

IRINA ATHANASIU

DIANA RAICIU RADU SION IRINA MOCANII

**Universitatea "Politehnica" din București
Catedra de Calculatoare**

LIMBAJE FORMALE

Și

AUTOMATE

(îndrumar pentru aplicații)

©MATRIX ROM
CP. 16 - 162
77500 - BUCUREȘTI
tel. 01.4113617, fax 01.4114280
e-mail: matrix@fx.ro

Despre autori,

Acest îndrumar are mulți autori cu o contribuție secvențială. A început cu niște foi scrise de mână care conțineau probleme și schițe de rezolvări, a mai existat și un text introductiv în lex în format electronic. În 1992 Diana Raiciu (în prezent Diana Mărculescu, profesor în Department of Electrical and Computer Engineering din Universitatea Carnegie Mellon, SUA) a realizat prima formă electronică a acestui îndrumar. Următorul autor (1999) a fost Radu Sion (în prezent doctorand în departamentul de calculatoare din Universitatea Purdue, SUA). A urmat în 2001 doamna asistenta Irina Mocanu. Fiecare dintre autori a corectat textul anterior și a mai adăugat probleme noi. Cu fiecare nou autor textul părea că devine gata. Am hotărât să public textul în formatul curent care desigur este încă departe de ce putea să fie, pentru că se aniversează anul acesta 10 ani de când credeam că îndrumarul va fi foarte curând gata de publicare.

Descrierea CIP a Bibliotecii Naționale a României

Limbaje formale și automate. îndrumar pentru aplicații/ Irina Athanasiu,
Diana Raiciu, Radu Sion, Irina Mocanu, București, Matrix Rom, 2002

98 pagini, 25 cm

Bibliogr.

ISBN 973-685-407-8

Athanasiu, Irina

Raiciu, Diana

Sion, Radu

Mocanu, Irina

IV.
004.43

Martie 2002

Irina Athanasiu

ISBN 973 - 685 - 407 - 8

/ Elemente de teoria limbajelor formale		2
1.1	Gramatici	2
1.1.1	Ierarhia Chomsky	3
	Probleme	4
1.1.2	Lema de pompare	16
	Probleme	17
1.1.3	Transformări asupra GIC	18
1.1.3.1	Eliminare recursivitate stânga	18
1.1.3.2	Eliminare X producții	18
1.1.3.3	Eliminare simbolii neutilizați	19
1.1.3.4	Eliminare simbolii inaccesibili	20
	Probleme	20
1.2	Mulțimi regulate. Expresii regulate	23
	Probleme	24
1.3	Acceptoare	27
1.3.1	Automate finite	27
	Probleme	28
1.3.2	Automate cu stivă (push-down)	44
	Probleme	45
1.3.3	Mașina Turing	52
	Probleme	54
2 Lex - generator de analizare lexicală		65
2.1	Expresii regulate. Structura unei specificații lex.	66
2.2	Elemente avansate	70
2.2.1	Funcționarea analizorului lexical generat de lex	71
2.2.2	Stări de start	72
2.2.3	Macrodefiniții, funcții și variabile predefinite	73
2.2.4	Fișiere de intrare multiple	73
2.3	Exemple comentate	74
3 Teme pentru acasă		87
4 Bibliografie		91

1 Elemente de teoria limbajelor formale

Definiția 1.1. Se numește *alfabet* orice mulțime finită T de simbolii.

Definiția 1.2. Se numește *limbaj peste un alfabet T* orice submulțime a mulțimii T^* .

1.1 Gramatici

Definiția 1.1.1. O *gramatică* este un cvadruplu $G = (N, T, P, S)$ unde:

- N este o mulțime finită de simbolii numiți simbolii neterminalii
- T este o mulțime finită de simbolii numiți terminale ($N \cap T = \emptyset$)
- P este o submulțime a mulțimii $(N \cup T)^* N (N \cup T)^* x (N \cup T)^*$ și reprezintă mulțimea producțiilor gramaticii G . Un element $p \in P$, $p = (a, P)$ se reprezintă sub forma: $oc \rightarrow p$
- $S \in N$ se numește simbolul de start al gramaticii

Definiția 1.1.2. Se numește *formă propozițională în gramatica G* orice șir din $(N \cup T)^*$ obținut prin aplicarea recursivă a următoarelor reguli:

1. S este o formă propozițională;
2. Dacă $ocpy$ este o formă propozițională și există o producție $p \rightarrow \delta$ atunci și $a\delta y$ este o formă propozițională.

Definiția 1.1.3. Se numește *propoziție generată de gramatica G* o formă propozițională în G care conține numai terminale.

Definiția 1.1.4. Se spune că între două forme propoziționale a și p în gramatica G există *relația de derivare* notată cu $a \Rightarrow p$ dacă există două șiruri w_1, w_2 și o producție $y \rightarrow \delta$ astfel încât $a = w_1 y w_2$ și $p = w_1 \delta w_2$. Se spune că între formele propoziționale a și P există relația de derivare notată $a \Rightarrow^* P$ dacă există formele propoziționale $\delta_0, \delta_1, \dots, \delta_k$ astfel încât $a = \delta_0 \Rightarrow \delta_1 \Rightarrow \dots \Rightarrow \delta_k = p$. Închiderea tranzitivă a relației de derivare se notează cu \Rightarrow^* .

Definiția 1.1.5. Fie G o gramatică. Se numește *limbaj generat de gramatică G* și se notează $L(G)$ mulțimea propozițiilor generate de G . Altfel spus:

$$L(G) = \{w \in T^* \mid S \Rightarrow^* w\}$$

Definiția 1.1.6. Două gramatici G_1, G_2 se numesc *echivalente* dacă generează același limbaj adică $L(G_1) = L(G_2)$.

1.1.1 Ierarhia Chomsky

Gramaticile pot să fie clasificate conform complexității producțiilor în următoarea ierarhie:

- gramatici de tip 0 (fără restricții) - au producțiile de forma:
 $a \rightarrow P$ cu $a \in (N \cup T)^* N (N \cup T)^*$, $P \in (N \cup T)^*$
- gramatici de tip 1 (dependente de context - GDC) - au producțiile de forma:
 $a A P \rightarrow a y p$, $a, P \in (N \cup T)^*$, $A \in N$, $y \in (N \cup T)^+$
sau de forma
 $S \rightarrow k$.

În al doilea caz S nu apare în membrul drept al nici unei producții.

- gramatici de tip 2 (independente de context - GIC) au producțiile de forma:
 $A \rightarrow a$, $A \in N$, $a \in (N \cup T)^*$.
- gramatici de tip 3 (regulate la dreapta - GR) - au producțiile de forma:
 $A \rightarrow H \delta c B c u A e N$, $B \in (N \cup \{A\})$ și $a \in T^*$.

Corespunzător mulțimii gramaticilor $G[k]$, $k = 0, 1, 2, 3$ există mulțimea $L[k]$ care reprezintă mulțimea limbajelor ce pot fi generate de o gramatică din $G[k]$. $L[0]$ este mulțimea limbajelor generale care pot fi generate de o gramatică de tip 0 (fără restricții), $L[1]$ este mulțimea limbajelor dependente de context (generate de o GDC), $L[2]$ mulțimea limbajelor independente de context (generate de o GIC), iar $L[3]$ mulțimea limbajelor regulate (generate de o GR). Se poate arăta că are loc incluziunea:

$$L[3] \subset L[2] \subset L[1] \subset L[0].$$

Termenul de "dependent de context" provine din faptul că producția $ocAp \rightarrow ayP$, poate să fie interpretată ca - neterminalul A poate să fie înlocuit cu șirul y doar în contextul dat de șirurile a și p . Similar termenul "independent de context" provine din faptul că o producție de forma $A \rightarrow a$ se poate interpreta ca - neterminalul A poate să fie înlocuit cu șirul a indiferent de context (indiferent de simbolii între care apare).

Se observă că într-o gramatică independentă de context prin aplicarea unei producții de formă $A \rightarrow a$ asupra unei forme propoziționale, lungimea acesteia poate să devină mai mică (dacă a este șirul vid), egală (dacă a este un terminal sau neterminal) sau mai mare (dacă a este un șir de lungime mai mare decât 1). Acest lucru nu este valabil pentru gramaticile dependente de context unde lungimea unei forme propoziționale în urma aplicării unei producții de forma

$$a A P \rightarrow a y P, a, P \in (N \cup T)^*, A \in N, y \in (N \cup T)^+$$

devine mai mare (dacă y este un șir de lungime mai mare decât 1), cel mult egală (dacă y este un șir format dintr-un neterminal sau un terminal). Prin urmare, nu este evident că orice gramatică independentă de context este și dependentă de context. Se poate demonstra că pentru orice gramatică independentă de context există o altă gramatică independentă de context echivalentă obținută prin transformări, astfel încât modul de creștere a lungimii formelor propoziționale să fie similară cu cea a gramaticilor dependente de context.

Intuitiv, o caracterizare a gramaticilor de mai sus poate fi făcută din punctul de vedere al "numărării" simbolilor. Astfel, gramaticile regulate pot "număra" până la o limită finită. De exemplu, limbajul specificat prin $L = \{0^i T \mid 0 < i < 5\}$ poate fi generat de o gramatică regulată

problema 1.1-5

$$L = \{ a^n b^m c^m d^n \mid n > 1 \text{ și } m \geq 1 \}$$

Soluție:

Orice șir w din L începe cu simbolul a , și se termină cu simbolul d . w se poate scrie ca $w = a u d$ unde u este un șir care fie este de tipul celui din mulțimea L , fie este de forma $b^m c^m$, $m \geq 1$. Acesta din urmă se poate scrie sub forma $b v c$ cu v de aceeași formă sau v este șirul vid. Astfel, ținând seama de observațiile anterioare, o gramatică ce generează limbajul L este: $G = (\{ S, A \}, \{ a, b, c, d \}, P, S)$ unde $P = \{ S \rightarrow a S d \mid a A d, A \rightarrow b A c \mid b c \}$.

Verificare:

Pentru șirul $a a b c d d$ ($n = 2$ și $m = 1$) se obține următoarea secvență de derivări: $a S \Rightarrow a a A d \Rightarrow a a b c d d$.

Comentarii:

Se observă că neterminalele A este utilizat pentru a genera subșiruri de forma $b^m c^m$. Gramatica G este independentă de context.

problema 1.1-6

$$L = \{ w \in \{0,1\}^* \mid w \text{ conține un număr egal de } 0 \text{ și } 1 \}$$

Soluție:

Fie un șir oarecare astfel încât $w \in \{0,1\}^*$ și $\#_0(w) = \#_1(w)$ unde prin $\#_a(w)$ am notat numărul de apariții ale simbolului a în șirul w . Dacă w începe cu simbolul 0 , atunci există u și v , $u, v \in \{0,1\}^*$ cu u și v având fiecare un număr egal de 0 și 1 (eventual zero) astfel încât șirul w se poate scrie $w = 0 u 1 v$, cu $u, v \in \{0,1\}^*$ cu număr egal de 0 și 1 astfel încât $w = 0 u 1 v$. Producția care generează acest tip de șiruri este: $S \rightarrow 0 S 1 S$. De asemenea, dacă $w \in L$ începe cu 1 , se arată în mod analog că $w = 1 u 0 v c u u$, v șiruri cu număr egal de 0 și 1 . Pe baza celor de mai sus, rezultă că o gramatică care generează limbajul L este $G = (\{ S \}, \{0,1\}, P, S)$ unde: $P = \{ S \rightarrow 0 S 1 S \mid 1 S 0 S \mid A \}$.

Verificare:

Pentru șirul $0 1 0 1 1 0$, se obține următoarea secvență de derivări: $S \Rightarrow 0 S 1 S \Rightarrow 0 1 S 0 S \Rightarrow 0 1 0 S 1 S \Rightarrow 0 1 0 1 S \Rightarrow 0 1 0 1 1 S 0 S \Rightarrow 0 1 0 1 1 0 S \Rightarrow 0 1 0 1 1 0$.

Comentarii:

Gramatica G este o gramatică de independentă de context.

problema 1.1-7

$$L = \{ w \in \{0,1\}^* \mid w \text{ nu conține subșirul } 011 \}$$

Soluție:

Presupunem că parcurgem un șir din limbajul L începând cu simbolul cel mai din stânga. Se observă că pot să existe mai multe situații în ceea ce privește simbolul curent:

1. Simbolul curent este 1 și cel anterior a fost tot 1 . Atâta timp cât simbolul curent este 1 nu suntem la un început al unui sub șir de forma 011 .
2. Simbolul analizat este 0 adică suntem la începutul unei posibile secvențe 011 care trebuie rejectată. Atâta timp cât simbolul curent este 0 este aceeași situație. Dacă simbolul curent este 1 se trece în situația următoare (3).
3. Simbolul curent este 1 și cel precedent a fost 0 , adică suntem în interiorul unei posibile secvențe 011 care trebuie rejectată. De aici, singurul simbol posibil este 0 (1 ar genera secvența 011), care ar putea să înceapă o nouă posibilă secvență 011 . Ca urmare, pentru un simbol 1 suntem în situația 2.

Vom codifica cele trei situații posibile prin neterminalele S, A, B . Corespunzător celor de mai sus o gramatică G care generează limbajul L este: $G = (\{ S, A, B \}, \{0,1\}, P, S)$, $P = \{ S \rightarrow \hat{1} S \mid 0 A \mid X, A \rightarrow 0 A 1 B \mid X, B \rightarrow 0 A \mid X \}$.

Verificare:

Pentru șirurile: $10 1$, respectiv $0 10 1$, secvențele de derivări sunt: $S \Rightarrow 1 S \Rightarrow 1 0 A \Rightarrow 1 0 1 B \Rightarrow 1 0 1$ și $S \Rightarrow 0 A \Rightarrow 0 1 B \Rightarrow 0 1 0 A \Rightarrow 0 1 0 1 B \Rightarrow 0 1 0 1$.

Comentarii:

Gramatică G este regulată, iar limbajul L este de un limbaj regulat.

pick-iis.i 1.1-K

$$L = \{ w \in \{0,1\}^* \mid |w| \text{ divizibilă cu } 3 \}$$

Soluție:

Vom presupune că parcurgem șirul w de la stânga la dreapta considerând un prefix al acestuia: $w = u a v$ unde u este prefixul, a este simbolul curent (0 sau 1), iar v este restul șirului ($u, v \in \{0,1\}^*$). Sunt posibile următoarele situații:

1. Numărul de simboluri din prefixul u este de forma $3 \cdot n$, $n \in \mathbb{N}$
2. Numărul de simboluri din prefixul u este de forma $3 \cdot n + 1$
3. Numărul de simboluri din prefixul u este de forma $3 \cdot n + 2$

Asociem neterminalele S primei stări (care este de altfel cea corespunzătoare generării unui șir din limbajul L), iar celorlalte două stări neterminalele A , respectiv B . Din prima stare se trece în a doua după ce se "citește" simbolul curent (numărul de simboluri citite devine $3 \cdot n + 1$), iar de aici, după "citirea" unui nou simbol, se trece în starea 3 și apoi din nou în starea 1 după "citirea" următorului simbol din șir, etc. Gramatica care generează limbajul L este: $G = (\{ S, A, B \}, \{0,1\}, P, S)$ unde $P = \{ S \rightarrow \bullet 0 A \mid 1 A \mid X, A \rightarrow 0 B 1 B, B \rightarrow 0 S \mid \hat{1} S \}$.

Verificare:

Pentru șirul 010011 , se obține secvența de derivări $S \Rightarrow 0 A \Rightarrow 0 1 B \Rightarrow 0 1 0 S \Rightarrow 0 1 0 0 A \Rightarrow 0 1 0 0 1 B \Rightarrow 0 1 0 0 1 1 S \Rightarrow 0 1 0 0 1 1$.

Comentarii:

Altă gramatică regulată echivalentă este $G_1 = (\{ S \}, \{0,1\}, P, S)$ unde: $P = \{ S \rightarrow 0 0 0 S \mid 0 0 1 S \mid 1 0 1 S \mid 1 0 1 1 S \mid 1 0 1 1 0 S \mid 1 0 1 1 1 S \mid X \}$.

1.1.1

Soluție:

Șirurile din L pot fi de două tipuri: cu $i > j$, respectiv $i < j$. Astfel, un șir $w \in L$ poate fi scris fie sub forma $w = 0^i 1^j$ dacă $i > j$ și $n = i - j$, fie sub forma $w = 0^i 1^j$ dacă $i < j$ și $n = j - i$. În ambele cazuri, $n > 0$. Pentru generarea primului șir, vom observa că acesta începe cu 0 și se termină cu 1 și poate fi scris sub forma $w = 0^i 1^j$ unde u este fie de aceeași formă, fie de formă 0^n , cu $n > 0$. În mod analog, cel de-al doilea șir începe cu 0 și se termină cu 1 și poate fi scris sub formă $w = 0^i 1^j$ unde u este fie un șir de aceeași formă, fie de forma 1^n , cu $n > 0$. Ținând cont de observațiile de mai sus, o gramatică ce generează limbajul L este $G = (\{ S, A, B \}, \{ 0, 1 \}, P, S)$ unde $P = \{ S \rightarrow 0S1 \mid 0A \mid B, A \rightarrow 0A \mid A., B \rightarrow IB \mid A. \}$.

Verificare:

Pentru șirul 0001 se obține următoarea secvență de derivări: $S \Rightarrow 0S1 \Rightarrow 00A1 \Rightarrow 000A1 \Rightarrow 00001$.

Comentarii:

G este o gramatică independentă de context.

Problema 1.1-10

$$L = \{ x^m y^n \mid n < m \text{ sau } 2^* m < n, n, m > 1 \}$$

Soluție:

Limbajul L conține două tipuri de șiruri: șiruri de forma $x^m y^n$ cu $2^* m < n$, respectiv cu $n < m$. Primul tip de șir poate să fie scris sub forma $x^m y^{2^* m + r}$ unde $n = 2^* m + r$, cu $r > 0$ și $m > 1$; iar pentru generarea lui pot să fie utilizate producțiile: $S_j \rightarrow x A y y, A \rightarrow x A y y \mid y C \text{ și } C_i \rightarrow y C \mid A$. Al doilea tip de șir poate să fie scris sub forma $x^n x^m$ unde $m = n + r$ cu $r > 0$ și $n > 1$. Pentru generarea acestui șir se pot utiliza producțiile: $S_2 \rightarrow x B y, B \rightarrow x B y \mid x D, D \rightarrow x D \mid X$. Prin urmare, o gramatică care generează limbajul L este: $G = (\{ S, S_1, S_2, A, B, C, D \}, \{ x, y \}, P, S)$ unde $P = \{ S \rightarrow S_1 \mid S_2, S_1 \rightarrow x A y y, A \rightarrow x A y y \mid y C, S_2 \rightarrow x B y, B \rightarrow x B y \mid x D, C \rightarrow y C \mid A, D \rightarrow x D \mid A \}$.

Verificare:

Pentru șirul $x x x y y$ (pentru care $m = 3$ și $n = 2$, adică $m > n$) se obține secvența de derivări: $S \Rightarrow S_1 \Rightarrow x B y \Rightarrow x x B y \Rightarrow x x x D y \Rightarrow x x x y y$. Fie șirul $x x y y y y$ (pentru care $m = 2$ și $n = 5$ adică $2^* m < n$). Pentru acest șir se obține secvența de derivări: $S \Rightarrow S_2 \Rightarrow x A y y \Rightarrow x x A y y y \Rightarrow x x y C y y y y \Rightarrow x x y y y y y$.

Comentarii:

G este o gramatică independentă de context

1.1.1

$$L = \{ a^m b^n c^p d^q \mid m + n = p + q, m, n, p, q > 0 \}$$

Soluție:

Relația din enunț poate fi scrisă sub forma $m - q = p - n$. Pot apare două cazuri:

- $m > q$. Notăm $r = m - q = p - n$, unde $r > 0$ și $m = q + r$ și $p = n + r$. Fie L_1 limbajul generat în acest caz. Atunci L_1 poate fi descris prin: $L_1 = \{ a^q a^r b^n c^n c^r d^q \mid n, q > 0, r > 0 \}$. Orice șir de aceasta formă, poate să fie scris ca $w = a u d$ unde u este de aceeași formă cu w sau de forma $a^r b^n c^n c^r$ și în acest caz poate fi scris sub forma: $u = a v c$ unde v este fie de aceeași formă cu u , fie este de forma $b^n c^n$. În acest caz, v poate fi scris sub formă $v = b z c$ unde șirul z este de aceeași formă cu v , fie șirul vid. Prin urmare, vor fi utilizate producțiile: $S \rightarrow a S d \mid A, A \rightarrow a A c \mid B, B \rightarrow b B c \mid A$.
- $m < q$. Notăm $r = q - m = n - p$, unde $r > 0$ și prin urmare $n = p + r$ și $q = m + r$. Fie L_2 limbajul generat în acest caz. Prin urmare $L_2 = \{ a^m b^p b^r c^p d^m \mid m, p > 0, r > 0 \}$. Un șir $w \in L_2$, poate fi scris sub formă $w = a u d$, unde u este de aceeași formă cu w , j sau w este de formă $b^r b^p c^p d^r$ și în acest caz u fie se poate scrie sub forma $u = b v d c u v$; un șir de aceeași formă cu u , sau u este de forma $x = b^p c^p$. Șirul nou obținut poate fi generat observând că se poate scrie sub forma $x = b z c$ unde z este un șir de aceeași formă cu x . Pot să fie utilizate, producțiile: $S \rightarrow a S d \mid b C d, C \rightarrow b C d \mid D, D \rightarrow b D c \mid X$. Se observă că $L = L_1 \cup L_2$ și că cele două neterminale B și D se comportă la fel (adică generează același limbaj), și putem renunța la neterminala D . Prin urmare, o gramatică G care generează limbajul L este: $G = (\{ S, A, B, C \}, \{ a, b, c, d \}, P, S)$ unde $P = \{ S \rightarrow a S d \mid A \mid b C d, A \rightarrow a A c \mid B, B \rightarrow b B c \mid X, C \rightarrow b C d \mid B \}$.

Verificare:

Vom considera două exemple de derivare: $a a b c c c$ (cazul $m > q$), respectiv $a b b c d d$ (cazul $m < q$). Pentru șirul $a a b c c c$ se obține secvența de derivări: $S \Rightarrow A \Rightarrow a A c \Rightarrow a a c c c$; $A c c c \Rightarrow a a b c c c \Rightarrow a a b b c c c$, respectiv pentru cel de-al doilea șir: $S \Rightarrow a S d \mid$

G este o gramatică independentă de context.

Problema li-42-7

$$L = \{ a^m b^n \mid n < m < 2^* n, n, m > 1 \}$$

Soluție:

Notăm $c u p = m - n$ și $q = 2^* n - m$ unde $p > 0, q > 0$. Prin urmare, $n = p + q$ și $m = 2^* p + i$; q și limbajul poate fi scris sub forma: $L = \{ a^{2^* p} a^q b^p \mid p, q > 0 \}$. Un șir care aparține acestui limbaj, poate să fie scris sub forma $w_1 = a a u b$ unde u este de aceeași formă cu w_1 sau de formă $a^q b^q$. Șirul $a^q b^q$ poate să fie generat în mod asemănător lui w_1 , observând că poate să fie scris sub forma $w_2 = a v b$ unde v este de aceeași formă cu w_2 sau este șirul ab . Prin urmare o gramatică este $G = (\{ S, A \}, \{ a, b \}, P, S)$ unde $P = \{ S \rightarrow a a S b \mid a a A b, A \rightarrow a A b \mid a b \}$.

Verificare:

Fie şirul: $a a a b b$ ($n = 2, m = 3$). Şirul poate să fie obţinut utilizând următoarea secvenţă de derivări: $S \Rightarrow a a A b \Rightarrow a a a b b$.

Comentarii:

G este o gramatică, independentă de context

problema 1.1-13

$L = \{ w c w^R \mid w \in \{a,b\}^* \}$ unde prin $u = w^R$ am notat şirul care are proprietatea $u_i = w_{n+1-i}$ pentru orice $i = 1, 2, \dots, n$ unde $n = |w|$ şi $w = w_1 w_2 \dots w_n, u = u_1 u_2 \dots u_n$ (şirul reflectat corespunzător lui w).

Soluție:

Orice şir din limbajul L, are proprietatea că începe şi se termină cu acelaşi simbol sau este doar simbolul c. În primul caz şirul poate să fie scris sub forma $w = x u x$ unde $x \in \{ a, b \}$, iar u are aceeaşi proprietate ca w. Luând în considerare cele de mai sus, gramatica G care generează limbajul L este $G = (\{ S \}, \{ a, b, c \}, P, S)$ unde $P = \{ S \rightarrow aSa \mid bSb \mid c \}$. Într-adevăr, orice şir u generat de G are proprietatea că există $w \in \{a,b\}^*$ astfel încât $u = w c w^R$ şi pentru orice şir u e L, există o derivare $S \Rightarrow u$.

Verificare:

Pentru şirul $a b a c a b a$, se obține următoarea secvenţă de derivări: $S \Rightarrow a S a \Rightarrow a b S b a \Rightarrow a b a S b a b a \Rightarrow a b a c a b a$

Comentarii:

G este o gramatică independentă de context

problema 1.1-14

$L = \{ w e \{a,b\}^* w = w^R \}$

Soluție:

Un şir w care respectă condiția de mai sus poate să aibă una din următoarele forme: $v^R, v a v^R, v b v^R$ unde $v \in \{ a, b \}^*$ este un şir oarecare. Pentru a genera aceste şiruri, putem să utilizăm gramatica $G = (\{ S \}, \{ a, b \}, P, S)$ cu $P = \{ S \rightarrow aSa \mid bSb \mid a \mid b \mid X \}$.

Verificare:

Fie şirul $a b a a b a b a$. Pentru obținerea acestui şir putem să considerăm următoarea secvenţă de derivări: $S \Rightarrow a S a \Rightarrow a b S b a \Rightarrow a b a S b a b a \Rightarrow a b a a S a a b a \Rightarrow a b a a b a b a$.

