

Tehnici de compilare

Mircea Drăgan

November 7, 2019

Cuprins

1	Introducere	3
1.1	De ce învățăm tehnici de compilare?	4
1.2	Limbaje de programare și traducerea programelor	5
1.3	Structura unui compilator	8
2	Analiza lexicală	11
2.1	Analiza lexicală	11
2.2	Minimizarea automatelor finite deterministe	20
2.2.1	Studiu de caz	24
2.3	Construcția expresiei regulate pentru un limbaj regulat .	28
3	Analiza sintactică	33
3.1	Algoritmi TOP-DOWN	34
3.1.1	Algoritmul general de analiză top-down	34
3.1.2	Programarea algoritmului top-down general	36
3.1.3	Analiza top-down fără reveniri	41
3.1.4	Programarea unui analizor sintactic. Studiu de caz	43
3.2	Algoritmi BOTTOM-UP	50
3.2.1	Gramatici cu precedență simplă	50
3.2.2	Relații de precedență	52
3.2.3	Proprietăți ale gramaticilor cu precedență simplă	53
3.2.4	Determinarea relațiilor de precedență pentru gramatici cu precedență simplă	55
3.2.5	Studiu de caz	56
3.2.6	Gramatici operatoriale	59
3.2.7	Gramatici operatoriale	61

3.2.8	Determinarea relațiilor de precedență pentru gra-	
	matici operatoriale	66
3.2.9	Studiu de caz	66
3.3	Analiza sintactică LR(k)	68
3.3.1	Construcția analizorului pentru gramatici LR(0) .	70
3.3.2	Analizoare SLR	75
3.3.3	Gramatici LR(1) și analizoare LALR(1)	78

Capitolul 1

Introducere

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

Există în prezent un număr mare de limbaje de programare de *nivel înalt* sau evolute care se caracterizează printr-o anumită naturalețe, în sensul că descrierea procesului de prelucrare cu ajutorul limbajului este apropiată de descrierea naturală a procesului respectiv. Dintre limbajele de acest tip cu o anumită răspândire în momentul de față menționăm limbajele PASCAL, FORTRAN, C, JAVA, etc. Un program redactat într-un limbaj de programare poartă denumirea de **program sursă**.

Fiecare sistem de calcul, în funcție de particularitățile sale, posedă un anumit limbaj propriu, numit *cod mașină*, acesta fiind singurul limbaj înțeles de procesorul sistemului. Un astfel de limbaj conține un număr mic de instrucțiuni elementare de manipulare a datelor (dar care se execută foarte rapid). Un program redactat în limbajul cod mașină al sistemului de calcul îl numim **program obiect**.

Procesul de transformare al unui program sursă în program obiect se numește **compilare** sau *translatare*, uneori chiar *traducere*.

Tehnicile de compilare sunt tehnici de programare specializate utilizate în primul rând pentru la scrierea programelor de translatare, aplicabile la realizarea unei game de programe similare translatoarelor.

Primul compilator a fost realizat în 1957-1958 pentru limbajul FOR-

TRAN, de către o echipă de cercetatori de la I.B.M., condusă de John Backus. Acum, scrierea unui compilator se face mult mai ușor și simplu datorită elaborării unei teorii a compilării fundamentate matematic de Teoria limbajelor formale și apariția unor unelte software specializate pentru generarea automată a unor componente ale translatorului.

1.1 De ce învățăm tehnici de compilare?

Desigur că foarte puțini informaticieni vor avea nevoie să scrie un compilator pentru un limbaj general de nivel înalt, ca de exemplu C, Pascal, Java. Totuși, iată câteva motive pentru care acest curs apare în planul de învățământ al majorității instituțiilor de învățământ superior din domeniul Computer Science:

- Se consideră că tehniciile de compilare împreună cu teoria limbajelor formale face parte din cunoștințele de bază (*cultura generală obligatorie*) ale unui informatician.
- Un bun meseriaș (*profesionist*) trebuie să își cunoască sculele, iar pentru un informatician limbajele și compilatoarele sunt printre principalele instrumente de lucru.
- Tehnicile de compilare folosite pentru verificarea corectitudinii programelor sunt utile și în alte ramuri ale informaticii (*Sisteme de operare, Baze de date, Prelucrarea textelor, etc.*).
- Este foarte posibil ca un informatician să scrie un minicompiletor sau interpreter pentru un limbaj specializat.

Primul motiv pare un pic *exagerat* deoarece este legat de istoricul dezvoltării informaticii ca știință și aceasta este un domeniu în continuă schimbare, deci este greu de definit ce înseamnă cultură generală în domeniu.

Al doilea motiv este mult mai intemeiat deoarece înțelegerea modului de construcție a compilatoarelor va ajuta programatorul în depanarea programelor scrise în limbi de nivel înalt și optimizarea programului ce se execută pe mașina de calcul.

1.2. LIMBAJE DE PROGRAMARE SI TRADUCEREA PROGRAMELOR

Cunoașterea mecanismelor de verificare a programelor ajută la înțelegerea și manipularea facilă a oricărui tip de text structurat, ca de exemplu XML. Ultimul motiv devine tot mai actual datorită folosirii limbajelor specializate (*DSL–Domain Specific Languages*) destinate unei clase de probleme, de exemplu: interogarea bazelor de date, formatarea textelor și a imaginilor, simularea economică, etc.

1.2 Limbaje de programare si traducerea programelor

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

Limbajele de programare de *nivel înalt*, gen PASCAL, FORTRAN, C, JAVA, sunt oarecum independente față de sistemul de calcul.

O altă clasă importantă de limbi, sunt *limbajele de asamblare*, sau de nivel inferior, ale căror caracteristici depind de sistemul de calcul considerat. În general, fiecare sistem de calcul (sau tip de sistem de calcul), are propriul său limbaj de asamblare; de exemplu, limbajul de asamblare ale sistemelor de calcul echipate cu procesoare de tip Intel Z-80 este denumit în mod curent ASSEMBLER. Instrucțiile unui limbaj de asamblare corespund cu operațiile simple ale sistemului de calcul iar stocarea datelor în memorie este realizată direct de utilizator la nivelul locațiilor elementare ale memoriei. Există de asemenea anumite *pseudo-instrucții* sau *directive* referitoare la alocarea memoriei, generaarea datelor, segmentarea programelor, etc., precum și macroinstructiile care permit generarea unor secvențe tipice de program sau accesul la bibliotecile de subprograme.

In afara de limbajele evoluate și de limbajele de asamblare, există numeroase *limbaje specializate*, numite uneori și de *comandă*, care se referă la o anumită clasă de aplicații. Menționăm de exemplu limbajul pentru prelucrarea listelor LISP, limbajele utilizate în cadrul softwarelui matematic (Mathematica, Maple, MATCAD, etc.) și multe altele. În

general însă astfel de limbaje nu sunt considerate limbaje de programare propriuizise.

Translatarea unui **program sursă** în **program obiect** se numește **compilare** în cazul limbajelor evolute, iar în cazul limbajelor de asamblare se utilizează termenul de *asamblare*.

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

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

Procesul de compilare este un proces relativ complex și comportă operații care au un anumit caracter de autonomie. Din aceste motive procesul de compilare este de obicei descompus în mai multe subprocese sau faze, fiecare fază fiind o operație coerentă, cu caracteristici bine definite. În principiu aceste faze sunt parcurse secvențial (pot să existe și anumite reveniri) iar programul sursă este transformat succesiv în formate intermediare. Se poate considera că fiecare fază primește de la faza precedentă un fișier cu programul prelucrat într-un mod corespunzător fazei respective, îl prelucrează și furnizează un fișier de ieșire, iarăși într-un format bine precizat, fișier utilizat în faza următoare. Există cinci faze de compilare principale: analiza lexicală, analiza sintactică, generarea formatului intermediar, generarea codului, optimizarea codului și două faze auxiliare, tratarea erorilor și tratarea tabelelor (vezi figura 1.1).

1.2. LIMBAJE DE PROGRAMARE SI TRADUCEREA PROGRAMELOR 7

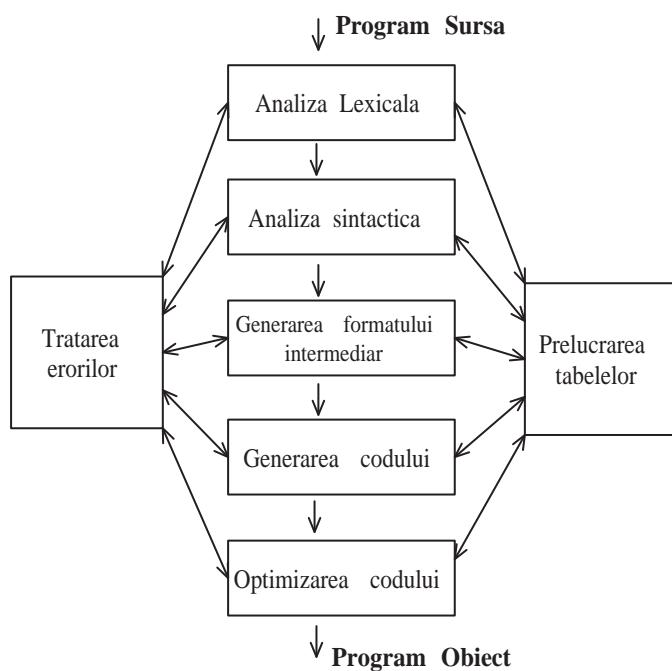


Figura 1.1: Fazele compilării

1.3 Structura unui compilator

Analiza lexicală are ca obiectiv principal determinarea unităților lexicale ale unui program (*secvențe de caractere cu înțeles logic*), furnizarea codurilor acestora și detectarea eventualelor erori lexice. De exemplu, în secvența de program $ion = a + 3 * c$ sunt găsite următoarele unități lexice, sau **atomii** (*tokens* în limba engleză):

- identificatorii ion, a și c;
- operatorii $=$, $+$, $*$;
- numărul întreg 3;

Pe lângă aceste operații de bază, la analiza lexicală se mai pot efectua anumite operații auxiliare precum: eliminarea spațiilor albe (dacă limbajul permite utilizarea fără restricții a acestora), ignorarea comentariilor, diferite conversiuni ale unor date (care se pot efectua la această fază), completarea tabelelor compilatorului.

Analiza sintactică grupează ierarhic atomii în colecții mai mari, numite unități sintactice ale programului (secvențe de text pentru care se poate genera format intermediu: expresiile, instrucțiunile, declarațiile, etc.) și verifică programul din punct de vedere sintactic. Este faza centrală a unui compilator, deseori toate celelalte faze sunt rutine apelate de analizorul sintactic pe măsură ce este posibilă efectuarea unei părți din faza respectivă. Deseori se folosește o variantă a arborelui de derivare pentru gramatici independente de context, numit **arbore sintactic**, în care fiecare nod reprezintă o operație iar fiile sunt argumentele operației (fig. 1.2). Tot la analiza sintactică se definitivează tabelele de informații (verificarea tipurilor, analiza domeniului de vizibilitate, etc) și se realizează prelucrarea erorilor sintactice.

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

Generarea codului este o fază în care se realizează codul obiect corepunzător programului. Practic în această fază se obține programul

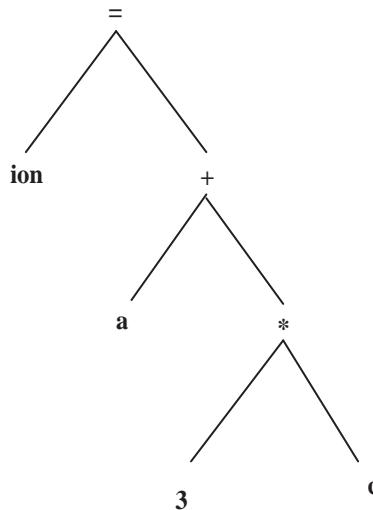


Figura 1.2: Arborele sintactic pentru $ion = a + 3 * c$

în limbaj de asamblare corespunzător programului sursă redactat într-un limbaj evoluat. În mod obișnuit generatorul de cod este constituit din rutinele generatoare de cod corespunzătoare diverselor unități ale formatului intermedian.

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

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

Delimitarea fazelor compilării ca în structura precedentă urmărește organizarea logică a unui compilator. În implementări concrete fazele

se grupează de obicei în **front-end** și **back-end**. Front end cuprinde de obicei fazele dependente de limbajul sursă și include analiza lexicală, analiza sintactică, crearea tabelei de simboluri și generarea formatului intermedian. Back end include acele porțiuni ce depind de calculatorul destinație și include generarea și optimizarea codului.

Capitolul 2

Analiza lexicală

2.1 Analiza lexicală

Procesul de analiză lexicală este o fază a procesului de compilare în care se determină unitățile lexicale (cuvintele, atomii) ale unui program sursă, se furnizează codurile interne ale acestora și se detectează eventualele erori lexicale. Analizorul lexical mai poate efectua o serie de operații auxiliare precum: eliminarea blank-urilor, ignorarea comentariilor, diferite conversiuni ale unor date, completarea tabelelor compilatorului, gestiunea liniilor textului sursă.

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

- Unități lexicale simple, sunt unități lexicale care nu comportă atrbute suplimentare, de exemplu, cuvintele cheie, operatorii;
- Unități lexicale compuse (atributive), sunt unități lexicale care comportă anumite atrbute suplimentare, de exemplu, identificatorii și constantele. Atributele sunt informații specifice, de exemplu, tipul identificatorului sau al constantei (întreg, real, etc.).

Este necesară specificarea tipului unui identificator sau a unei constante din startul programului deoarece structura programului obiect sau reprezentarea internă a constantelor depinde de acest tip.

