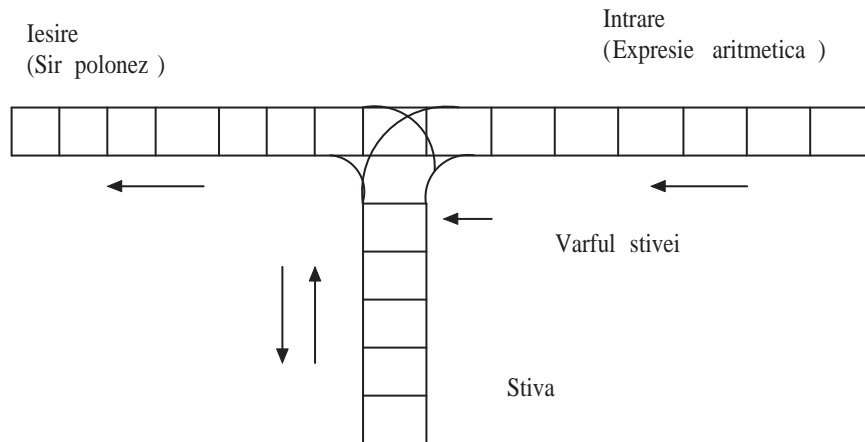


Figura 5.9: Structura dispozitivului

(3) Citirea unei paranteze închise provoacă extragerea operatorilor din stivă până la prima paranteză deschisă; aceasta este extrasă și ea din stivă (dar paranteza extrasă nu trece în șirul polonez) iar în continuare cele două paranteze sunt ignorate.

(4) În finalul procesului, deci după ce ultimul operand a fost trecut în șirul polonez, se extrag toți operatorii din stivă și se scriu în șirul polonez.

Exemplu. Pentru expresia $a * b + (c - d * e) * f$, stare dinamică a stivei de operatori este prezentată în figura 5.10.

5.3 Generarea formatului intermediar pentru instrucții clasice

5.3.1 Instrucțiunea de atribuire

to fill

5.3.2 Instrucțiunea If

to fill

5.3.3 Variabile indexate

to fill

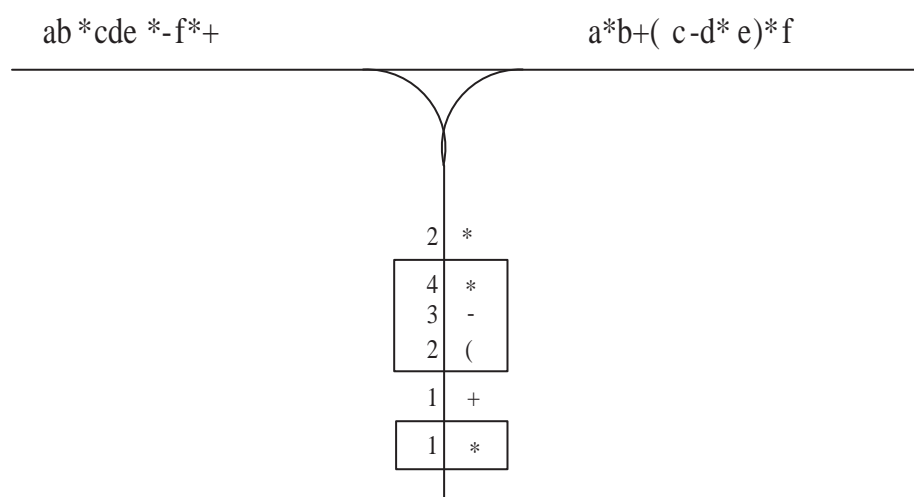


Figura 5.10: Starea dinamică a stivei

Capitolul 6

Masina Turing

6.1 Limbaje de tipul zero

Limbajele de tipul 0 sunt generate de gramatici Chomsky fără restricții asupra regulilor de derivare. Se poate arăta că orice gramatică de tipul zero acceptă următoarea formă normală

$$\begin{aligned} S &\rightarrow SA, \\ B &\rightarrow C|i|\lambda, \\ DE &\rightarrow FH, \end{aligned}$$

unde $A, C, F, H \neq S$. Se observă că singura deosebire față de gramaticile dependente de context este aceea că putem avea și λ -reguli.