Comentarii:

G este o gramatică independentă de context

IMLIII.1 I.1 15

L =

Soluție:

Orice şir nevid w e L conține la începutul şirului toate simbolurile a, după care urmează toate simbolurile b şi se termină cu simbolurile c. Pentru a genera toate a-urile se va folosi neterminalul S_1 , pentru a genera toate b-urile se va folosi neterminalul S_2 , iar pentru a genera toate c-urile se va folosi neterminalul S_3 . Prin urmare, gramatica G care generează limbajul L poate să fie: $G = (\{ S, S_1, S_2, S_3 \}, \{ a, b, c \}, P, S)$ unde: $P = \{ S \rightarrow S_1 S_2 S_3, S_1 \rightarrow a S_1, S_2 \rightarrow b S_2, S_3 \rightarrow c S_3 \mid c \}$.

Verificare:

Fie şirul $a a b b c c c c$. Pentru obținerea acestui şir putem să considerăm următoarea secvenţă de derivări: $S \Rightarrow S_1 S_2 S_3 \Rightarrow a S_1 S_2 S_3 \Rightarrow a a S_2 S_3 \Rightarrow a a b S_2 S_3 \Rightarrow a a b b S_2 S_3 \Rightarrow a a b b b S_3 \Rightarrow a a b b c S_3 \Rightarrow a a b b c c S_3 \Rightarrow a a b b c c c S_3 \Rightarrow a a b b c c c c$

Comentarii:

G este o gramatică independentă de context

problema 1.1-16

$L = \{ a^i b^j c^k \mid i \neq j \text{ sau } j \neq k, i, j, k > 0 \}$

Soluție:

Pot apare 4 cazuri:

1. Dacă $i < j$, notăm $p = j - i$, rezultă $j = p + i, p > 0$. Fie L1 limbajul generat în acest caz. Atunci L1 poate fi descris prin: $L1 = \{ a^i b^{i+p} c^k \mid p > 0, i, k > 0 \}$. Vom folosi neterminalul $S1$ pentru generarea subşirului $a^i b^i$, pentru a genera subşirul b^p se va folosi neterminalul $S2$, iar pentru a genera subşirul c^k se va folosi neterminalul $S3$. Prin urmare vor fi utilizate producțiile $S \rightarrow S1 S2 S3, S1 \rightarrow a S1 b \mid a b, S2 \rightarrow b S2, S3 \rightarrow c S3 \mid c$.
2. Dacă $i > j$, notăm $p = i - j$, rezultă $i = p + j, p > 0$. Fie L2 limbajul generat în acest caz. Atunci L2 poate fi descris prin: $L2 = \{ a^{i+p} b^j c^k \mid p > 0, j, k > 0 \}$. Vom folosi neterminalul $S1$ pentru generarea subşirului a^{i+p} , pentru a genera subşirul $a^j b^j$ se va folosi neterminalul $S2$, iar pentru a genera subşirul c^k se va folosi neterminalul $S3$. Prin urmare vor fi utilizate producțiile $S \rightarrow S1 S2 S3, S1 \rightarrow a S1, S2 \rightarrow a S2 b \mid a b, S3 \rightarrow c S3 \mid c$.
3. Dacă $j < k$, notăm $p = k - j$, adică $k = j + p, p > 0$. Fie L3 limbajul generat în acest caz. Atunci L3 poate fi descris prin: $L3 = \{ a^i b^j c^{j+p} \mid p > 0, i, j > 0 \}$. Vom folosi neterminalul $S1$ pentru generarea subşirului a^i , pentru a genera subşirul $b^j c^j$ se va folosi neterminalul $S2$, iar pentru a genera subşirul c^p se va folosi neterminalul $S3$. Prin urmare vor fi utilizate producțiile $S \rightarrow S1 S2 S3, S1 \rightarrow a S1 \mid a, S2 \rightarrow b S2 c \mid b c, S3 \rightarrow c S3 \mid c$.

4. Dacă $j > k$, notăm $p = j - k$, și $j = p + k$, $p > 0$. Fie L_4 limbajul generat în acest caz Atunci L_4 poate fi descris prin: $L_4 = \{ a^i b^p c^k \mid p > 0, i, k > 0 \}$. Vom folosi neterminatul S_1 pentru generarea subșirului a^i , pentru a genera subșirul b^p se va folosi neterminatul S_2 , iar pentru a genera subșirul c^k se va folosi neterminatul S_3 . Prin urmare vor fi utilizate producțiile $S \rightarrow S, S_2 S_3, S, - * a \text{ Si } \mid a, S_2 \rightarrow \bullet b S_2, S_3 \rightarrow b S_3, c \mid b c$.

: Corespunzător celor 4 cazuri notăm neterminalele S_i, S_2, S_3 cu $S_{i_a}, S_{j_b}, S_{i_c}, S_{i_d}, S_{2a}, S_{2b}, S_{2c}, S_{2d}, S_{3a}, S_{3b}, S_{3c}, S_{3d}$. Prin urmare gramatica ce generează limbajul este: $G = (\{ S, S_{i_a}, S_{2a}, S_{3a}, S_{i_b}, S_{2b}, S_{3b}, S_{i_c}, S_{2c}, S_{3c}, S_{i_d}, S_{2d}, S_{3d} \}, \{ a, b, c \}, P, S)$ unde: $P = \{ S \rightarrow S_{i_a} S_{2a} S_{3a} \mid S_{i_b} S_{2b} S_{3b} \mid S_{i_c} S_{2c} S_{3c} \mid S_{i_d} S_{2d} S_{3d}, S_{i_a} \rightarrow \bullet a S_{j_a} b \mid a b, S_{2a} \rightarrow \bullet b S_{2a} \mid b, S_{3a} \rightarrow \bullet c S_{3a} \mid c, S_{j_b} \rightarrow a S_{i_b} \mid a, S_{2b} \rightarrow a S_{2b} \mid a b, S_{3b} \rightarrow c S_{3b} \mid c, S_{i_c} \rightarrow a S_{i_c} \mid a, S_{2c} \rightarrow b S_{2c} \mid b c, S_{3c} \rightarrow c S_{3c} \mid c, S_{i_d} \rightarrow a S_{i_d} \mid a, S_{2d} \rightarrow b S_{2d} \mid b, S_{3d} \rightarrow b S_{3d} \mid b c \}$.

G este o gramatică independentă de context

$$L = \{ a^i b^j c^k \mid i < j, i, j, k > 0 \}$$

Notăm $p = j - i$, adică $j = i + p$, $p > 0$. Atunci L poate fi descris prin: $L = \{ a^i b^p c^k \mid p > 0, i, k > 0 \}$. Vom folosi neterminatul A pentru generarea subșirului a^i , pentru a genera subșirul b^p se va folosi neterminatul B , iar pentru a genera subșirul c^k se va folosi neterminatul C . Prin urmare, gramatica G care generează limbajul L poate să fie: $G = (\{ S, A, B, C \}, \{ a, b, c \}, P, S)$ unde: $P = \{ S \rightarrow ABC, A \rightarrow a A b \mid A, B \rightarrow b B \mid b, C \rightarrow c C \mid A \}$.

; G este o gramatică independentă de context

$$L = \{ a^i b^j c^k \mid i + j = k, i, j, k > 0 \}$$

Dacă $k = i + j$ atunci L poate fi descris prin: $L = \{ a^i w c^j \mid i, j > 0 \}$. Orice șir w din limbaj începe cu simbolul a , și se termină cu simbolul c . Deci, w se poate scrie ca $w = a u c$ unde u este un șir care fie este de tipul celui din mulțimea L , fie este de forma $V d, j > 0$. Acesta din urmă se poate scrie sub forma $b v c$ cu v de aceeași formă sau v este șirul bc . Prin urmare gramatica G care generează limbajul L poate să fie: $G = (\{ S, A \}, \{ a, b, c \}, P, S)$ unde: $P = \{ S \rightarrow a S c \mid a A c, A \rightarrow b A c \mid b c \}$.

; G este o gramatică independentă de context

1.1-1')

I \ \ i \ i m i

Soluție:

Vom presupune că parcurgem $-inul \ \ \$ de la aiânga ia Jivapia considerând un prefix al acestuia: $w = u a v$ unde u este prefixul, a este simbolul curent (a sau b), iar v este restul șirului ($u, v \in \{ a, b \}^*$). Sunt posibile următoarele situații:

1. Numărul de simboluri a din prefixul u este par
 2. Numărul de simboluri a din prefixul u este impar
- ; Asociem neterminatul S primei stări (care este de altfel ceea ce corespunde generării unui șir din limbajul L), iar celelalte stări neterminale A . Din prima stare se trece în a doua dacă simbolul curent a a fost un a și se rămâne în starea S dacă simbolul curent a a fost un b . Din starea A se trece în starea S dacă simbolul curent a a fost un a și se rămâne în starea A dacă simbolul curent a a fost un b . Deci gramatica care generează limbajul L este: $G = (\{ S, A \}, \{ a, b \}, P, S)$ unde $P = \{ S \rightarrow a A \mid b S \mid \hat{A}, A \rightarrow b A \mid a S \}$.

Comentarii:

G este o gramatică regulată

jiiuillii.i 1.1 IU ; " - - "-

$$L = \{ w \in \{ a, b \}^* \mid \#_a(w) = 2 \#_b(w) \}$$

Soluție:

Pentru a genera șirurile din limbaj de câte ori apare un simbol b trebuie să apară doi simboluri a . Nu este importantă ordinea în care apar simbolii. Deci gramatica care generează limbajul L este: $G = (\{ S \}, \{ a, b \}, P, S)$ unde $P = \{ S \rightarrow a S a S b S \mid a S b S a S \mid b S a S a S \mid \hat{A} \}$.

Comentarii:

G este o gramatică independentă de context

Limbajul parantezelor bine formate

Soluție:

Gramatica care generează limbajul L este: $G = (\{ S \}, \{ (,) \}, P, S)$ unde $P = \{ S \rightarrow S S \mid (S) ! () \}$ sau $G = (\{ S \}, \{ (,) \}, P, S)$ unde $P = \{ S \rightarrow S S \mid (S) \mid X \}$ dacă este acceptat șirul vid.

Comentarii:

G este o gramatică independentă de context

problema 1.1-22

$$L = \{ wwjwe \{0, 1\}^* \}$$

Soluție:

Analizând forma generală a șirurilor ce fac parte din limbajul L se observă că în generarea acestor șiruri este necesară "memorarea" primei jumătăți a șirurilor pe măsură ce se face această generare, pentru a putea utiliza apoi informația memorată în generarea celei de-a doua jumătăți a șirului. În mod evident, această "memorare" nu se poate realiza utilizând mecanismele oferite de gramaticile regulate sau independente de context. Am văzut că se pot genera șiruri de formă w / c u ajutorul unor gramatici independente de context. Vom genera un șir în care unui simbol 1 din prima jumătate a șirului îi va corespunde un simbol notat U în a doua jumătate a șirului. În mod corespunzător, unui simbol 0 din prima jumătate a șirului îi va corespunde un simbol notat Z în a doua jumătate a șirului. Pentru a marca sfârșitul șirului vom utiliza un neterminat B. Se observă că în acest mod s-a memorat prima jumătate a șirului într-o formă inversată. În continuare, din modul în care se utilizează următoarele derivări ar trebui să rezulte și restul șirului. Astfel, transformările pe care le vom opera asupra unui șir generat ca mai sus, pot să fie exprimate în termeni intuitivi astfel: orice Z sau U se transformă în 0, respectiv 1 dacă sunt urmate de terminatorul B, iar orice simbol 0 sau 1 se deplasează către capătul din stânga, respectiv orice simbol Z sau U se deplasează către dreapta (pentru a ajunge la terminatorul B).

De exemplu: $001UZZB \Rightarrow 001UZ0B \Rightarrow 001UOZB \Rightarrow 001U00B \Rightarrow 0010U0B \Rightarrow 00100UB \Rightarrow 001001B \Rightarrow 001001$. Astfel, mulțimea producțiilor gramaticii va conține $S \rightarrow AB$ cu neterminatul A utilizat pentru crearea unui șir de formă $w_1 w_2$ unde $w_1 \in \{0, 1\}^*$, $w_2 \in \{U, Z\}^*$, iar între w_1 și w_2 există relația: $w_1 = w_2^R$ în care $U = 1, Z = 0$. Acest șir este generat prin producțiile: $A \rightarrow 0AZ \mid 1AU \mid X$. Pentru generarea limbajului L mai sunt necesare producțiile care realizează înlocuirea neterminalelor U, Z cu 1, respectiv 0. Dacă sunt urmate de neterminatul B și dubla deplasare a simbolilor: $ZB \rightarrow 0B, UB \rightarrow 1B, Z0 \rightarrow 0Z, U0 \rightarrow 0U, Z1 \rightarrow 1Z, U1 \rightarrow 1U$. De asemenea, este necesară producția $B \rightarrow A$ pentru a elimina simbolul utilizat ca terminator. Deci o gramatică G care generează limbajul L este $G = (\{S, A, B, U, Z\}, \{0, 1\}, P, S)$ unde: $P = \{S \rightarrow AB, A \rightarrow 0AZ \mid 1AU \mid A., ZB \rightarrow 0B, UB \rightarrow 1B, Z0 \rightarrow 0Z, U0 \rightarrow 0U, Z1 \rightarrow 1Z, U1 \rightarrow 1U, B \rightarrow A.\}$.

criticare:

Pentru șirul 0101 se obține șirul de derivări: $S \Rightarrow AB: 0AZB \Rightarrow 01AUZB \Rightarrow 01ZB \Rightarrow 010UB \Rightarrow 010UB \Rightarrow 0101B \Rightarrow 0101$

Comentarii:

Se observă că G este o gramatică de tip 0 (tară restricții).

Dacă se utilizează producția $B \rightarrow A$ într-un șir de derivări înaintea transformării tuturor neterminalelor U și Z, se poate obține o formă propozițională din care nu se mai poate obține un șir din limbajul L.

De exemplu: $S \Rightarrow AB \Rightarrow 0AZB \Rightarrow 0AZ$ și nu mai există nici o producție care să poată să fie aplicată. Din forma propozițională obținută nu se mai poate obține un șir care să aparțină limbajului L. Astfel de rezultate țin de aspectul nedeterminist al gramaticii respective.

1.1-23

$$L = \{a^n b^n \mid n > 1\}$$

Soluție:

Problema poate să fie abordată într-o manieră similară cu cea anterioară adică se „memorează” prima parte a șirului (a^n), însă generarea nu poate continua verificând că există un număr corespunzător de simboluri b și de simboluri c. Prin urmare, va trebui să „memorăm” prima parte a șirului (a^n) și apoi să generăm același număr de perechi de forma (BC) unde B, C sunt neterminale asociate cu terminalele b, c. În continuare, din derivările ulterioare, va rezulta

transformarea celei de-a doua părți a șirului în șirul $b^n c^n$. Vom considera deci, șirurile de forma $a^n (BC)^n$ unde B, C sunt neterminale prin intermediul cărora se va face generarea subșirului $b^n c^n$. Pentru aceasta este necesară substituția oricărui subșir de forma CB cu șirul BC, respectiv înlocuirea neterminalelor B, C cu b, respectiv c. Astfel, o gramatică G care generează limbajul L este $G = (\{S, B, C\}, \{a, b, c\}, P, S)$ unde: $P = \{S \rightarrow aSBC \mid aBC, CB \rightarrow BC, bB \rightarrow bb, bC \rightarrow bc, cB \rightarrow Bc, cC \rightarrow ce\}$.

Verificare:

Pentru șirul $a^3 b^3 c^3$ de exemplu, se obține șirul de derivări: $S \Rightarrow aSBC \Rightarrow a a S B C B C \Rightarrow a a a b C B C B C \Rightarrow a a a b B C C B C \Rightarrow a a a b B C B C C \Rightarrow a a a b B B C C C \Rightarrow a a a b b B C C C \Rightarrow a a a b b b C C C \Rightarrow a a a b b b c C C \Rightarrow a a a b b b c c C \Rightarrow a a a b b b c c c = a^3 b^3 c^3$.

Comentarii:

G este o gramatică fără restricții

1.1-24

$$L = \{a^m b^n c^m d^n \mid m, n > 1\}$$

Soluție:

Vom genera întâi forme propoziționale de tipul $a^m b^n D^n C^m$ unde C, D sunt neterminale asociate terminalelor c, d. Într-adevăr, o formă propozițională ca mai sus poate fi generată de producțiile: $S \rightarrow aSCj a C, A \rightarrow bAD \mid bD$.

Următorul pas este deplasarea simultană a simbolilor C spre stânga și respectiv D spre dreapta, înlocuirea cu terminalele corespunzătoare c, d făcându-se de la stânga la dreapta: $DC \rightarrow CD, bC \rightarrow bc, cC \rightarrow ce, cD \rightarrow cd, dC \rightarrow Cd, dD \rightarrow dd$. Astfel, o gramatică G care generează limbajul L este $G = (\{S, A, C, D\}, \{a, b, c, d\}, P, S)$ unde: $P = \{S \rightarrow aSC \mid a^A Q A \rightarrow bAD \mid bD, DC \rightarrow CD, bC \rightarrow bc, cC \rightarrow ce, cD \rightarrow cd, dC \rightarrow Cd, dD \rightarrow dd\}$.

1.1.3 Transformări asupra GIC

Gramaticile independente de context (GIC) sunt utilizate pentru analiza sintactică în cadrul compilatoarelor. Pentru că procesul de analiza sintactică se simplifică mult dacă gramaticile l utilizate satisfac anumite condiții, este util să se facă transformări asupra acestui tip de gramatici, transformări care să producă gramatici echivalente (care generează același limbaj).

1.1.3.1 Eliminare recursivitate stânga

Spunem că o gramatică este *recursivă stânga* dacă există un neterminal A astfel încât existai o derivare $A \Rightarrow^+ Ax$ pentru $A \in N, x \in (N \cup T)^*$. Dacă pentru un neterminal A există] producțiile

$$AB_2f... | AB_m|y_1| \dots | y_n$$

unde Y_i nu începe cu A, $1 < i < n$, se spune că avem *recursivitate stânga imediată* și] producțiile anterioare se pot înlocui cu:

$$A \rightarrow Y, A' | y_2 A' | \dots | y_n A'$$

$$A' \rightarrow B_1 A' | B_2 A' | \dots | B_m A' | X$$

Această construcție elimină recursivitatea stângă imediată. Dacă gramatica nu permite i derivări de tipul $A \Rightarrow^+ A$ (nu are cicluri) și nu conține X producții poate să fie transformată în: vederea eliminării recursivității stânga utilizând următorul algoritm.

Intrare: o gramatică fără cicluri și X producții
Ieșire; o gramatică echivalentă fără recursivitate stânga.
 Se alege o ordine a neterminalelor, fie ea: A_i, A_j, \dots, A_n ,
pentru $i = 2 \dots n$ **execută**

pentru $j = 1 \dots i-1$ **execută**

înlocuiește fiecare producție de forma $A_i \rightarrow r$ cu producțiile-
 $A_i \rightarrow U_i r | u_1 r | \dots | u_k r$ unde $A_i \rightarrow u_i$ $u_1 \dots | u_k$
 sunt toate producțiile pentru A_j .

elimină recursivitatea stângă imediată între producțiile A_i .

1.1.3.2 Eliminare A producții

O gramatică nu are X producții dacă satisface una din-următoarele condiții:

- Nu există nici o producție care să aibă în partea dreaptă șirul vid sau
- Există o singură producție care are în partea dreaptă șirul vid și anume producția $S \rightarrow X$.
 Simbolul S nu apare în acest caz în partea dreaptă a nici unei producții.

Algoritmul de transformare este:

Intrare: o gramatică $G = (N, T, P, S)$
T,șire: o gramatică $G' = (N', T, P', S')$ care satisface condițiile $L(G) = L(G')$ și G' nu are X producții

$i = 0$ *
 $N_0 = \{ A | A \in X \}$
repetă
 $N_{i+1} = \{ A | A \Rightarrow^+ a e, a \in EN, \dots \} \cup N_i$
până când $N_i = N_{i+1}$

dacă $S \in N_i$
 $N' = N \cup \{ S' \}$
 $P' = \{ S' \rightarrow A, S' \rightarrow \dots \}$
altfel
 $N' = N$
 $S' = S$
 $P' = P$

pentru fiecare $p \in P$ **execută**
dacă p este de forma
 $A \rightarrow a_1 B_1 a_2 \dots B_k a_k, k > 0, B_i \in N, 1 < i < k, a_j \in (N - N_i) \cup T, 0 < j < k$
 $P' = P \cup (\{ A \rightarrow a_1 | \dots | a_k \} \cup \{ A \rightarrow \hat{A} \}) - \{ A \rightarrow \hat{A} \}$
altfel
 $P' = P \cup P$

1.1.3.3 Eliminare simbolii neutilizați

Un simbol neterminal este nefinalizat dacă din el nu se poate deriva un șir format numai din simbolii terminali.

Algoritmul de transformare este:

Intrare: o gramatică $G = (N, T, P, S)$
Ieșire: o gramatică $G' = (N', T, P', S)$ care satisface condițiile $L(G) = L(G')$ și $\forall A \in N, A \Rightarrow^+ w, w \in T^*$.
 $i = 0$;
 $N_0 = \{ A \in N | A \in T \}$
repetă
 $i++$
 $N_i = \{ A | A \rightarrow \dots \in P, \dots \in (N_i \cup T)^* \} \cup N^{i-1}$
Până când $N_i = N_{i+1}$
 $N' = N_i$.
 P' c: P și este format numai producțiile din P care au în partea stângă simbolii din N' .

1.1.3.4 Eliminare simbolii inaccesibili

Un simbol neterminal este inaccesibil dac  nu poate apare  ntr-o form  propozi ional . Algoritmul de transformare este:

```
Intrare: o gramatic  G = (N, T, P, S)
Ieşire: o gramatic  G' = (N', T, P', S) care satisface condi iile L(G) = L(G')
         i VA eN' exist  w e (N U T)*, s =>w  i A apare  n w.

i = 0;
N0 = {S}
repet 

    N* = { AeN | A apare  n partea dreapt  a produ iilor pentru un neterminal
           din Nj-j} UN,
p n  c nd Ni = N*
N' = Ni
P' c: P con ine numai ;produ iile corespunz toare neterminalelor din N'
```

Probleme

problema 1.1-28

S  se elimine recursivitatea-st ng  pentru urm toarea gramatic : G = ({S, L}, { a, ,, (,) }, P, S) unde P = {S -> (L) | a, L -> L, S | a}

Solu ie:

Se alege o ordine pentru neterminale, fie ea: S < L . Pentru produ ia S -> (L) | a nu se face nici o modificare. La fel pentru S -> a.

Pentru produ ia L -> L,S | a se elimin  recursivitatea imediat   i rezult :

```
L -> aL'
L' -> ,SL' | X
```

In final dup  eliminarea recursivit  ii st nga gramatica este:

```
G = ({ S, L, L'}, { a,, G}), P, S) unde: P = { S-> (L) | a, L,-> aL', L' -> ,SL' |X}
```

pro bit mu 1.1-29

S  se elimine recursivitatea st ng  pentru urm toarea gramatic :

```
G = ({ S, A, B }, { a, b }, P, S) unde P = {S -> A | B, A -> Sa | Bb, B -> Sb | Aa}
```

Solu ie:

Se alege o ordine pentru neterminale S < A < B

Pentru produ ia S -> A | B nu se face nici o modificare.

Pentru produ ia A -> Sa | Bb avem S < A deci se va folosi S -> A | B | X  i rezult 

A -> Aa | Ba | a |Bb din care se elimin  recursivitatea imediat   i rezult :

```
AA -> BaA' | aA' | BbA', A' -> aA' | X
```

Pentru produ ia B -> Sb | Aa se substituie cu S -> Sb | Aa | X  i rezult  B H> Ab | Bb | b| Aa iar apoi se  nlocuieşte A -> BaA' | aA' | BbA'  i rezult  B -> BaA'b | aA'b j BbA'b | Bb | b | B a A' a | a A' a | B b A' a. Din care se elimin  recursivitatea imediat   i rezult :

```
; BB -> a A' b B' | bB' | a A' aB', B' -> aA'b B' | b A' b B' | b B' | a A' a B' | b A'a B' | X
```

 n final dup  eliminarea recursivit  ii st nga gramatica este:

```
; G = ({ S, A, A', B, B' }, { a, b }, P, S) unde: P = { S-> A | B, A -> BaA' | aA' | BbA', A' -> aA' i X, B -> a A' b B' | bB' | a A' aB', B' -> aA'b B' | b A' b B' | b B' | a A' a B' | b A'a B' }
```

probi-Pia 1.1-30

S  se elimine recursivitatea st ng  pentru urm toarea gramatic : G = ({ A, B, C }, { a, b, c }, P, A) unde P = {A -> BC | a, B -> Ca | Ab, C -> Ab | cC | b}

Solu ie:

Se alege o ordine pentru neterminale A < B < C. Pentru produ ia A -> BC | a nu se face nici o modificare. Pentru produ ia B -> Ca | Ab deoarece A < B se va folosi A -> BC j a  i rezult  B -> Ca | BCb | ab din care se elimin  recursivitatea imediat   i rezult : B -> CaB' | jabB', B' -> CbB'|L

• Pentru produ ia C -> AB j cC | b se substituie cu A -> BC | a  i rezult  C -> BCb | aB | cC | b iar apoi se  nlocuieşte B -> CaB' | abB'  i rezult  C -> CAB'CB | abB'Cb | aB | cC | b

: Din care se elimin  recursivitatea imediat   i se ob ine: C -> abB'CbC | aBC | cCC | bC, C -> AB'CBC | X.  n final dup  eliminarea recursivit  ii st nga gramatica este:

```
G = ({ A, B, B\ C, C }, { a, b, c }, P, A) unde: P = { A-> BC | a, B -> CaB' | abB', B'' -+ CbB' | X, C -> abB'CbC | aBC j cCC | bC, C -> AB'CBC j X}
```

1.1-31

* i ->• elimine ' piiHluciiii.- Jiu y.nn iik\i li i[s. \. li. l ',. / J. l' ; ' P. M unde P = {S -> ABC, A -> BB | X, B -> CC | a, C -> AA | b}

Solu ie:

```
N0 = {A}, Ni = {A, C}, N2 = {A, B, C}, N3 = {S, A, B, C} = Nf
```

: S e N_f deci se introduc produ iile S' -> S  i S' -> X

```
\ Produ ia S -> ABC devine S '-> ABC | AB | AC | BC | A |B | C
```

```
, Produ ia A -> BB devine A -> BB | B
```

```
: Produ ia B -> CC \ a devine B -> CC | C | a
```

```
! Produ ia CC -> AA | b devine C -> AA | A | b
```

; Dup  eliminarea X produ iilor gramatica este G = ({ S', S, A, B, C }, { a, b }, P, S') unde: P = { S' -> S , S' -> X, S -> ABC | AB | AC | BC | A |B | C, A -> BB | B, B -> CC | C | a, C -> ,AA |A |b }

problema 1.1-32

Să se elimine simbolii nefinalizați, iar apoi cei inaccesibili pentru gramatica $G = (\{S, A, B\}, \{a, b\}, P, S)$ unde $P = \{S \rightarrow A \mid a, A \rightarrow AB, B \rightarrow b\}$

Soluție:

Eliminare simbolii nefinalizați: $N_0 = \{S, B\}, N_1 = \{S, B\}, N_f = N_1$.