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

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

Un exemplu de unități lexicale și coduri asociate ar putea fi cele din figura 2.1.

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

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

va produce sirul de unități lexicale

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

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

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

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

Figura 2.1: Coduri asociate unitătilor lexicale

```

comentarii          /* text */

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

```

Înaintea analizei lexicale (sau ca subrutină) se preprocesează textul și abia apoi se introduce rezultatul în analizorul lexical.

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

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

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

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

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

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

unde l -literă, c -cifra.

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

Descrieri echivalente ale unităților lexicale cu ajutorul expresiilor regulate sunt

cuvinte cheie = if | do | else | for

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

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

Operatori aritmetici = + | - | * | /

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

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

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

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

De obicei textul sursă conține și secvențe ce se pot descrie cu ajutorul expresiilor regulate, dar care nu sunt unități lexicale (de exemplu

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

Există posibilitatea ca mai multe secvențe cu aceeași origine să corespundă la diferite descrieri ale unităților lexicale. Se consideră unitate lexicală cel mai lung sir ce se potrivește unei descrieri (**longest match rule**). Dacă sunt două reguli care se potrivesc la același sir de lungime maximă atunci se consideră o prioritate asupra descrierilor (**rule priority**).

De exemplu, în textul următor,

```
i if if8
```

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

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

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

Etapă I. Descriem secvențele cu ajutorul expresiilor regulate

$$\text{IF} = \text{"if"}$$

$$\text{ID} = (\text{a}|\text{b}|\text{c}|\dots|\text{z})(\text{a}|\text{b}|\text{c}|\dots|\text{z}|\text{0}|\text{1}|\dots|\text{9})*$$

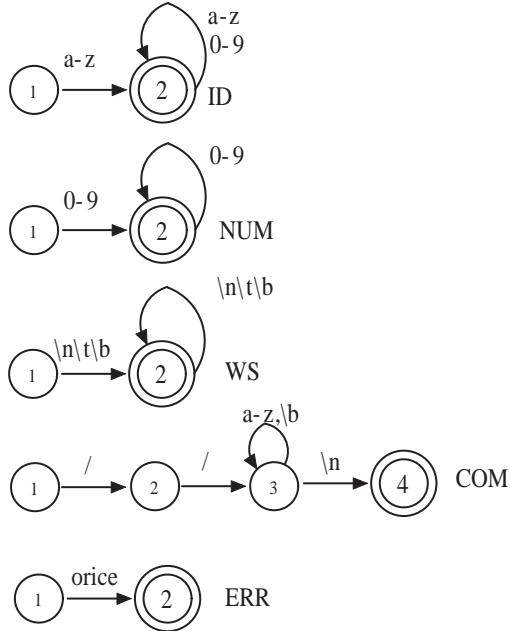


Figura 2.2: Automatele finite corespunzătoare expresiilor regulate

$$\text{NUM} = (0|1|\dots|9)(0|1|\dots|9)^* = (0|1|\dots|9)^+$$

$$\text{WS} = (\backslash n|\backslash t|" ")^+$$

$$\text{COMMENT} = //"(a|b|c|\dots|z|0|1|\dots|9|" ")^*\backslash n$$

$$\text{ALL} = a|b|\dots|z|0|\dots|9|\backslash t| \dots \text{ toate caracterele ASCII}$$

Etapa II. Corespunzător expresiilor avem următoarele automate finite deterministe echivalente, prezentate în figura 2.2 (stările au fost notate prin numere întregi):

Etapa III. Se construiește sistemul tranzițional (vezi figura 2.3) ce recunoaște limbajul reuniune, adăugând o nouă stare inițială (notată cu 1), pe care o conectăm prin arce punctate (ce corespund λ -tranzițiilor). Construcția provine din *legarea* sistemelor tranziționale *în paralel*. S-

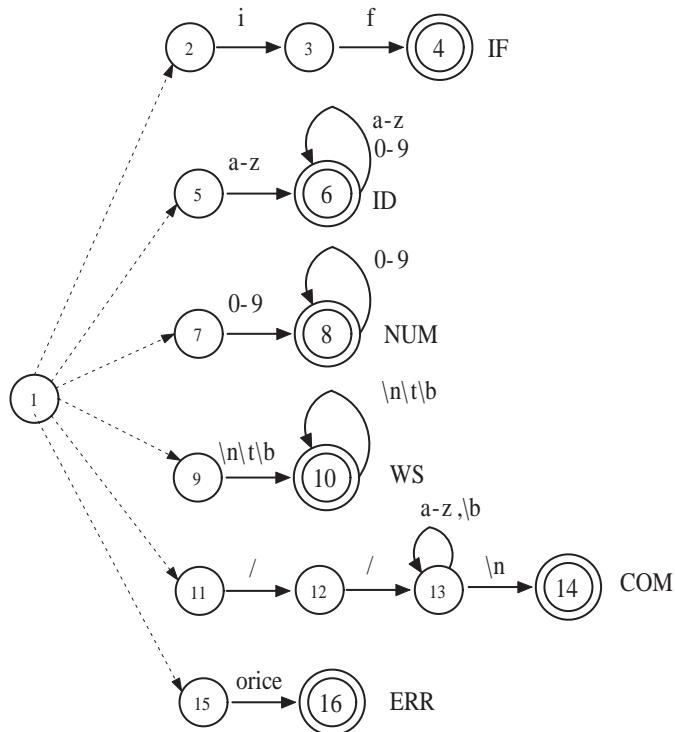


Figura 2.3: Sistemul tranzitonal ce recunoaște reuniunea limbajelor

au renumerotat stările sistemului tranzitonal, asignând nume simbolice stărilor finale.

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

În cazul particular al automatului nostru, diagrama de stări este dată în figura 2.4.

Etapa IV. Programarea analizorului lexical.

Automatul finit obținut are stări finale asociate cu clasele de cuvinte recunoscute. Se asociază acțiuni stărilor finale, corespunzător

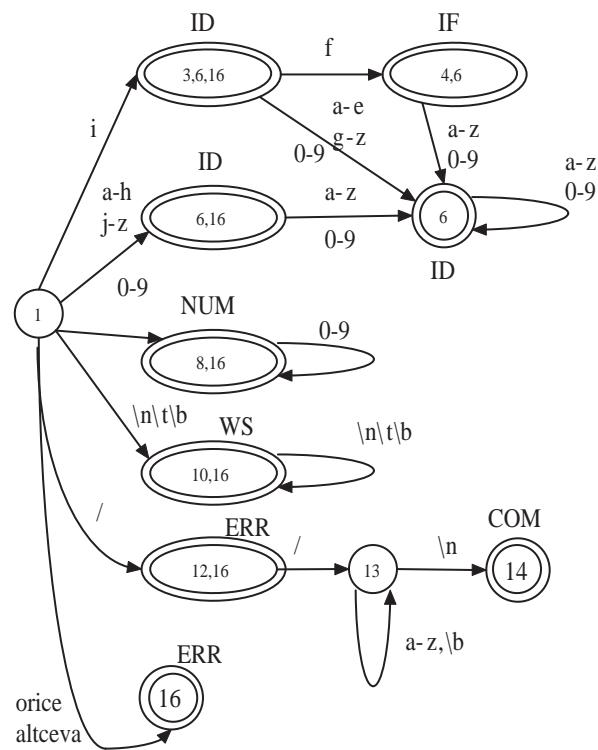


Figura 2.4: Automatul finit determinist ce recunoaște reuniunea limbajelor

definițiilor unităților lexicale (de exemplu pentru constante numerice se generează reprezentarea internă, se memorează în tabelul de constante și se returnează codul unității lexice NUM). Pentru programarea analizorului se folosesc trei variabile de tip pointer în textul sursă: **FirstSymbol**, **CurrentSymbol**, **LastFinalSymbol**, ce rețin indicele caracterului de început al secvenței, indicele caracterului ce urmează la citire, indicele ultimului caracter ce corespunde atingerii unei stări finale. Cei trei pointeri au fost reprezentați prin semnele grafice |, ⊥ respectiv ⊤. De asemenea considerăm o variabilă *State*, ce reține starea curentă a automatului.

In tabelul 2.5 este indicată evoluția analizorului lexical (inclusiv acțiunile asociate) pentru cazul analizei următorului text sursă.

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

Pentru simplificarea codificării, stările automatului finit determinist au fost redenumite prin numere întregi începând cu starea inițială 1, starea {3, 6, 16} = 2, {4, 6} = 3, {6, 16} = 4, și.a.m.d. de la stânga la dreapta și de sus în jos. De obicei funcția de evoluție asociată automatului finit determinist se memorează sub forma unui tablou bidimensional de întregi, ca în figura 2.6. Starea 0 este asociată cu blocarea automatului. Ajungerea în această stare echivalează cu găsirea ultimei unități lexice, între pointerii | și ⊤. Se execută acțiunea asociată stării finale și se reia căutarea (*resume*) următoarei unități lexice începând cu caracterul imediat următor pointerului ⊤.

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

2.2 Minimizarea automatelor finite deterministe

In procesul de generare automată a unui analizor lexical apare uneori o fază de optimizare a automatului finit determinist în sensul determinării unui automat finit determinist echivalent care să aibă un număr minim

Last Final	Current State	Current Input	Accept Action
0	1	_⊥ <i>f</i> i f 8 % // h a \n	
2	2	<i>i</i> ^T _⊥ <i>f</i> i f 8 % // h a \n	
3	3	<i>i f</i> ^T _⊥ i f 8 % // h a \n	
	0	<i>i f</i> ^T _⊥ i f 8 % // h a \n	return <i>cc</i> = < <i>if</i> > resume
0	1	<i>i f</i> _⊥ ^T i f 8 % // h a \n	
7	7	<i>i f</i> _⊥ ^T <i>i f</i> 8 % // h a \n	
	0	<i>i f</i> _⊥ ^T <i>i f</i> 8 % // h a \n	resume
0	1	<i>i f</i> _⊥ ^T <i>i f</i> 8 % // h a \n	
2	2	<i>i f</i> <i>i</i> ^T _⊥ <i>f</i> 8 % // h a \n	
3	3	<i>i f</i> <i>i f</i> ^T _⊥ 8 % // h a \n	
5	5	<i>i f</i> <i>i f</i> 8 _⊥ ^T % // h a \n	
	0	<i>i f</i> <i>i f</i> 8 _⊥ ^T % // h a \n	return <i>id</i> = < <i>if8</i> > resume
0	1	<i>i f</i> i f 8 _⊥ ^T % // h a \n	
11	11	<i>i f</i> i f 8 % _⊥ ^T // h a \n	print ("illegal character: %"); resume
	0	<i>i f</i> i f 8 % _⊥ ^T // h a \n	
0	1	<i>i f</i> i f 8 % _⊥ ^T // h a \n	
0	8	<i>i f</i> i f 8 % _⊥ ^T // h a \n	
0	9	<i>i f</i> i f 8 % _⊥ ^T // h a \n	
...	

 Figura 2.5: Evoluția analizorului lexical pentru textul *if if8%// ha\n*

```

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

```

Figura 2.6: Reprezentarea funcției de evoluție a automatului finit

de stări, deci o reprezentare internă de dimensiune mai mică. Procedeul folosit pentru optimizare se bazează pe teorema de caracterizare algebraică a limbajelor regulate ce oferă o variantă de automat finit minimal.

Teorema 2.1 Caracterizarea limbajelor regulate

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

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

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

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

Pentru demonstrația completă se poate revedea cursul de Limbaje formale și teoria automatelor din anul I. Folosind clasele de echivalență determinate de congruență μ se poate construi un automat finit ce recunoaște limbajul E astfel":

$AF_\mu = (I^*/\mu, I, f, C_\lambda, \Sigma_f)$,
unde funcția de evoluție f și mulțimea de stări finale sunt

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

Automatul astfel definit are proprietatea de minimalitate din punctul de vedere al numărului de stări. Deoarece construcția automatului presupune identificarea claselor de echivalență a unei relații definite pe o multime infinită, procedeul nu oferă un algoritm acceptabil de minimizare a automatului finit. Din acest motiv, minimizarea se face prin gruparea stărilor unui automat finit dat folosind o relație de echivalență pe multimea finită a stărilor și definirea unui automat finit echivalent ce are proprietatea că este *izomorf* cu AF_μ .

Întreaga construcție se bazează pe proprietatea de *separabilitate* a două stări, care provine de la următoarea observație. Pentru oricare două cuvinte dintr-o clasă de echivalență definită de μ , $p, q \in \mu$ avem $\chi_E(pr) = \chi_E(qr), \forall r \in I^*$, adică $f(s_0, pr) \in \Sigma_f$ dacă și numai dacă $f(s_0, qr) \in \Sigma_f$ pentru orice cuvânt r . Deci, stările $s = f(s_0, p)$, $s' = f(s_0, q)$ au proprietatea că $f(s, r) \in \Sigma_f$ dacă și numai dacă $f(s', r) \in \Sigma_f$ pentru orice cuvânt arbitrar r . Se spune că stările s și s' sunt *inseparabile* (prin cuvinte de orice lungime). Cum orice automat definește o partiționare a lui I^* în clase de echivalență incluse în totalitate în clasele determinante de μ , adică o clasă a lui μ este formată prin reuniunea unui număr finit de clase determinante de automat, seincearcă o partiționare a multimii stărilor în submulțimi pentru care are loc proprietatea inseparabilitate caracteristică tuturor claselor din reuniune.

Vom da o descriere a procedeului după modelul prezentat de T. Jucan în [4]. Fie $E \subset I^*$ un limbaj regulat peste alfabetul I și $AFD = (\Sigma, I, f, s_0, \Sigma_f)$ un automat finit determinist ce recunoaște E .