Demonstrația existenței formei normale se face prin transformări succesive ale regulilor. În prealabil gramatica se completează cu neterminale noi astfel încât orice regulă $u \rightarrow v$ care nu este de ștergere să respecte condiția de monotonie, adică $|u| \leq |v|$. Acest lucru se poate realiza după cum urmează. Fie

$$X_1 \dots X_n \rightarrow Y_1 \dots Y_m, n > m,$$

o regulă care nu respectă condiția de monotonie. Vom pune

$$X_1 \dots X_n \rightarrow Y_1 \dots Y_m Z_{m+1} \dots Z_n, Z_{m+1} \rightarrow \lambda, \dots, Z_n \rightarrow \lambda.$$

Este clar că orice derivare directă $u \Rightarrow v$ în gramatica inițială se poate obține și în gramatica modificată și reciproc (neterminalele Z nu pot apărea în stânga altor relații). În acest fel, gramatica se va transforma astfel încât să conțină două categorii de reguli: reguli care respectă condiția de monotonie și λ -reguli.

În continuare reducem lungimea părții drepte a regulilor folosind următoarea transformare. Fie $u \rightarrow v \in P, |v| > 2$; atunci $v = Y_1 Y_2 Y_3 v'$.

1. dacă $u = X_1$ atunci regula se înlocuiește cu

$$\begin{aligned} X_1 &\rightarrow Z_1 Z_2, \\ Z_1 &\rightarrow Y_1, \\ Z_2 &\rightarrow Y_2 Y_3 v'. \end{aligned}$$

2. dacă $u = X_1X_2u'$ atunci regula se înlocuiește cu

$$\begin{aligned} X_1X_2 &\rightarrow Z_1Z_2, \\ Z_1 &\rightarrow Y_1, \\ Z_2u' &\rightarrow Y_2Y_3v'. \end{aligned}$$

În ambele cazuri Z_1, Z_2 sunt neterminale noi.

Este clar că noile lungimea părții drepte a regulilor se reduce cu o unitate pentru cazul în care există minim 3 simboluri în dreapta. repetând procedura se poate obține o gramatică echivalentă în care regulile au una din formele:

1. $B \rightarrow C|i|\lambda,$
2. $DE \rightarrow FH,$
3. $X \rightarrow YZ.$

Regulile (1) și (2) satisfac condițiile de la forma normală, iar cele de forma (3) presupunem că nu satisfac aceste condiții, deci că $X \neq S$ sau $Y \neq S$. În general, putem avea mai multe astfel de reguli; pentru simplificare presupunem că există două reguli de forma (3):

$$(3'); \quad X_1 \rightarrow Y_1Z_1, \quad X_2 \rightarrow Y_2Z_2.$$

Construim o gramatică echivalentă $G' = (V_N \cup \{S', \bar{X}_1, \bar{X}_2\}, V_T, S', P')$, unde P' conține toate regulile care convin (de formele (1) și (2)) precum și

$$\begin{aligned} S' &\rightarrow S, \\ S' &\rightarrow S'\bar{X}_1|S'\bar{X}_2, \\ X_1\bar{X}_1 &\rightarrow Y_1Z_1, X_2\bar{X}_2 \rightarrow Y_2Z_2. \end{aligned}$$

De asemenea vom include în P' următoarele reguli de comutare (relativ la simbolurile nou introduse \bar{X}_1, \bar{X}_2):

$$\begin{aligned} A\bar{X}_1 &\rightarrow \bar{X}_1A, \\ A\bar{X}_2 &\rightarrow \bar{X}_2A, \end{aligned} \quad , \forall A \in V_N.$$

Se vede că G' este o gramatică în formă normală. Este trivial să arătăm că $L(G) = L(G')$.

Fie $p \in L(G)$, deci $S \xRightarrow[G]{*} p$. Presupunem că în această derivare există o singură secvență $u \xRightarrow{*} v$ în care se aplică reguli de forma (3'), adică

$$G : \quad S \xRightarrow{*} u \xRightarrow{*} v \xRightarrow{*} p.$$

Presupunem că secvența $u \xRightarrow{*} v$ are forma

$$u = u_1 X_1 u_2 X_2 u_3 X_2 u_4 \xRightarrow{*} u_1 Y_1 Z_1 u_2 Y_2 Z_2 u_2 Y_2 Z_3 u_4 = v,$$