Rezultă că A este simbol nefinalizat, se vor elimina producțiile corespunzătoare neterminalelor A și gramatica va deveni: $G = (\{S, B\}, \{a, b\}, P, S)$ unde $P = \{S \rightarrow a, B \rightarrow b\}$.

Eliminare simbolii inaccesibili: $N_0 = \{S\}, N_1 = \{S\}, N_f = N_1$

B este inaccesibil, rămâne producția $S \rightarrow a$. După eliminarea simbolilor nefinalizați și inaccesibili $G = (\{S\}, \{a\}, P, S)$ unde $P = \{S \rightarrow a\}$

Coaerentă:

Ordinea corectă de aplicare a celor doi algoritmi este eliminare simbolii nefinalizați și apoi eliminare simbolii inaccesibili.

problema 1.1-33

Să se elimine simbolii nefinalizați, iar apoi cei inaccesibili pentru gramatica $G = (\{S, A, B, C, D\}, \{a, b, d\}, P, S)$ unde $P = \{S \rightarrow bA \mid aB, A \rightarrow a \mid aC \mid aD \mid aS \mid bAA, B \rightarrow b \mid Cb \mid Db \mid bS \mid aBB, C \rightarrow Cb \mid Db, D \rightarrow Cb \mid Db, D \rightarrow Cd \mid dAC\}$

Soluție:

Eliminare simbolii nefinalizați:

$N_0 = \{A, B\}, N_1 = \{S, A, B\}, N_f = N_1$

Rezultă că C, D ∈ N_f, adică sunt simbolii nefinalizați, se vor elimina producțiile corespunzătoare lor și gramatica va deveni: $G = (\{S, A, B\}, \{a, b, d\}, P, S)$ unde $P = \{S \rightarrow bA \mid aB, A \rightarrow a \mid aS \mid bAA, B \rightarrow b \mid bS \mid aBB\}$

Eliminare simbolii inaccesibili:

$N_0 = \{S\}, N_1 = \{S, A, B\}, N_f = N_1$

Nu există simbolii inaccesibili.

În final gramatica este: $G = (\{S, A, B\}, \{a, b, d\}, P, S)$ unde $P = \{S \rightarrow bA \mid aB, A \rightarrow a \mid aS \mid bAA, B \rightarrow b \mid bS \mid aBB\}$

problema 1.1-34

Să se elimine recursivitatea stângă, X producțiile, simbolii neutilizați și inaccesibili pentru gramatica $G = (\{S, B, C, D, E\}, \{a, b\}, P, S)$ unde $P = \{S \rightarrow aBa, B \rightarrow Sb \mid bCC \mid DaB, C \rightarrow abb \mid DD, E \rightarrow aC, D \rightarrow aDB\}$

Soluție:

Eliminare recursivitate stângă. Se alege ordinea neterminalelor $S < B < C < E < D$

Pomind de la $B \rightarrow \bullet Sb \mid bCC \mid DaB, S \rightarrow aBa$. Rezultă $B \rightarrow aBab \mid bCC \mid DaB$ care nu este recursivă stângă-

Eliminare Xproducții. Gramatica nu are X producții.

Eliminare simbolii nefinalizați. $N_0 = \{C\}, N_1 = \{B, C, E\}, N_2 = \{S, B, C, E\} = N_f$

$D \notin N_f$ se elimină toate producțiile corespunzătoare neterminalelor D.

Gramatica devine $G = (\{S, B, C, E\}, \{a, b\}, P, S)$ unde $P = \{S \rightarrow aBa, B \rightarrow Sb \mid bCC, C \rightarrow abb, E$

Eliminare simbolii inaccesibili. $N_0 = \{S\}, N_1 = \{S, B\}, N_2 = \{S, B, C\} = N_f \in N_f$, rezultă E simbol inaccesibil.

În final rezultă: $G = (\{S, B, C\}, \{a, b\}, P, S)$ unde $P = \{S \rightarrow aBa, B \rightarrow Sb \mid bCC, C \rightarrow abb\}$

problema 1.1-35

Să se elimine recursivitatea stângă, X producțiile, simbolii neutilizați și inaccesibili pentru gramatica $G = (\{S, T\}, \{a, b, c\}, P, S)$ unde $P = \{S \rightarrow TbT, T \rightarrow TAT \mid ca\}$

Soluție:

Eliminare recursivitate stângă. Alegem ordinea neterminalelor $S < T$. Se consideră producția $T \rightarrow TaT$ se elimină recursivitatea stângă și rezultă $T \rightarrow caT'$ și $T' \rightarrow aTT' \mid X$

Rezultă $G = (\{S, T, T'\}, \{a, b, c\}, P, S)$ unde $P = \{S \rightarrow TbT, T \rightarrow caT', T' \rightarrow aTT' \mid X\}$

Eliminare Xproducții. $N_0 = \{T'\}, N_1 = \{T'\}, N_f = N_1$, Rezultă $G = (\{S, T, T'\}, \{a, b, c\}, P, S)$ unde $P = \{S \rightarrow TbT, T \rightarrow caT' \mid ca, T' \rightarrow aTT' \mid aT\}$

Eliminare simbolii nefinalizați. $N_0 = \{T\}, N_1 = \{S, T, T'\}$. Nu există simbolii nefinalizați.

Eliminare simbolii inaccesibili. $N_0 = \{S\}, N_1 = \{S, T\}, N_2 = \{S, T, T'\}, N_f = N_2$. Nu există simbolii inaccesibili.

În final se obține $G = (\{S, T, T'\}, \{a, b, c\}, P, S)$ unde $P = \{S \rightarrow TbT, T \rightarrow caT' \mid ca, T' \rightarrow aTT' \mid aT\}$.

1.1 Mulțimi regulate. Expresii regulate

Definiția 1.2.1. Fie T un alfabet finit. Se numește mulțime regulată asupra alfabetului T mulțimea definită recursiv astfel:

- 1- 0 este o mulțime regulată asupra alfabetului T.
2. Dacă a ∈ T, atunci { a } este o mulțime regulată asupra alfabetului T.
3. Dacă P, Q sunt mulțimi regulate, atunci P ∪ Q, PQ, P* sunt mulțimi regulate asupra alfabetului T.
4. O mulțime regulată asupra alfabetului T nu se poate obține decât aplicând regulile 1-3. Am notat prin P ∪ Q, PQ, respectiv P* reuniunea, concatenarea a două mulțimi, respectiv închiderea tranzitivă a unei mulțimi.

Definiția 1.2.2. O expresie regulată este definită recursiv în modul următor:

1. X este o expresie regulată care generează mulțimea 0.
2. Dacă a ∈ T, atunci a este o expresie regulată care generează mulțimea { a }.
3. Dacă p, q sunt expresii regulate care generează mulțimile P, Q, atunci: (p + q) sau (p | q) este o expresie regulată care generează mulțimea P ∪ Q.

- (pq) este o expresie regulată care generează mulțimea PQ.
 (p)* este o expresie regulată care generează mulțimea P*.
 4. O expresie regulată nu se poate obține decât prin aplicarea regulilor 1-3.

Proprietăți ale expresiilor regulate:

1. $a|P = B|oc$ (comutativitate reuniune)
2. $a|(B|y) = (a|P)|y$ (asociativitate reuniune)
3. $ct(py) = (a p) y$ (asociativitate concatenare)
4. $a(P|y) = aP|a y$ (distributivitate la stânga a concatenării față de reuniune)
5. $(a|P)y = ay|p y$ (distributivitate la dreapta a concatenării față de reuniune)
6. $a^{\wedge} = A, a = a$ (element neutru pentru concatenare)
7. $a^* = a I a^*$
8. $(a^*)^* = a^*$
9. $a|a = a$
10. $(a^*p^*)^* = (a + P)^*$

Utilizând expresii regulate se pot construi ecuații regulate. Soluția generală a ecuației de forma: $X = aX + b$ unde a, b, X sunt expresii regulate, este: $X = a^*(b + y)$ unde y este o expresie regulată oarecare. Soluția minimală (punctul fix al ecuației) este: $X = a^*b$.

Propoziție Fie G o gramatică regulată. L(G) este o mulțime regulată.

Definiția 1.2.3. Două expresii regulate se numesc echivalente dacă descriu aceeași mulțime regulată.

Probleme

problema 1.12-1

Să se rezolve sistemul de ecuații:

$$\begin{aligned} X &= a_1X + a_2Y + a_3 \\ Y &= b_1X + b_2Y + b_3 \end{aligned}$$

Soluție:

Pentru prima ecuație, soluția este $X = a_1^*(a_2Y + a_3)$, înlocuind în cea de-a doua ecuație, obținem $Y = b_1a_1^*(a_2Y + a_3) + b_2Y + b_3$ sau echivalent, folosind proprietățile expresiilor regulate: $Y = b_1a_1^*a_2Y + b_1a_1^*a_3 + b_2Y + b_3$ sau $Y = (b_1a_1^*a_2 + b_2)Y + (b_1a_1^*a_3 + b_3)$. Deci, pentru Y se obține soluția $Y = (b_1a_1^*a_2 + b_2)^*(b_1a_1^*a_3 + b_3)$. În mod corespunzător, înlocuind în prima ecuație, se obține următoarea soluție pentru X: $X = a_1^*a_2Y + a_1^*a_3$ sau $X = a_1^*a_2(b_1a_1^*a_2 + b_2)^*(b_1a_1^*a_3 + b_3) + a_1^*a_3$.

problema 1.12-2

Să se rezolve sistemul de ecuații:

$$\begin{aligned} X_1 &= OX_2 + 1X_1 + \hat{A}, \\ X_2 &= OX_3 + 1X_2 \\ X_3 &= OX_1 + 1X_3 \end{aligned}$$

Soluție:

Din ultima ecuație obținem $X_3 = 1^*0X_1$. Înlocuind în a doua ecuație obținem $X_2 = 1X_2 + 01^*0X_1$, de unde rezultă $X_2 = 1^*01^*0X_1$. Dacă se înlocuiește în prima ecuație rezultă $X_1 = 01^*01^*0X_1 + 1X_1 + X$. Deci $X_1 = (01^*01^*0 + 1)^*$. Dacă se înlocuiește rezultatul obținut pentru X_1 în formulele corespunzătoare lui X_2 și X_3 obținem $X_2 = 1^*01^*0(01^*01^*0 + 1)^*$ și $X_3 = 1^*0(01^*01^*0 + 1)^*$.

Comentarii:

Se poate demonstra că în expresia regulată X_1 numărul de simboluri 0 este divizibil cu 3. Prin urmare, X_1 poate fi scris sub formă echivalentă: $X_1 = (1^*01^*01^*0)^*1^*$.

problema 1.12-3

Să se construiască expresia regulată care generează mulțimea regulată egală cu limbajul generat de gramatica regulată cu producțiile: $P \rightarrow 0Q11P$, $Q \rightarrow 0R11P$, $R \rightarrow 0R11R10$

Soluție:

Asociem fiecărui netcrninal o expresie regulată și fiecărei producții de forma $A \rightarrow a|p$, o ecuație de forma $A = a + p$ unde a și p sunt tot expresii regulate. Corespunzător setului de producții de mai sus, se obține următorul sistem de ecuații:

$$\begin{aligned} P &= 0Q + 1P \\ Q &= 0R + 1P \\ R &= 0R + 1R + 0 \end{aligned}$$

Din ultima ecuație, se obține folosind proprietatea de distributivitate a concatenării față de reuniune $R = (0 + 1)R + 0$ care are soluția $R = (0 + 1)^*0$. Înlocuind în a doua ecuație, obținem: $Q = 0(0 + 1)^*0 + 1P$. Înlocuind Q în prima ecuație: $P = 00(0 + 1)^*0 + 01P + 1P$ adică $P = (01 + 1)P + 00(0 + 1)^*0$. Această ecuație are soluția $P = (01 + 1)^*00(0 + 1)^*0$.

Să se construiască expresia regulată care generează mulțimea regulată egală cu limbajul regulat generat de gramatica regulată:

$$\begin{aligned} S &\rightarrow 0A[1S|X \\ A &\rightarrow 0B|1A \\ B &\rightarrow 0S|1B \end{aligned}$$

Soluție:

Corespunzător setului de producții de mai sus, se obține următorul sistem de ecuații:

$$\begin{aligned} S &= 0A + 1S + X \\ A &= 0B + 1A \\ B &= 0S + 1B \end{aligned}$$

! Din ultima ecuație, obținem $B = 1^* 0 S$. Înlocuind în a doua ecuație, aceasta devine: $A = 0 1^* : 0 S + I A$ care are soluția $A = 1^* 0 1^* 0 S$. Prima ecuație devine în urma înlocuirii expresiei : regulate $A, S = 0 1^* 0 1^* 0 S + 1 S + X$ sau, folosind distributivitatea concatenării față de reuniune. $S = (0 1^* 0 1^* 0 + 1) S + X$ care are soluția $S = (0 1^* 0 1^* 0 + 1)^* X = (0 1^* 0 1^* 0 + 1)^* \dots$

problema 1.2-5

Să se construiască gramatica regulată care generează limbajul generat de expresia regulată: $(a | b)^* a^* b^+ a^*$

Soluție:

Notăm $S = (a | b)^* a^* b^+ a^*$ și $A = a^* b^+ a^*$. Corespunzător, se poate scrie ecuația ce are soluția $S: S = (a + b) S + A$. De asemenea, dacă notăm $B' = b^+ a^*$ și $B = b^* a^*$, putem scrie următoarea relație $B' = b B$ (deoarece $b^+ = b b^*$).

Pe de altă parte, $A = a^* B'$ și deci se poate scrie relația $A = aA + B'$ sau $A = aA + bB$ unde $B = b^* a^*$. Notând $C = a^*$, B devine $B = b^* C$ care este soluția ecuației $B = b B + C$. De asemenea, $C = a^*$ este soluția ecuației $C = a C + X$. Prin urmare, corespunzător expresiei regulate de mai sus se obține următorul sistem de ecuații:

$$S = aS + bS + A$$

$$A = aA + bB$$

$$B = bB + C$$

$$C = aC + A$$

și respectiv următorul set de producții:

$$S \rightarrow aS \quad bSjA$$

$$A \rightarrow aA | bB$$

$$B \rightarrow bB$$

Gramatica G care generează limbajul descris de expresia de mai sus, este $G = (\{ S, A, B, C \}, \{ a, b \}, P, S)$ cu P mulțimea de producții de mai sus.

problema 1.2-6

Să se construiască gramatica regulată care generează limbajul descris de expresia regulată: $(a|b)^* a (a|b)$

Soluție:

Notăm $S = (a | b)^* a (a | b)$ și $A' = a (a | b)$. Deci, S este soluția ecuației $S = (a + b) S + A'$ iar A' poate fi scris sub forma $A' = a A$ unde $A = a + b$. Deci, corespunzător expresiei regulate S , poate să fie scris următorul sistem de ecuații: $S = a S + b S + A', A' = a A, A = a + b$ sau $S = aS + bS + aA, A = a + b$. Gramatica ce generează limbajul generat de expresia S este $G = (\{ S, A \}, \{ a, b \}, P, S)$ unde P conține producțiile: $\{ S \rightarrow a S | b S | a A, A \rightarrow a | b \}$.

problema 1.2-7

Șă se construiască gramatica care generează limbajul descris de expresia regulată: $(a | b)^* a (a | b) (a | b)$

Soluție:

Notăm $S = (a | b)^* a (a | b) (a | b)$ și respectiv $A' = a (a | b) (a | b)$. S este soluția ecuației: $S = (a + b) S + A'$ care poate să fie scrisă echivalent $S = a S + b S + A' (1)$. Dar $A' = a A$ unde $A = (a | b) (a | b)$. Prin urmare, relația (1) devine $S = aS + bS + aA (2)$. Dacă notăm $B = (a | b)$, atunci A devine $A = a B + b B (3)$ și $B = a + b (4)$ Corespunzător relațiilor (2) - (4), setul de producții al gramaticii G care generează limbajul descris de expresia regulată S este :

$$S \rightarrow a S | b S | a A$$

$$A \rightarrow a B | b B$$

$$B \rightarrow a | b$$

Gramatica ce generează limbajul descris de expresia regulată dată este $G = (\{ S, A, B \}, \{ a, b \}, P, S)$ unde P conține producțiile $S \rightarrow a S | b S | a A, A \rightarrow a B | b B, B \rightarrow a | b$

1.3 Acceptoare

Spre deosebire de gramatici și expresii regulate care generează limbaje formale acceptoarele sunt dispozitive care sunt în stare să „recunoască” dacă un șir de simbolii face parte dintr-un limbaj pentru care mecanismul este acceptor.

1.3.1 Automate finite

Definiția 1.3.1. Se numește automat finit un obiect matematic $AF = (Q, T, m, s_0, F)$ unde:

- Q reprezintă mulțimea finită a stărilor
- T este o mulțime finită de elemente numită alfabet de intrare
- m este funcția parțială a stării următoare $m: Q \times (T \cup \{ X \}) \rightarrow P(Q)$ unde prin $P(Q)$ s-a notat mulțimea părților lui Q
- $s_0 \in Q$ este o stare numită starea de start
- $F \subset Q$ este o mulțime de stări numită mulțimea stărilor finale sau de acceptare

Definiția 1.3.2. Se numește graf de tranziție pentru automatul finit $AF = (Q, T, m, s_0, F)$ un graf orientat cu arce etichetate $G = (N, A)$ în care nodurile (mulțimea N) reprezintă stările automatului finit, iar arcele (mulțimea $A \subset N \times N$) sunt definite astfel: $(s_i, s_j) \in A$ dacă există $a \in T$ astfel încât $s_j \in m(s_i, a)$. În acest caz, arcul (S_j, s_j) este etichetat cu simbolul a .

Definiția 1.3.3. Se spune că un șir $w \in T^*$ este acceptat de automatul finit AF dacă există o cale în graful de tranziție între starea de start și o stare finală, astfel încât șirul simbolilor care etichetează arcele este șirul w . Mulțimea șirurilor acceptate de un automat finit formează limbajul acceptat de automatul finit respectiv.

Definiția 1.3.4. Se numește automat finit determinist un automat finit $AF = (Q, T, m, s_0, F)$ pentru care funcția de tranziție $m: Q \times T \rightarrow Q$. Se observă că în acest caz: nu există X -tranziții

pentru orice $s \in Q$ și orice $a \in T$ există o unică stare s' e S astfel încât $m(s, a) = s'$.

Propoziție. Pentru orice automat finit nedeterminist (AFN) există un automat finit determinist (AFD) care acceptă același limbaj.

Propoziție Limbajele acceptate de automate finite sunt limbaje regulate (generate de gramatici regulate).

Având în vedere că limbajele regulate sunt generate și de expresii regulate, există o echivalență ca putere între gramatici regulate, expresii regulate și automate finite. Automatele finite deterministe sunt utilizate pentru implementarea analizoarelor lexicale. Expresiile regulate sunt utilizate pentru specificare atomilor lexicali recunoscuți de un analizor lexical în mod corespunzător au fost elaborați algoritmi pentru construirea de automate finite deterministe direct din expresii regulate.

Probleme

problema 1.3-1

i Să se construiască automatul finit care accepta limbajul generat de gramatica:

- $P \rightarrow 0 Q \mid 1 P$
- $Q \rightarrow 0 R \mid 1 P$
- $R \rightarrow 0 R \mid 1 U$

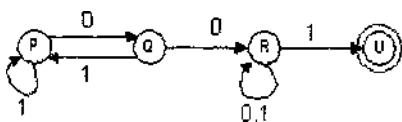
Să se reprezinte graful de tranziție corespunzător.

Soluție;

Gramatica este în mod evident regulată, prin urmare există un automat finit care acceptă limbajul generat de aceasta. Vom construi acest automat asociind fiecărui neterminat o stare și fiecărei producții o tranziție. Rezultă că putem construi următoarea funcție de tranziție:

- $m(P,0) = \{Q\}, m(P, 1) = \{P\}$
- $m(Q,0) = \{R\}, m(Q,1) = \{P\}$
- $m(R,0) = \{R\}, m(R, 1) = \{R, U\}$

unde U este o stare nouă introdusă, în care se va produce acceptarea (deci va fi singura stare finală). Prin urmare, automatul finit care acceptă limbajul generat de gramatica de mai sus este: $AF = (\{ P, Q, R, U \}, \{ 0, 1 \}, m, P, \{ U \})$. Corespunzător, graful de tranziție asociat este:



Comentarii:

Se observă că în acest caz, automatul finit obținut este nedeterminist (pentru starea R și simbolul 1 există două stări succesoare posibile: R și U).

problema 1.3-2

ga se construiască automatul finit care acceptă limbajul generat de gramatica G:

- $S \rightarrow 0 A \mid 1 S \quad X, A \rightarrow 0 B \mid 1 A, B \rightarrow 0 S \mid 1 B$

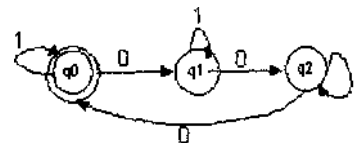
Să se reprezinte graful de tranziție asociat.

Soluție:

Fie următoarea mulțime de stări $Q = \{ q_0, q_1, q_2 \}$ astfel încât asociem fiecărui neterminat o stare din Q: lui S i se asociază q_0 , lui A i se asociază q_1 , lui B i se asociază q_2 . Corespunzător construim funcția de tranziție în modul următor:

- $m(q_0, 0) = q_1, m(q_0, 1) = q_0, m(q_1, 0) = q_2, m(q_1, 1) = q_1, m(q_2, 0) = q_0, m(q_2, 1) = q_2$

Automatul finit este $AF = (\{ q_0, q_1, q_2 \}, \{ 0, 1 \}, m, q_0, \{ q_0 \})$. Se observă că AF este determinist. Graful de tranziție corespunzător este:



Comentam;

Șirurile acceptate conțin un număr de simbolii 0 divizibil cu 3. Automatul obținut este determinist.

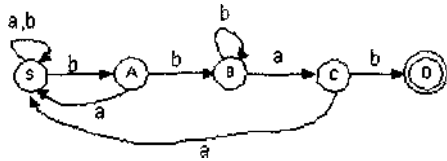
problema 1.3-1

Să se construiască gramatica regulată pentru limbajul $L = \{ w \in \{a,b\}^* \mid w = ubbab, u \in \{a, j b\}^* \}$ construind mai întâi graful de tranziție asociat.

Soluție:

Fie S starea inițială a automatului. Trebuie să accepte șirurile care se termină cu bbab. Din starea S trecem în starea A în momentul în care apare un b (am recunoscut primul b din subșirul bbab). Din starea A se trece în starea B dacă s-a întâlnit un b (s-a recunoscut șirul bb) și se trece în starea S dacă s-a primit un a (în continuare se încearcă să se identifice subșirul bbab).

Din starea B se trece în starea C dacă s-a întâlnit un a (s-a recunoscut bba) și se rămâne în B pentru b. Din starea C se trece în starea D dacă s-a întâlnit un a, adică în starea D arii găsit \ bbab, deci D este stare finală. Din starea C se revine în S dacă s-a întâlnit un a, se reîncearcă să se identifice subșirul bbab. în starea S se poate rămâne dacă s-a întâlnit a sau b (pentru a se putea accepta șirul u).



Pentru a construi gramatica fiecărei stări i se asociază un neterminal. Rezultă gramatica $G = \langle \{ S, A, B, C \}, \{ a, b \}, P, S \rangle$ unde: $P = \{ S \rightarrow a S \mid b S \mid b A, A \rightarrow b B \mid a S, B \rightarrow b B \mid a C, C \rightarrow a S \}$

! .3-4

Să se construiască gramatica regulată care generează limbajul $L = \{ a_1 a_2 \dots a_n \mid n > 3, a, e \in \{ x, y \}, a_{n-2} = y \}$ construind mai întâi graful de tranziție asociat.

Soluție:

Fie S starea inițială a automatului. În starea S se acceptă x și y pentru a se genera orice prefix. De asemenea din S se trece în starea A dacă s-a întâlnit y (s-a recunoscut a_{n-2}). Din starea A se trece în starea B dacă s-a întâlnit x sau y. Din starea B se trece în starea C dacă s-a întâlnit x sau y (C este și stare finală).

; Pentru a construi gramatica, fiecărei stări i se asociază un neterminal. Rezultă gramatica $G = \langle \{ S, A, B, C \}, \{ x, y \}, P, S \rangle$ unde: $P = \{ S \rightarrow x S \mid y S \mid y A, A \rightarrow x B \mid y B, B \rightarrow x y \}$.

problema !.3-5

i Să se construiască automatul finit care acceptă limbajul descris de expresia regulată:

$(a \mid b)^* a^* b^+ a^*$

Să se reprezinte graful de tranziție corespunzător.