Definiție 2.1 Spunem că un cuvânt $p \in I^*$ separă stările $s_1, s_2 \in \Sigma$ în raport cu Σ_f dacă una și numai una din stările $f(s_1, p)$ și $f(s_2, p)$ se află în Σ_f , adică $\chi_{\Sigma_f}(f(s_1, p)) \neq \chi_{\Sigma_f}(f(s_2, p))$.

Definiție 2.2 Stările s_1 și s_2 se zic k -inseparabile în raport cu Σ_f , notat $s_1 \stackrel{k}{=} s_2$, dacă pentru orice cuvânt $p \in I^*$ cu lungimea $|p| \leq k$ avem $\chi_{\Sigma_f}(f(s_1, p)) = \chi_{\Sigma_f}(f(s_2, p))$.

Definiție 2.3 Stările s_1 și s_2 se zic inseparabile în raport cu Σ_f , notat simplu $s_1 \equiv s_2$ dacă sunt k -inseparabile în raport cu Σ_f pentru orice număr natural k .

Este evident că relațiile \equiv^k și \equiv sunt relații de echivalență pe mulțimea stărilor. Pentru determinarea relațiilor între stări sunt utile următoarele două leme, ușor de demonstrat.

Lema 2.1 *Pentru orice $k \geq 1$, $s_1, s_2 \in \Sigma$ avem $s_1 \stackrel{k}{\equiv} s_2$ dacă și numai dacă $s_1 \stackrel{k-1}{\equiv} s_2$ și $f(s_1, i) \stackrel{k-1}{\equiv} f(s_2, i)$, $\forall i \in I$.*

Lema 2.2 *Dacă $|\Sigma| = n$ atunci stările s_1 și s_2 sunt inseparabile dacă și numai dacă sunt $(n - 2)$ -inseparabile.*

Pentru calculul relațiilor se pornește cu $\stackrel{0}{\equiv}$, inseparabilitate prin cîvinte de lungime zero, care împarte mulțimea stărilor în două clase de echivalență Σ_f și $\Sigma - \Sigma_f$, apoi se aplică de cel mult $n - 2$ ori lema 2.1.

Următorul rezultat pe care îl enunțăm fără demonstrație, dă un mod de construcție al automatului finit minimal.

Teorema 2.2 *Fie E un limbaj regulat acceptat de un automat finit deterministic $AFD = (\Sigma, I, f, s_0, \Sigma_f)$ neblocabil (adică $f(s, i) \neq \emptyset, \forall s \in \Sigma, i \in I$) și fără stări inaccessible (adică fiecare stare este extremitatea unei traectorii ce pornește din starea inițială). Automatul $AFD' = (\Sigma/\equiv, I, f', C_{s_0}, \Sigma'_f)$, unde $f'(C_s, i) = C_{f(s,i)}$, $\forall i \in I$ și $\Sigma'_f = \{C_s | s \in \Sigma_f\}$ recunoaște limbajul E și este izomorf cu automatul minimal AF_μ definit de relația μ din teorema de structură algebraică a limbajelor regulate.*

2.2.1 Studiu de caz

Considerăm un automatul finit deterministic descris prin diagrama de stări din figura 2.7. Să observăm că automatul este **complet** (din fiecare nod avem exact un arc pentru fiecare simbol de intrare) dar conține stări inaccessible. Starea D nu poate să apară pe nici o traectorie ce incepe cu starea inițială A , deci starea D se elimină din graf împreună cu arcele ce pornesc din ea.

La construcția automatului finit minimal echivalent se folosesc lemele 2.1 și 2.2 pentru identificarea claselor de echivalență determinate de relația \equiv pe mulțimea stărilor. O posibilitate de organizare a calculelor

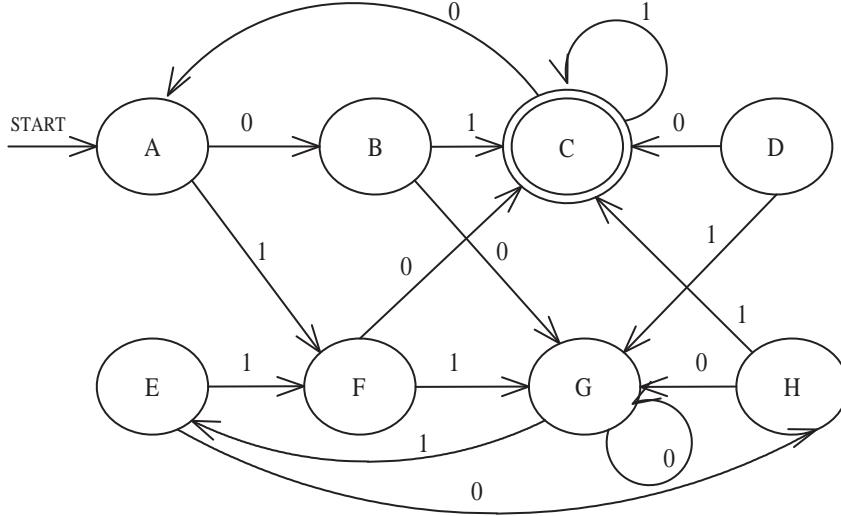


Figura 2.7: Automatul finit determinist

este cea prezentată în figura 2.8 unde se încearcă marcarea perechilor de stări separabile prin cuvinte de lungime $0, 1, 2, \dots, n - 2$.

Inițial se marchează în tabel perechile de stări s_i, s_j separabile prin cuvinte de lungime 0, adică perechile de stări ce conțin o stare finală și una nefinală, exact ca în figura 2.8.

Pentru calculul perechilor de stări separabile prin cuvinte de lungime 1 se aplică lema 2.1, adică se parcurge tabelul și pentru fiecare pereche de stări nemarcată se marchează perechea dacă există un simbol de intrare a astfel încât perechea $f(s_i, a), f(s_j, a)$ este deja marcată. În exemplul nostru vom marca perechea A, H deoarece $f(A, 1) = F$ și $f(H, 1) = C$ iar în tabel perechea de stări F, C este deja marcată, apoi marcăm perechea A, F și.m.d. Prin parcurgerea tabelului de la stânga la dreapta și de sus în jos vom obține tabelul descris în figura 2.9

Se repetă procedeul descris la pasul precedent (adică marcarea perechilor nemarcate încă dacă există simboluri de intrare ce ne conduc prin aplicarea funcției de evoluție la o pereche deja marcată) până când la o parcurgere nu se mai pot adăuga marcaje. Numărul de iterări nu

	H	G	F	E	C	B
A					*	
B					*	
C	*	*	*	*		
E						
F						
G						

Figura 2.8: Tabelul inițial

	H	G	F	E	C	B
A	*		*		*	*
B		*	*	*	*	
C	*	*	*	*		
E	*		*			
F	*	*				
G	*					

Figura 2.9: Tabel ce conține marcajele după prima parcurgere

depășește $n - 2$ conform lemei 2.2.

Tabelul final este prezentat în figura 2.10 unde se observă că avem doar două perechi de stări echivalente A, E și B, H . Diagrama de stări pentru automatul finit minimal echivalent se obține imediat pornind de la diagrama 2.7 alegând pentru construcția arcelor orice reprezentant al clasei de echivalență. Rezultatul este descris de figura 2.11.

2.3 Construcția expresiei regulate pentru un limbaj regulat

Expresiile regulate sunt notații pentru limbaje regulate (se poate construi un *sistem tranzițional* ce recunoaște limbajul notat de expresie folosind legarea sistemelor tranziționale în serie, paralel sau scurtcircuit vezi [2]). Problema care rămâne de rezolvat este cum se poate construi o expresie regulată pentru un limbaj recunoscut de un automat finit.

Vom considera un automat finit determinist $AFD = (S, I, f, s_0, S_f)$ cu n stări $S = \{s_1, s_2, \dots, s_n\}$ unde starea inițială este $s_0 = s_1$. Putem presupune că automatul este complet și fără stări inaccesibile, ca în algoritmul de minimizare a unui AFD .

Ideea construcției provine din asocierea unui cuvânt recunoscut cu o traекторie ce pornește cu starea inițială și se termină cu o stare finală din S_f . Se consideră mulțimile de cuvinte ce definesc traectorii de la starea s_i la starea s_j , iar stările intermediare au indici cel mult k . Formal,

$$\mathbf{R}_{i,j}^k = \{p \in I^* | f(s_i, p) = s_j\}, i, j \in \{1, \dots, n\}, k \geq 0$$

și pe traекторie toate stările intermediare au indici cel mult k .

Pentru $k = 0$ se consideră mulțimile de cuvinte ce definesc traectorii fără stări intermediare, adică

$$\mathbf{R}_{i,j}^0 = \begin{cases} \{a \in I | f(s_i, a) = s_j\} & , i \neq j \\ \{a \in I | f(s_i, a) = s_j\} \cup \{\lambda\} & , i = j \end{cases}$$

Este evident că pentru mulțimile $\mathbf{R}_{i,j}^0$ avem notații cu expresii regulate $\mathbf{r}_{i,j}^0$ deoarece sunt mulțimi finite. De exemplu pentru mulțimea $\{a, b, c\}$

2.3. CONSTRUCTIA EXPRESIEI REGULATE PENTRU UN LIMBAJ REGULAT 29

	H	G	F	E	C	B
A	*	*	*		*	*
B		*	*	*	*	
C	*	*	*	*		
E	*	*	*			
F	*	*				
G	*					

Figura 2.10: Tabel cu stările separabile marcate

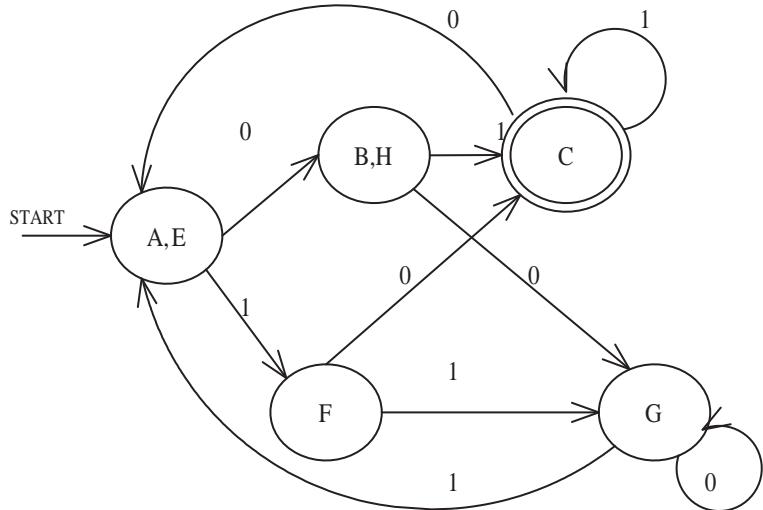


Figura 2.11: Automatul finit determinist minimal

notația este $a|b|c$. Pentru calculul mulțimilor $\mathbf{R}_{i,j}^k$, $k > 0$ se folosește formula

$$\mathbf{R}_{i,j}^k = \mathbf{R}_{i,j}^{k-1} \cup \mathbf{R}_{i,k}^{k-1} (\mathbf{R}_{k,k}^{k-1})^* \mathbf{R}_{k,j}^{k-1}$$

Justificarea formulei este următoarea: Pentru un cuvânt din $\mathbf{R}_{i,j}^k$ traectoria pornește de la starea s_i și apoi parcurge stări intermediare cu indici maxim k până ajunge la starea s_j . Dacă toți indicii stărilor intermediare sunt mai mici decât k atunci cuvântul face parte din $\mathbf{R}_{i,j}^{k-1}$. Dacă există pe traекторie și indicele k atunci putem parta cuvântul în subcuvinte folosind stările s_k într-un sir de subcuvinte astfel: de la s_i la prima apariție a stării s_k (adică un cuvânt din $\mathbf{R}_{i,k}^{k-1}$), apoi un subcuvânt cuprins între prima și a doua apariție pe traectorie a stării k (adică un cuvânt din $\mathbf{R}_{k,k}^{k-1}$), apoi un subcuvânt cuprins între a doua și a treia apariție pe traectorie a stării k (adică un cuvânt din $\mathbf{R}_{k,k}^{k-1}$), și a. m. d., iar după ultima apariție a stării s_k avem un cuvânt ce ne duce în starea s_j (adică un cuvânt din $\mathbf{R}_{k,j}^{k-1}$).

Formula de calcul poate fi transformată într-o formulă cu expresii regulate, adică

2.3. CONSTRUCTIA EXPRESIEI REGULATE PENTRU UN LIMBAJ REGULAT 31

$$\mathbf{r}_{i,j}^k = \mathbf{r}_{i,j}^{k-1} | \mathbf{r}_{i,k}^{k-1} (\mathbf{r}_{k,k}^{k-1})^* \mathbf{r}_{k,j}^{k-1}$$

Este evident că limbajul recunoscut de automatul finit este

$$L(afd) = \bigcup_{s_i \in S_f} \mathbf{R}_{1,i}^n$$

care poate fi notată de expresia regulată $|_{s_i \in S_f} \mathbf{r}_{1,i}^n$.

Capitolul 3

Analiza sintactică

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

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

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

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

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

3.1 Algoritmi TOP-DOWN

3.1.1 Algoritmul general de analiză top-down

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

Pentru descrierea riguroasă a acestui algoritm am urmat modelul din cursul *Limbaje Formale și Tehnici de Compilare* predat la UVT în anii 80 de către Dl. Prof. Ștefan Mărușter, unul dintre pionierii informaticii în România.