unde $u_1, u_2, u_3, u_4 \in V_N^*$. În G' putem scrie

$$\begin{aligned} S' &\Rightarrow S' \bar{X}_2 \Rightarrow S' \bar{X}_2 \bar{X}_2 \Rightarrow S' \bar{X}_1 \bar{X}_2 \bar{X}_2 \Rightarrow S \bar{X}_1 \bar{X}_2 \bar{X}_2 \\ &\xRightarrow{*} u \bar{X}_1 \bar{X}_2 \bar{X}_2 = u_1 X_1 u_2 X_2 u_3 X_2 u_4 \bar{X}_1 \bar{X}_2 \bar{X}_2 \xRightarrow{*} u_1 X_1 \bar{X}_1 u_2 X_2 \bar{X}_2 u_3 X_2 \bar{X}_2 u_4 \\ &\xRightarrow{*} u_1 Y_1 Z_1 u_2 Y_2 Z_2 u_3 Y_2 Z_2 u_4 = v \xRightarrow{*} p \end{aligned}$$

Prin urmare, $S' \xRightarrow[G']{*} p$ și $p \in L(G')$.

Invers, fie $p \in L(G')$, deci $S' \xRightarrow[G']{*} p$. Dacă în derivarea $S' \xRightarrow[G']{*} p$ apar simbolurile \bar{X}_1, \bar{X}_2 ele nu pot apărea decât în urma aplicării regulilor $S' \rightarrow S' \bar{X}_1$ și $S' \rightarrow S' \bar{X}_2$, deci la începutul derivării, apoi singura posibilitate de continuare este $S' \rightarrow S$; de asemenea, aceste simboluri nu pot fi eliminate decât cu regulile $X_1 \bar{X}_1 \rightarrow Y_1 Z_1$ și $X_2 \bar{X}_2 \rightarrow Y_2 Z_2$, etc. Astfel derivarea noastră trebuie să aibă forma descrisă mai sus, de unde rezultă $S \xRightarrow[G]{*} p, p \in L(G)$. \square

Relativ la puterea generativă a gramaticilor de tipul zero se poate arăta că $\mathcal{L}_1 \subset \mathcal{L}_0$ (strict). Demonstrația se face pe o cale indirectă.

Menționăm și următoarea "teoremă a spațiului de lucru". Fie G o gramatică de tipul zero. Pentru o derivare

$$D : S = u_0 \Rightarrow u_1 \Rightarrow \dots \Rightarrow u_n = p \in V_T^*$$

definim

$$WS_D(p) = \max_{i=1, \dots, n} |u_i|$$

și

$$WS(p) = \min_D \{WS_D(p)\}.$$

Observație. WS este prescurtare de la *working space*.

Spunem că o gramatică de tipul zero G are spațiul de lucru mărginit dacă există $k \in \mathbb{N}$ astfel încât pentru $\forall p \in L(G)$ să avem $WS(p) \leq k|p|$. Să observăm că orice gramatică care nu are reguli de ștergere (excepție $S \rightarrow \lambda$ și atunci S nu apare în dreapta) satisface această condiție cu $k = 1$.

Teorema 6.1 (teorema spațiului de lucru) *dacă o gramatică G are spațiul de lucru mărginit, atunci $L(G) \in \mathcal{L}_1$.*