Soluție:

i Expresia regulată $(a \mid b)^* a^* b^+ a^*$ generează limbajul generat și de gramatica G determinată în Problema 1.2.5, adică: $S \rightarrow a S \mid b S \mid A, A \rightarrow a A \mid b B, B \rightarrow b B \mid C, C \rightarrow a C \mid X$

Construim automatul finit care acceptă limbajul descris de expresia regulată de mai sus.

Funcția de tranziție parțială pentru automatul finit care acceptă limbajul descris de gramatica G este: $m(S, a) = \{ S \}, m(S, b) = \{ S \}, m(S, X) = \{ A \}, m(A, a) = \{ A \}, m(A, b) = \{ B \}, m(B, b) = \{ B \}, m(B, X) = \{ C \}, m(C, a) = \{ C \}$

! Deoarece există producția $C \rightarrow X$, starea asociată neterminalului C este de acceptare (sau finală). Același lucru s-ar fi întâmplat pentru orice producție de forma $A \rightarrow a$, cu a e T.

Automatul finit care acceptă limbajul descris de expresia $(a \mid b)^* a^* b^+ a^*$ este AF = $\langle \{ S, A, B, C \}, \{ a, b, c \}, m, S, \{ C \} \rangle$

Comentarii:

Automatul AF astfel construit este nedeterminist (are X-tranziții).

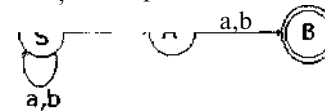
problema 1,3-6

Să se construiască automatul finit care acceptă limbajul descris de expresia regulată:

$(a \mid b)^* a (a \mid b)$

Soluție:

Gramatica G care generează limbajul descris de expresia regulată $(a \mid b)^* a (a \mid b)$ este (vezi Problema 1.2.6): $S \rightarrow a S \mid b S \mid a A, A \rightarrow a \mid b$. Asociem fiecărui neterminal o stare și fiecărei producții o tranziție. Corespunzător, obținem următoarea funcție de tranziție parțială: $m(S, a) = \{ S, \hat{A} \}, m(S, b) = \{ S \}, m(A, a) = \{ B \}, m(A, b) = \{ B \}$ unde B este o stare de acceptare (finală) nou introdusă. Se observă ca automatul finit este nedeterminist. Graful de tranziție corespunzător este:

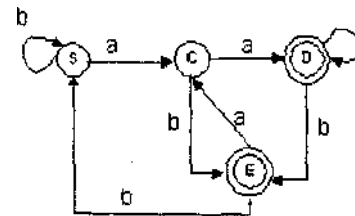


Dacă dorim să construim automatul finit determinist care acceptă limbajul descris de $(a \mid b)^* a (a \mid b)$, putem să efectuăm transformări asupra gramaticii G astfel încât din mulțimea producțiilor să rezulte o funcție de tranziție de tipul: $m: Q' \times T \rightarrow Q'$. Pentru aceasta, transformăm gramatică G prin factorizare stânga (terminalul a începe și în producția $S \rightarrow a$). Obținem: $S \rightarrow a C \mid b S, A \rightarrow a \mid b, C \rightarrow S \mid A$.

Prin substituție de începuturi ("corner substitution"), obținem: $S \rightarrow a C \mid b S, A \rightarrow a \mid b, C \rightarrow a C \mid b S \mid a \mid b$, unde observăm că neterminalul A este neutilizat (deci îl putem elimina). Aplicând din nou factorizarea stânga urmată de substituție de începuturi, obținem:

$SH \rightarrow a C \mid b S, C \rightarrow a D \mid b E, D \rightarrow C \mid X, E \rightarrow S \mid X$, respectiv, $S \rightarrow a C \mid b S, C \rightarrow a D \mid b E, D \rightarrow a D \mid b E \mid X, E \rightarrow a C \mid b S \mid X$. Se observă că în acest caz se poate construi o funcție de tranziție totală de forma: $m(S, a) = C, m(S, b) = S, m(C, a) = D, m(C, b) = E, m(D, a) = D, m(D, b) = E, m(E, a) = C, m(E, b) = S$

; Graful de tranziție pentru acest automat finit determinist este:



Stările asociate neterminalelor D, respectiv E sunt stări de acceptare datorită existenței producțiilor: $DH \rightarrow X$ și $E \rightarrow X$.

Comentarii:

Se poate construi automatul finit determinist care acceptă limbajul de mai sus pornind direct de la expresia regulată.

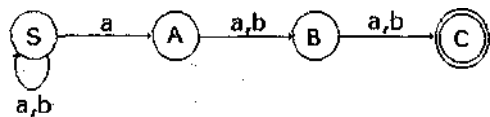
proMwii i I '- "

Să se construiască automatul finit care acceptă limbajul generat de expresia regulată $(a \mid b)^* a (a \setminus b) (a \mid b)$

Soluție.

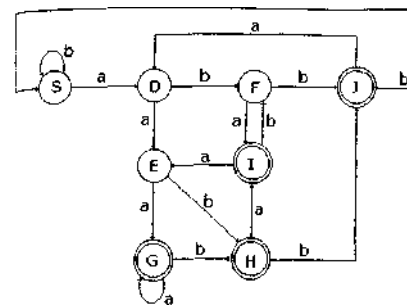
Fie gramatica G care generează limbajul descris de expresia regulată $(a \mid b)^* a (a \mid b) (a \mid b)$ (vezi Problema 1.2.7.): $S \rightarrow a S \mid b S \mid a A$, $A \rightarrow a B \mid b B$, $B \rightarrow a \mid b$. Construim pe baza mulțimii de producții de mai sus următoarea funcție de tranziție: $m(S, a) = \{ S, A \}$, $m(S, b) = \{ S \}$, $m(A, a) = \{ B \}$, $m(A, b) = \{ B \}$, $m(B, a) = \{ C \}$, $m(B, b) = \{ C \}$ unde C este o stare de acceptare nou introdusă.

Se observă că automatul finit obținut: $AF = (\{ S, A, B, C \}, \{ a, b \}, m, S, \{ C \})$ este nedeterminist. Graful, corespunzător este:



Pentru a obține automatul finit determinist care acceptă limbajul generat de gramatica G de mai sus, se pot face transformări asupra gramaticii. Prin *factorizare stânga* (terminalul a începe și producția $S \rightarrow a S$ și producția $S \rightarrow a A$) se obține: $S \rightarrow a D \mid b S$, $A \rightarrow a B \mid b B$, $B \rightarrow a \mid b$, $D \rightarrow S \mid A$. Prin *substituție de începuturi* ("corner substitution"), obținem: $S \rightarrow a D \mid b S$, $B \rightarrow a \mid b$, $D \rightarrow a E \mid b F$, $E \rightarrow D \mid B$, $F \rightarrow S \mid B$, respectiv, $S \rightarrow a D \mid b S$, $D \rightarrow a E \mid b F$, $E \rightarrow a E \mid b F \mid a \mid b$, $F \rightarrow a D \mid b S \mid a \mid b$. Aplicând din nou factorizarea stânga urmată de substituție de începuturi, obținem: $S \rightarrow a D \mid b S$, $D \rightarrow a E \mid b F$, $E \rightarrow a G \mid b H$, $F \rightarrow a I \mid b J$, $G \rightarrow E \mid A$, $H \rightarrow F \mid A$, $I \rightarrow D \mid A$, $J \rightarrow S \mid A$, respectiv, $S \rightarrow a D \mid b S$, $D \rightarrow a E \mid b F$, $E \rightarrow a G \mid b H$, $F \rightarrow a I \mid b J$, $G \rightarrow a G \mid b H \mid X$, $H \rightarrow a I \mid b J \mid A$, $I \rightarrow a E \mid b F \mid A$, $J \rightarrow a D \mid b S \mid A$.

Se observă că în acest caz se poate construi o funcție de tranziție totală de forma: $m(S, a) = D$, $m(S, b) = S$, $m(D, a) = E$, $m(D, b) = F$, $m(E, a) = G$, $m(E, b) = H$, $m(F, a) = I$, $m(F, b) = J$, $m(G, a) = G$, $m(G, b) = H$, $m(H, a) = I$, $m(H, b) = J$, $m(I, a) = E$, $m(I, b) = F$, $m(J, a) = D$, $m(J, b) = S$. Graful de tranziție al automatului determinist este:

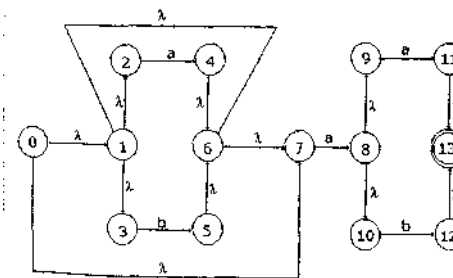


problema 1.3>#

Să se construiască AFN (automatul finit nedeterminist) care generează același limbaj ca și următoarea expresie regulată și apoi să se construiască AFD-ul (automatul finit determinist) corespunzător. Să se facă și construcția directă: expresie regulată-AFD $(a \mid b)^* a (a \mid b)$

Soluție:

AFN corespunzător expresiei regulate $(a \mid b)^* a (a \mid b)$ este:



Construcția AFD se face pe baza algoritmului de construire a unui AFD echivalent unui AFN q_0 (vezi curs pentru notații și algoritmul). Starea inițială a AFD va fi: $q_0 = X$ închidere($\{ 0 \}$) = $\{ 0, 1, 2, 3, 7 \}$.

în continuare, celelalte stări precum și tranzițiile corespunzătoare se vor determina astfel:

$(x) = q_j$ unde $q_j = \text{inchiidere}(\text{move}(q_i, x))$, obținem:

$q_1 = A_inchiidere(\text{move}(q_0, a)) = A_inchiidere(\{ 4, 8 \}) = \{ 1, 2, 3, 4, 6, 7, 8, 9, 10 \}$

$q_2 = _inchiidere(\text{move}(q_0, b)) = _inchiidere(\{ 5 \}) = \{ 1, 2, 3, 5, 6, 7 \}$

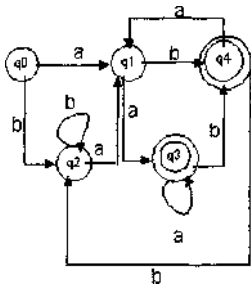
$q_3 = M_inchiidere(\text{move}(q_1, a)) = A_inchiidere(\{ 4, 8, 11 \}) = \{ 1, 2, 3, 4, 6, 7, 8, 9, 10, 11, 13 \}$

$q_3^b = U_inchiidere(\text{move}(q_1, b)) = I_inchiidere(\{ 5, 12 \}) = \{ 1, 2, 3, 5, 6, 7, 12, 13 \}$

Se obține următoarea tabelă de tranziții:

stare	stare următoare	
	a	b
q0	q1	q2
q1	q3	q4
q2	q1	q2
q3	q3	q4
q4	q1	q2

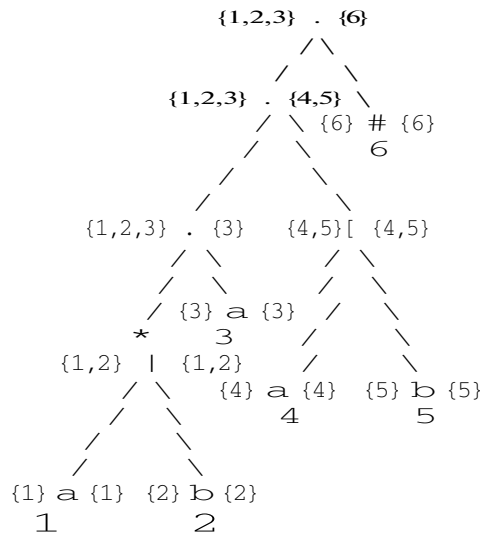
și corespunzător, graful de tranziție:



Dacă dorim acum să construim direct automatul finit determinist care acceptă expresia de mai sus, trebuie să utilizăm algoritmul de construire a AFD direct din expresia regulată (vezi curs). Considerăm arborele corespunzător expresiei regulate, terminate cu terminatorul #:

$(a|b)^* a(a|b) \#$
 1 2 3 4 5 6

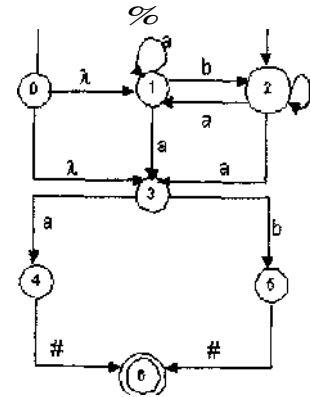
Arborele este:



, Am notat în stânga fiecărui nod firstpos(nod), iar în dreapta lastpos(nod). Calculând pentru i fiecare nod frunză followpos(i) (i codul unui nod frunză), se obține:

nod	followpos
a 1	{ 1, 2, 3 }
b 2	{ 1, 2, 3 }
a 3	{ 4, 5 }
a 4	{ 6 }
b 5	{ 6 }
# 6	-

Corespunzător tabelii de mai sus, se obține automatul finit nedeterminist reprezentat prin următorul graf de tranziție (vom asocia fiecărui nod frunză o stare în graful de tranziție), arcele sunt stabilite după următoarea regulă: există un arc între nodul i și nodul j dacă $j \in \text{followpos}(i)$. Arcul se etichetează cu simbolul corespunzător codului j. De asemenea, se introduce o stare inițială 0 din care există λ -tranziții în stările din firstpos(rad) (rad este rădăcina arborelui corespunzător expresiei regulate). Vom obține următorul graf de tranziție:

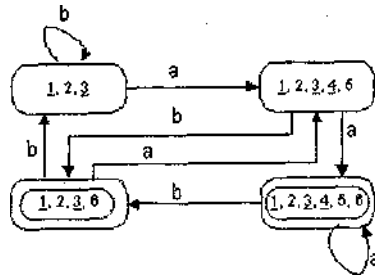


Se observă că în stările 4, respectiv 5 este acceptat sufixul aa, respectiv ab. Automatul finit determinist corespunzător expresiei regulate se obține considerând că stare inițială firstpos(rad) = { 1, 2, 3}. Tranzițiile se determină astfel:

$$m(q_i, x) = \cup \{ \text{followpos}(i) \} = q_j$$

$$i = \text{cod}(x)$$

Se obține următorul graf de tranziție pentru AFD:



Am subliniat în fiecare mulțime nodurile i pentru care este îndeplinită condiția $i = \text{cod}(a)$.

Comentarii:

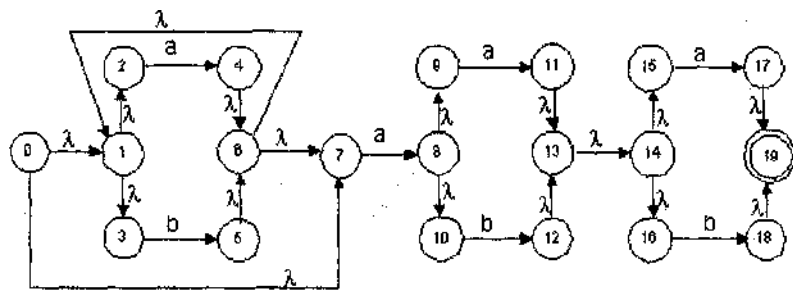
Se observă că în starea $\{1,2,3,4,5,6\}$ se acceptă sufixul aa, iar în starea $\{1,2,3,6\}$ se acceptă sufixul ab. Ambele sunt stări terminale deoarece conțin codul pentru # (terminatorul expresiei regulate).

Să se construiască AFN pentru următoarea expresie regulată și apoi să se construiască AFD-ul corespunzător. Să se facă și construcția directă: expresie regulată-AFD:

$$(a|b)^*a(a|b)(a|b)$$

Soluție:

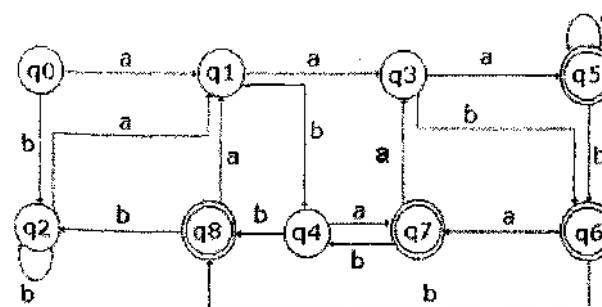
AFN corespunzător expresiei regulate de mai sus este:



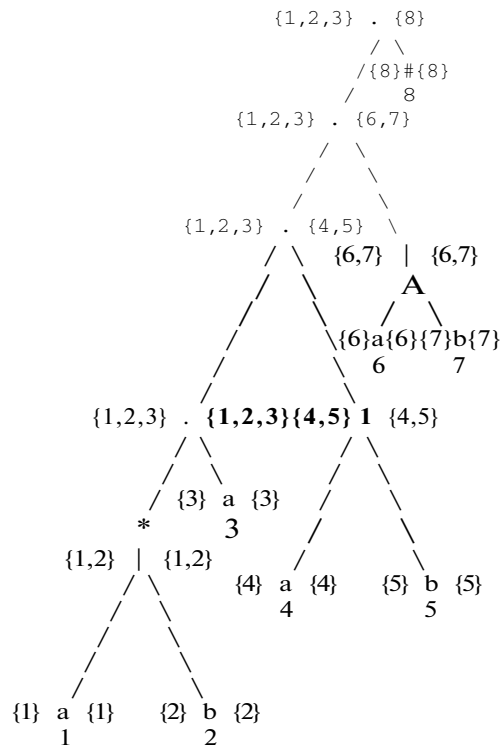
Pentru a obține AFD din AFN-ul de mai sus, starea inițială a noului automat determinist va fi: $q_0 = A_inchidere(\{0\}) = \{0, 1, 2, 3, 7\}$ deoarece starea inițială a AFN este starea 0. în continuare, vom determina stările următoare și tranzițiile utilizând:
 $m(q_i, x) = \hat{_}inchidere(\text{move}(q_i, x))$

Vom obține: $m(q_0, a) = A_inchidere(\text{move}(q_0, a)) = X_inchidere(\{4, 8\}) = \{1, 2, 3, 4, 6, 7, 8, 9, 10\} = q_1$, $m(q_0, b) = A_inchidere(\text{move}(q_0, b)) = A_inchidere(\{5\}) = \{1, 2, 3, 5, 6, 7\} = q_2$, $m(q_1, a) = X_inchidere(\text{move}(q_1, a)) = X_inchidere(\{4, 8, 11\}) = \{1, 2, 3, 4, 6, 7, 8, 9, 10, 11, 13, 14, 15, 16\} = q_3$, $m(q_1, b) = _inchidere(\text{move}(q_1, b)) = \hat{_}inchidere(\{5^{\wedge} : 12\}) = \{1, 2, 3, 5, 6, 7, 12, 13, 14, 15, 16\} = q_4$, $m(q_2, a) = A_inchidere(\text{move}(q_2, a)) = X_inchidere(\{4, 8\}) = q_1$, $m(q_2, b) = A_inchidere(\text{move}(q_2, b)) = A_inchidere(\{1, 2, 3, 5, 6, 7\}) = q_2$, $m(q_3, a) = X_inchidere(\text{move}(q_3, a)) = X_inchidere(\{4, 8, 11, 17\}) = \{1, 2, 3, 4, 6, 7, 8, 9, 10, 11, 13, 14, 15, 16, 17, 19\} = q_5$, $m(q_3, b) = \hat{_}inchidere(\text{move}(q_3, b)) = X_inchidere(\{5, 12, 18\}) = \{1, 2, 3, 5, 6, 7, 12, 13, 14, 15, 16, 18, 19\} = q_6$, $m(q_4, a) = X_inchidere(\text{move}(q_4, a)) = X_inchidere(\{4, 8, 17\}) = \{1, 2, 3, 4, 6, 7, 8, 9, 10, 17, 19\} = q_7$, $m(q_4, b) = A_inchidere(\text{move}(q_4, b)) = A_inchidere(\{5, 18\}) = \{1, 2, 3, 5, 6, 7, 18, 19\} = q_8$, $m(q_5, a) = A_inchidere(\text{move}(q_5, a)) = A_inchidere(\{4, 8, 11, 17\}) = \{1, 2, 3, 4, 6, 7, 8, 9, 10, 11, 13, 14, 15, 16, 17, 19\} = q_5$, $m(q_5, b) = X_inchidere(\text{move}(q_5, b)) = X_inchidere(\{5, 12, 18\}) = \{1, 2, 3, 5, 6, 7, 12, 13, 14, 15, 16, 18, 19\} = q_6$, $m(q_6, a) = X_inchidere(\text{move}(q_6, a)) = _inchidere(\{4, 8, 17\}) = \{1, 2, 3, 4, 6, 7, 8, 9, 10, 17, 19\} = q_7$, $m(q_6, b) = \hat{_}inchidere(\text{move}(q_6, b)) = A_inchidere(\{5, 18\}) = \{1, 2, 3, 5, 6, 7, 18, 19\} = q_8$, $m(q_7, a) = X_inchidere(\text{move}(q_7, a)) = \hat{_}inchidere(\{4, 8, 11\}) = \{1, 2, 3, 4, 6, 7, 8, 9, 10, 11, 13, 14, 15, 16\} = q_3$, $m(q_7, b) = _inchidere(\text{move}(q_7, b)) = X_inchidere(\{5, 12\}) = \{1, 2, 3, 5, 6, 7, 12, 13, 14, 15, 16\} = q_4$, $m(q_8, a) = X_inchidere(\text{move}(q_8, a)) = i$, $X_inchidere(\{4, 8\}) = \{1, 2, 3, 4, 6, 7, 8, 9, 10\} = q_1$, $m(q_8, b) = A_inchidere(\text{move}(q_8, b)) = i$, $= A_inchidere(\{5\}) = \{1, 2, 3, 5, 6, 7\} = q_2$.

O altă soluție se obține pornind direct de la expresia regulată construind AFN, respectiv AFD ! corespunzătoare. în mod corespunzător funcției de tranziție de mai sus, obținem următorul graf de tranziție asociat AFD. Stările finale vor fi toate stările care conțin starea finală 19 din AFN inițial:



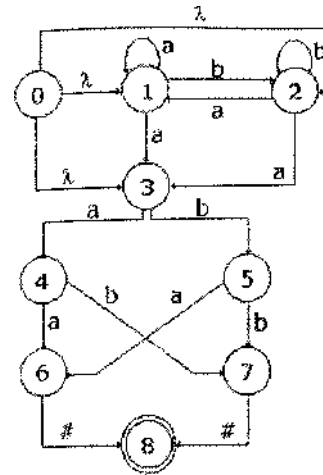
Arborele corespunzător expresiei regulate este:



în stânga fiecărui nod se află firstpos(nod), iar în dreapta lastpos(nod). Calculând pentru fiecare frunză followpos(i) (i este codul frunzei), obținem:

nod	followpos
a 1	{ 1, 2, 3 }
b 2	{ 1, 2, 3 }
a 3	{ 4, 5 }
a 4	{ 6, 7 }
b 5	{ 6, 7 }
a 6	{ 8 }
b 7	{ 8 }
# 8	-

Procedând în mod analog problemei precedente obținem următorul AFN:

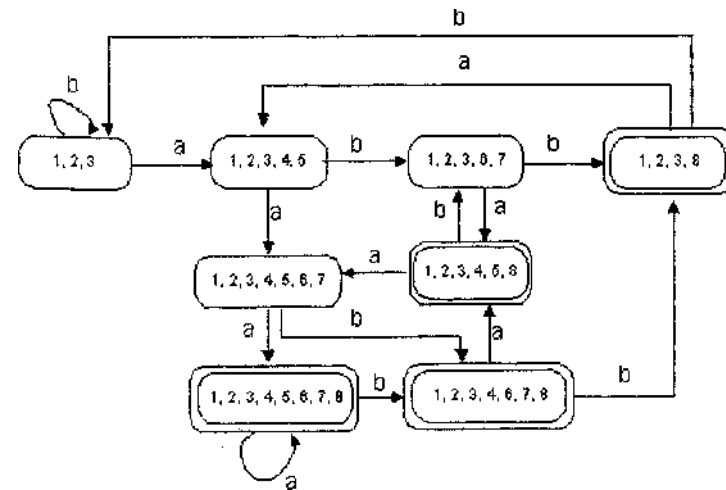


Dacă dorim să obținem AFD corespunzător automatului de mai sus, considerăm ca stare inițială firstpos(rad) = { 1, 2, 3 }. Tranzițiile se obțin astfel:

$$m(q_i, x) = v \text{ followpos}(j) = q_k$$

$$j = \text{cod}(x)$$

Obținem următorul automat finit determinist (în care sunt marcate ca finale stările ce conțin codul 8 corespunzător simbolului terminator #):



problem;! 1.3-10

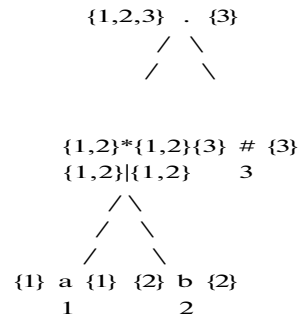
Să se construiască automatul finit determinist pentru expresia regulată: $(a|b)^*$.