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

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

Considerăm toate regulile care au X_0 în stânga:

$$x_0 \rightarrow x_{11} \dots x_{1n_1} | x_{21} \dots x_{2n_2} | \dots | x_{k1} \dots x_{kn_k}$$

Initial x_0 devine *activ* și alege prima regulă $x_0 \rightarrow x_{11} \dots x_{1n_1}$. Dacă această alegere este corectă trebuie să avem $x_0 \Rightarrow x_{11} \dots x_{1n_1} \xrightarrow{*} p$ și în conformitate cu *lema de localizare* a gramaticilor de tipul 2, cuvântul

p se poate descompune în forma $p = p_1 p_2 \dots p_{n_1}$, unde $x_{1j} \xrightarrow{*} p_j$, $j = 1, \dots, n_1$.

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

Observații

- Fiecare simbol devine activ, procedează exact ca și părintele său, alege prima regulă, activează primul simbol, etc.
- Nu se cunoaște anticipat descompunerea $p = p_1 p_2 \dots p_{n_1}$. Deci x_{1j} transmite *succes* dacă reușește să găsească o derivare $x_{1j} \xrightarrow{*} p_j$, unde p_j este un subcuvânt oarecare al lui p , cu singura condiție ca p_{1j} să înceapă din punctul unde s-a terminat p_{1j-1} . De exemplu, dacă $p = i_1 i_2 \dots i_8 \dots$ și $p_1 = i_1 i_2 i_3 i_4$, $p_2 = i_5 i_6$, atunci x_{13} trebuie să găsească o derivare de forma $x_{13} \xrightarrow{*} i_7 i_8 \dots$. În particular, dacă $x_{1j} \in V_T$ decizia de *succes* sau *eșec* depinde de faptul dacă x_{1j} coincide sau nu cu simbolul din p care urmează (În exemplul de mai sus, dacă x_{13} coincide sau nu cu i_7).
- Când se reactivează un simbol și i se cere o nouă derivare, acesta reactivează ultimul simbol fiu și îi cere același lucru.

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

$$\left\{ \begin{array}{l} E \rightarrow T + E|T \\ T \rightarrow F * T|F \\ F \rightarrow (E)|a \end{array} \right.$$

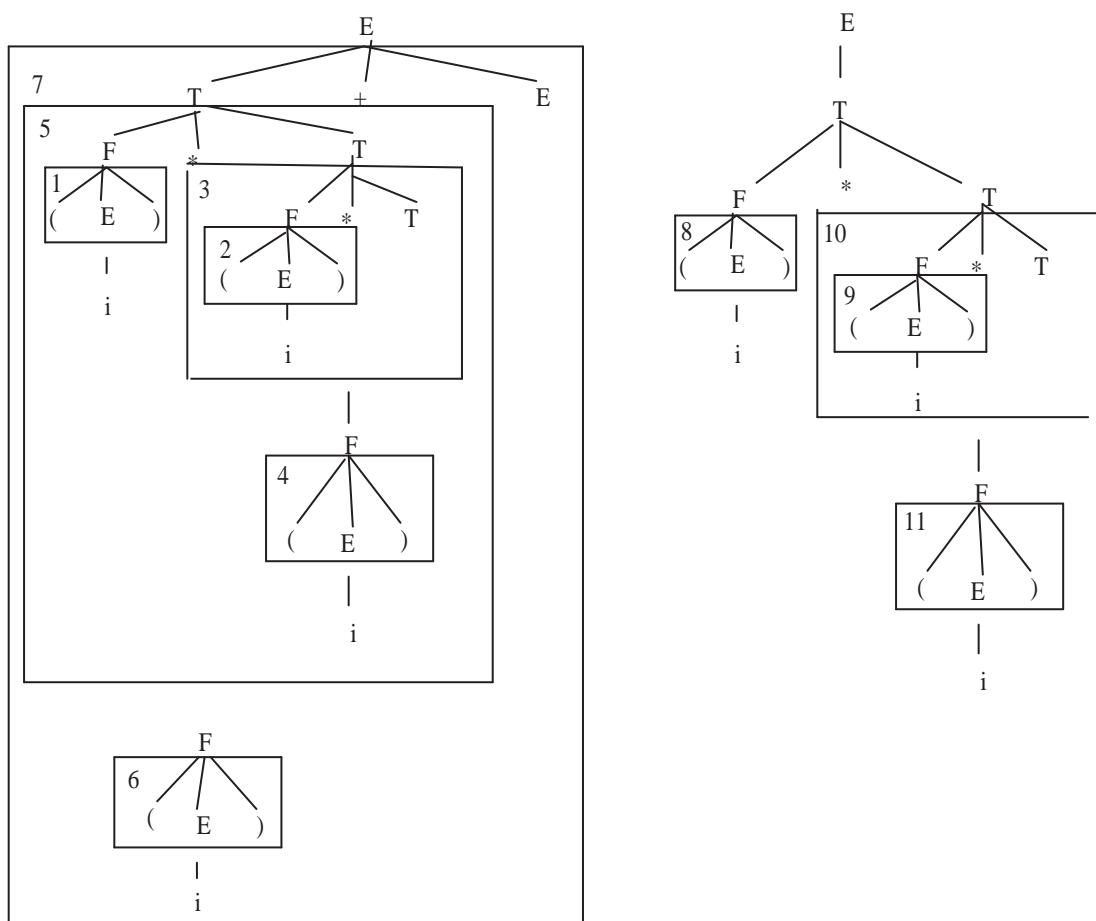
Procesul de analiza sintactică top-down pentru cuvântul $p = a * a$ este prezentat în figura 3.1.

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

3.1.2 Programarea algoritmului top-down general

Descrierea formală a algoritmului este făcută după Aho și Ullman, *The theory of Parsing, Translation, and Compiling*. Vom considera că regulile gramaticii sunt numerotate, iar pentru regulile cu neterminalul A în partea stangă, $A \rightarrow \alpha_1|\alpha_2|\dots|\alpha_k$ vom nota A_i indexul regulii i . Algoritmul folosește două stive $L1$ și $L2$, precum și un index ce reprezintă caracterul curent în sirul de intrare $w = a_1a_2\dots a_n, n >= 0$. Pentru descrierea formală vom defini o configurație printr-un cvadruplu (s, i, α, β) , unde s este *starea algoritmului*, i reprezintă indexul locației curente (pointer pe următorul caracter din sirul de intrare), α este conținutul stivei $L1$ iar β este conținutul stivei $L2$. Stiva $L2$ va conține la fiecare pas forma propozițională din derivarea extrem stângă, iar simbolul aflat în vârful stivei va fi simbolul activ din arborele de derivare. Stiva $L1$ va conține un istoric al alegerii regulilor și simbolurile de intrare care au fost deja generate (*prefixele viabile* ale sirului de intrare). Pentru ușurință prezentării vom considera că se folosește un marcat de sfârșit pentru cuvântul examinat, notat cu $\$$, acesta fiind pe poziția $n + 1$.

Algoritmul folosește trei stări q , b și t cu semnificațiile urmatoare: q operație normală, b revenire (*backtracking*) și t terminare. Configurația inițială a algoritmului va fi $(q, 1, \lambda, S\$)$. Funcționarea algoritmului va fi

Figura 3.1: Arborele sintactic pentru $p = i * i$

o secvență de pași discrete ce pot fi descriși printr-o relație pe mulțimea configurațiilor, notată cu \rightarrow ca și în cazul evoluției automatului push down.

Conform cu descrierea informală a algoritmului, avem următoarele tipuri de pași:

- **Expandarea arborelui**

$$(q, i, \alpha, A\beta) \rightarrow (q, i, \alpha A_1, \gamma_1 \beta)$$

unde, $A \rightarrow \gamma_1$ este prima regula din listă, aplicată celui mai din stânga neterminál în forma propozițională curentă.

- **Potrivirea simbolului terminal din sirul de intrare cu terminalul activat de algoritm**

$$(q, i, \alpha, a\beta) \rightarrow (q, i + 1, \alpha a, \beta)$$

Dacă simbolul din poziția curentă i din sirul de intrare coincide cu următorul terminal derivat, $a_i = a$, atunci terminalul a aflat în vârful stivei $L2$ se extrage și se scrie pe stiva $L1$. Totodată se avansează indexul i în sirul de intrare.

- **Terminare cu succes**

$$(q, n + 1, \alpha, \$) \rightarrow (t, n + 1, \alpha, \lambda)$$

S-a găsit marcajul de sfârșit al sirului examinat și avem o formă propozițională ce coincide cu sirul examinat. Lista regulilor aplicate pentru obținerea derivării extrem stângi se poate obține prin eliminarea terminalelor de pe stiva $L1$.

- **Nepotrivirea simbolului terminal activ cu următorul simbol din sirul de intrare**

$$(q, i, \alpha, a\beta) \rightarrow (b, i, \alpha, a\beta) \quad a_i \neq a$$

Se intră în modul revenire deoarece forma propozițională actuală nu va genera cuvântul de examinat.

- **Revenire pe șirul examinat**

$$(b, i, \alpha a, \beta) \xrightarrow{} (b, i - 1, \alpha, a\beta), \forall a \in V_T$$

În modul revenire se mută (**shift**) simbolurile de intrare înapoi de pe stiva $L1$ pe $L2$.

- **Încercarea unei reguli noi pentru un ultimul neterminal activat**

$$(b, i, \alpha A_j, \gamma_j \beta) \xrightarrow{}$$

- $(q, i, \alpha A_{j+1}, \gamma_{j+1} \beta)$, dacă γ_{j+1} este următoarea regulă din lista neterminalului A . Adică se înlocuiește partea dreaptă a ultimei reguli încercate cu partea dreaptă a următoarei reguli.
- blocarea algoritmului, atunci când $i = 1, A = S$ și există doar j reguli pentru simbolul de start. În această situație șirul de intrare w nu aparține limbajului generat de gramatica dată.
- $(b, i, \alpha, A\beta)$ pentru orice altă situație. Este cazul când toate regulile lui A au fost încercate și revenirea se face prin extragerea numărului regulii A_j de pe stiva $L1$ și înlocuirea părții drepte α_j cu neterminalul A în stiva $L2$.

EXEMPLU: Considerăm gramatica simplificată ce generează expresiile aritmetice fără paranteze, pentru care regulile sunt numerotate astfel:

$$\left\{ \begin{array}{ll} (1) & E \rightarrow T + E \\ (2) & E \rightarrow T \\ (3) & T \rightarrow F * T \\ (4) & T \rightarrow F \\ (5) & F \rightarrow a \end{array} \right.$$

În descriere vom utiliza E_1 pentru $T + E$, E_2 pentru T , samd. Examinarea șirului de intrare $w = a + a$ va fi descrisă de următoarea secvență de configurații.

$$\begin{aligned}
(q, 1, \lambda, E\$) &\xrightarrow{} (q, 1, E_1, T + E\$) \\
&\xrightarrow{} (q, 1, E_1 T_1, F * T + E\$) \\
&\xrightarrow{} (q, 1, E_1 T_1 F_1, a * T + E\$) \\
&\xrightarrow{} (q, 2, E_1 T_1 F_1 a, *T + E\$) \\
&\xrightarrow{} (b, 2, E_1 T_1 F_1 a, *T + E\$) \\
&\xrightarrow{} (b, 1, E_1 T_1 F_1, a * T + E\$) \\
&\xrightarrow{} (b, 1, E_1 T_1, F * T + E\$) \\
&\xrightarrow{} (q, 1, E_1 T_2, F + E\$) \\
&\xrightarrow{} (q, 1, E_1 T_2 F_1, a + E\$) \\
&\xrightarrow{} (q, 2, E_1 T_2 F_1 a, +E\$) \\
&\xrightarrow{} (q, 3, E_1 T_1 F_1 a +, E\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_1, T + E\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_1 T_1, F * T + E\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_1 T_1 F_1, a * T + E\$) \\
&\xrightarrow{} (q, 4, E_1 T_2 F_1 a + E_1 T_1 F_1 a, *T + E\$) \\
&\xrightarrow{} (b, 4, E_1 T_2 F_1 a + E_1 T_1 F_1 a, *T + E\$) \\
&\xrightarrow{} (b, 3, E_1 T_2 F_1 a + E_1 T_1 F_1, a * T + E\$) \\
&\xrightarrow{} (b, 3, E_1 T_2 F_1 a + E_1 T_1, F * T + E\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_1 T_2, F + E\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_1 T_2 F_1, a + E\$) \\
&\xrightarrow{} (q, 4, E_1 T_2 F_1 a + E_1 T_2 F_1 a, +E\$) \\
&\xrightarrow{} (b, 4, E_1 T_2 F_1 a + E_1 T_2 F_1 a, +E\$) \\
&\xrightarrow{} (b, 3, E_1 T_2 F_1 a + E_1 T_2 F_1, a + E\$) \\
&\xrightarrow{} (b, 3, E_1 T_2 F_1 a + E_1 T_2, F + E\$) \\
&\xrightarrow{} (b, 3, E_1 T_2 F_1 a + E_1, T + E\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_2, T\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_2 T_1, F * T\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_2 T_1 F_1, a * T\$) \\
&\xrightarrow{} (q, 4, E_1 T_2 F_1 a + E_2 T_1 F_1 a, *T\$) \\
&\xrightarrow{} (q, 4, E_1 T_2 F_1 a + E_2 T_1 F_1 a, *T\$) \\
&\xrightarrow{} (b, 4, E_1 T_2 F_1 a + E_2 T_1 F_1 a, *T\$) \\
&\xrightarrow{} (b, 3, E_1 T_2 F_1 a + E_2 T_1 F_1, a * T\$) \\
&\xrightarrow{} (b, 3, E_1 T_2 F_1 a + E_2 T_1, F * T\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_2 T_2, F\$) \\
&\xrightarrow{} (q, 3, E_1 T_2 F_1 a + E_2 T_2 F_1, a\$) \\
&\xrightarrow{} (q, 4, E_1 T_2 F_1 a + E_2 T_2 F_1 a, \$) \\
&\xrightarrow{} (t, 4, E_1 T_2 F_1 a + E_2 T_2 F_1 a, \lambda)
\end{aligned}$$

O derivare extrem stângă se obține prin eliminarea terminalelor și asocierea numărului regulii cu notația rămasă pe stivă, adică ordinea de aplicare a regulilor este 145245.