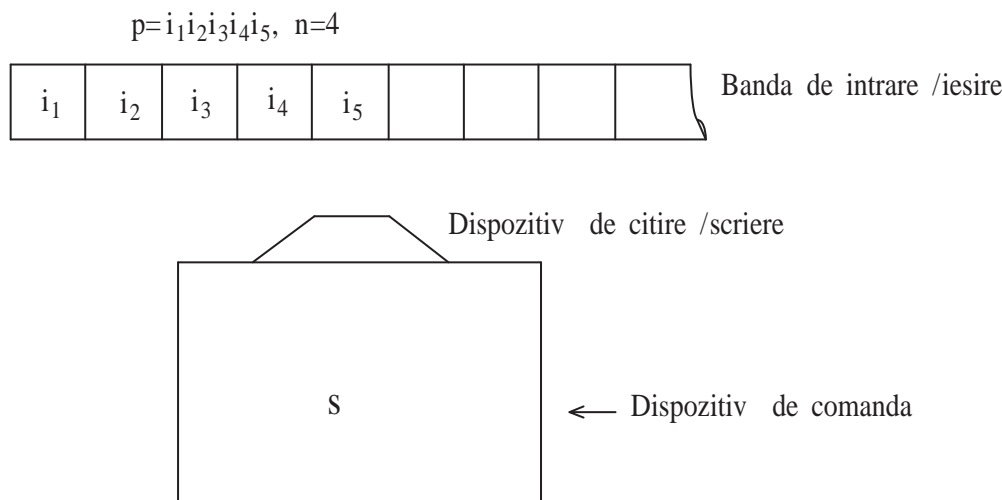


Figura 6.1: Reprezentarea schematică a mașinii Turing

6.2 Mașina Turing

Conceptul de mașină Turing. Mașina Turing este un mecanism de recunoaștere a limbajelor de tipul zero. În felul acesta, familiilor de limbaje din clasificarea Chomsky le corespund automate specifice de recunoaștere.

- Familia \mathcal{L}_3 - automate finite,
- Familia \mathcal{L}_2 - automate push-down,
- Familia \mathcal{L}_1 - automate liniar mărginite,
- Familia \mathcal{L}_0 - mașina Turing.

O mașină Turing (prescurtat MT) se compune dintr-un dispozitiv de comandă care posedă un dispozitiv de scriere-citire și o bandă de intrare. Banda de intrare se consideră mărginită la stânga și nemărginită la dreapta; ea conține simboluri ale unui alfabet de intrare, începând din prima poziție. Alfabetul are un simbol special numit blank și notat b (sau spațiu), care se consideră înregistrat în toată partea neocupată a benzii. Indicele simbolului din dreptul dispozitivului de scriere-citire îl notăm cu n și el va constitui un element al stării mașinii Turing. Dispozitivul de comandă se află într-o anumită stare internă, element al unei mulțimi finite de stări.

Schema unei mașini Turing este dată în figura 6.1. O mașină Turing funcționează în pași discreți. Un pas de funcționare constă din:

Dispozitivul de comandă citește simbolul aflat în dreptul dispozitivului de scriere-citire; în funcție de simbolul citit, mașina Turing trece într-o nouă stare internă, scrie pe bandă un nou simbol (întotdeauna diferit de blank) în locul simbolului citit și mută banda cu o poziție (spre stânga sau dreapta) sau o lasă

pe loc. Convențional, vom nota cu $+1$ o mișcare spre stânga, cu -1 spre dreapta și cu 0 pe loc.

Din punct de vedere matematic, o mașină Turing este un sistem

$$MT = (\Sigma, I, f, s_0, \Sigma_f)$$

unde:

Σ este mulțimea de stări;

I este alfabetul de intrare;

$f : (\Sigma \times I)' \longrightarrow \Sigma \times (I - \{b\}) \times \{1, 0, -1\}$, $(\Sigma \times I)' \subseteq \Sigma \times I$,

este funcția de evoluție;

$s_0 \in \Sigma$ este simbolul de start;

$\Sigma_f \subseteq \Sigma$ este mulțimea de stări finale.

Observație. Funcția f nu este definită pe întregul produs $\Sigma \times I$; dacă mașina ajunge într-o stare s iar în dreptul dispozitivului de scriere-citire se află simbolul i și $(s, i) \notin (\Sigma \times I)'$ spunem că mașina se blochează. Există și alte cazuri de blocare, de exemplu situația în care dispozitivul de scriere-citire se află în dreptul primului simbol, iar pasul de funcționare prevede o mișcare a benzii spre dreapta.

O stare (sau configurație) a unei mașini Turing este un triplet de forma $\sigma = (s, p, n)$ unde s este starea internă, p este cuvântul scris pe bandă iar n este indicele simbolului din dreptul dispozitivului de scriere-citire.

Vom spune că starea $\sigma_1 = (s_1, p_1, n_1)$ *evoluează direct* în starea $\sigma_2 = (s_2, p_2, n_2)$ și vom scrie $\sigma_1 \mapsto \sigma_2$ dacă se efectuează un pas de evoluție. În termenii funcției de evoluție, avem

(1) $f(s_1, i_{n_1}) = (s_2, i, +1)$, $p_2 = i_1 \dots i_{n_1-1} i i_{n_1+1} \dots i_m$, $n_2 = n_1 + 1$;

(2) $f(s_1, i_{n_1}) = (s_2, i, -1)$, $p_2 = i_1 \dots i_{n_1-1} i i_{n_1+1} \dots i_m$, $n_2 = n_1 - 1$;

(3) $f(s_1, i_{n_1}) = (s_2, i, 0)$, $p_2 = i_1 \dots i_{n_1-1} i i_{n_1+1} \dots i_m$, $n_2 = n_1$;