Soluție;

Considerăm expresia regulată terminată cu simbolul # și codificăm în mod corespunzător simbolii:

$(a|b)^* \#$
 1 2 3

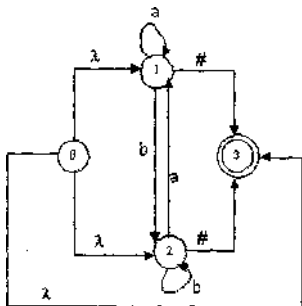
Pentru expresia regulată de mai sus, arborele corespunzător este:



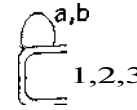
Determinând pentru fiecare cod i followpos(i) pe baza mulțimilor firstpos(nod) (specificată în stânga nodurilor), respectiv lastpos(nod) (specificată în dreapta nodurilor), se obține următorul tabel:

nod	followpos
a 1	{ 1, 2, 3 }
b 2	{ 1, 2, 3 }
# 3	-

Graful de tranziție corespunzător automatului finit se determină asociind fiecărui cod o stare, iar arcele se determină astfel: există un arc între starea i și starea j dacă j e followpos(i), iar eticheta arcului este simbolul cu codul j . Obținem:

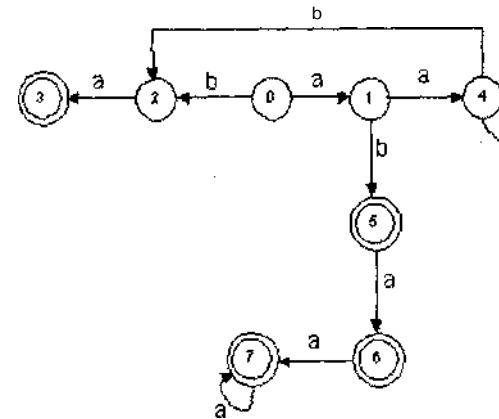


Am introdus o stare de start (0) din care există \wedge -tranziții în stările corespunzătoare elementelor din mulțimea firstpos(rad) = { 1, 2, 3 }. De asemenea, starea 3 corespunzătoare simbolului terminal # este stare finală. Automatul finit obținut astfel este nedeterminist. Dacă dorim să obținem automatul finit determinist, vom considera starea de start firstpos(rad) = { 1, 2, 3 }, iar tranzițiile se determină ca în problema anterioară. Obținem următorul automat finit determinist:



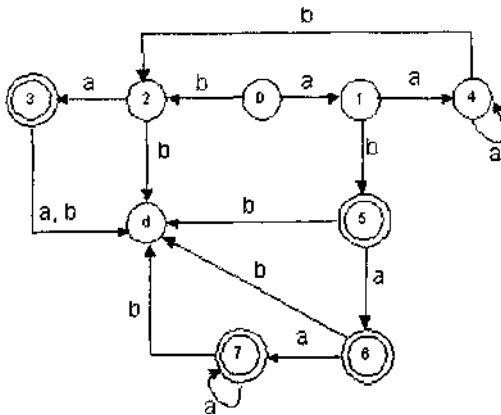
Să se minimizeze automatele specificate prin grafurile de tranziție:

prolik'ma 1.3-11



Soluție;

Introducem o nouă stare d astfel încât funcția de tranziție să fie totală, adică pentru orice stare q pentru care $m(q, x)$ nu este definită ($x \in \{ a, b \}$), adăugăm $m(q, x) = d$. Noul graf de tranziție devine:



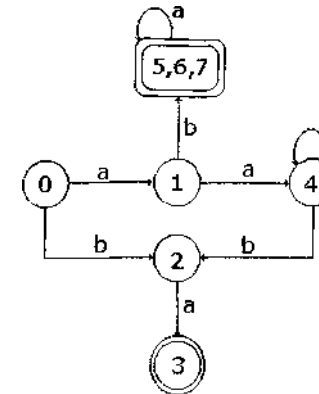
Fie partiția disjunctă pe mulțimea stărilor: $P = \{ \{ 0, 1, 2, 4 \}, \{ 3, 5, 6, 7 \}, \{ d \} \}$ unde prima mulțime este mulțimea stărilor nefmale, iar a doua este mulțimea stărilor finale. în continuare, pentru fiecare mulțime $Q \in P$ și toate stările $q \in Q$ pentru care: $m(q, x) \in Q'$ astfel încât există q' cu $m(q', y) \in Q'$ (1). Se construiește o nouă partiție: $P = (P \setminus \{ Q \}) \cup \{ Q \setminus new, new \}$ unde new este mulțimea stărilor care îndeplinesc condiția (1). într-adevăr, în cazul nostru obținem. Pentru $Q = \{ 0, 1, 2, 4 \}$, tranzițiile stărilor componente sunt:

$m(0, a) = 1 \in Q$ $m(0, b) = 2 \in Q$, $m(1, a) = 4 \in Q$ $m(1, b) = 5 \notin Q$ deci 1 e new
 $m(2, a) = 3 \notin Q$ deci 2 e new, $m(4, a) = 4 \in Q$ $m(4, b) = 2 \in Q$

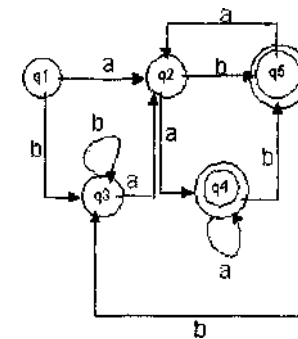
Prin urmare, noua partiție pe mulțimea stărilor este: $P = \{ \{0,4\}, \{1,2\}, \{3,5,6,7\}, \{d\} \}$.
 Pentru $Q = \{ 0, 4 \}$, tranzițiile stărilor sunt: $m(0, a) = 1 \in \{ 1, 2 \}$ $m(0, b) = 2 \in \{ 1, 2 \}$,
 $m(4, a) = 4 \in \{ 1, 2 \}$ deci 4 e new. Partiția obținută pe mulțimea stărilor devine: $P = \{ \{0\}, \{ 4 \}, \{ 1, 2 \}, \{ 3, 5, 6, 7 \}, \{ d \} \}$. Pentru $Q = \{ 1, 2 \}$ (este prima mulțime din partiție care conține mai mult de un element) rezultă: $m(1, a) = 4 \in \{ 4 \}$ $m(1, b) = 5 \in \{ 3, 5, 6, 7 \}$,
 $m(2, a) = 3 \notin \{ 4 \}$ deci 3 e new.

Partiția astfel obținută este: $P = \{ \{ 0 \}, \{ 1 \}, \{ 2 \}, \{ 4 \}, \{ 3, 5, 6, 7 \}, \{ d \} \}$
 Pentru $Q = \{ 3, 5, 6, 7 \}$ tranzițiile stărilor componente sunt: $m(3, a) = d \in \{ d \}$ $m(3, b) = d \in \{ d \}$,
 $m(5, a) = 6 \in \{ d \}$ deci 5 e new, $m(6, a) = 7 \in \{ d \}$ deci 6 e new, $m(7, a) = 7 \in \{ d \}$ deci 7 e new.

Partiția obținută în acest caz este: $P = \{ \{ 0 \}, \{ 1 \}, \{ 2 \}, \{ 3 \}, \{ 4 \}, \{ 5, 6, 7 \}, \{ d \} \}$. Fie acum $Q = \{ 5, 6, 7 \}$. Tranzițiile stărilor componente sunt: $m(5, a) = 6 \in \{ 5, 6, 7 \}$ $m(5, b) = d \in \{ d \}$,
 $m(6, a) = 7 \in \{ 5, 6, 7 \}$ $m(6, b) = d \in \{ d \}$
 $m(7, a) = 7 \in \{ 5, 6, 7 \}$ $m(7, b) = d \in \{ d \}$, adică nu mai apare o nouă partiție pe mulțimea stărilor. Prin urmare, automatul finit minimizat este:



.i.ihluu.il.*-12

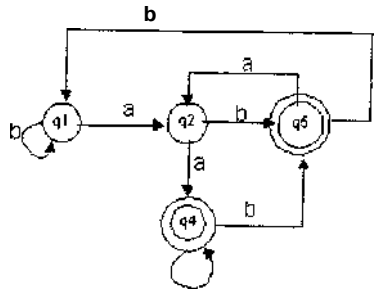


Soluție;

Funcția de tranziție este totală, deci putem considera următoarea partiție pe mulțimea stărilor:
 $P = \{ \{q1, q2, q3\}, \{q4, q5\} \}$

Aplicând același algoritm de la punctul a), considerăm mulțimea: $Q = \{ q1, q2, q3 \}$ pentru care tranzițiile stărilor componente sunt: $m(q1, a) = q2 \in \{ q2 \}$, $m(q1, b) = q3 \in \{ q1, q3 \}$,
 $m(q2, a) = q4 \notin \{ q1, q2, q3 \}$ deci $q2$ e new, $m(q3, a) = q2 \in \{ q2 \}$, $m(q3, b) = q3 \in \{ q1, q3 \}$. Noua partiție pe mulțimea stărilor devine: $P = \{ \{q1, q3\}, \{q2\}, \{q4, q5\} \}$.
 Fie $Q = \{ q1, q3 \}$. Tranzițiile pentru această mulțime de stări: $m(q1, a) = q2 \notin \{ q2 \}$,
 $m(q1, b) = q3 \in \{ q1, q3 \}$, $m(q3, a) = q2 \in \{ q2 \}$, $m(q3, b) = q3 \in \{ q1, q3 \}$

Fie acum $Q = \{ q4, q5 \}$. Tranzițiile pentru stările componente sunt: $m(q4, a) = q4 \in \{ q4, q5 \}$
 $m(q4, b) = q5 \in \{ q4, q5 \}$, $m(q5, a) = q2 \notin \{ q4, q5 \}$ deci $q5$ e new. Noua partiție pe mulțimea stărilor este: $P = \{ \{ q1, q3 \}, \{ q2 \}, \{ q4 \}, \{ q5 \} \}$. Automatul finit minimizat, P^{en} tru care starea corespunzătoare mulțimii $\{ q1, q3 \}$ este notată cu $q1$, este:



1.3.2 Automate cu stivă (push-down)

Definiția 1.3.5. Se numește automat cu stivă un obiect matematic definit în modul următor: $P = (Q, T, Z, m, q_0, z_0, F)$ unde:

- Q - este o mulțime finită de simboluri ce reprezintă stările posibile pentru unitatea de control a automatului;
- T - este mulțimea finită a simbolurilor de intrare;
- Z - este mulțimea finită a simbolurilor utilizați pentru stivă;
- m - este o funcție, $m : Q \times (T \cup \{X\}) \times Z \rightarrow P(S \times Z^*)$ (prin $P(Q)$ s-a notat mulțimea părților lui Q) este funcția care descrie modul în care se obține starea următoare și informația care se introduce în stivă pentru o combinație (stare, intrare, conținut stivă) | dată;
- $q_0 \in Q$ este starea inițială a unității de control;
- $z_0 \in Z$ este simbolul aflat în vârful stivei în starea inițială;
- $F \subseteq Q$ reprezintă mulțimea finită a stărilor finale.

Definiția 1.3.6. O configurație de stare a automatului este un triplet $(q, w, a) \in Q \times T^* \times Z^*$ unde:

- q - reprezintă starea curentă a unității de control;
- w - reprezintă partea din banda de intrare care nu a fost încă citită. Dacă $w = X$ înseamnă că s-a ajuns la sfârșitul benzii de intrare;
- a - reprezintă conținutul stivei.

automatului este relația $|$ - definită asupra mulțimii [următor: $(q, aw, za) |$ - (q', w, Pa) unde $(q', P) \in m(q, a, z), q' \in Q, a \in T \cup \{1\}, w \in T^*, z \in Z, a \in Z^*, fie Z^*$.

Dacă $a = X$ înseamnă că, dacă unitatea de control este în starea q , capul de citire este pe simbolul a iar simbolul din vârful stivei este z atunci automatul poate să își schimbe configurația în modul următor: starea unității de control devine q' , capul de citire se deplasează cu o poziție la dreapta iar simbolul din vârful stivei se înlocuiește cu p .

Dacă $a = X$ înseamnă că avem o A.-tranzitie pentru care simbolul aflat în dreptul capului de citire pe banda de intrare nu contează (capul de citire nu se va deplasa), însă starea unității de control și conținutul memoriei se pot modifica. O astfel de tranziție poate să aibă loc și după ce s-a parcurs întreaga banda de intrare.

Dacă se ajunge într-o configurație pentru care stiva este goală nu se mai pot executa tranziții.

Relația $|$ - se poate generaliza la $|$ - \setminus $|$ - \setminus $|$ - \setminus , într-o manieră similară relației de derivare pentru forme prepoziționale.

Definiția 1.3.8. O configurație inițială pentru un automat cu stivă este o configurație de forma (q_0, w, z_0) unde $w \in T^*$. O configurație finală este o configurație de forma (q, X, a) cu $q \in F, a \in Z^*$.

Definiția 1.3.9. Spunem că un șir w este acceptat de un automat cu stivă prin stări finale dacă $(q_0, w, z_0) |$ - \setminus (q, X, a) pentru $q \in F$ și $a \in Z^*$. Mulțimea șirurilor acceptate de un automat cu stivă se notează cu $L(P)$.

Definiția 1.3.10. Spunem că un șir w este acceptat de un automat cu stivă prin stivă goală dacă $(q_0, w, z_0) |$ - \setminus (q, X, X) pentru $q \in Q$. Limbajul acceptat de un automat pushdown P în acest mod se notează cu $Le(P)$.

Propoziție. Dacă $L(P)$ este limbajul acceptat de automatul P prin stări finale atunci se poate construi un automat pushdown P' astfel încât $L(P) = Le(P')$.

Propoziție. Dacă $Le(P)$ este limbajul acceptat de automatul P prin stivă goală atunci se poate construi un automat pushdown P' astfel încât $Le(P) = L(P')$.

Propoziție. Mulțimea limbajelor acceptate de automate cu stivă este mulțimea limbajelor independente de context.

Probleme

Să se determine limbajul acceptat de următorul automat cu stivă:

- $PD = (\{q_0, q_1, q_2, q_3\}, \{a, b, c\}, \{a, b, c, z\}, m, q_0, z, \{q_3\})$
- $m(q_0, a, z) = \{(q_0, az)\}$
- $m(q_0, a, a) = \{(q_0, aa)\}$
- $m(q_0, c, a) = \{(q_1, a)\}$
- $m(q_1, b, a) = \{(q_1, A.)\}$
- $m(q_1, b, z) = \{(q_1, bz)\}$
- $m(q_1, b, b) = \{(q_1, bb)\}$
- $m(q_1, c, b) = \{(q_2, b)\}$
- $m(q_2, b, b) = \{(q_2, X)\}$
- $m(q_2, X, z) = \{(q_3, X)\}$

Să;

Soluție:

Evoluția automatului PD pentru un anumit șir de intrare poate fi caracterizată în modul următor: din starea inițială q_0 se poate face o tranziție dacă șirul începe fie cu simbolul a fie cu simbolul c . Se observă că simbolii a sunt introduși în stivă. Fie n numărul de simboluri a din prefixul șirului de intrare, deci șirul de intrare poate fi de forma: $w = a^n ev$ unde $n > 0$, iar v este un sufix al șirului de intrare. În continuare, se trece în starea q_1 în care poate urma doar simbolul b în șirul de intrare și pentru fiecare simbol b de pe banda de intrare, se descarcă un simbol a din vârful stivei.

Deci, numărul de simbolii b care poate urma după c este tot n, iar w este de forma; $w = a^n c b^n u$, u este un sufix al șirului de intrare. Simbolul care poate urma pe banda de intrare în momentul în care stiva este vidă (s-au descărcat toți simbolii a din stivă) este tot b.

În continuare, pot apare doar simbolii b care se introduc în stivă. Fie m numărul acestora. Configurația în acest moment a prefixului analizat din șirul de intrare este: $w = a^n c b^n b^m x$, unde x este restul din șir rămas neanalizat. Se observă că din starea q1, dacă pe banda de intrare nu urmează un b, poate urma doar un c. Se trece în starea q2 în care sunt acceptați doar simbolii b pe banda de intrare. Ori de câte ori apare un b pe banda de intrare, se descarcă un b din vârful stivei. În momentul în care stiva s-a golit, se trece în starea finală q3. Deci șirurile acceptate de automat sunt de forma: $w = a^n c b^n b^m c b^m u$, $n > 0$.

Pentru șirul de intrare aacbbcb, mișcările efectuate de automat sunt:

$(q0, aacbbcb, z) \rightarrow (q0, acbbcb, az) \rightarrow (q0, cbbcb, aaz) \rightarrow (q1, bbbcb, aaz)$
 $\rightarrow (q1, bbcb, az) \rightarrow (q1, bcb, z) \rightarrow (q1, cb, bz) \rightarrow (q2, b, bz)$

$Hq2, A, z) \rightarrow (q3, Aa)$

deci șirul este acceptat.

problema 1.3-14

Să se construiască automatul cu stivă care acceptă următorul limbaj:

$L = \{ w c w^R \mid w \in \{a, b\}^* \}$

Soluție:

Vom folosi în construcția automatului PD care acceptă limbajul L faptul că șirul de intrare este simetric și are număr impar de simbolii, simbolul din centru fiind c. Prin urmare, pentru a putea stabili dacă șirul de intrare poate fi acceptat sau nu, trebuie memorată prima jumătate a șirului (până se întâlnește simbolul c) în stivă și apoi comparată cu a doua jumătate a șirului (vom compara simbolul analizat în mod curent de pe banda de intrare cu simbolul din vârful stivei; dacă sunt egali, se descarcă stiva). Prin urmare, automatul cu stivă care acceptă limbajul L, este: $PD = (\{q0, q1, q2\}, \{a, b, c\}, \{a, b, c, z\}, m, q0, z, \{q2\})$

$m(q0, d, e) = \{(q0, de)\} \forall d \in \{a, b\}, \forall e \in \{a, b, z\}$

$m(q0, c, d) = \{(q1, d)\} \forall d \in \{a, b, z\}$

$m(q1, d, d) = \{(q1, A)\} \forall d \in \{a, b\}$

$m(qU, z) = \{(q2, A.)\}$

Comentarii:

Se observă ca automatul este determinist și acceptarea este prin stări finale.

Să se construiască automatul cu stivă care acceptă următorul limbaj:

$L = \{ 0^i 1^j \mid 0 < j \leq i \}$

Soluție:

Numărul de simbolii 0 este mai mare sau egal decât numărul de simbolii 1. Vom proceda în mod analog exemplului precedent introducând inițial în stivă toți simbolii 0 din prima parte a șirului și ulterior descărcând stiva pentru fiecare simbol 1 întâlnit pe banda de intrare. Condiția ca un șir să fie acceptat este ca stiva să nu se golească înainte de a ajunge la sfârșitul șirului. Rezultă deci automatul:

$= (\{q0, q1, q2\}, \{0, 1\}, \{0, z\}, m, q0, z, \{q2\})$

$m(q0, 0, a) = \{(q0, 0a)\} \forall a \in \{0, z\}$

$m(q0, 1, 0) = \{(q1, X)\}$

$m(q1, 1, 0) = \{(q1, X)\}$

$m(q1, ^a) = \{(q2, X)\} \forall a \in \{0, z\}$

Verificare:

Pentru șirul de intrare 00011, evoluția automatului este:

$(q0, 00011, z) \rightarrow (q0, 0011, Oz)$

$\rightarrow (q0, 011, 00z)$

$\rightarrow (q0, 11, 000z)$

$\rightarrow (q1, 1, 00z)$

$\rightarrow (q1, A, 0z)$

$\rightarrow (q2, A, z)$

prin urmare șirul este acceptat. De asemenea, pentru șirul de intrare 011, evoluția automatului este:

funcția de tranziție nu este definită în acest caz, prin urmare șirul nu este acceptat (starea q1 nu este de acceptare).

Comentarii:

Automatul PD este nedeterminist deoarece există două tranziții din starea q1 nedistinguite:

$m(qU, 0) = \{(q2, A.)\}$

problema 1.3-16

Să se construiască automatul cu stivă pentru limbajul generat de următoarea gramatică:

$S \rightarrow a S b \mid a A b$

$A \rightarrow c$

Soluție:

Putem construi un automat cu stivă care acceptă limbajul generat de gramatică prin stivă goală: $PD' = (\{q_0\}, \{a, b, c\}, \{a, b, c, S, A\}, m, q_0, S, \{\})$, $m(q_0, X, S) = \{(q_0, aSb), (q_0, aAb)\}$, $m(q_0, \hat{A}, A) = \{(q_0, c)\}$, $m(q_0, x, x) = \{(q(U))\} \forall x \in \{a, b, c\}$

O altă soluție se poate obține considerând un automat cu stivă extins (în acest caz considerăm că stiva crește spre dreapta): $PD'' = (\{q_0, q_1\}, \{a, b, c\}, \{a, b, c, z, S, A\}, m, q_0, z, \{q_1\})$, $m(q_0, x, X) = \{(q_0, x)\} \forall x \in \{a, b, c\}$, $m(q_0, X, aSb) = \{(q_0, S)\}$

$m(q_0, X, aAb) = \{(q_0, S)\}$, $m(q_0, X, c) = \{(q_0, A)\}$, $m(q_0, X, zS) = \{(q_1, X)\}$

Verificare;

Pentru șirul de intrare aacbb, se obțin evoluțiile celor două automate:

$(q_0, aacbb, S) \mid (q_0, aacbb, aSb) \mid (q_0, acbb, St) \mid (q_0, acbb, aAbb) \mid (q_0, cbb, Abb)$
 $\mid (q_0, cbb, cbb) \mid (q_0, bb, bb) \mid (q_0, b, b) \mid (q_0, X, X)$

iar pentru a doua variantă:

$(q_0, aacbb, z) \mid (q_0, acbb, za) \mid (q_0, cbb, zaa) \mid (q_0, bb, zaac) \mid (q_0, bb, zaaA)$
 $\mid (q_0, b, zaaAb) \mid (q_0, b, zaS) \mid (q_0, X, zaSb) \mid (q_0, X, zS) \mid (q_1, X, z)$

problema 1.3-17

Să se construiască automatul cu stivă care acceptă următorul limbaj:

$L = \{w \mid w \text{ conține un număr egal de } a \text{ și } b\}$

Soluție:

Vom utiliza aceeași tehnică ca și în problemele precedente și anume: deoarece numărul de simbolii a este egal cu numărul de simbolii b din șirul de intrare, putem introduce în stivă toți simbolii a, dacă șirul începe cu a, respectiv toți simbolii b, dacă șirul începe cu b. În continuare, de fiecare dată când pe banda de intrare se află un simbol diferit de cel din vârful stivei se descarcă stiva, altfel se introduce în stivă. Șirul este acceptat dacă în final stiva este vidă. Rezultă: $PD = (\{q_0\}, \{a, b\}, \{a, b, z\}, m, q_0, z, \{q_0\})$, $m(q_0, x, z) = \{(q_0, xz)\} \forall x \in \{a, b\}$, $m(q_0, x, x) = \{(q_0, xx)\} \forall x \in \{a, b\}$, $m(q_0, x, y) = \{(q(U))\} \forall x, y \in \{a, b\}, x \neq y$, $m(q_0, X, z) = \{(q_0, z)\}$

O altă soluție se poate construi observând că limbajul L este generat de gramatica: $S \rightarrow aSbS \mid bSaS \mid X$ (vezi, de exemplu, Problema 1.1.6 din paragraful 1.1.1). Prin urmare putem construi un automat cu stivă care acceptă limbajul L prin stivă goală: $PD' = (\{q_0\}, \{a, b\}, \{a, b, S\}, m, q_0, S, \{\})$, $m(q_0, X, S) = \{(q_0, aSbS), (q_0, bSaS), (q_0, X)\}$, $m(q_0, c, c) = \{(q_0, X)\} \forall c \in \{a, b\}$

În fine, o a treia soluție se poate obține considerând un automat cu stivă extins (în acest caz considerăm stiva crește spre dreapta): $PD'' = (\{q_0, q_1\}, \{a, b\}, \{a, b, z, S\}, m, q_0, z, \{q_1\})$, $m(q_0, x, X) = \{(q_0, x)\} \forall x \in \{a, b\}$, $n \neq X$, $X = \{(q_0, S)\}$, $m(q_0, X, aSbS) = \{(q_0, S)\}$, $m(q_0, X, bSaS) = \{(q_0, S)\}$, $m(q_0, X, zS) = \{(q_1, X)\}$

Verificare:

Pentru șirul de intrare abbaa, se obțin evoluțiile celor trei automate:

$(q_0, abbaa, z) \mid (q_1, bbaa, az) \mid (q_1, bbaa, z) \mid (q_1, baa, bz) \mid (q_1, aa, bbz)$
 \mid

sau:

$(q_0, abbaa, S) \mid (q_0, abbaa, aSbS) \mid (q_0, bbaa, SbS) \mid (q_0, bbaa, bS)$
 $\mid (q_0, bbaa, S) \mid (q_0, bbaa, bSaS) \mid (q_0, baa, SaS) \mid (q_0, baa, bSaSaS)$
 $\mid (q_0, aa, SaSaS) \mid (q_0, aa, aSaS) \mid (q_0, a, SaS) \mid (q_0, a, aS) \mid (q_0, X, S)$

iar pentru ultima variantă:

$(q_0, abbaa, z) \mid (q_0, bbaa, za) \mid (q_0, bbaa, zaS) \mid (q_0, bbaa, zaSb)$
 $\mid (q_0, baa, zaSbb) \mid (q_0, aa, zaSbbb) \mid (q_0, aa, zaSbbbS)$
 $\mid (q_0, a, zaSbbbSa) \mid (q_0, a, zaSbbbSaS) \mid (q_0, a, zaSbbS)$
 $\mid (q_0, X, zaSbbSa) \mid (q_0, X, zaSbbSaS) \mid (q_0, X, zaSbS)$
 $\mid (q(U, zS) \mid (q_1, A, z)$

Comentarii:

Considerând noțiunea de derivare putem să introducem două noi notații. *Derivarea stânga* este o derivare în care se înlocuiește întotdeauna cel mai din stânga neterminat în timp ce într-o *derivare dreapta* se înlocuiește întotdeauna cel mai din dreapta neterminat. Ultimele două variante corespund derivărilor stânga, respectiv dreapta, ^le șirului de intrare.