3.1.3 Analiza top-down fără reveniri

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

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

$$G \left\{ \begin{array}{l} < \text{program} > \rightarrow \text{Program } D \ I \ \text{EndProgram} \\ D \rightarrow d; X \\ X \rightarrow d; X | \lambda \\ I \rightarrow iY \\ Y \rightarrow ; \ iY | \lambda \end{array} \right.$$

Să considerăm urmatorul cuvânt de analizat

Program

d;
d;
i;
i;
i

EndProgram

O derivare extrem stângă pentru acest cuvânt este

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

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

Dacă se consideră gramatica ce generează expresii aritmetice simple, în forma considerată la algoritm general top-down (3.2.9), atunci la primul pas al unei derivări extrem stângi pentru cuvântul $(a + a * a) + a$ nu se poate decide regula de ales prin citirea primului terminal (deoarece ar fi necesar să consultăm sirul rezultat până la întâlnirea operatorului $+$ aflat după paranteza închisă).

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

$$\left\{ \begin{array}{l} E \rightarrow TE' \\ E' \rightarrow +TE' | -TE' | \lambda \\ T \rightarrow FT' \\ T' \rightarrow *FT' | /FT' | \lambda \\ F \rightarrow (E) | id | num \end{array} \right.$$

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

Figura 3.2: Gramatica pentru generarea blocurilor de instrucțiuni.

3.1.4 Programarea unui analizor sintactic. Studiu de caz

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

Să considerăm următoarea gramatică (vezi figura 3.2) ce generează un bloc de instrucțiuni de atribuire, condiționale de tip **if** (expresie aritmetică) **then** instrucțiune sau repetitive de tip **while**. Instrucțiune poate fi o instrucție simplă sau bloc de instrucțiuni separate prin delimitatorul ;(SEMICOLON).

Un text sursă ce poate fi analizat de această gramatică este cel din figura 3.3.

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

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

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

Figura 3.3: Program sursa de analizat sintactic

Un program pentru analiza sintactică top-down fără reveniri corespunzător gramaticii din figura 3.2 este reprezentat parțial în figura 3.4. Lista codurilor de unități lexicale (terminalele gramaticii) conține SEMI LPAR RPAR §.a.m.d. Presupunem că analizorul lexical funcționează ca funcție ce returnează codul următoarei unități lexicale din textul sursă. Variabila de tip întreg *token* conține codul returnat de analizorul lexical *ALEX()*. S-au mai definit două funcții: *err()* pentru tratarea erorilor de sintaxă depistate și *eat(int tok)* ce consumă din textul sursă unitatea lexicală *tok* pe care ne așteptăm să o găsim în text.

Rularea programului pentru exemplul din figura 3.3 va produce o secvență de apeluri recursive ca în figura 3.5.

Un pic de algebră

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

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

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

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

```

final int if=1, then=2, LPAR=3, RPAR=4, LBRACE=5, RBRACE=6,
PLUS=7, MINUS=8, ORI=9, DIV=10, SEMI=11, id=12,
while=13, do=14, LET=15;
void err();
int ALEX();
int token = ALEX();
void eat(int tok){
    if (tok==token) token=ALEX() else err()
};

void bloc(){
    eat(LBRACE); lista(); eat(RBRACE)
};

void lista(){
    instr();L()
};

void L(){
    switch(token){
        case SEMI: eat(SEMI); instr(); L(); break;
        default: break }
};

void instr(){
    switch(token){
        case if: eat(if); eat(LPAR); E(); eat(RPAR); eat (then); instr(); break;
        case id: eat(id); eat(LET); E(); break;
        case while: eat(while); eat(LPAR); E(); eat(RPAR); eat(do); instr(); break;
        case LBRACE: eat(LBRACE); instr(); eat(RBRACE); break;
        default: printf("syntax error: if, identif, while or left brace expected");
                  err()
    };
    ...
};

void E(){
    T(); Eprime();
};

void Eprime(){
switch(token) {
    case PLUS: eat(PLUS);T();Eprime();break;
    case MINUS: eat(MINUS);T();Eprime();break;
    default: break;
}
};

...

```

```

bloc()
    eat(LBRACE);
    lista()
        instr()
            eat(id);
            eat(LET);
            E()
                T()
                    F()
                        eat(num);
                        return_F;
                        Tprime()
                        return_Tprime;
                        return_T;
                        Eprime()
                        return_Eprime;
                        return_E;
                    return_instr;
                L()
                    eat(SEMI)
                    instr()
                        .... aici se recunoaste ; b = b + 1
                    return_instr;
                L()
                    eat(SEMI);
                    instr()
                        .... aici se recunoaste ; if ( ... )
                    return_instr;
                L()
                    .... aici se recunoaste ; c = b * 72
                    return_L();
                return_L;
            return_Lista;
            eat(RBRACE);
        return_bloc;
    
```

Figura 3.5: Execuția analizei lexicale pentru textul sursă

Pentru $k = 1$ se obține cazul gramaticilor din secțiunile precedente. Restricția asupra numărului de simboluri citite în avans restrânește drastic mulțimea limbajelor ce pot fi analizate cu astfel de gramatici. Un inconvenient major pentru $k > 1$ este creșterea exponențială a dimensiunii tabelei de predicție (câte o intrare în tabel pentru fiecare combinație de k simboluri). O variantă mai eficientă de analiză se obține cu gramatici $LR(1)$ (*Left to right parsing, Rightmost derivation, 1 symbol lookahead*).

Pentru orice $\alpha \in V_G^+, X \in V_N, X \rightarrow \alpha \in \mathcal{P}$ vom considera următoarele mulțimi:

$$\begin{aligned} Prim(\alpha) &= \{a \in V_T \mid \alpha \xrightarrow{*} a\beta, \beta \in V_G^*\} \\ Urm(X) &= \{a \in V_T \mid S \xrightarrow{*} uXv, a \in Prim(v)\} \\ SD(X, \alpha) &= \{a \in V_T \mid a \in Prim(\alpha) \text{ sau } (\alpha \xrightarrow{*} \lambda \text{ și } a \in Urm(X))\} \\ NULL(G) &= \{X \in V_N \mid X \xrightarrow{*} \lambda\} \end{aligned}$$

Dacă un terminal $a \in Prim(\alpha)$ atunci a poate apărea pe prima poziție într-un sir derivat din α . Mulțimea $SD(X, \alpha)$ poartă denumirea de *mulțimea simbolurilor directoare* asociate regulii $X \rightarrow \alpha$, iar $NULL(G)$ conține neterminalele ce se pot șterge (în unul sau mai mulți pași) cu reguli din G .

Teoremă.

$$G \in LL(1) \Leftrightarrow SD(X, \alpha) \cap SD(X, \beta) = \emptyset, \forall X \rightarrow \alpha, X \rightarrow \beta \in \mathcal{P}, \beta \neq \alpha$$

Determinarea mulțimilor definite anterior se poate face ușor. Pentru $NULL(G)$, se poate aplica algoritmul definit la eliminarea regulilor de ștergere pentru o gramatică independentă de context.

```

 $NULL = \{X \in V_N \mid X \rightarrow \lambda \in \mathcal{P}\}$ 
repeat
  (1)  $AUX = NULL;$ 
  (2)  $NULL = NULL \cup \{X \in V_N \mid X \rightarrow p, p \in AUX^+\};$ 
until  $NULL = AUX \square$ 

```

Calculul mulțimilor $Prim(X)$ și $Urm(X)$ se poate face cu algoritmul următor:

Pas 1: { Inițializare }

$$\text{Prim}(a) = \{a\}; \forall a \in V_T$$

$$\text{Prim}(X) = \text{Urm}(Y) = \Phi; \forall X, Y \in V_N$$

Pas 2: **Repeat**

Pentru fiecare regulă $X \rightarrow Y_1 Y_2 \dots Y_k$ **Do**

For($i = 1; i \leq k; i++$)

(2.1)**if** ($i = 1$ sau $Y_1 Y_2 \dots Y_{i-1} \in \text{NULL}^+$)

$$\text{Prim}(X) = \text{Prim}(X) \cup \text{Prim}(Y_i);$$

(2.2)**if** ($i = k$ sau $Y_{i+1} \dots Y_k \in \text{NULL}^+$)

$$\text{Urm}(Y_i) = \text{Urm}(Y_i) \cup \text{Urm}(X);$$

(2.3)**For** ($j = i + 1; j \leq k; j++$)

if ($j = i + 1$ sau $Y_{i+1} \dots Y_{j-1} \in \text{NULL}^+$)

$$\text{Urm}(Y_i) = \text{Urm}(Y_i) \cup \text{Prim}(Y_j);$$

end for

end for

end do

Until $\text{Prim}(X), \text{Urm}(X)$ nu se schimbă la o iterație. \square

Calculul mulțimilor $\text{Prim}(\alpha)$ revine la definiția recursivă

$$\text{Prim}(X\alpha) = \begin{cases} \text{Prim}(X) & \text{daca } X \notin \text{NULL}(G) \\ \text{Prim}(X) \cup \text{Prim}(\alpha) & \text{daca } X \in \text{NULL}(G) \end{cases}$$

Dacă se consideră gramatica ce generează expresii aritmetice 3.1.3, atunci rezultatul aplicării algoritmilor este prezentat în tabelul următor:

	NULL	Prim	Urm
E		$(, i$	$)$
E'	DA	$+ , -$	$)$
T		$(, i$	$+ , - ,)$
T'	DA	$* , /$	$+ , - ,)$
F		$(, i$	$* , / ,)$

iar mulțimile simbolurilor directoare sunt

$$SD(E, TE') = \{(, i\}$$

$$SD(E', +TE') = \{+\}$$

$$SD(E', \lambda) = \{\}\}$$

$$SD(T, FT') = \{(, i\}$$

$$\begin{aligned}
 SD(T', *FT') &= \{*\} \\
 SD(T', \lambda) &= \{+)\} \\
 SD(F, (E)) &= \{()\} \\
 SD(F, id) &= \{id\} \\
 SD(F, num) &= \{num\}
 \end{aligned}$$

Funcțiile C sau *Java* corespunzătoare neterminalelor se modifică puțin, în sensul că vom avea clauze corespunzătoare terminalelor ce se regăsesc în multimile simbolurilor directoare asociate regulilor. Dacă la reconstrucția derivării trebuie aplicată o regulă pentru X și următorul terminal din textul de analizat nu face parte din multimea simbolurilor directoare asociate vreunei reguli cu X în stânga, atunci se apelează modulul de tratare a erorilor sintactice.

Pentru cazul gramaticii anterioare, câteva din secvențele de cod asociate neterminalelor sunt prezentate în continuare.

```

...
void E(){
switch(token) {
    case LPAR:
    case id:
    case num: T(); Eprime(); break;
    default: printf("syntax error: (, identifier or number expected");
              err()
}
}

void Tprime(){
switch(token) {
    case ORI: eat(ORI); F(); Tprime(); break;
    case DIV: eat(DIV); F(); Tprime(); break;
    case PLUS:
    case RPAR: break;
    default: printf("syntax error: *, /, +, ) expected");
              err()
}
}

```

Tratarea erorilor este dependentă de proiectantul compilatorului

și limbaj. Principal, există trei modalități de tratare.

- Se oprește analiza sintactică în punctul unde s-a depistat o eroare, cu transmiterea unui mesaj prietenos către utilizator, de exemplu: *Eroare sintactică: așteptam să urmeze delimitator de instrucțiune. Mai învață sintaxa limbajului! BYE!*
- Se încearcă repararea greșelii de sintaxă inserând în textul sursă un simbol din mulțimea de simboluri directoare asociate regulii ce s-a aplicat. Desigur că este suficient să presupunem că am inserat în text, astfel încât analiza poate continua (desigur nu vom uita să anunțăm utilizatorul că nu cunoaște regulile de sintaxă și l-am corectat noi!). Inserarea poate conduce la cicluri infinite, astfel că nu este recomandabilă totdeauna.
- Se încearcă găsirea unui simbol ce se potrivește ignorând toate termenalele textului sursă până se întâlnește un simbol din mulțimea simbolurilor directoare. Se transmite același mesaj prietenos către autorul textului sursă, apoi se continuă analiza.

3.2 Algoritmi BOTTOM-UP

3.2.1 Gramatici cu precedență simplă

Fie $G = (V_N, V_T, x_0, \mathcal{P})$ o gramatică de tipul doi și $p \in L(G)$ o formă propozițională, adică un cuvânt peste alfabetul general V_G astfel încât $x_0 \xrightarrow{*} p$. Vom nota cu $\mathcal{A}_{x_0,p}$ arborele de derivare care are rădăcina x_0 și frontiera p .

Definiția 1. Vom spune că f este o frază simplă a lui p dacă f este un subcuvânt al lui p și în plus:

- $\exists x \in V_N$, cu $x \rightarrow f \in \mathcal{P}$;
- $\mathcal{A}_{x,f} \subset \mathcal{A}_{x_0,p}$.

Exemplu. Considerăm gramatica G_E 3.2.9 și forma propozițională $p = T * F + a * (E + T)$. Arborele de derivare este prezentat în figura 3.6. Se poate observa că frazele simple sunt $T * F, a, E + T$.