(4) $f(s_1, b) = (s_2, i, +1)$, $p_2 = p_1 i$, $n_2 = n_1 + 1$;

(5) $f(s_1, b) = (s_2, i, -1)$, $p_2 = p_1 i$, $n_2 = n_1 - 1$;

(6) $f(s_1, b) = (s_2, i, 0)$, $p_2 = p_1 i$, $n_2 = n_1$;

Vom spune că σ' *evoluează* (fără specificația direct) în σ'' și vom nota $\sigma' \xrightarrow{*} \sigma''$ dacă $\sigma' = \sigma''$ sau dacă există $\sigma_1, \dots, \sigma_n$ astfel încât

$$\sigma' = \sigma_1 \mapsto \sigma_2 \mapsto \dots \mapsto \sigma_n = \sigma''.$$

Limbajul recunoscut de o mașină Turing este prin definiție

$$L(MT) = \{p | p \in I^*, (s_0, p, 1) \xrightarrow{*} (s, q, \epsilon), s \in \Sigma_f\}.$$

Observație. Este posibil ca mașina să ajungă într-o stare finală înainte de citirea integrală a lui p ; analiza stării finale trebuie făcută numai după parcurgerea cuvântului.

Exemplu. Considerăm mașina turing $MT = (\Sigma, I, f, s_0, \Sigma_f)$ unde $\Sigma = \{s_0, s_1, s_2\}$, $I = \{0, 1\}$, $\Sigma_f = \{s_1\}$ iar funcția de evoluție este dată de

f	s_0	s_1	s_2
b	$(s_2, 1, 1)$		$(s_0, 0, -1)$
0	$(s_1, 0, 1)$		
1	$(s_0, 1, 1)$	$(s_0, 0, 1)$	

Evoluția mașinii pentru $p = 001$ este

$$(s_0, 011, 1) \mapsto (s_1, 011, 2) \mapsto (s_0, 001, 3) \mapsto (s_0, 001, 4) \mapsto \\ (s_2, 0011, 5) \mapsto (s_0, 0011, 4) \mapsto (s_0, 00110, 5) \mapsto (s_1, 00110, 6)$$

Se poate observa că după citirea întregului cuvânt mașina poate să efectueze un număr de pași suplimentari până la ajungerea într-o stare finală, deci este posibil ca $|p| < |q|$.

Situații în care o mașină Turing se blochează (în aceste situații cuvântul scris pe bandă nu este recunoscut):

1. MT ajunge într-o stare s , în dreptul dispozitivului de citire-scriere se află simbolul i și $(s, i) \notin (\Sigma \times I)'$;
2. MT este în starea s , a citit simbolul din prima poziție i și $f(s, i) = (s', i', -1)$;
3. MT efectuează un ciclu infinit în interiorul cuvântului.

Definiție 6.1 Vom spune că o mașină Turing este nestaționară dacă

$$f : (\Sigma \times I)' \rightarrow \Sigma \times (I \setminus \{b\}) \times \{-1, +1\}.$$

Prin urmare, o mașină Turing nestaționară nu lasă în nici o situație banda pe loc.

Lema 6.1 Orice mașină Turing este echivalentă cu o mașină Turing nestaționară.

Demonstrație. Pornind de la o MT dată construim o MT' nestaționară astfel:

Dacă $f(s, i) = (s', i', \pm 1)$ vom pune $f'(s, i) = f(s, i)$;

Dacă $f(s, i) = (s', i', 0)$ vom pune $f'(s, i) = (s'', i', -1)$ și $f'(s'', j) = (s', j, +1), \forall j \in I$, unde s'' este o stare nouă introdusă.

Astfel, în cazul unei rămânări pe loc a lui MT, noua mașină va face un pas spre stânga și unul spre dreapta, fără să modifice nici starea și nici conținutul benzii. Evident, cele două mașini sunt echivalente.

Limbajele recunoscute de mașini Turing.

Teorema 6.2 Un limbaj este recunoscut de o mașină Turing dacă și numai dacă este de tipul zero.

Demonstrație. Partea I. $E = L(MT) \Rightarrow E \in \mathcal{L}_0$.

Fie $MT = (\Sigma, I, f, s_0, \Sigma_f)$ o mașină Turing astfel încât $E = L(MT)$. Putem presupune că MT este nestaționară. Fie $I_\lambda = (I \setminus \{b\}) \cup \{\lambda\}$.