$S \Rightarrow aSbS \Rightarrow abS \Rightarrow abbSaS \Rightarrow abbbSaSaS \Rightarrow abbaSaS \Rightarrow abbaaSaS \Rightarrow abbaa$

Respectiv:

$S \Rightarrow aSbS \Rightarrow aSbbSaS \Rightarrow aSbbSa \Rightarrow aSbbbSaSa \Rightarrow aSbbbSaa \Rightarrow aSbbbaa \Rightarrow abbaa$

1.3-18

Să se construiască automatul cu stivă care acceptă următorul limbaj:

$L = \{a^n b^m \mid n \geq 1, m \geq 1, n \leq m \text{ sau } m < 2 * n\}$

Soluție:

Deoarece între numărul de simbolii a (n), respectiv numărul de simbolii b (m) există fie relația $n < m$ (cazul 1), fie relația $m < 2 * n$ (cazul 2), vom trata cele două cazuri în același mod: se depun în stivă toți simbolii a de pe banda de intrare prin simpla copiere în stivă (cazul 1), deci vom avea în stivă cei n simbolii a sau prin duplicare (cazul 2), caz în care vor exista în stiva $2 * n$ simbolii a. În continuare, atât în cazul 1 cât și în cazul 2 se descarcă stiva pentru fiecare simbol b de pe banda de intrare. În cazul 1 se trece în starea de acceptare dacă nu s-a ajuns la sfârșitul șirului înainte de a se goli stiva, iar în cazul 2 dacă la terminarea șirului, stiva este nevidă. Prin urmare, un automat PD care acceptă limbajul L este:

$PD = (\{q_0, q_1, q_2, q_3, q_4, q_5, q_6\}, \{a, b\}, \{a, z\}, m, q_0, z, \{q_3, q_6\})$
 $m(q_0, a, z) = \{(q_1, az), (q_4, aaz)\}$

(tratare cazul 1):

$| m(q1, a, a) = \{ (q1, aa) \}, m(q1, b, a) = \{ (q2, A) \}, m(q2, b, a) = \{ (q2, A) \}$
 $: m(q2, X, z) = \{ (q3, z) \}$ cazul $n = m, m(q2, b, z) = \{ (q3, z) \}$ cazul $n < m$
 $' m(q3, b, z) = \{ (q3, z) \}$

(tratare cazul 2):

$j m(q4, a, a) = \{ (q4, aaa) \}, m(q4, b, a) = \{ (q5, X) \}, m(q5, b, a) = \{ (q5, A) \}$
 $m(q5, X, a) = \{ (q6, X) \}$

Verificare:

Pentru șirul de intrare aabbb se obține următoarea evoluție a automatului cu stivă:

$(q0, aabbb, z) \mid (q1, abbb, az) \mid (q1, bbb, aaz) \mid (q2, bb, az) \mid (q2, b, z) \mid (q3, X, z)$
sau $(q0, aabbb, z) \mid (q4, abbb, aaz) \mid (q4, bbb, aaaaz) \mid (q5, bb, aaz) \mid (q5, b, aaz)$
 $\mid (q5, X, az) \mid (q6, X, z)$

Deci șirul aabbb este acceptat de automat. Pentru șirul de intrare aabb, evoluția automatului PD este: $(q0, aabb, z) \mid (q1, abb, az) \mid (q1, bb, aaz) \mid (q2, b, az) \mid (q2, X, z) \mid (q3, X, z)$
sau

$(q0, aabb, z) \mid (q4, abb, aaz) \mid (q4, bb, aaaaz) \mid (q5, b, aaaz) \mid (q5, X, aaz) \mid (q6, X, az)$

Deci șirul aabb este acceptat de automat conform ambelor condiții. Pentru șirul aabbbb evoluția automatului este: $(q0, aabbbb, z) \mid (q1, abbbb, az) \mid (q1, bbbb, aaz) \mid (q2, bbb, az) \mid (q2, bb, z) \mid (q3, b, z) \mid (q3, A, z)$

Deci șirul este acceptat de automat (conform primei condiții). De asemenea, pentru șirul de intrare aab, evoluția automatului este: $(q0, aab, z) \mid (q1, af, az) \mid (q1, b, aaz) \mid (q2, A, az)$ și în acest caz nu mai este aplicabilă nici o tranziție sau: $(q0, aab, z) \mid (q4, ab, aaz) \mid (q4, b, aaaaz) \mid (q5, X, aaaaz) \mid (q6, X, aaz)$. Adică șirul este acceptat de automat (conform condiției a doua).

problema 1.3-19

Să se construiască automatul cu stivă care acceptă următorul limbaj:

$L = \{ a^n b^m a^n \mid n, m \geq 1 \} \cup \{ a^n b^m c^m \mid n, m \geq 1 \}$

Soluție:

Deoarece limbajul specificat de L este de fapt reuniunea a doua limbaje, va trebui ca în specificarea automatului PD să realizăm niște mecanisme prin care să fie acceptate șiruri din ambele limbaje. Astfel, șirurile nu pot începe decât cu simbolul a. Toți simbolii a din prefixul șirului se introduc în stivă. De asemenea, toți simbolii b care urmează sunt introduși în stivă. În continuare: Dacă următorul simbol de pe banda de intrare (după ce s-au epuizat toți simbolii b) este a, atunci se încearcă acceptarea unui șir de forma: $a^n b^m a^n$, caz în care se descarcă toți simbolii b din stivă și se verifică dacă numărul de simbolii a din stivă este egal cu numărul de simbolii a rămași pe banda de intrare.

Dacă următorul simbol de pe banda de intrare este c, se încearcă acceptarea unui șir de forma: $a^n b^m c^m$, caz în care se verifică dacă numărul de simbolii c din sufixul șirului este egal cu numărul de simbolii b din stivă. În final, se descarcă din stivă toți simbolii a.

Rezultă automatul:

$PD = (\{ q0, q1, q2, q3, q4, q5, q6 \}, \{ a, b, c \}, \{ a, b, z \}, m, q0, z, \{ q5 \})$

i (încarcăsimbolii a și b în stivă):

$m(q0, a, z) = \{ (q1, az) \}, m(q1, a, a) = \{ (q1, aa) \}, m(q1, b, a) = \{ (q1, ba) \}$
 $m(q1, b, b) = \{ (q1, bb) \}$

(varianta 1):

$m(q1, a, b) = \{ (i^{\wedge}, A) \}$ am consumat un simbol a de la intrare

$m(q2, X, b) = \{ (q2, X) \}$ descarcă simbolii b din stivă

$m(q2, a, a) = \{ (q3, X) \}$

$m(q3, A, a) = \{ (q4, A) \}$ scot un a din stivă ce corespunde simbolului a consumat în starea q1

$m(q4, a, a) = \{ (q4, A) \}$

$m(q4, X, z) = \{ (q5, z) \}$ stare finală

(varianta 2):

$m(q1, c, b) = \{ (q6, A) \}, m(q6, c, b) = \{ (q6, A) \}, m(q6, X, z) = \{ (q5, z) \}$

(descarcă simbolii a din stivă):

$m(q4, A, a) = \{ (q4, A) \}, m(q4, A, z) = \{ (q5, z) \}$ (stare finală)

problem: 1.3-20

Să se construiască automatul cu stivă care acceptă limbajul expresiilor aritmetice cu paranteze

Soluție:

Vom considera că există două stări q1, respectiv q2 pentru care sunt valabile următoarele semnificații: q1 este starea de start. Atât timp cât pe banda de intrare simbolul analizat este "(", se copiază în stivă. Dacă simbolul de pe banda de intrare este a, se trece în starea q2. În starea q2, dacă simbolul de pe banda de intrare este operator (+, *), se trece în starea q1. Dacă simbolul este ")", se descarcă din stivă un simbol "(".

Rezultă automatul $PD = (\{ q1, q2, q3 \}, \{ a, +, *, (,) \}, \{ z, (, m, q1, z, \{ q3 \} \})$,

$m(q1, (, x) = \{ (q1, (x) \} \mid x \in \{ (, z \}, m(q1, a, x) = \{ (q2, x) \} \mid x \in \{ (, z \}$

$m(q2,), () = \{ (q2, A) \}, m(q2, y, x) = \{ (q1, x) \} \mid y \in \{ +, * \}, m(q2, X, z) = \{ (q3, X) \}$

Verificări::

Pentru șirul de intrare $a+a*(a+(a)*a)$, evoluția automatului este:

$(q1, a+a*(a+(a)*a), z) \mid (q2, +a*(a+(a)*a), z) \mid (q1, a*(a+(a)*a), z)$
 $;$ $\mid (q2, *(a+(a)*a), z) \mid (q1, (a+(a)*a), z) \mid (q1, a+(a)*a), (z)$
 $;$ $\mid (q2, +(a)*a), (z) \mid (q1, (a)*a), (z) \mid (q1, a)*a), ((z) \mid (q2,)*(a), ((z)$
 $\mid (q2, *a), (z) \mid (q1, a), (z) \mid (q2,), (z) \mid (q2, A, z) \mid (q3, A, X)$

Comentarii:

Șirurile sunt acceptate prin stivă goală.

1.3.3 Mașina Turing

Definiția 1.3.11. Se numește mașină Turing obiectul matematic: $MT = (Q, T, m, q_0)$

unde

- Q este o mulțime finită a stărilor mașinii Turing, $hg Q$;
- T este alfabetul finit de intrare care conține simbolul # dar nu conține simbolii L și R ;
- m este funcția de tranziție $m : Q \times T \rightarrow (Q \cup \{h\}) \times (T \cup \{L, R\})$.
- $q_0 \in Q$ este starea inițială a mașinii Turing.

Dacă $q \in Q$, $a \in T$ și $m(q, a) = (p, b)$ înseamnă că fiind în starea q , având simbolul a sub capul de citire mașina trece în starea p . Dacă $b \in T$ atunci simbolul a va fi înlocuit cu b , altfel capul se va deplasa cu o poziție la stânga sau la dreapta în funcție de valoarea simbolului $b \in \{L, R\}$. Se observă că mașina Turing a fost definită ca un acceptor determinist.

Definiția 1.3.12. O configurație pentru o mașină Turing este $(q, a, a, /?)$ notată și $(q, aay?)$ unde $q \in (Q \cup \{h\})$, $a \in T^*$, $a \in T$, $\exists e (T - \{#\}) \cup \{X\}$ în care

- q este starea curentă
- a este șirul aflat pe banda de intrare la stânga capului de citire / scriere
- a este simbolul aflat sub capul de citire / scriere
- $/?$ este șirul aflat la dreapta capului de citire / scriere.

Definiția 1.3.13. Relația de tranziție \vdash se definește asupra configurațiilor în modul următor. Fie (q_1, w_1, a_1, u_1) și (q_2, w_2, a_2, u_2) două configurații. Atunci:

$(q_1, w_1, a_1, u_1) \vdash (q_2, w_2, a_2, u_2)$ dacă și numai dacă există $b \in T \cup \{L, R\}$ astfel încât $m(q_1, a_1) = (q_2, b)$ și este îndeplinită una din următoarele condiții:

1. $b \in T, w_1 = w_2, u_1 = u_2, a_2 = b;$

2. $b = L, w_1 = w_2 a_2$
 dacă $a_1 \neq \#$ sau $u_1 \neq X$
 $u_2 = a_1 u_1$
 altfel $u_2 = X$

3. $b = R, w_2 = w_1 a_1$
 dacă $u_1 = X$
 $u_2 = X$
 $a_2 = \#$
 altfel $u_1 = a_2 u_2$

Un calcul efectuat de o mașină Turing este o secvență de configurații c_0, c_1, \dots, c_n astfel încât $n > 0$ și $c_0 \vdash c_1 \vdash \dots \vdash c_n$. Se spune despre calcul că are loc în n pași.

O mașină Turing poate să fie utilizată și ca *acceptor de limbaj*. Și anume să considerăm un limbaj L definit asupra alfabetului T care nu conține simbolii #, D și N . Fie $dL : T^* \rightarrow \{D, N\}$ o funcție definită în modul următor - pentru $\forall w \in T^*$

D dacă $w \in L$

$dL(w) =$
 \backslash
 N dacă $w \notin L$

Se spune ca limbajul L este *decidabil în sens Turing* (Turing decidable) dacă și numai dacă funcția dL este calculabilă Turing. Dacă dL este calculată de o mașina Turing MT spunem ca MT decide L .

Noțiunea de acceptare este mai largă decât noțiunea de decidabilitate în legătura cu limbajele. Și anume dacă L este un limbaj asupra alfabetului T , spunem ca limbajul L este *Turing acceptabil* dacă există o mașina Turing care se oprește dacă și numai dacă $w \in L$.

Propoziție. Limbajele Turing decidabile sunt și Turing acceptabile. Reciproca nu este însă adevărată.

Definiția 1.3.14. O schemă de mașină Turing este tripletul $*F = (M, r, MO)$ unde

- M este o mulțime finită de mașini Turing care au toate un alfabet comun T și mulțimi de stări distincte;
- $MO \in M$ este mașina inițială;
- T este o funcție parțială, $r : M \times T \rightarrow M$.

O schemă de mașină Turing reprezintă o mașina Turing T compusă din mașinile care formează mulțimea M . Funcționarea acesteia începe cu funcționarea mașinii MO . Dacă MO se oprește atunci \forall poate continua eventual funcționarea conform altei mașini din M . Dacă MO se oprește cu capul de citire / scriere pe caracterul a atunci există următoarele situații:

- $t|(MO, a)$ este nedefinită, în acest caz \forall se oprește;
- $-n(MO, a) = M'$, în acest caz funcționarea continuă cu starea inițială a mașinii M' .

$4 > = (M, r, MO)$ o schemă de mașină Turing reprezintă mașina $MT = (Q, T, m, s)$ unde

- $Q = Q_0 \cup \dots \cup Q_k \cup \{q_0, \dots, q_k\}$
- $s = s_0$
- m este definită în modul următor:
 - dacă $q \in Q_j, 0 < i < k, a \in T$ și $n_{ij}(q, a) = (p, b), p \neq h$ atunci $m(q, a) = m_{ij}(q, a) = (p, b)$;
 - dacă $q \in Q_0, 0 < i < k, a \in T$ și $m_i(q, a) = (h, b)$ atunci $m(q, a) = (q_i, b)$;
 - dacă $r_i(M_i, a) (0 < i < k, a \in T)$ este nedefinită atunci $m(q_i, a) = (h, a)$;
 - dacă $r_i(M_i, a) = M_j (0 < i < k, a \in T)$ și $m_j(s_j, a) = (p, b)$ atunci

(p, b)

$m(q_i, a) =$
 \backslash
 $(q_j, b)p = h$

Propoziție. Adăugarea unor facilități la definiția Mașinii Turing ca de exemplu :

- banda de intrare infinită la ambele capete;
- mai multe capete de citire / scriere;
- mai multe benzi de intrare;
- bandă de intrare organizată în două sau mai multe dimensiuni (eventual ca o memorie);

nu crește puterea de calcul oferită. Nu se schimbă nici clasa funcțiilor care pot să fie calculate și nici a limbajelor acceptate sau decise.

Probleme

problema 1.3-21

Fie mașina Turing $MT = (\{q_0, q_1, q_2\}, \{a, \#\}, m, q_0)$ cu

$m(q_0, a) = (q_1, L)$
 $m(q_0, \#) = (q_0, \#)$
 $m(q_1, a) = (q_2, \#)$
 $m(q_1, \#) = (h, \#)$
 $m(q_2, a) = (q_1, a)$
 $m(q_2, \#) = (q_0, L)$

care este evoluția mașinii dacă pleacă din configurația inițială: $(q_0, \#a^n a) \quad n > 0$?

Soluție:

Să presupunem că n este un număr impar de exemplu $n = 1$. în acest caz evoluția mașinii

- Turing este:

$(q_0, \#aa) \rightarrow (q_1, \#aa) \rightarrow (q_2, \#\#a) \rightarrow (q_0, \#\#a)$ se agață

: Dacă n este par, de exemplu 2 atunci:

$(q_0, \#a^2a) \rightarrow (q_1, \#\#aa) \rightarrow (q_1, \#\#aa) \rightarrow (q_0, \#\#aa) \rightarrow (q_1, \#\#aa) \rightarrow (h, \#\#aa)$

Deci mașina se oprește pentru șiruri de forma: $\#a^{2k}a$ și le transformă în $\#\#a^{2k}a$. Pentru șiruri de forma $\#a^{2k+1}a$ mașina nu se oprește.

problema 1.3-22

Să se construiască mașina Turing care acceptă limbajul $L = \{ w \in \{a, b\}^* \mid w \text{ conține doi simboluri consecutive } a \}$.

Soluție:

Se pleacă din starea inițială q_0 și se merge spre dreapta, dacă se întâlnește un simbol a se va trece într-o stare nouă și se continuă deplasarea la dreapta. Mașina se va opri pe al doilea simbol a întâlnit.

$MT = (\{q_0, q_1\}, \{a, b, \#\}, m, q_0)$ cu

$m(q_0, a) = (q_1, R)$
 $m(q_0, x) = (q_0, R) \quad x \in \{b, \#\}$
 $m(q_1, a) = (h, a)$
 $m(q_1, x) = (q_0, R) \quad x \in \{b, \#\}$

Descrieți în cuvinte transformarea efectuată de următoarea mașina Turing.

$MT = (\{q_0, q_1, q_2\}, \{a, b, \#\}, m, q_0)$ cu

$m(q_0, x) = (q_0, x) \quad x \in \{a, b\}$ arbitrar
 $m(q_0, \#) = (q_1, L)$
 $m(q_1, a) = (q_1, L)$
 $m(q_1, b) = (q_1, L)$
 $m(q_1, \#) = (q_2, b)$
 $m(q_2, a) = (q_2, R)$
 $m(q_2, b) = (q_2, R)$
 $m(q_2, \#) = (h, \#)$

Soluție:

Să considerăm de exemplu comportarea pentru configurația inițială $(q_0, \#\#ab\#)$.

$(q_0, \#\#ab\#) \rightarrow (q_1, \#\#ab\#) \rightarrow (q_1, \#\#ab\#) \rightarrow (q_2, \#ab\#) \rightarrow (q_2, \#ab\#) \rightarrow (h, \#ab\#)$

Mașina transformă o bandă de intrare de forma $\#w\#$ în $bw\#$.

Considerând reprezentarea unară a numerelor naturale, care este funcția pe numere naturale calculată de următoarea mașină Turing:

$MT = (\{q_0, q_1, q_2, q_3\}, \{1, \#\}, m, q_0)$ cu

$m(q_0, \#) = (q_1, L)$
 $m(q_0, 1) = (q_0, 1)$ arbitrar
 $m(q_1, \#) = (h, R)$
 $m(q_1, 1) = (q_2, \#)$
 $m(q_2, \#) = (q_3, L)$
 $m(q_2, 1) = (q_2, 1)$ arbitrar
 $m(q_3, \#) = (h, R)$
 $m(q_3, 1) = (h, \#)$

Soluție:

Să considerăm câteva exemple de evoluție.

$(q_0, \#\#) \rightarrow (q_1, \#\#) \rightarrow (h, \#\#)$, $(q_0, \#\#1) \rightarrow (q_1, \#\#1) \rightarrow (q_2, \#\#\#) \rightarrow (q_3, \#\#\#) \rightarrow (h, \#\#\#)$
 $(q_0, \#\#1\#) \rightarrow (q_1, \#\#1\#) \rightarrow (q_2, \#\#\#\#) \rightarrow (q_3, \#\#\#\#\#) \rightarrow (h, \#\#\#\#\#\#)$

;

Se observă că $f(n) = \lfloor \frac{n}{2} \rfloor$

;

Soluție:

Pentru a se verifica $\#_a(w) = \#_b(w) = \#_c(w)$ se caută de la stânga la dreapta primul a, după care tot de la stânga la dreapta primul b, respectiv primul c. La fiecare căutare (un simbol a, b, respectiv c) dacă simbolul respectiv este găsit, atunci acesta va fi șters. Pentru a șterge fiecare simbol găsit a, b, respectiv c se va folosi un simbol d g {a, b, c} (dacă ștergerea s-ar realiza cu # nu am mai ști care este capătul din dreapta al șirului inițial).

Mașina Turing se oprește dacă șirul de pe banda de intrare aparține limbajului L. Pentru aceasta înainte de a căuta a, b, respectiv c trebuie să vedem dacă primul simbol căutat la dreapta (diferit de d) nu este #, caz în care șirul aparține limbajului, deci mașina Turing se oprește. Pentru șiruri care nu aparțin limbajului mașina va cicla căutând la dreapta un simbol a sau b sau c.

$x \neq \#$

Halt

Să se construiască mașina Turing care calculează funcția $f(x, y) = x * y$, unde x și y sunt numere naturale

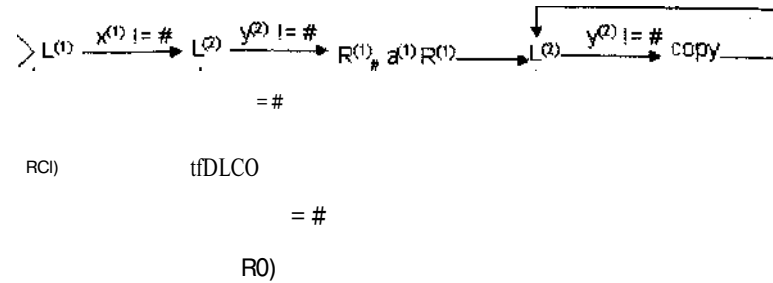
Soluție:

Considerăm că pentru x, y și rezultat se utilizează reprezentarea unară utilizând simbolul 1. Vom construi mașina Turing cu două benzi: b1 conține #x# și b2 conține #y#, rezultatul va fi pe b1 #x*y#. Sunt posibile următoarele situații:

1. Dacă x = 0 și y = 0 rezultatul înmulțirii este 0.
2. Dacă 0 și y = 0 trebuie ștersă banda 1.

Dacă x = 0 și y = 0 atunci va trebuie copiat de y - 1 ori conținutul inițial al benzii 1 la sfârșitul benzii 1. Pentru a păstra valoarea inițială de pe banda 1 se marchează sfârșitul șirului cu a, iar în final după terminarea calculării înmulțirii se înlocuiește a cu I și se mută cu o poziție la stânga sfârșitul șirului de pe banda 1.

Astfel mașina Turing testează mai întâi dacă x = 0 sau y = 0, caz în care rezultatul este 0.



Pentru a copia o singură dată pe x, după simbolul a pe banda b1 vom scrie mașina Turing compusă copy, care lucrează numai cu banda b1. Inițial b1 conține #wa#, iar în final #waw# sau inițial b1 conține #waw#, iar în final #waw#w# (w' reprezintă șirul w copiat de un număr de ori).

Mașina Turing compusă care realizează copierea este:

problem» 1.3-33

Să se construiască mașina Turing care decide limbajul: $L = \{w \mid \exists j \in \mathbb{N} \text{ } w \in \{a, b\}^j\}^*$

Soluție:

Pentru a determina în mod determinist mijlocul șirului se folosește un caracter c care se propagă spre stânga șirului. Pentru a verifica dacă cele două "jumătăți" ale șirului, delimitate de simbolul c, sunt egale se vor folosi două benzi.

1. se propagă simbolul c pe banda b1 spre stânga cu o poziție și apoi se copiază pe banda b1 pe banda b2 (copierea se face parcurgând banda b1 de la stânga la dreapta, iar banda b2 de la dreapta la stânga)
2. se compară cele două jumătăți găsite (banda b1 este parcursă de la stânga la dreapta pornind inițial de la simbolul c, banda b2 este parcursă de la dreapta la stânga). Compararea are loc până se întâlnește pe banda 1 #. Dacă cele două jumătăți astfel determinate sunt egale trebuie ștersă banda 1 și se scrie D. Dacă cele două jumătăți determinate nu sunt egale se propagă cu încă o poziție la stânga pe banda 1 simbolul j c, după care se repetă copierea pe banda 1 și apoi compararea. Propagarea la stânga cu o poziție a simbolului c pe banda 1 are loc până se întâlnește # (marginea din stânga a șirului), caz în care banda 1 este ștersă și trebuie scris N.

Pentru aceasta se vor construi mașinile Turing compuse:

copy pentru a realiza copierea pe banda 2
cmp pentru a compara cele două jumătăți găsite

copy (mașina Turing compusă de copiere)

De exemplu:

Inițial

b1: #abbabcb#

b2: ##

după copiere

b1: #abbabcb#

b2: #bcbabba#

| $\langle \rangle$ =

cmp (mașina Turing compusă de comparare)

Fie **cmp** mașina Turing compusă care compară cele două jumătăți delimitate de simbolul c.

Cele două jumătăți se află pe banda b1, respectiv banda b2.

Inițial:

b1: #abbabcb#

b2: #bcbabba#

Banda b1 este parcursă de la stânga la dreapta, iar banda b2 de la dreapta la stânga.

în final:

b1: #abbabcb#

b2: #bcbabba#

Dacă cele două jumătăți sunt egale:

Inițial

b1: #abbcabb#

t>2: #bbacba#

în final:

b1: #abbcabb#

b2: #bbacba#



I = e

Mașina Turing care decide limbajul L este:

R(1) NO) RtD

R(1) 00) RC)

cmp

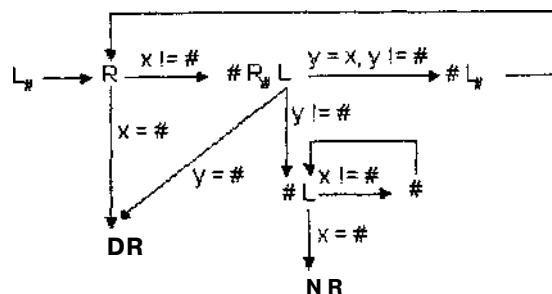
problema 1.3-34

Să se construiască mașina Turing compusă care decide limbajul:

$$L = \{ w \in \{a, b\}^* \mid w = w^R \}$$

Soluție:

Inițial banda de intrare conține #w#, iar în Iku! \wedge sau #N#, după cum w aparține sau nu limbajului L. Se verifică egalitatea simbolilor de pe poziții egale față de capetele șirului. Pe măsură ce se parcurge șirul caracterele se șterg și se înlocuiesc cu #.



2 Lex - generator de analizoare lexicale

lex este un generator de analizoare lexicale prezent în orice versiune UNIX. Față de varianta standard (anii 70 - '80) există numeroase variante mai noi. Cea mai cunoscută variantă se numește *flex* (Fast lexical analyzer generator) și face parte dintre instrumentele realizate de FSF (Free Software Foundation) în cadrul proiectului GNU (GNU is Not UNIX). Spre deosebire de *lex* pentru *flex* (ca și pentru alte programe din domeniu public) programul sursă este disponibil, și în același timp performanțele exprimate în timp de execuție și memorie ocupată sunt mai bune. Din acest motiv în cele ce urmează vom prezenta de fapt varianta *flex*.

lex poate să fie utilizat în programare și pentru alte aplicații decât scrierea de compilatoare. Astfel, ori de câte ori un program trebuie să identifice în fișierul de intrare șiruri de caractere având o structură de tip "atom lexical" (adică numere întregi sau reale în diferite formate, cuvinte cheie, identificatori, caractere speciale, etc.) pentru care se execută acțiuni speciale, putem să folosim o funcție generată de *lex*.