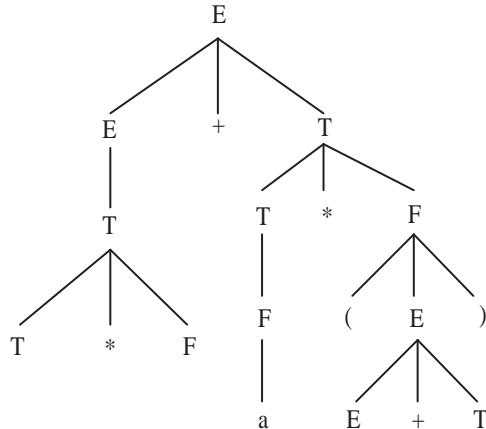


Figura 3.6: Arborele de derivare corespunzător derivării lui $p = T * F + a * (E + T)$

- (1) Initializare p ;
- (2) Se determină f_p , fraza simplă stângă a lui p ;
- (3) Se determină regula $x \rightarrow f_p$;
- (4) Se **reduce** fraza simplă stângă, adică dacă $p = r f_p s$, se pune $p = r x s$;
GOTO(2).

Figura 3.7: Algoritmul bottom-up

Observație. Subcuvântul F respectă prima condiție dar nu este frază simplă, întrucât nu este satisfăcută condiția a doua din definiție.

Definiția 2. Fraza simplă cea mai din stânga poartă denumirea de *frază simplă stângă*.

Vom nota fraza simplă stângă corespunzătoare lui p cu f_p în exemplul nostru, $f_p = T * F$. După cum vom vedea în continuare, fraza simplă stângă are un rol important în algoritmii de analiză sintactică bottom-up. În principiu, acești algoritmi prevăd următorii pași (descriși de figura 3.7):

Regulile determinate la pasul (3), aplicate în ordine inversă, ne vor

da derivarea corespunzătoare lui p . Să mai observăm că în acest mod arborele de derivare se construiește de jos în sus (bottom-up). Problema cea mai dificilă este desigur determinarea frazei simple stângi de la pasul (2). În principiu ar trebui cunoscut anticipat arborele de derivare corespunzător lui p , dar tocmai acesta este scopul analizei sintactice. După cum vom vedea, această dificultate se poate elimina, deci putem determina fraza simplă stângă fără a cunoaște acest arbore. De fapt, fraza simplă stângă este legată de anumite proprietăți ale gramaticii pe care le tratăm în continuare.

3.2.2 Relații de precedență

Definiția 3. Fie $x, y \in V_G$ două simboluri oarecare ale gramaticii.

Vom spune că:

- (1) $x \prec y$ (x precede pe y) dacă $\exists p$ astfel încât $p = rxys$, $x \notin f_p$, $y \in f_p$;
- (2) $x \pm y$ (x este egal cu y) dacă $\exists p$ astfel încât $p = rxys$, $x \in f_p$, $y \in f_p$;
- (3) $x \succ y$ (y succede lui x) dacă $\exists p$ astfel încât $p = rxys$, $x \in f_p$, $y \notin f_p$;

Relațiile $\prec \pm \succ$ se numesc *relații de precedență simple* (atunci când nu există posibilitatea de confuzie se folosesc denumirile *mai mic*, *egal*, *mai mare* pentru relațiile de precedență). Să observăm că aceste relații nu se exclud reciproc, deci că putem avea, de exemplu, $x \prec y$ și în același timp $x \succ y$, întrucât pot exista cazuri de gramatici în care putem găsi două forme propoziționale p_1, p_2 astfel încât să avem simultan cele două relații. Să mai facem de asemenea observația că există mai multe tipuri de astfel de relații, cu definiții asemănătoare, întrucât relațiile trebuie să satisfacă condiții suplimentare care depind și de gramatică; pentru toate aceste tipuri de relații vom utiliza aceleași notății.

Definiția 4. O gramatică de tipul 2 spunem că este cu *precedență simplă* dacă între oricare două simboluri ale gramaticii **există cel mult o relație** de precedență.

Pentru algoritmul de analiză sintactică ascendentă este important ca să nu existe reguli cu părțile drepte identice; altminteri, în cadrul pasului (4) din algoritmul de analiză sintactică nu se poate decide la cine să se facă reducerea frazei simple stângi. De asemenea regulile de ștergere nu sunt admise, pentru acest caz noțiunea de frază simplă neavând sens. Aceste condiții presupunem în mod sistematic că sunt satisfăcute.

Fie acum $V_G = \{x_1, x_2, \dots, x_n\}$ alfabetul general, unde am adoptat o anumită ordine a simbolurilor. Definim matricea de precedență a gramaticii,

$$M = (a_{ij}), \text{ unde } a_{ij} = \begin{cases} \prec & \text{daca } x_i \prec x_j \\ \pm & \text{daca } x_i \pm x_j \\ \succ & \text{daca } x_i \succ x_j \end{cases}$$

Exemplu. Considerăm următoarea gramatică

$$G = (\{A, B\}, \{0, 1\}, A, \{A \rightarrow A0|1B, B \rightarrow 1\}).$$

Matricea de precedență va fi

	A	B	0	1
A			⊲	
B			⊷	
0			⊷	
1		±	⊷	⊲

3.2.3 Proprietăți ale gramaticilor cu precedență simplă

Proprietatea 1. Fie $p = a_1 \dots a_n$ o formă propozițională într-o gramatică cu precedență simplă (a_i sunt simboluri neterminale sau terminale ale gramaticii). Atunci $f_p = a_i \dots a_j$ dacă și numai dacă

$$a_1 \stackrel{\pm}{\prec} a_2 \dots \stackrel{\pm}{\prec} a_{i-1} \prec a_i \pm a_{i+1} \dots \pm a_j \succ a_{j+1} \quad (*)$$

Demonstrația acestei proprietăți se efectuează considerând toate structurile de arbori de derivare posibile într-o situație dată. Vom exemplifica această idee pentru cuvântul $p = a_1 \dots a_7$. Dacă presupunem că $f_p = a_3a_4a_5$, atunci, conform definiției frazei simple stângi, avem

$$a_2 \prec a_4, \quad a_4 \pm a_5, \quad a_5 \succ a_6$$

Mai trebuie arătat că $a_1 \stackrel{\pm}{\prec} a_2$. În acest scop vom considera toate posibilitățile de structuri arborescente; o parte din ele sunt prezentate în figura 3.8.

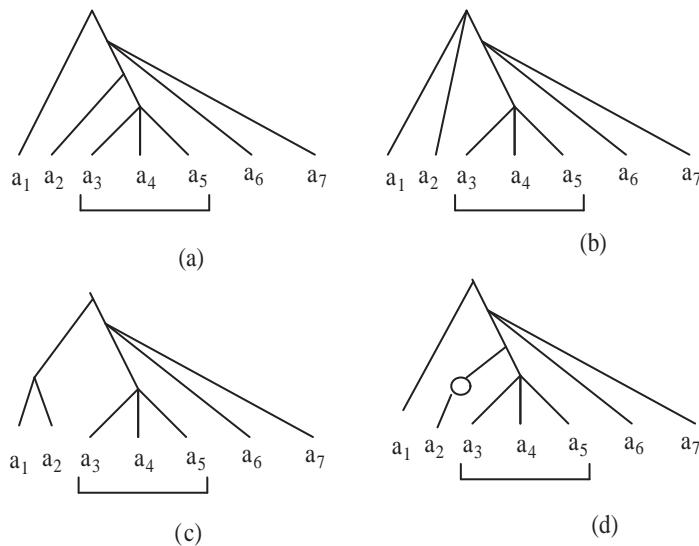


Figura 3.8: Arbori de derivare posibili

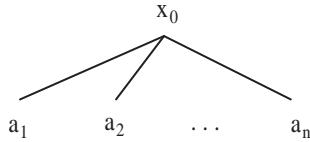
Se poate observa că în ipoteza că $f_p = a_3a_4a_5$, singurele situații posibile sunt (a) și (b). Celelalte situații contrazic această ipoteză, de exemplu, în cazurile (c) și (d) fraza simplă stângă ar fi , respectiv, a_1a_2 sau a_2 . Acum, în cazul (a), prin reducerea succesivă a frazei simple stângi, vom obține un arbore în care fraza simplă stângă va fi a_1a_2 , adică $a_1 \pm a_2$ iar în cazul (b), fraza simplă stângă va fi a_2 și atunci $a_1 \prec a_2$.

Implicația inversă se bazează pe cea directă; se presupune că $f_p = a_l \dots a_k$ cu $l \neq i$ și $k \neq j$ și se consideră în continuare toate posibilitățile de poziționare a indicilor l, k față de i, j . De exemplu, dacă $f_p = a_2a_3$ atunci, conform implicației directe, avem

$$a_1 \prec a_2 \pm a_3 \succ a_4$$

ceea ce contrazice relația (*), conform căreia $a_3 \pm a_4$ iar gramatica are precedență simplă.

Observație La delimitarea frazei simple stângi se procedează astfel: parcurgem sirul de simboluri de la stânga la dreapta până se găsește

Figura 3.9: Arborele de derivare pentru $ni = 1$

relația \succ (sau marginea dreaptă a formei propoziționale) pentru determinarea ultimului simbol din f_p , apoi se parcurge sirul spre stânga, trecând peste relații \pm , până la găsirea unei relații \prec (sau marginea dreaptă a formei propoziționale) pentru determinarea primului simbol din f_p .

Proprietatea 2. O gramatică cu precedență simplă este neambiguă.

Schită de demonstrație. Fie p o formă propozițională și $\mathcal{A}_{x_0,p}$ arborele de derivare corespunzător. Vom arăta că acest arbore este unicul arbore de derivare care are rădăcina x_0 și frontieră p . Procedăm prin inducție asupra numărului de noduri interne ni . Dacă $ni = 1$ atunci arborele de derivare va avea forma din figura 3.9 și unicitatea acestuia este evidentă. Presupunem acum că proprietatea este adevărată pentru un ni oarecare și considerăm cazul $ni + 1$. Să presupunem că ar exista doi arbori diferenți cu rădăcina x_0 și frontieră p ; fie aceștia $\mathcal{A}'_{x_0,p}$ și $\mathcal{A}_{x_0,p}$.

Deoarece gramatica este cu precedență simplă, rezultă că fraza simplă stângă, care este aceeași în cei doi arbori, este situată în aceeași poziție, $p = rf_ps$. Efectuăm reducerea frazei simple f_p (conform regulii aplicate $x \rightarrow f_p$) și vom obține arborii $\mathcal{A}'_{x_0,q}$ și $\mathcal{A}_{x_0,q}$ unde $q = rXs$ iar cei doi arbori trebuie să fie diferenți. Dar numărul de noduri interne este ni , în contrazicere cu ipoteza.

3.2.4 Determinarea relațiilor de precedență pentru gramatici cu precedență simplă

Așa cum s-a precizat deja, relațiile de precedență sunt proprietăți intrinseci ale gramaticii, nu depind de contextul în care se află cele

două simboluri. Prin urmare, cunoașterea anticipată a acestor relații împreună cu proprietatea 1, rezolvă complet problema pasului 2 de la algoritmul de analiză bottom-up, adică determinarea frazei simple stângi. În acest subparagraf vom prezenta o teoremă de caracterizare pe baza căreia se pot determina relațiile de precedență simple și implicit faptul dacă gramatica este sau nu cu precedență simplă.

Vom defini două relații specifice gramaticilor și care vor fi utilizate în continuare, numite, respectiv, **F** (First) și **L** (Last). Fie $x \in V_N$ și $y \in V_G$ două simboluri ale gramaticii. Vom spune că

$$\begin{aligned} xFy &\text{ daca } \exists x \rightarrow yu \in \mathcal{P}, u \in V_G^* \\ xLy &\text{ daca } \exists x \rightarrow uy \in \mathcal{P}, u \in V_G^* \end{aligned}$$

Închiderile tranzitive, respectiv, tranzitive și reflexive ale acestor relații le vom nota cu F^+ , L^+ și respectiv, F^* , L^* . Există numeroși algoritmi direcți care permit calcularea acestei închideri, cu complexitate redusă.

Teorema 2. Avem

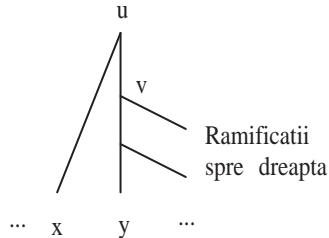
- (1) $x \prec y \Leftrightarrow \exists u \rightarrow \dots xv \dots \in \mathcal{P}, vF^+y;$
- (2) $x \pm y \Leftrightarrow \exists u \rightarrow \dots xy \dots \in \mathcal{P};$
- (3) $x \succ y \Leftrightarrow \exists u \rightarrow \dots vw \dots \in \mathcal{P}, vL^+x, wF^*y;$

Schită de demonstrație. Se analizează posibilitățile de arbori de derivare care răspund cerințelor teoremei. De exemplu, pentru cazul (1), trebuie să existe structura de arbore ca în figura 3.10.

Este clar că dacă $x \prec y$, atunci conform definiției frazei simple stângi, trebuie să existe p astfel încât $p = rxys$, $x \notin f_p$, $y \in f_p$, adică să existe structura din 3.10. Invers, dacă există o astfel de structură, atunci fie o formă propozițională care conține u . Efectuăm reduceri succesive până când vom obține o formă propozițională pentru care u aparține frazei simple stângi; la arborele astfel obținut, adăugăm subarborele de mai sus. Vom avea în mod evident $x \notin f_p$, $y \in f_p$.