Construim o gramatică de tipul zero astfel: $G = (V_N, V_T, S, P)$ unde

$$V_N = \Sigma \cup \{(i, j) | i \in I_\lambda, j \in I\} \cup \{S, X_1, X_2\}; \quad V_T = I \setminus \{b\}.$$

Definiția lui P :

$$S \rightarrow s_0 X_1, \quad X_1 \rightarrow (i, i) X_1, \quad i \in I \setminus \{b\},$$

$$X_1 \rightarrow X_2, \quad X_2 \rightarrow (\lambda, b) X_2, \quad X_2 \rightarrow \lambda,$$

$$\text{dacă } f(s, i) = (s', i', 1) \text{ atunci } s(i_1, i) \rightarrow (i_1, i') s', \quad \forall i_1 \in I_\lambda,$$

$$\text{dacă } f(s, i) = (s', i', -1) \text{ atunci } (i_1, i_2) s(i_3, i) \rightarrow s'(i_1, i_2) (i_3, i'), \quad i_1, i_3 \in I_\lambda, \quad i_2 \in I,$$

$$\text{dacă } s \in \Sigma_f, \text{ atunci } s(i_1, i_2) \rightarrow s i_1 s \text{ și } (i_1, i_2) s \rightarrow s i_1 s.$$

Fie $p \in L(MT)$, $p = i_1 \dots i_n$. Presupunem că MT utilizează în recunoaștere m poziții de pe bandă situate la dreapta cuvântului p . În G avem

$$S \Rightarrow s_0 X_1 \xrightarrow{*} s_0(i_1, i_1) \dots (i_n, i_n) X_2 \xrightarrow{*} s_0(i_1, i_1) \dots (i_n, i_n) (\lambda, b)^m.$$

Cuvântul p fiind recunoscut de MT avem

$$(1) \quad (s_0, i_1 \dots i_n, 1) \xrightarrow{*} (s_f, i'_1 \dots i'_h, k), \quad h \geq n.$$

Vom arăta că (1) implică existența unei derivări de forma

$$2) s_0(i_1, i_1) \dots (i_n, i_n) (\lambda, b)^m \xrightarrow{*} (i_1, i'_1) \dots (i_{k-1}, i'_{k-1}) s_f(i_k, i'_k) \dots (i_{n+m}, i'_{n+m}),$$

unde

$$\begin{aligned} i_1, \dots, i_n &\in I \setminus \{b\}, \quad i_{n+1}, \dots, i_{n+m} = \lambda, \\ i'_1, \dots, i'_h &\in I \setminus \{b\}, \quad i'_{h+1}, \dots, i'_{n+m} = \lambda. \end{aligned}$$

Demonstrație prin inducție asupra lungimii l a evoluției.

Pentru $l = 0$ avem

$$(s_0, i_1 \dots i_n, 1) \xrightarrow{*} (s_0, i_1 \dots i_n, 1), \quad k = 1, h = n, i'_j = i_j.$$

Partea dreaptă a lui 2) va avea forma $s_0(i_1, i_1) \dots (i_n, i_n) (\lambda, b)^m$ și deci 2) este adevărată.

Presupunem că implicația este adevărată pentru l oarecare și considerăm o evoluție de lungime $l + 1$. Punem în evidență ultima evoluție directă

$$(s_0, i_1 \dots i_n, 1) \xrightarrow{*} (s, i''_1 \dots i''_g, j) \mapsto (s_f, i'_1 \dots i'_h, k).$$

În general $h = g$ sau $h = g + 1$ în conformitate cu următoarele două cazuri (figura 6.2).

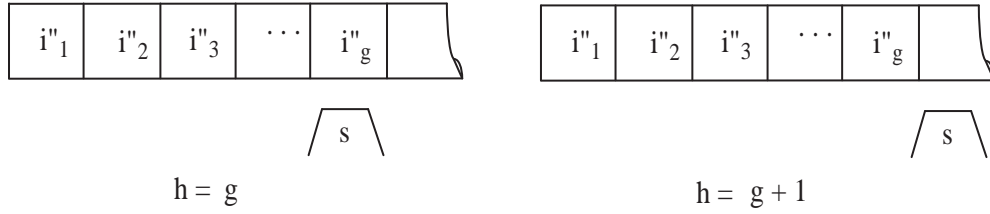


Figura 6.2: Configurații posibile ale mașinii Turing

Întotdeauna $k = j \pm 1$. apoi $i''_t = i'_t$, $t = 1, \dots, g$, $t \neq j$, adică

$$\begin{array}{c} i''_1, i''_2, \dots, i''_{j-1}, i''_j, i''_{j+1}, \dots, i''_g; \\ i'_1, i'_2, \dots, i'_{j-1}, i'_j, i'_{j+1}, \dots, i'_g. \end{array}$$

În urma ultimei evoluții directe vor diferi numai simbolurile i''_j, i'_j .