Un generator de programe este de fapt un translator care realizează traducerea unui text care reprezintă specificațiile programului ce trebuie generat (sursa) într-un text scris într-un limbaj de programare (obiect). În cazul *lex*-ului generarea analizorului lexical utilizează și un cod fix care reprezintă emulatorul automatului determinist corespunzător. Cu alte cuvinte algoritmul de analiză propriu-zis. În cazul *lex*-ului acest cod este scris în C, și programele generate vor fi programe C.

Formularea "execuția programului *lex*" se referă printr-un abuz de limbaj la execuția programului generat pe baza specificațiilor și nu la execuția programului care realizează traducerea.

Pentru a obține un program executabil din textul generat de către *lex*, textul C generat trebuie să fie compilat (utilizând un compilator de C) obținându-se un modul obiect.

Specificațiile sursă pentru *lex* constau dintr-o descriere a unităților lexicale ce urmează să fie recunoscute de către programul generat. Pentru fiecare unitate lexicală se specifică eventual și o acțiune reprezentată de o secvență de cod C care urmează să se execute la recunoașterea unui șir ce corespunde formei generale. Dacă specificația este corectă atunci se va genera un fișier *lex.yy.c* (*lexyy.c* sub MS-DOS). În cadrul acestui fișier apare cel puțin o funcție numită *yylexO^{c a r c}* reprezintă de fapt analizorul lexical. Execuția acesteia realizează căutarea în șirul de intrare a unui subșir care să "se potrivească" cu unul din modelele descrise în specificații. Dacă un astfel de șir este găsit atunci se va executa acțiunea asociată modelului respectiv.

Pentru a înțelege modul în care se utilizează programul *lex* trebuie să precizăm câteva lucruri. Și anume, funcționarea oricărui analizor lexical are o parte fixă care nu depinde de modelele de atomi lexicali căutate în șirul de intrare. Algoritmul corespunzător este de fapt simularea unui automat finit determinist. Această parte constituie scheletul analizorului lexical. Acest schelet este un text scris în limbajul C. Partea care variază de la un analizor lexical la altul are de fapt două componente și anume modelele de atomi lexicali căutate în șirul de intrare și acțiunile pe care eventual analizorul lexical trebuie să le execute la identificarea unui șir care se "potrivește" cu un model de atom lexical. Specificațiile *lex* conțin descrierea părții variabile. Deoarece ce se generează în final este un text C, acțiunile sunt descrise ca secvențe de instrucțiuni C. În aceste secvențe pot să apară orice instrucțiuni C inclusiv apeluri de funcții. Secvențele respective vor "îmbrăca" scheletul părții fixe. Specificarea modelelor de atomi

lexicali se face sub formă de expresii regulate. Pornind de la ele programul lex va genera tabelele care descriu funcția de tranziție a automatului determinist corespunzător. Aceste tabele reprezintă o parte din structurile de date pe care lucrează "scheletul". Utilizând "scheletul" și secvențele de instrucțiuni oferite de către programator programul lex va genera un text care reprezintă analizorul lexical. Programul lex "nu înțelege" nici scheletul și nici acțiunile specificate de către programator, aceste elemente reprezintă numai șiruri de caractere pe care lex știe să le combine. În mod corespunzător cel care scrie specificațiile lex trebuie să aibă în vedere aspecte ca domeniul de valabilitate al variabilelor utilizate în programul generat. Și asta ținând cont că există de fapt două tipuri de variabile, variabile pe care le controlează numai cel care scrie specificațiile lex și pentru ele se face atât declararea cât și utilizarea în cadrul secvențelor specificate de către programator și variabile care sunt declarate în cadrul scheletului (sunt predefinite din punct de vedere al celui care scrie specificațiile). Declararea acestor variabile apare în textul generat înainte de începutul funcției yylex() Referirea acestora se poate face atât în secvențele de cod care se execută asociat diferitelor modele de atomi (aceste secvențe vor fi interne funcției yylex()) cât și în funcții care vor fi incluse în textul generat după textul funcției yylex().

2.1 Expresii regulate. Structura unei specificații lex.

În cele ce urmează vom numi program lex (flex) textul ce conține descrierea unităților lexicale și a acțiunilor corespunzătoare. Să considerăm un exemplu foarte simplu pentru a oferi o idee despre cum arată specificațiile programele lex.

rcma 2-1

```
/* declarații de variabile utilizate in
programul generat
*/
int numar_cifre = 0, numar_litere = 0;
%}
/* declarații de macrodefiniții utilizate in specificarea regulilor */

cifra [0-9]
litera [A-Za-z]

{cifra} { numar_cifre++; }
{litera} { număr litere++;}

main() {
yylex();
printf("\n număr cifre = %i, număr litere = %i\n",
numar_cifre, numar_litere);
```

Descriere;

Programul numără literele și cifrele care apar în fișierul de intrare.

Nu este cel mai simplu program lex posibil. Conține însă majoritatea elementelor Unui program lex. Se observă că programul este format din trei secțiuni separate de caracterele %%. Prima secțiune poate conține text de program C care va fi copiat nemodificat în textul programului generat și macrodefiniții. Textul care se copiază conține declarațiile de variabile și funcții ce vor fi utilizate în acțiunile asociate șirurilor recunoscute. Acest text trebuie inclus între caracterele %{ și respectiv %} (fiecare apărând pe un rând separat). Macrodefinițiile sunt utilizate pentru a scrie mai compact descrierea modelelor de atomi lexicali. Pentru exemplul considerat au fost definite două nume cifra și litera. Primul nume are asociată o descriere a unei liste de caractere care conține numai cifre, în timp ce al doilea nume are asociată o listă care conține orice literă (mare sau mică).

Pentru caractere se consideră ordonarea lexicografică obișnuită. Conform acestei ordonări cifrele sunt caractere care au valori într-un interval în care '0' este cel mai mic iar '9' este cel mai mare, literele mici sunt caractere care au valori mai mari decât literele mari, literele mari și respectiv mici au valori în câte un interval astfel încât 'A' este caracterul cu valoarea cea mai mică dintre literele mari iar 'Z' este caracterul cu valoarea cea mai mare dintre literele mari, etc. Corespunzător definiția unei cifre este un caracter având codul cuprins în intervalul închis [0-9] iar o litera este un caracter având codul cuprins în unul dintre intervalele [A-Z] și respectiv fa-zi, în programele lex utilizarea caracterelor este destul de rigidă. Astfel dacă definiția literei s-ar fi făcut sub forma:

litera [A-Z a-z]

atunci și caracterul blank ar fi fost considerat literă.

A doua secțiune conține specificarea modelelor și a prelucrărilor asociate. Fiecare model este scris pe o linie separată începând din prima coloană. Pe aceeași linie separată de un blank sau un caracter de tabulare este specificată acțiunea corespunzătoare.

A treia secțiune conține text C care va fi copiat în programul generat. De obicei aici apar funcțiile referite în acțiunile din secțiunea a doua și care referă variabile predefinite. Deoarece în cazul nostru generăm un program care se va executa independent (adică nu generăm o funcție care va fi referită din alte module compilate) în secțiunea a treia este conținut codul pentru programul principal. În programul principal se apelează funcția yylex() care va fi generată conform specificațiilor din a doua secțiune. În exemplul considerat execuția funcției yylex() se termină când se ajunge la sfârșitul fișierului de intrare.

În general dacă funcțiile apelate în cadrul acțiunilor sunt semnificative ca dimensiuni este de preferat ca acestea să apară în module de program compilate separat. În acest mod se simplifică depanarea programului generat.

Pentru descrierea unităților lexicale ce urmează să fie recunoscute se utilizează expresii regulate. În acest caz o expresie regulată este interpretată ca un șablon (model) care descrie o mulțime de șiruri. Un astfel de șablon poate să fie utilizat pentru clasificarea șirurilor de caractere în șiruri care corespund sau nu șablonului. Pentru lex se utilizează a doua interpretare a expresiilor regulate. Adică, programul generat va identifica într-un șir de intrare subșirurile care "se potrivesc" cu descrierile făcute în specificații.

În notația utilizată de lex pentru reprezentarea expresiilor regulate se utilizează caracterele obișnuite utilizate pentru limbajul C. O serie de caractere ca de exemplu: *, +, |, ? au semnificație specială în scrierea expresiilor regulate. Aceste caractere se numesc metacaractere. În afară de : *, +, |, ? se mai utilizează ca metacaractere: ", \, {, }, ^, <, >, \$, /, (,)•• în tabelul 2.1 sunt conținute construcțiile lex în care se utilizează metacaractere.

notație	Semnificație
"șir"	șirul de caractere șir; chiar dacă conține numai un metacarakter, de exemplu "*", permite specificarea blanurilor și metacaracterelor în expresii regulate (în mod normal un blank încheie șirul de caractere care descrie o expresie regulată).
\x	dacă "x" este unul din caracterele speciale din limbajul C ca de exemplu t sau n atunci notația reprezintă caracterul respectiv, altfel notația reprezintă caracterul "x" chiar dacă este un metacarakter.
[lista]	unul dintre caracterele din listă. Lista este formată din caractere individuale sau din intervale. De exemplu [ABC] reprezintă unul dintre caracterele A, B, C, iar [A-Za-z] reprezintă orice literă. Dacă primul caracter din listă este] sau - atunci el este tratat ca atare (și pierde semnificația de metacarakter).
[^listă]	orice caracter mai puțin cele din listă. De exemplu [^0-9] înseamnă orice caracter, cu excepția cifrelor.
	orice caracter mai puțin caracterul linie nouă (\n).
^e	un șir reprezentat de expresia regulată e dar numai la început de linie.
e\$	un șir reprezentat de expresia regulată e dar numai la sfârșit de linie.
<y>e	un șir reprezentat de expresia regulată e dacă analizorul este în starea y (noțiunea de stare va fi explicată ulterior).
e/f	un șir reprezentat de expresia regulată e dacă urmează un șir reprezentat de expresia regulată f. Într-o expresie regulată poate să apară o singură dată această condiție (la sfârșitul expresiei).
e?	un șir reprezentat de expresia regulată e sau șirul vid.
e*	0, 1, 2, ... apariții ale unor șiruri reprezentate de expresia regulată e.
e+	1, 2, 3, ... apariții ale unor șiruri reprezentate de expresia regulată e.
e{m,n}	între m și n apariții ale unor șiruri reprezentate de expresia regulată e.
e{m,}	cel puțin m apariții ale unor șiruri reprezentate de expresia regulată e.
e{m}	exact m apariții ale unor șiruri reprezentate de expresia regulată e.
e f	un șir de caractere provenind din expresia regulată e sau din expresia regulată f
(x)	dacă x un caracter atunci are același efect ca și includerea caracterului între ghilimele. Dacă x este o expresie regulată atunci notația poate să fie utilizată pentru schimbarea priorității operatorilor utilizați.
{XX}	dacă xx este un nume definit în prima secțiune a unui program lex atunci notația reprezintă rezultatul substituirii numelui xx cu definiția asociată.

Tabelul 2-1

Mai există o construcție specială «EOF», care apare numai în Qex și care reprezintă condiția de sfârșit de fișier. Această construcție nu poate să intre în compunerea unei expresii regulate. În Tabelul 2-1 au apărut o serie de operații care sunt utilizate pentru formarea expresiilor regulate. Se pune problema care sunt regulile de precedență pentru aceste operații. Din păcate aici apar diferențe între lex și flex. În general ordinea operațiilor este în ordine descrescătoare:

{}	
* + ?	
concatenare	

De exemplu expresia:

[a-z][4-5]+/123+9 reprezintă șirul a4445555555123333339, dar nu reprezintă șirul a44455555551231231239, deoarece interpretarea expresiei este: (([a-z]([4-5]+))/(12(3+9)). Pentru a ușura citirea textului programului și pentru a obține programe portabile între diferitele variante de lex trebuie să se utilizeze paranteze pentru stabilirea ordinii de considerare a operatorilor.

Din Tabelul 2-1 mai rezultă faptul că semnificația unui metacarakter poate să depindă de contextul în care acesta apare. Astfel, în expresia:

cele patru puncte care apar nu au toate aceeași semnificație. Primul și ultimul reprezintă un caracter oarecare (diferit de linie nouă, dar care poate însă să fie și un punct). Al doilea punct indică faptul că pe a doua poziție în șirul căutat poate să apară și un punct. Al treilea caracter precedat de metacarakterul \ indică faptul că pe poziția 3 trebuie să apară un punct. Rezultă că atât șirul "a!.x" cât și șirul "...." "se potrivesc" cu expresia regulată.

După cum am mai amintit un program lex este format din maxim 3 secțiuni separate de o linie care conține caracterele %%. Structura generală pentru un program este :

secvențe de cod care se copiază în programul generat
tracodefiniții

reguli
aa-

secvențe de cod care se copiază în programul generat

l: i

Împărțirea secvenței de cod care se copiază în programul generat între cele două secțiuni se face în așa fel încât să se permită definirea variabilelor înainte de utilizare. Astfel toate variabilele definite de către programator și care urmează să fie referite în acțiuni trebuie să fie definite în prima secțiune. Dacă în funcțiile care sunt apelate din acțiuni apar referiri la variabile care vor fi definite în mod implicit în funcția yylex (de exemplu variabilele yytext sau yleng) atunci aceste funcții trebuie să apară în a treia secțiune a programului. Existența unor secvențe de cod este opțională. Secvențele de cod care apar în prima secțiune sunt cuprinse între caracterele "%{" și respectiv "%}" conținute fiecare pe câte o linie separată. Este foarte important să se facă distincția între codul care este "înțeles" și prelucrat de către lex și cel care va fi înțeles de către compilatorul de C. Astfel lex nu face nici o verificare asupra textului cuprins între caracterele %{} și respectiv %}. Abia când se face compilarea programului generat de către lex se va face și verificarea acestor secvențe.

Prima secțiune conține macrodefiniții utilizate pentru simplificarea scrierii modelelor. Forma generală a unei macrodefiniții este:

nome expresie regulată

unde nome este un cuvânt care începe cu o literă sau caracterul "_" și continuă cu litere, cifre sau caractere "-", "_"; expresieregulată este construită pe baza regulilor din Tabelul 2-1. Utilizarea unei macrodefiniții se poate face în orice expresie regulată care apare ulterior prin utilizarea numelui macrodefiniției între acolade.

Deoarece notațiile utilizate în expresiile regulate sunt destul de greoaie, este indicată utilizarea unor macrodefiniții pentru a ușura citirea acestora. De exemplu se poate considera pentru construirea expresiei regulate care reprezintă un număr real secvența de macrodefiniții:

```
cifra [0-9]
semn [+ -]?
IntregFaraSemn {cifra}+
Exponent ([Ee]{semn}{IntregFaraSemn})
Real ({semn}{IntregFaraSemn}\.{IntregFaraSemn}?{Exponent}?)
```

sau direct:

```
Real [+ -]?[0-9]+\.( [0-9]+)?([Ee] [+ -]?[0-9]+) ?
```

Cele două definiții sunt echivalente ca efect, dar prima soluție este mai ușor de urmărit. Regulile au forma generală:

expresie regulată acțiune

Expresia regulată este scrisă începând din prima coloană a liniei pe care apare. Sfârșitul expresiei regulate este reprezentat de un caracter delimitator: blank, linie nouă, un caracter de tabulare, etc. Acțiunea are forma unei instrucțiuni în limbajul C. Dacă pentru o acțiune sunt necesare mai multe instrucțiuni atunci ele vor fi grupate între caracterele "{" și "}". De fapt pentru claritate și pentru evitarea unor efecte laterale neplăcute este de preferat să se încadreze întotdeauna acțiunile între acolade. În orice caz acțiunea trebuie să înceapă pe linia pe care este scrisă expresia regulată.

Dacă pentru mai multe expresii regulate corespunde o aceeași acțiune atunci se poate utiliza o notație de forma:

```
expresie_regulatã1 |
expresie_regulatã2 |
... |
expresie_regulatãn acțiune
```

Pentru toate cele n expresii regulate se va executa aceeași acțiune.

Dintre cele trei secțiuni este obligatorie numai secțiunea de reguli, care trebuie însă să fie precedată de perechea de caractere "%%". Cel mai simplu text de specificații conține numai aceste caractere. În acest caz efectul funcției generate constă din tipărirea șirului de intrare.

2.2 Elemente avansate

Pentru a înțelege modul în care se utilizează programul lex trebuie să precizăm câteva lucruri. Și anume, funcționarea oricărui analizor lexical are o parte fixă care nu depinde de modelele de atomi lexicali căutate în șirul de intrare. Algoritmii corespunzător este de fapt simularea unui *automat finit determinist*. Această parte constituie scheletul analizorului lexical. Acest schelet este un text scris în limbajul C la care se adaugă o serie de variabile și structuri de date. Partea care variază de la un analizor lexical la altul are de fapt două componente și anume

modelele de atomi lexicali căutate în șirul de intrare și acțiunile pe care eventual analizorul lexical trebuie să le execute la identificarea unui șir care se "potrivește" cu un model de atom lexical. Specificațiile lex conțin descrierea părții variabile. Deoarece ce se generează în final este un text C, acțiunile sunt descrise ca secvențe de instrucțiuni C. În aceste secvențe pot să apară orice instrucțiuni C inclusiv apeluri de funcții. Secvențele respective vor "îmbrăca" scheletul părții fixe. Specificarea modelelor de atomi lexicali se face sub formă de *expresii regulate*. Pornind de la ele programul lex va genera tabelele care descriu funcția de tranziție a automatului determinist corespunzător. Aceste tabele reprezintă o parte din structurile de date pe care lucrează "scheletul". Utilizând "scheletul" și secvențele de instrucțiuni oferite de către programator, programul lex va genera un text care reprezintă analizorul lexical. Programul lex "nu înțelege" nici scheletul și nici acțiunile specificate de către programator, aceste elemente reprezintă numai șiruri de caractere pe care lex știe să le combine. În mod corespunzător cel care scrie specificațiile lex trebuie să aibă în vedere aspecte ca domeniul de valabilitate al variabilelor utilizate în programul generat, și asta ținând cont că există de fapt două tipuri de variabile, variabile pe care le controlează numai cel care scrie specificațiile lex și pentru ele se face atât declararea cât și utilizarea în cadrul secvențelor specificate de către programator și variabile care sunt declarate în cadrul scheletului (sunt predefinite din punct de vedere al celui care scrie specificațiile). Declararea acestor variabile apare în textul generat înainte de începutul funcției yylexQ. Referirea acestora se poate face atât în secvențele de cod care se execută asociat diferitelor modele de atomi (aceste secvențe vor fi interne funcției yylex()) cât și în funcții care vor fi incluse în textul generat după textul funcției yylexQ-

2.2.1 Funcționarea analizorului lexical generat de lex

Analizorul lexical generat este de fapt o funcție yylex(). În execuția acestei funcții se parcurge textul de intrare și se caută un subșir care începe cu primul caracter și "se potrivește" cu o expresie regulată specificată în reguli. Dacă există mai multe soluții se va lua în considerare cel mai lung subșir care "se potrivește". Pentru mai multe soluții având aceeași lungime se consideră soluția care a fost găsită prima (în ordinea de parcurgere secvențială a regulilor). În cazul expresiilor regulate care descriu și contextul dreapta al șirului analizat (de exemplu expresiile regulate care termină cu \$ sau care conțin condiții de tip-/) în determinarea lungimii "de potrivire" participă și contextul dreapta. De exemplu pentru expresiile:

```
"ab" {printf(" s-a recunoscut ab");}
"abc" {printf(" s-a recunoscut abc");}
"ab'V'cc" {printf(" s-a recunoscut ab urmat de ce");}
```

pentru șirul abec se va afișa textul:

```
s-a recunoscut ab urmat de ce
```

Variabila globală yytext este un pointer spre începutul șirului care "s-a potrivit". Lungimea șirului este memorată în yyleng. După ce se actualizează valorile variabilelor yytext și yyleft, conform identificării unei expresii regulate se va executa acțiunea asociată expresiei, care poate să utilizeze variabilele yytext și yyleng.

După execuția acțiunii asociate unui subșir găsit se va continua parcurgerea șirului de intrare începând cu caracterul care urmează subșirului selectat. Dacă nu se găsește un subșir care să se potrivească cu o regulă se execută acțiunea implicită, adică se va copia caracterul curent în fișierul de ieșire. Din acest motiv dacă acest efect nu este dorit trebuie să se prevadă o ultimă regulă de forma:

`\n /* se potrivește cu orice caracter */`

care va "înghiți" orice caracter care nu "s-a potrivit". În continuare se va încerca găsirea unei "potriviri" începând cu următorul caracter, etc. A se vedea programele 2.2. - 2.4.

2.2.2 Stări de start

Pentru toate exemplele considerate până acum analizorul generat încearcă în ordine toate regulile specificate pentru a determina cel mai lung subșir care "se potrivește" cu o regulă. Uneori în funcție de context anumite reguli trebuie să fie ignorate. În acest scop într-un program lex se pot utiliza stări de start. Dacă nu se specifică altfel, analizorul este în starea 0 numită simbolic `INITIAL`. Se pot declara și alte nume de stări în prima secțiune sub forma:

```
hs numei      numen

sau

%x          numen
```

În primul caz dacă o acțiune este prefixată de `<numep` atunci acțiunea respectivă se va executa numai dacă starea curentă este `numej`. În același caz toate acțiunile care nu au nici un prefix vor fi luate în considerare indiferent de starea curentă. A doua notație nu este recunoscută de lex dar este recunoscută de flex și de alte variante ale programului lex. Dacă declararea numelui stării s-a făcut cu `x` (starea `nume`; este denumită în acest caz stare exclusivă) atunci acțiunile care nu au prefixul `<numep` nu vor fi executate dacă starea curentă este `nume-j`. Se poate prefixa o regulă și cu o listă de stări. În acest caz regula va fi considerată dacă starea curentă este una dintre stările din listă. Trecerea dintr-o stare în alta se face prin apelarea unei macrodefiniții predefinite `BEGIN`. Forma de apel a acestei macrodefiniții este:

```
BEGIN(nume);
```

unde `nume` este numele stării în care trebuie să treacă analizorul generat. A se vedea programele 2.5. - 2.6.

Utilizarea stărilor exclusive este de preferat celor obișnuite pentru că măresc claritatea programului. Stările exclusive există însă numai în flex. Pentru lex se poate obține o comportare de tip stare exclusivă utilizând o variabilă din programul generat care va fi testată înainte de fiecare verificare.

2.2.3 Macrodefiniții, funcții și variabile predefinite

În cadrul acțiunilor pot să apară orice instrucțiuni C valide, deoarece generatorul de programe lex nu face nici o analiză a acestora. Există însă câteva macrodefiniții lex predefinite care pot să apară în orice acțiune. În discuția despre stări am întâlnit deja una dintre aceste macrodefiniții (`BEGIN`). Mai există alte două macrodefiniții: `ECHO` și `REJECT`.

Efectul apelului macroinstrucțiunii `ECHO` constă din copierea valorii `yytext` la ieșire. Cu alte cuvinte, efectul acestei macroinstrucțiuni constă din copierea la ieșire a subșirului pentru care se execută acțiunea.

Macroinstrucțiunea `REJECT` este utilizată atunci când se dorește modificarea algoritmului de alegere a subșirului care "se potrivește". Și anume s-a precizat că dacă există mai multe soluții se alege cel mai lung subșir. În caz de egalitate, dictează ordinea în care au fost scrise expresiile regulate. Dacă se utilizează macroinstrucțiunea `REJECT` se renunță la cea mai bună soluție și se consideră a doua. A se vedea programele 2.7. - 2.9.

Se recomandă însă să se evite utilizarea macrodefiniției `REJECT` deoarece reduce viteza programului generat și în același timp conduce la programe greu de citit.

În acțiuni se pot utiliza și o serie de funcții predefinite:

- `yymoreO` adaugă ca prefix la următorul șir care "se potrivește" șirul pentru care se execută acțiunea curentă (șirul indicat de `yytext`). Vezi programul 2.10.
- `yyles(n)` se renunță la ultimele `yyleng - n` caractere din `yytext` care vor fi reutilizate pentru formarea subșirului următor. Vezi programul 2.11.
- `unput(c)` caracterul `c` este "pus" în bufferul de intrare devenind următorul caracter de intrare. Vezi programul 2.12.
- `inputO` citește un caracter din șirul de intrare. Această funcție permite amestecarea recunoașterilor de subșiruri făcută de către analizor cu cea făcută în cadrul acțiunilor. Vezi programul 2.13.
- `yyterminate()` - este de fapt o macroinstrucțiune care poate să fie redefinită, reprezintă o instrucțiune `return`. Ca efect se abandonează execuția funcției `yylex()`, adică se produce o revenire forțată din analizor. În mod implicit la întâlnirea sfârșitului de fișier se produce execuția acestei macroinstrucțiuni.

2.2.4 Fișiere de intrare multiple

În toate exemplele considerate până acum fișierul de intrare a fost considerat în mod implicit fișierul standard de intrare. Dar, pentru situațiile reale, fișierul de intrare este un fișier disc. Din acest motiv este necesar un mecanism de specificare a fișierului pe care analizorul generat îl va prelucra. Variabila `yyin` conține întotdeauna numele fișierului din care se citește textul prelucrat de către analizorul generat. Vezi programul 2.14.

Tratarea fișierelor multiple depinde de modul în care acestea sunt utilizate în timp. Evident există două situații. În prima situație este vorba de o înlănțuire de fișiere (adică după ce este tratat un fișier se trece la al doilea apoi la al treilea etc. în orice caz nu este necesară o revenire la un fișier anterior). În a doua situație trecerea la un nou fișier se face înainte de a se termina prelucrarea fișierului curent. La terminarea prelucrării unui fișier se face revenirea la fișierul

anterior (ca de exemplu dacă se tratează directiva #include).