3.2.5 Studiu de caz

Să considerăm cazul gramaticii pentru generarea expresiilor aritmetice simple G_E , cu regulile obișnuite

Figura 3.10: Arborele de derivare corespunzator cazului $x \prec y$

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

Relațiile First, Last și închiderile acestora sunt:

$$\begin{array}{lll}
 E F \{E, T\} & E F^+ \{E, T, F, (, a\} & E F^* \{E, T, F, (, a\} \\
 T F \{T, F\} & T F^+ \{T, F, (, a\} & T F^* \{T, F, (, a\} \\
 F F \{(, a\} & F F^+ \{(, a\} & F F^* \{F, (, a\} \\
 \\
 E L \{T\} & E L^+ \{T, F, a,)\} & E F^* \{E, T, F, a,)\} \\
 T L \{F\} & T L^+ \{F, a,)\} & T F^* \{T, F, a,)\} \\
 F L \{a,)\} & F L^+ \{a,)\} & F F^* \{F, a,)\}
 \end{array}$$

Conform teoremei de calcul a relațiilor de precedență vom obține matricea de precedență (vezi figura 3.11) asociată simbolurilor gramaticii G_E .

Deci gramatica G_E nu are proprietatea de precedență simplă, astfel că nu poate fi aplicat direct algoritmul de analiză sintactică bottom up.

Vom exemplifica algoritmul de analiză pentru gramatica din figura 3.12, matricea de precedență asociată este cea alăturată.

Reconstituirea derivării cuvântului $p = aaabaababbabb$ este prezentată în tabelul 3.13, unde fiecare linie corespunde formei propoziționale, în care fraza simplă stângă este subliniată, precedată de regula identificată.

	E	T	F	$+$	$*$	$($	$)$	a
E				\pm			\pm	
T				γ	γ	\pm		
F				γ	γ			
$+$		$\prec \pm$	\prec	\prec			\prec	\prec
$*$			\pm			\prec	\prec	\prec
$($	$\prec \pm$	\prec	\prec			\prec		\prec
$)$			γ	γ	γ		γ	γ
a				γ	γ		γ	

Figura 3.11: Matricea de precedență pentru G_E

$G : S \rightarrow aSSb \mid ab,$	S	S	a	b
	\pm	\prec	\pm	
S	\pm	\prec	\pm	
a	\pm	\prec	\pm	
b	γ	γ	γ	γ

Figura 3.12: Exemplu de gramatică cu precedență simplă

<i>Regula</i>	<i>Forma propositională</i>
$S \rightarrow ab$	$a \prec a \prec \underline{a \pm b} \succ aababbabb$
$S \rightarrow ab$	$a \prec a \pm S \prec a \prec \underline{a \pm b} \succ abbbabb$
$S \rightarrow ab$	$a \prec a \pm S \prec a \pm S \prec \underline{a \pm b} \succ bbabb$
$S \rightarrow ab$	$a \prec a \pm S \prec \underline{a \pm S \pm S \pm b} \succ babb$
$S \rightarrow aSSb$	$a \prec \underline{a \pm S \pm S \pm b} \succ abb$
$S \rightarrow aSSb$	$a \pm \underline{S \prec a \pm b} \succ b$
$S \rightarrow ab$	$a \pm S \pm S \pm b$
$S \rightarrow aSSb$	S

Figura 3.13: Reconstrucția derivării cuvântului $p = aaabaababbabb$

3.2.6 Gramatici operatoriale

Gramaticile operatoriale sunt gramatici cu o structură specială a regușilor de rescriere, în care anumite simboluri sunt considerate operatori iar celealte operanzi, operatorii trebuind să separe operanzii. Această structură este preluată de la gramaticile care generează expresiile aritmetice iar principiul de analiză este de asemenea preluat de la algoritmul de evaluare a expresiilor aritmetice. În continuare vom prezenta acest principiu pentru evaluarea expresiilor aritmetice fără paranteze, cazul expresiilor cu paranteze putându-se reduce la cel fără paranteze, de exemplu utilizând tehnica de modificare a ponderilor operatorilor la întâlnirea parantezelor.

În primul rând se atribuie operatorilor anumite ponderi care vor defini ordinea de efectuare a operațiilor. De obicei, se consideră $p(+)=p(-)=1 < 2 = p(*) = p(/)$, ceea ce înseamnă că mai întâi se fac adunările și scăderile, apoi înmulțirile și împărțirile, etc. Efectuarea unei operații depinde de contextul în care se află operatorul respectiv, și anume, dacă ponderea operatorului precedent este mai mare sau egală cu ponderea operatorului curent, atunci se poate efectua operația definită de operatorul precedent. Din acest motiv, acest tip de gramatici se numesc *cu operator de precedență* (operator precedence grammars).

Putem realiza un algoritm simplu de evaluare utilizând două stive

- **P-** stiva operatorilor;
- **S-** stiva operanzilor.

Vom utiliza indicii i și k pentru a indica nivelele celor două stive. Prin urmare, notația p_i , respectiv s_k au sensul de operatorul, respectiv operandul din stivă aflat pe nivelul i , respectiv k . Pentru simplificare, vom utiliza aceeași notație p_i atât pentru operator cât și pentru ponderea operatorului respectiv.

Algoritmul de evaluare:

PAS 1: *Inițializări:* Se introduce primul operator și primul operand în stivele P și S și se inițializează nivelele celor două stive, $i = k = 2$;

PAS 2: Se introduce următorul operator în stiva P și se actualizează nivelul, $i = i + 1$;

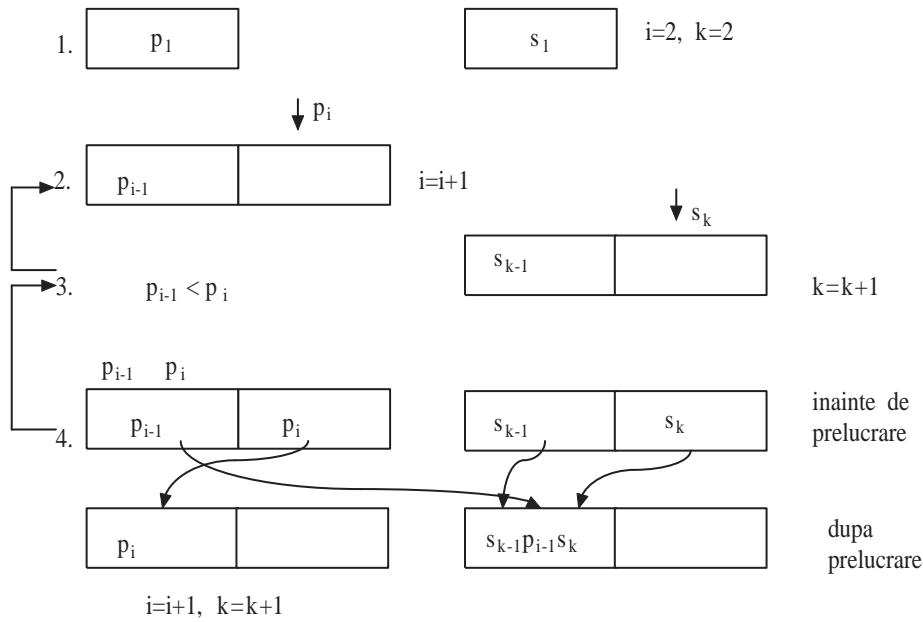


Figura 3.14: Prelucrările efectuate asupra celor două stive

PAS 3: Dacă $p_{i-1} < p_i$ atunci se introduce următorul operand în lista S , se actualizează nivelul stivei, $k = k + 1$, și se revine la pasul 2;

PAS 4: Dacă $p_{i-1} \geq p_i$ atunci:

în stiva P : p_{i-1} se înlocuiește cu p_i ;

în stiva S : s_{k-1} se înlocuiește cu $s_{k-1}p_{i-1}s_k$;

se actualizează nivelele celor două stive, $i = i - 1; k = k - 1$ și se revine la pasul 3.

Observație. Pentru uniformitatea algoritmului se introduce operatorul marcaj, notat cu $\#$, cu ponderea cea mai mică; orice expresie se încadrează între două marcaje iar oprirea algoritmului are loc la pasul 4 în cazul în care $p_1 = p_2 = \#$. Schematic, acest algoritm este prezentat în figura 3.14.

Exemplu. Tabelul 3.15 conține starea dinamică a celor două stive pentru expresia $\#a + b * c/d - f\#$. Menționăm că în cadrul pașilor 2 și 4 s-a completat câte un rând nou în tabel.

Observație. Ponderile operatorilor sunt legate de ierarhia lor în gramatică, cu cât un operator este mai jos cu atât ponderea lui este

	1	2	3	4	1	2	3	4
1	# ₀				<i>a</i>			
2	# ₀	+ ₁			<i>a</i>	<i>b</i>		
3	# ₀	+ ₁	- ₁		<i>a</i>	<i>b</i>		
4	# ₀	- ₁			<i>a</i> + <i>b</i>	<i>c</i>		
5	# ₀	- ₁	* ₂		<i>a</i> + <i>b</i>	<i>c</i>	<i>d</i>	
6	# ₀	- ₁	* ₂	/ ₂	<i>a</i> + <i>b</i>	<i>c</i>	<i>d</i>	
7	# ₀	- ₁	/ ₂		<i>a</i> + <i>b</i>	<i>c</i> * <i>d</i>	<i>f</i>	
8	# ₀	- ₁	/ ₂	# ₀	<i>a</i> + <i>b</i>	<i>c</i> * <i>d</i>	<i>f</i>	
9	# ₀	- ₁	# ₀		<i>a</i> + <i>b</i>	<i>c</i> * <i>d</i> / <i>f</i>		
10	# ₀	# ₀			<i>a</i> + <i>b</i> - <i>c</i> * <i>d</i> / <i>f</i>			

Figura 3.15: Evaluarea expresiei $\#a + b - c * d/f\#$

mai mare. Acest fapt este ilustrat în figura 3.16.

Pentru cazul expresiilor aritmetice cu paranteze algoritmul se păstrează, dar este aplicat după preprocesarea expresiei prin modificarea dinamică a ponderilor operatorilor și renunțarea la paranteze. Se parcurge expresia de la stânga la dreapta alocând ponderi operatorilor; la întâlnirea unei paranteze deschise ponderile se incrementează cu 10, iar la întâlnirea unei paranteze închise se decrementează ponderile operatorilor cu 10. În continuare parantezele sunt ignorate. De exemplu, expresia

$$a * (b - c * d) - (x + y/(z - t)) * h$$

devine prin preprocesare

$$a *_{2} b -_{11} c *_{12} d -_{1} x +_{11} y /_{12} z -_{21} t *_{2}$$

ceea ce va asigura ordinea corectă de evaluare.

3.2.7 Gramatici operatoriale

Fie $G = (V_N, V_T, x_0, \mathcal{P})$ o gramatică de tipul 2. Vom considera că simbolurile neterminale V_N sunt operanții limbajului iar simbolurile terminale V_T sunt operatorii. Să observăm că în cazul gramaticii G_E

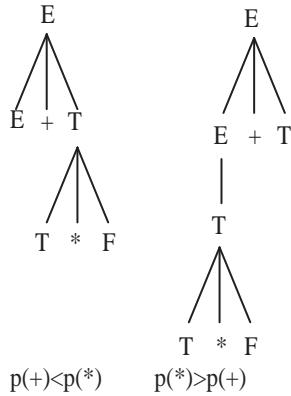


Figura 3.16: Ierarhia operatorilor

care generează expresii aritmetice simple o asemenea convenție este justificată, deoarece oricare din neterminalele E, T, F derivează în terminalul a (generic, acest terminal reprezintă operanzi), sau într-o expresie (eventual între paranteze) care are rolul de operand.

Definiția 1. Vom spune că o gramatică este operatorială dacă regulile de rescriere au forma

$$X \rightarrow N_i T_i N_{i+1} T_{i+1} \dots T_j N_{j+1},$$

unde N_k sunt neterminale (*operanzi*), inclusiv cuvântul vid, iar T_k sunt terminale (*operatori*).

Observație. Într-o gramatică operatorială orice formă propozițională are structura:

$$p = N_1 T_1 N_2 T_2 \dots T_n N_{n+1}$$

Aceasta înseamnă că într-o formă propozițională a unei gramici operatoriale între oricare doi operanzi există cel puțin un operator; pot însă exista mai mulți operatori alăturați, de exemplu, $\dots ((a \dots$

Definiția 2. Fie $p = N_1 T_1 N_2 T_2 \dots T_n N_{n+1}$ o formă propozițională într-o gramatică operatorială. Vom spune că f este o *frază simplă* a lui p dacă $f = N_i T_i N_{i+1} T_{i+1} \dots T_j N_{j+1}$ este un subcuvânt al lui p care conține cel puțin un simbol terminal și în plus:

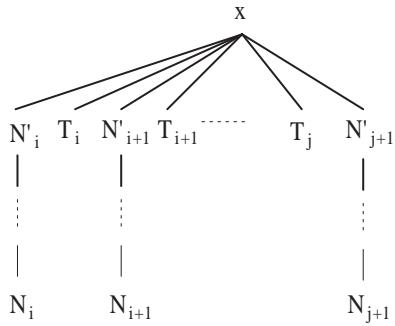


Figura 3.17: Structura subarborelui

- (1) $\exists x \rightarrow N'_i T_i N'_{i+1} T_{i+1} \dots T_j N'_{j+1} \in \mathcal{P}$, și $N'_k \xrightarrow{*} N_k \forall k = i, \dots, j$.
- (2) $\mathcal{A}_{x,f} \subset \mathcal{A}_{x_0,p}$.

Structura de arbore corespunzătoare unei fraze simple este prezentată în figura 3.17.