Din ipoteza inductivă rezultă

$$s_0(i_1, i_1) \dots (i_n, i_n)(\lambda, b)^m \xrightarrow{*} (i_1, i''_1) \dots (i_{j-1}, i''_{j-1}) s(i_j, i''_j) \dots (i_{n+m}, i''_{n+m}).$$

În conformitate cu definiția evoluției directe avem

$$\begin{array}{l} f(s, i''_j) = (s_f, i'_j, 1), \quad k = j + 1, \quad s(i_j, i''_j) \rightarrow (i_j, i'_j) s_f; \\ f(s, i'_j) = (s_f, i'_j, -1), \quad k = j - 1, \quad (i_{j-1}, i''_{j-1}) s(i_j, i'_j) \rightarrow s_f(i_j, i''_{j-1})(i_j, i'_j). \end{array}$$

În ambele cazuri

$$s_0(i_1, i_1) \dots (i_n, i_n)(\lambda, b)^m \xrightarrow{*} (i_1, i'_1) \dots (i_{k-1}, i'_{k-1}) s_f(i_k, i'_k) \dots (i_{n+m}, i'_{n+m}).$$

Acum, deoarece $s_f \in \Sigma_f$, putem scrie

$$\begin{array}{l} S \xrightarrow{*} s_0(i_1, i_1) \dots (i_n, i_n)(\lambda, b)^m \\ \xrightarrow{*} (i_1, i'_1) \dots (i_{k-1}, i'_{k-1}) s_f(i_k, i'_k) \dots (i_{m+n}, i'_{m+n}) \\ \Rightarrow (i_1, i'_1) \dots (i_{k-1}, i'_{k-1}) s_f i_k s_f(i_{k+1}, i'_{k+1}) \dots (i_{m+n}, i'_{m+n}) \\ \Rightarrow (i_1, i'_1) \dots (i_{k-2}, i'_{k-2}) s_f i_{k-1} s_f i_k s_f i_{k+1} s_f(i_{k+2}, i'_{k+2}) \dots (i_{m+n}, i'_{m+n}) \\ \xrightarrow{*} s_f i_1 s_f i_2 s_f \dots s_f i_n s_f \xrightarrow{*} i_1 \dots i_n = p. \end{array}$$

Prin urmare $p \in L(G)$ și $L(MT) \subseteq L(G)$. Analog se arată și incluziunea inversă și deci $L(MT) = L(G)$. \square

Bibliografie

1. Octavian C. Dogaru, *Bazele informaticii. Limbaje formale*, Tipografia Universității din Timișoara, 1989.

2. Gheorghe Grigoraș, *Limbae formale și tehnici de compilare*, Tipografia Universității "Alexandru Ioan Cuza", Iași, 1984.
3. J. E. Hopcroft și J. D. Ullman, *Introduction to Automata Theory, Languages and Computation*, Reading Mass., 1979.
4. Solomon Marcus, *Gramatici și automate finite*, Editura Academiei, București, 1964.
5. Ștefan Mărușter, *Curs de Limbae formale și tehnici de compilare*, Tipografia Universității din Timișoara, 1980.
6. Gheorghe Orman, *Limbae formale*, Editura tehnică, București, 1982.
7. Gheorghe Păun, *Probleme actuale în teoria Limbajelor formale*, Editura Științifică și Enciclopedică, București, 1984.
8. Teodor Rus, *Mecanisme formale pentru specificarea limbajelor*, Editura Academiei, București, 1983.
9. Arto Salomaa, *Formal languages*, Academic Press, New York, 1973.
10. Dan Simovici, *Limbae formale și tehnici de compilare*, Editura didactică și pedagogică, București, 1978.
11. Luca Dan Șerbănați, *Limbae de programare și compilatoare*, Editura Academiei, București, 1987.