Să considerăm situația în care se face o înlănțuire a fișierelor. In acest caz ceea mai bună soluție este utilizarea funcției yywrap. Această funcție este apelată în mod automat de către analizorul generat în momentul în care se detectează un sfârșit de fișier. Dacă rezultatul apelului funcției este unu înseamnă că nu există un alt fișier de tratat și execuția funcției yylex s-a terminat. Dacă rezultatul apelului funcției este zero înseamnă că prelucrarea nu s-a terminat. În acest ultim caz funcția yywrap realizează și comutarea yyin pe noul fișier de intrare. Vezi programul 2.15.

2.3 Exemple comentate

11:111:1 2-2

```
blanc [ \t]+
%%
{blanc} { printf (" "); }.
{blanc}$ { printf("#"); /* ignoră blancurile de la sfârșit */ }
```

```
main()
{ yylex();}
```

Descriere:

Programul înlocuiește secvențele de blankuri și caractere de tabulare ("t") cu câte un singur blank. De asemenea șterge blankurile de la sfârșitul liniilor înlocuindu-le cu un caracter #. Adică pentru un fișier de intrare de forma:

```
a   b c
  a  b d
```

se va obține un fișier de ieșire de forma:

```
a b c#
abd#
```

Comentarii:

Se observă că o secvență de blankuri aflată la sfârșitul unei linii satisface ambele reguli, dar este luată în considerare cea de a doua regulă care se potrivește subșirului mai lung (cel care conține și caracterul linie nouă, care nu este însă luat în considerare).

```
șterge
"alfa"   { printf("beta"); }
gama     { printf("delta"); }
```

```
itiain() {
  yylex();
```

Descriere:

Șterge aparițiile șirului de caractere "șterge" din text și înlocuiește șirul "alfa" cu "beta", respectiv "gama" cu "delta". Dacă se utilizează un fișier de intrare de forma:

```
alfa șterge alfa
alfa beta gama delta
```

Se va obține un fișier de ieșire de forma:

```
beta beta
beta beta delta delta
```

Comentarii:

În program se utilizează cele două forme echivalente de reprezentare a șirurilor de caractere. Să considerăm și o modificare a programului anterior:

```
%%
șterge ;
"alfa"+   printf('beta');
gama+     printf('delta');
main(){
  yylex ();
} ' ..
```

De data aceasta dacă vom executa programul generat asupra fișierului de intrare:

```
alfaalfa șterge alfaaaaaaaaaaaaaaaaa
gamagama șterge gamaaaaaaaaaaaaaaaaa
```

Se va obține fișierul de ieșire:

```
beta betaaaaaaaaaaaaaaaaa
deltadelta delta
```

Se observă că operatorul "+" se aplică în primul caz asupra întregului șir ("alfa"+) respectiv asupra ultimului caracter, în al doilea caz (gama+).

Acțiunile asociate expresiilor regulate pot să conțină instrucțiuni de forma return expresie. În cazul în care se produce o "potrivire" între un șir de caractere și un astfel de model execuția funcției yylex se termină și se "întoarce" valoarea expresiei. În cazul în care nu se execută o astfel de instrucțiune se continuă parcurgerea șirului de intrare până când acesta se termină. Să considerăm tema 2-4. pentru care fiecare acțiune se încheie cu o instrucțiune return.....


```

/* declarații de inacrodefiniții utilizate în specificarea
regulilor
*/

cifra [0-9]
semn [+]?
IntregFaraSemn {cifra}+
Exponent ([Ee]{semn}{IntregFaraSemn})
Real ({semn}{IntregFaraSemn}\.{IntregFaraSemn}?{Exponent}?)

```

```

{Real} { printf("real\n"); return yyleng;}
{IntregFaraSemn} { printf("intreg\n"); ; return yyleng;}
. { return -1; }
\n { return -2; }
%%
main() {
    printf("lungime = %i\n", yylex());
}

```

Descriere:

Executând programul generat asupra șirului de intrare:

123.56

Se va obține fișierul de ieșire:

```

real
lungime = 6

```

Comentarii:

În programul 2.4 sunt tratate toate cazurile. Adică dacă primul caracter care apare în fișierul de intrare nu reprezintă un început de număr întreg sau real atunci analiza se oprește cu rezultat -1 iar pentru caracterul linie nouă rezultatul este -2.

#ina 2-5

```

char * a;
int numar_aparitii = 0, numar_cifre = 0;
%}
litera [A-Za-z]
cifra [0-9]
%s URMĂTOR
%%
<INITIAL>{litera}+ { a = malloc(yyleng);
strcpy(a, yytext);
numar_aparitii = 1;
BEGIN (URMĂTOR);

```

```

<OKMATOR>{litera}+ { if (!strcmp(a, yytext))
numar_aparitii++;

{cifra} {numar_cifre++;}

main() {
yylex();
printf(„primul cuvânt %s apare de %i ori\n au fost %i cifre”,
a, numar_aparitii, numar_cifre);
free(a);
}

```

Descriere:

Programul numără aparițiile primului cuvânt (șir care conține numai litere) dintr-un text în textul care urmează. De asemenea, programul contorizează și cifrele care apar în text.

Comentarii:

A fost declarată o stare numită URMĂTOR TIVIVIWI în jLvi si.iiv c^iv comandată prin execuția macroinstrucțiunii:

```
BEGIN (URMĂTOR) ;
```

apelată pentru acțiunea corespunzătoare găsirii primului cuvânt. Dacă se execută analizorul generat asupra fișierului de intrare: **••**

```

123$ % &456 abbb bribnbnb abb nbnbnb abbb nbnbnb abbbbbbbbb
789 abbb hjhjhjhjh abb

```

se va obține fișierul de ieșire:

```

primul cuvânt abbb apare de 3 ori
au fost 9 cifre

```

Dacă programul se schimbă prin declararea stării URMĂTOR ca stare exclusivă execuția programului generat pentru același fișier de intrare va produce:

```
789
```

```

primul cuvânt abbb apare de 3 ori
au fost 6 cifre

```

Cifrele nu au fost luate în considerare decât atâta timp cât starea curentă a fost starea INIȚIAL.

ti 111.1 2-(i)

```

x COMENTARIU
{
int NumarLinieComentariu = 0;
int NumarLinieProgram = 0;

```

```

„/*” { BEGIN(COMENTARIU); ++NumarLinieComentariu; }

<CXMENTARIU>[^\n]* /* se ignora orice nu este ,*'sau \n */
<CCMENTARIU>"""*'+P*An)* { /* se ignora orice * care nu este
urmat de ,/' sau linie noua */ }

<CCMENTARIU>\n ++NumarLinieComentariu;
<CCMENTARIU>"""*'+"/" BEGIN (INIȚIAL) ;
\n ++NumarLinieProgram;

main() {
  yylex();
  printf(„au fost %i linii de comentariu \n”,
  NumarLinieComentariu);
  printf(„au fost %i linii de program\n”, NumarLinieProgram);
}

```

Beserie:

Programul generat numără liniile ce conțin comentarii respectiv text „obișnuit” într-un program. Comentariile încep și se termină cu caracterele „/*” respectiv „*/”.

Comentarii:

La începutul comentariului și la sfârșitul său se face comutarea între stările INIȚIAL și COMENTARIU. În starea COMENTARIU se ignoră orice caracter mai puțin caracterele „*” și respectiv „\n”. Pentru caracterul „*” se verifică ce caracter urmează. Se ignoră șirurile de caractere „*” care nu sunt urmate de „/” sau de „\n”. Cazul în care urmează caracterul „/” corespunde sfârșitului de comentariu, caracterele „\n” sunt tratate separat. Pentru caracterele „\n” se face incrementarea numărului de linii din comentariu. Se observă că o linie este contorizată fie ca linie program fie ca linie comentariu după cum caracterul \n a fost întâlnit în starea INIȚIAL sau în starea URMĂTOR. Dacă comentariul este inclus într-o linie care conține și text de program după el atunci linia respectivă va fi contorizată atât ca linie program cât și ca linie comentariu.

```

%%
„acesta” 1
„acesta este” I
„acesta este un” !
„acesta este un test” |
„acesta este un test ciudat” { ECHO; REJECT;}
/* sg ignoră ce nu se potrivește */

main() {
  yylex 0 r
}

```

Descriere:

Dacă se execută programul generat pentru programul asupra unui fișier de intrare care conține textul:

```
acesta este un test ciudat
```

se va obține textul:

```
acesta este un test ciudatacesta este un testacesta este unacesta esteacesta
```

C&mmKarik

Se observă că întâi s-a "potrivit" cel mai lung subșir pentru care s-a executat acțiunea ECHO care l-a tipărit apoi s-a executat acțiunea REJECT prin care s-a renunțat la soluția respectivă. Se va considera următorul subșir (mai scurt). Și pentru acesta se execută acțiunea { ECHO; REJECT} deci se va face afișarea, etc. Macrodefiniția REJECT este utilă dacă pentru un același șir trebuie să se execute două acțiuni diferite. De exemplu în programul 2.8. se recunosc câteva cuvinte cheie, dar în același timp se face și contorizarea tuturor cuvintelor care apar în text (inclusiv cuvintele cheie).

tenia 2-8

```

int NumarCuvinte = 0;
)
cuvânt [A-Za-z][A-Za-z0-9]*
%%
bun |
rau |
urit |
frumos { ECHO; printf " este adjectiv\n"; REJECT;}
maninca |
doarme I
scrie I
citește { ECHO; printf " este verb\n"; REJECT;}
{cuvânt} { NumarCuvinte++;}
/* ignoră orice altceva */

```

```

main () {
  yylex();
  printf("\n au fost %i cuvinte\n", NumarCuvinte);
}

```

Descriere:

Dacă se execută programul generat pentru programul asupra unui fișier de intrare care conține textul:

```
aaa
```

```
bun
```

Se va obține textul:

```
bun este adjectiv
```

```
au fost 2 cuvinte
```

Comentarii:

Subșirul "aaa" se potrivește cu a treia regula, deci el este numărat. Pentru subșirul "bun", acesta se potrivește cu prima regula, deci se afișează mesajul corespunzător și se execută ! macrodefiniția REJECT, adică se renunță la potrivirea găsită. Următoarea potrivire aleasă este regula a treia în care se numără cuvintele. Deci în final numărul de cuvinte este 2.

```
int NumarCuvinte = 0;
%
cuvint [A-Za-z][A-Za-zO-9]*
%%
{cuvint} { NumarCuvinte++; REJECT;}
bun I
rau I
urit I
frumos {
    ECHO;
    printfC este adjectiv\n" );
}
maninca I
doarme I
scrie I
citește { ECHO; printf " este verb\n";}
/* ignora orice altceva */
%%
main (){
    yylex();
    printf("\nau fost %i cuvinte\n", NumarCuvinte);
}
```

Deseriere:

Să presupunem că schimbăm ordinea regulilor și poziția macroinstrucțiunii REJECT din 2. ca in programul curent. Pentru același text de intrare ca în tema 2-8:

```
aaa
bun
se va obține un text de ieșire ca:
bun este adjectiv
au fost 7 cuvinte
```

Comentarii:

Răspunsul pare ciudat, dar ceea ce se întâmplă de fapt este că întâi subșirul "aaa" se potrivește cu prima regulă, deci el este numărat, apoi se renunță la această soluție și se caută altă "potrivire". Se va găsi o potrivire pentru un șir mai scurt "aa", deci și acest cuvânt se numără. Iar se renunță, urmează o potrivire cu subșirul "a", deci contorul ajunge la 3. În acest caz cu prima regulă nu mai există nici o potrivire, nici cu a doua, dar a treia regulă reușește I să "înghită" un "a". Se continuă de la al doilea caracter "a". La terminarea parcurgerii ; subșirului "aaa" contorul este 6. în cazul subșirului "bun" după ce prima regulă a executat macrodefiniția

REJECT se va încerca cu altă variantă și aceasta este potrivirea cu a doua regulă pentru care se consideră aceeași lungime. Această regulă "înghite" subșirul "bun" care deci este numărat o singură dată.

tema 2-10

```
%S DUPĂ
Cuvânt [A-Za-z]*
CuvintSpecial {Cuvânt}"$
CuvintObisnuit {Cuvânt}[\t \n]+
%%
<INITIAL>{CuvintSpecial} { BEGIN(DUPĂ); putchar('(');}
<INITIAL>{CuvintObisnuit} /* se ignora */
<DUPA>{CuvintObisnuit} { yymoreO;}
<DUPA>{CuvintSpecial} { yytext[yyleng -1] = ')';
    ECHO; BEGIN(INIȚIAL);
```

```
main(){
    yylex();
}
```

Descriere:

Programul generat extrage șirurile de caractere cuprinse între caractere "\$" din șirul de intrare și le copiază la ieșire între paranteze simple. Astfel pentru șirul:

```
a$a b$c
a $a b$c
aa$a b$c
aa$a b $c
aa$a b $ c
```

se va obține:

```
( a b ) ( a b ) ( a b ) ( a b )
```

Comentarii:

Soluția din program nu este unică, nici măcar nu este cea mai simplă, dar înțelegerea ei reprezintă un bun exercițiu.

Unu 2-11

```
%s DUPĂ
Cuvânt [A-Za-z]*
CuvintSpecial {Cuvânt}"$
CuvintObisnuit {Cuvânt}[\t \n]+
```

```

<INITIAL>"$" { BEGIN(DUPA); putchar ('('); }
<INITIAL>{CuvintObisnuit} /* se ignora ce a mai rainas */
<INITIAL>{CuvintSpecial} { yyles(yyleng - 1); }
<DUPA>"$"      { BEGIN(INITIAL);
  yytext[yyleng -1] = '(';
  ECHO;
}
<DUPA>{CuvintObisnuit} { yymore(); }
%%
main(){
  yylex();
}

```

Descriere:

O soluție mai complicată cu efect similar cu 2.10.

Comentarii:

Se folosește în acest caz funcția yyles pentru a extrage caracterul "\$

```

Cuvânt ([a-zA-Z][a-zA-ZO-9]*)
CuvintSpecial (\ ({Cuvânt}\))

```

```

{Cuvânt} {
  int i;
  printf("1 leng = %i, text = %s\n", yyleng,
  yytext);
  unput ('(');
  for (i = yyleng - 1; i >= 0; -i)
    unput (yytext[i]);
  unput ('(');
  printf("2 leng = %i, text = %s\n", yyleng, yytext);
}

```

```

{CuvintSpecial} |
\(\) ECHO;

```

```

main() {
  yylex();
}

```

Descriere:

Dacă se execută pentru șirul de intrare

abcde

se va obține:

```

1 leng = 5, text = abcde
2 leng = -1, text = abcd)

```

(abcd)

Comentarii:

Se observă că subșirul recunoscut abcde a fost "prelungit" dar numai la un capăt, primul caracter care este "dat înapoi" ("") înlocuiește ultimul caracter din subșir "e". Apelul yyles(yyleng -1) este echivalent cu unput(yytext[yyleng - 1]).

```

/*" { int c;
  for(;;){
  while((c = inputO) != '*' && c != EOF)
  ; /* se ignora ce nu poate sa fie sfârșit */
  if (c = '*') {
    while ((c = inputO) = '*')

  if (c = '/') {
    break; /* s-a găsit sfârșit */
    printf(" a fost comentariu\n");

  if (c = EOF) {
    printf(" comentariu neterminat \n");
    break;
}

```

Descriere:

Un exemplu de utilizare al funcției input() în tratarea comentariilor cuprinse între caractere duble cum sunt cele din limbajul C.

2-14

```

%{
#define LT 0
#define LE 1
#define EQ 2
#define NE 3
#define GT 4
#define GE 5
#define IF 10
#define THEN 11
#define ELSE 12
#define ID 13
#define NUMĂR 14
}

```

```
#define OPREL 15
#define ATRIB 16

int yyval;
char x[30]; /* tabela de simbolii ( o intrare) */
```

```
/* macroinstrucțiuni */
```

```
DELIMITATOR [ \t\n]
BLANC {DELIMITATOR}+
litera [A-Za-z]
ifra [0-9]
id {litera}({litera}|{cifra})*
număr {cifra}+(\.{cifra}+)?(E[+-]?{cifra}+)?
```

```
{BLANC} { /* nici o acțiune */}
if { printf(" IF"); }
then { printf(" THEN"); }
else { printf(" ELSE"); }
{id} { yyval = Tratareld0;
printf("%i, %i", yyval, ID);
```

```
{număr} { yyval = TratareNumar();
printf("%i, %i", yyval, NUMĂR);
```

```
{ yyval = LT; printf("%i, %j yyval, OPREL);}
{ yyval = LE; printf("%i, %J yyval, OPREL);}
{ yyval = EQ; printf(" li, %i yyval, OPREL);}
"O" '{ yyval = NE; printf("%i, %: yyval, OPREL);}
{ yyval = GT; printf("%i, %: yyval, OPREL);}
{ yyval = GE; printf("%i, li yyval, OPREL);}
{ yyval = ATRIB; printf(" %i, ", yyval, ATRIB),
```

```
int TratareNumar(void){
printf("\n număr");
return 1;
```

```
main( int argc, char **argv ){
++argv, --argc;
if ( argc > 0 )
yyin = fopen( argv[0], "r" );
else
yyin = stdin;
yylex ();
```

Descriere:

Programul reprezintă specificațiile flex pentru un subset al limbajului Pascal. Analizorul va recunoaște identificatori și numere, pentru care va apela funcții corespunzătoare și o parte din operatorii și cuvintele cheie care pot să apară într-un program Pascal.

Comentarii:

Funcția TratareNumar ar fi putut să fie inclusă și în prima secțiune de program deoarece nu referă variabile ce nu au fost încă definite. În schimb funcția Tratareld utilizează variabilele yytext și yyleng, deci nu poate să apară decât în a treia secțiune de program. Programul principal este mai complicat față de variantele anterioare și anume, înainte de apelul funcției yylex se va stabili valoarea variabilei yyin care poate să fie stdin sau un nume specificat de utilizator.

linia 2-15

```
fundef yywrap
int NumarCaractere = 0, NumarCuvinte = 0, NumarLinii= 0;
```

```
BLANC [ \t]
NOBLANC [^ \t\n]
```

```
{NOBLANC}+NumarCaractere += yyleng; ++NumarCuvinte;
{BLANCJ+ NumarCaractere += yyleng;
\n ++NumarLinii; ++NumarCaractere;
<< EOF >> yyterminate ();
```

```
char **ListaFisiere;
unsigned int FisierCurent = 0;
unsigned int NumarFisiere;
```

```
main( int argc, char **argv ) {
FILE *fișier; ListaFisiere = argv + 1;
NumarFisiere = argc -1;
if (argc > 1) {
FisierCurent = 0;
fișier = fopen(ListaFisiere[FisierCurent], "r");
if (!fișier){
printf("!!!eroare!!!");
exit(1);
}
yyin = fișier;
yylex ();
```

```
/* Noua definiție yywrap */
```

```
yywrap(){
    FILE *fisier;
    fisier = 0;
    fclose(yyin);
    printf("\nFisier %s", ListaFisiere[FisierCurent++]);
    printf("\nNumar caractere = %8d", NumarCaractere);
    printf("\nNumar cuvinte = %8d", NumarCuvinte);
    printf("\nNumar linii = %8d", NumarLinii);
    NumarCaractere = 0;
    NumarCuvinte = 0;
    NumarLinii = 0;
    if (.FisierCurent > NumarFisiere) return 1;
    fisier = fopen(ListaFisiere[FisierCurent], "r");
    if (!fisier){
        printf("!!! eroare!!!");
        exit(1);
    }
    yyin = fisier;
    return(fisier ? 0:1);
}
```

Descriere:

Programul calculează numărul de caractere, cuvinte și linii conținute într-o listă de fișiere transmise prin linia de comandă prin care este apelat programul generat. În programul principal se face deschiderea primului fișier. La fiecare apel al funcției yywrap se va face închiderea fișierului precedent, afișarea rezultatelor obținute pentru acesta și deschiderea fișierului următor. După fiecare deschidere de fișier variabila globală yyin este poziționată pe indicatorul fișierului deschis. În momentul în care s-a ajuns la sfârșitul listei funcția yywrap va avea rezultatul 1 și execuția funcției yylex se termină.

Comentarii:

Se observă că a fost necesară utilizarea directivei:

```
#undef yywrap
```

pentru a se putea defini o nouă funcție yywrap deoarece există o macrodefiniție yywrap definită ca:

```
#define yywrap() 1
```

3 Teme pentru acasă

tema 3-1

Sase scrie gramaticile ce generează limbajele:

- $L = \{ (ab)^n a^i b^j c^k \mid n \geq 0, i, j, k > 0 \}$
- $L = \{ a^i (ab)^n b^j c^k \mid i \geq 0, n, j, k > 0 \}$
- $L = \{ a^i v c^j v^b \mid j \geq 0, n, i, k > 0 \}$
- $L = \{ a^i b^j (ab)^n c^k \mid k \geq 0, n, j, i > 0 \}$
- $L = \{ w e \{a, b\} \#_a(w) \text{ par si } \#b(w) \text{ impar} \}$
- $L = \{ \mathbf{we} \{a, b\}^* \#_a(w) \text{ par si } \#b(w) \text{ par} \}$
- $L = \{ \mathbf{we} \{a, b\}^* \#_a(w) \text{ impar si } \#_b(w) \text{ par} \}$
- $L = \{ \mathbf{wg} \{a, b\}^* \#_a(w) \text{ impar si } \#b(w) \text{ impar} \}$
- $L = \{ w e \{a, b\} \text{ w nu conține șirul baab } \}$
- $L = \{ w e \{a, b\} \text{ w nu conține șirul abba } \}$
- $L = \{ w e \{a, b\} \text{ w nu conține șirul bbaa } \}$
- $L = \{ w e \{a, b\}^* \text{ w nu conține șirul abbb } \}$

\ Să se demonstreze folosind lema de pompare că următorul limbaj nu este independent de context:

- $\{j\}$
- $L = \{ a^i b^j c^k \mid i > j > k \}$
- $iL = \{ a^i b^j c^k \mid 0 \leq i \leq j \leq k \}$
- $^L = \{ a^i b^j c^k \mid i, j, k \geq 0, i! = j, i! = k, j! = k \}$

Pentru următoarele gramatici să se elimine recursivitatea stângă. În gramaticile rezultate să se elimine X producțiile, după care să se elimine simbolii nefinalizați și inaccesibili.

$G = (\{S, A, B\}, \{a, b\}, P, S)$ unde P este:

- $S \rightarrow Ba \mid Ab$
- $A \rightarrow Sa \mid Aab \mid a$
- $B \rightarrow Sb \mid Bba \mid b$

$G = (\{S, A, B, C\}, \{a, b\}, P, S)$ unde P este:
 $S \rightarrow A|B$
 $A \rightarrow Ba|Sb|b$
 $B \rightarrow AB|Ba$
 $C \rightarrow AS|b$

$G = (\{S, A, B, C\}, \{a, b\}, P, S)$ unde P este:
 $S \rightarrow AB|AC$
 $B \rightarrow Bb|b$
 $C \rightarrow Ac|BC$
 $A \rightarrow a|C$

$G = (\{S, A, B, C\}, \{a, b\}, P, S)$ unde P este:
 $S \rightarrow AB|AC$
 $A \rightarrow Sa|a$
 $B \rightarrow BC|AB$
 $C \rightarrow aB|b$

tema 4

: Pentru următoarele gramatici, să se elimine simbolii nefinalizați și inaccesibili.

$I G = (\{S, A, B\}, \{a, b\}, P, S)$ unde P este
 $S \rightarrow SS|Aaa$
 $A \rightarrow aAB|aS|b$
 $B \rightarrow bB$

$G = (\{S, A, B, C\}, \{a, b\}, P, S)$ unde P este
 $S \rightarrow AB|CA$
 $A \rightarrow a$
 $B \rightarrow BC|AB$
 $C \rightarrow aB|b$

$G = (\{S, A, B\}, \{a, b\}, P, S)$ unde P este
 $S \rightarrow a|A$
 $A \rightarrow AB$

$G = (\{S, A, B\}, \{a, b\}, P, S)$ unde P este
 $S \rightarrow A$
 $A \rightarrow bS|b$
 $C \rightarrow AS|b$

tema 3-5

Pentru expresiile regulate:

$(a|b)^*ab(a|b)^*$
 $(a|b)^*ab^*(a|b)$
 $(a|b)^*a(a|b)b^*$
 $(a|b)^+ab^*$

Să se construiască automatele finit nedeterminist și determinist și să se minimizeze stările.
 Să se facă și construcția directă pornind de la expresia regulată la automatul finit determinist.

Uina 3-6

Să se construiască automatele cu stivă care acceptă limbajele generate de gramaticile:

$G = (\{S, A, B\}, \{a, b\}, P, S)$ unde P este
 $S \rightarrow aB|bA$
 $A \rightarrow a|aS|bAA$
 $B \rightarrow b|bS|aBB$
 $G = (\{S, A\}, \{a, b\}, P, S)$ unde P este
 $S \rightarrow aAA$
 $A \rightarrow aS|bS|a$

$G = (\{S, A, B, C\}, \{a, b\}, P, S)$ unde P este
 $S \rightarrow ABC$
 $A \rightarrow BB|X$
 $B \rightarrow CC|a$
 $C \rightarrow AA|b$

u-niii '- "

Să se construiască mașinile Turing compuse care decide limbajele:

$$L = \{w \in \{a, b\}^* \mid \#_a(w) \text{ par și } \#_b(w) \text{ impar}\}$$

$$L = \{w \in \{a, b\}^* \mid \#_a(w) \text{ par și } \#_b(w) \text{ par}\}$$

$$L = \{w \in \{a, b\}^* \mid \#_a(w) \text{ impar și } \#_b(w) \text{ par}\}$$

$$L = \{w \in \{a, b\}^* \mid \#_a(w) \text{ impar și } \#_b(w) \text{ impar}\}$$

4 Bibliografie

1. Irina Athanasiu, Limbaje Formale si Compilatoare, Centrul de multiplicare, IPB, 1992
2. Alfred Aho, Jeffrey Ullman, The Theory of Parsing, Translation and Compiling, voi. 1
Parsing
3. Harry R. Lewis, Christos H. Papadimitriou, Elements of the theory of computation