Ca și în cazul gramaticilor cu precedență simplă, fraza simplă cea mai din stânga poartă denumirea de *frază simplă stângă*. Principiul de analiză sintactică este identic cu cel de la gramatici cu precedență simplă. Problema principală este și aici determinarea frazei simple stângi. Aceasta operație poate fi făcută utilizând relațiile de precedență dintre simbolurile gramaticii care pentru gramatici operatoriale au o definiție puțin modificată.

Definiția 3. Fie $x, y \in V_T$. Vom spune că:

- (1) $x \overset{\circ}{<} y$ (x precede pe y) dacă $\exists p$ astfel încât $p = rxNys$, $N \in V_N$, $x \notin f_p$, $y \in f_p$;
- (2) $x \overset{\circ}{=} y$ (x este egal cu y) dacă $\exists p$ astfel încât $p = rxNys$, $N \in V_N$, $x \notin f_p$, $y \in f_p$;
- (2) $x \overset{\circ}{>} y$ (y succede lui x) dacă $\exists p$ astfel încât $p = rxy$, $N \in V_N$, $x \in f_p$, $y \notin f_p$;

Matricea de precedență se definește numai pe mulțimea simbolurilor terminale, ceea ce conduce la o simplificare a algoritmului de analiză.

Definiția 4. O gramatică operatorială se spune că are *precedență simplă* dacă între două simboluri ale gramaticii există cel mult o relație

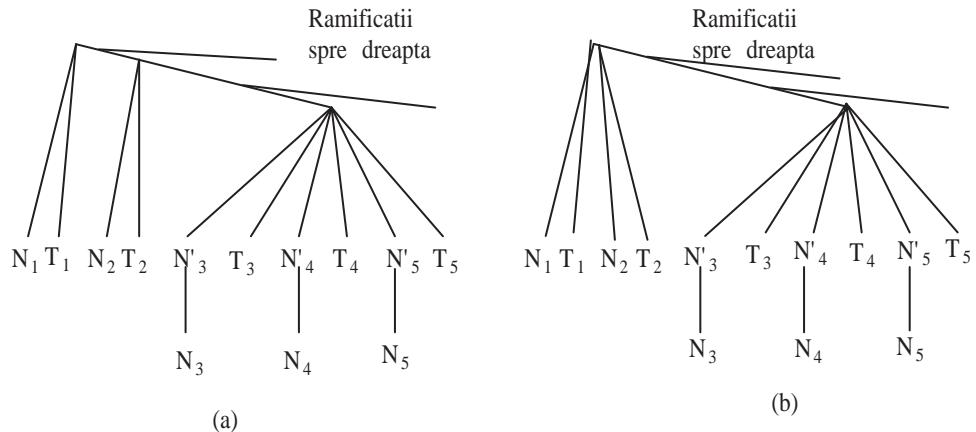


Figura 3.18: Structurile posibile de arbori

de precedentă.

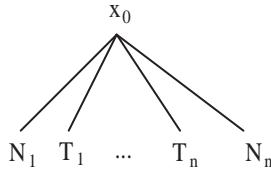
Ca și la celealte tipuri de gramatici, în vederea efectuării pasului 3 din algoritmul de analiză, vom presupune că nu există reguli cu părțile drepte identice. De asemenea, sunt valabile proprietățile de caracterizare și de neambiguitate.

Proprietatea 1. Fie $p = N_1 T_1 N_2 T_2 \dots T_n N_{n+1}$ o formă propozițională într-o gramatică cu precedență simplă. Atunci $f_p = N_i T_i N_{i+1} T_{i+1} \dots T_j N_{j+1}$ dacă și numai dacă

$$T_1 \stackrel{\circ}{\leq} T_2 \dots \stackrel{\circ}{\leq} T_{i-1} \stackrel{\circ}{<} T_i \stackrel{\circ}{=} T_{i+1} \dots \stackrel{\circ}{=} T_j \stackrel{\circ}{>} T_{j+1} \quad (*)$$

Schiță de demonstrație. Presupunem că $p = N_1 T_1 N_2 T_2 \dots T_n N_{n+1}$ și că fraza simplă stângă este $f_p = N_3 T_3 N_4 T_4 N_5 T_5$. Atunci din definiția frazei simple stângi, rezultă $T_2 < \overset{\circ}{T_3} = \overset{\circ}{T_4} = \overset{\circ}{T_5} > T_6$. Mai trebuie să arătăm că $T_1 < \overset{\circ}{T_2}$ sau $T_1 = \overset{\circ}{T_2}$. În acest scop analizăm structurile posibile de arbori (figura 3.18).

Efectuăm reduceri succesive până când $T_2 \in f_p$. În acest moment avem $T_1 \overset{\circ}{<} T_2$ pentru cazul (a) și $T_1 \overset{\circ}{=} T_2$ pentru cazul (b). Pentru implicația inversă presupunem că $f_p = N_k T_k N_{k+1} T_{k+1} \dots T_l N_{l+1}$, $l \neq$

Figura 3.19: Arborele de derivare pentru $ni = 1$

$j, k \neq i$ și considerăm toate cazurile de poziționare a indicilor k, l în raport cu i, j . De exemplu, dacă $l < j$ atunci, în conformitate cu implicația directă, rezultă $T_l \overset{\circ}{>} T_{l+1}$ iar conform relației (*) din ipoteză, avem $T_l \overset{\circ}{<} T_{l+1}$ sau $T_l \overset{\circ}{=} T_{l+1}$ ceea ce contrazice ipoteza inițială conform căreia gramatica este cu precedență simplă.

Proprietatea 2. O gramatică operatorială cu precedență simplă este neambiguă.

Schiță de demonstrație. Fie p o formă propozițională și $\mathcal{A}_{x_0,p}$ arborele de derivare corespunzător. Vom arăta că acest arbore este unicul arbore de derivare care are rădăcina x_0 și frontiera p . Procedăm prin inducție asupra numărului de noduri interne ni . Dacă $ni = 1$ atunci arborele de derivare va avea forma din figura 3.19 și unicitatea acestuia este evidentă. Presupunem acum că proprietatea este adevărată pentru un ni oarecare și considerăm cazul $ni + 1$. Să presupunem că ar exista doi arbori diferenți cu rădăcina x_0 și frontiera p ; fie aceștia $\mathcal{A}'_{x_0,p}$ și $\mathcal{A}_{x_0,p}$.

Deoarece gramatica este cu precedență simplă, rezultă că fraza simplă stângă, care este aceeași în cele două arbori, este situată în aceeași poziție, $p = rf_ps$. Efectuăm reducerea frazei simple f_p (conform regulii aplicate $x \rightarrow f_p$) și vom obține arborii $\mathcal{A}'_{x_0,q}$ și $\mathcal{A}_{x_0,q}$ unde $q = rXs$ iar cele două arbori trebuie să fie diferenți. Dar numărul de noduri interne este cel mult ni , ceea ce contrazice ipoteza inductivă.

3.2.8 Determinarea relațiilor de precedență pentru gramatici operatoriale

Vom defini mai întâi relațiile $F_{1,2}$ și $L_{1,2}$, analoage cu F și L de la gramaticile cu precedență simplă.

Definiție Fie $x \in V_N$ și $y \in V_G$. Vom spune că

$$\begin{aligned} xF_{1,2}y &\text{ daca } \exists x \rightarrow yu \in \mathcal{P}, \text{ sau } x \rightarrow Nyu \in \mathcal{P}, u \in V_G^*, N \in V_N \\ xL_{1,2}y &\text{ daca } \exists x \rightarrow uy \in \mathcal{P}, \text{ sau } x \rightarrow uyN \in \mathcal{P}, u \in V_G^*, N \in V_N \end{aligned}$$

Închiderile tranzitive, respectiv, tranzitive și reflexive ale acestor relații le vom nota cu $F_{1,2}^+$, $L_{1,2}^+$ și respectiv, $F_{1,2}^*$, $L_{1,2}^*$. Acum, determinarea relațiilor de precedență operatoriale se poate face cu ajutorul următoarei teoreme.

Teorema (*determinarea relațiilor de precedență*). Avem

- (1) $x \overset{\circ}{<} y \Leftrightarrow \exists u \rightarrow \dots xv \dots \in \mathcal{P}, vF_{1,2}^+y;$
- (2) $x \overset{\circ}{=} y \Leftrightarrow \exists u \rightarrow \dots xNy \dots \in \mathcal{P};$
- (3) $x \overset{\circ}{>} y \Leftrightarrow \exists u \rightarrow \dots vy \dots \in \mathcal{P}, vL_{1,2}^+x;$

Schiță de demonstrație. Demonstrația se poate face prin analiza structurilor posibile ale arborilor de derivare. Pentru cazul (1), trebuie să existe structura de arbore ca în figura 3.20.

Această teoremă ne permite să determinăm matricea de precedență a gramaticii pe baza căreia putem apoi aplica algoritmul de analiză bottom-up. Problema neterminalelor de la capetele frazei simple stângi se poate rezolva analizând părțile drepte ale regulilor de scriere. Într-adevăr, teorema ne dă numai terminalele care intră în fraza simplă stângă și neterminalele interioare; cele două eventuale neterminale de la capete vor apartine sau nu frazei simple stângi depinzând de forma părții drepte ale regulii respective. Evident, se pot face și ipoteze de neambiguitate suplimentare pentru gramatică, similare cu condiția de a nu exista reguli cu părțile drepte identice.

3.2.9 Studiu de caz

Să considerăm cazul gramaticii pentru generarea expresiilor aritmetice simple G_E , cu regulile obișnuite

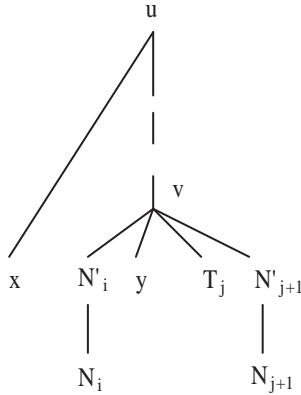


Figura 3.20: Structura posibilă de arbore

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

Relațiile $F_{1,2}, L_{1,2}$ și închiderile acestora sunt:

$$\begin{array}{lll} E F_{1,2} \{E, T, +\} & E F_{1,2}^+ \{E, T, F, +, *, (, a)\} & E F_{1,2}^* \{E, T, F, +, *, (, a)\} \\ T F_{1,2} \{T, F, *\} & T F_{1,2}^+ \{T, F, *, (, a)\} & T F_{1,2}^* \{T, F, *, (, a)\} \\ F F_{1,2} \{(, a)\} & F F_{1,2}^+ \{(, a)\} & F F_{1,2}^* \{F, (, a)\} \end{array}$$

$$\begin{array}{lll} E L_{1,2} \{T, +\} & E L_{1,2}^+ \{T, +, F, *, a,)\} & E F_{1,2}^* \{E, T, +, F, *, a,)\} \\ T L_{1,2} \{F, *\} & T L_{1,2}^+ \{F, *, a,)\} & T F_{1,2}^* \{T, F, *, a,)\} \\ F L_{1,2} \{a,)\} & F L_{1,2}^+ \{a,)\} & F F_{1,2}^* \{F, a,)\} \end{array}$$

Conform teoremei de calcul a relațiilor de precedență vom obține matricea de precedență (vezi figura 3.21) asociată simbolurilor gramiciei G_E .

Deci gramatica G_E are proprietatea de precedență operatorială simplă, astfel că se poate aplica direct algoritmul de analiză sintactică bottom up.

	+	*	()	<i>a</i>
+	o	o	o	o	o
*	>	>	<	>	<
(o	o	o	o	o
)	<	<	<	=	<
<i>a</i>	o	o	o	o	o
	>	>	>	>	>

Figura 3.21: Matricea de precedență operatorială pentru G_E

Reconstituirea derivării cuvântului $p = a + (a+a)*a*a$ este prezentată în tabelul 3.22, unde fiecare linie corespunde formei propoziționale, în care fraza simplă stângă este subliniată, precedată de regula identificată.

3.3 Analiza sintactică LR(k)

Denumirea LR(k) provine de la *Left to right parsing, Rightmost derivation, k-token lookahead* și reprezintă o categorie de algoritmi de analiză sintactică ascendentă în care se încearcă reconstituirea derivării extrem drepte pentru un sir de intrare prin identificarea la fiecare pas a frazei simple stângi prin examinarea formei propoziționale curente și a maxim k simboluri citite în avans din sirul de intrare rămas de examinat. Formalizarea riguroasă a acestor algoritmi a fost facută de către D. Knuth în 1965, iar o variantă pentru $k = 1$ a devenit populară odată cu sistemul operare Unix prin apariția utilitarului YACC (Yet Another Compiler Compiler) dezvoltat de AT&T Bell Laboratories în anii 70.

Prezentarea algoritmului urmăreste în mare măsură cea folosită de A. Appel în lucrarea *Modern compiler implementation in Java*.

Analizorul folosește o **stivă** și un **sir de intrare** care prin concatenare reprezintă forma propozițională curentă din derivarea extrem dreaptă corespunzătoare cuvântului de examinat. În funcție de primele k simboluri citite în avans din sirul de intrare (*the lookahead*) și conținutul

<i>Regula</i>	<i>Forma propositională</i>
$F \rightarrow a$	$\underline{a} > + \overset{\circ}{<}(\overset{\circ}{a} + a) * a * a$
$F \rightarrow a$	$F + \overset{\circ}{<}(\overset{\circ}{\underline{a}} + a) * a * a$
$F \rightarrow a$	$F + \overset{\circ}{<}(\overset{\circ}{F} + \overset{\circ}{\underline{a}}) * a * a$
$F \rightarrow a$	$F + \overset{\circ}{<}(\overset{\circ}{\underline{F}} + \overset{\circ}{F}) * a * a$
$T \rightarrow T * F, T \rightarrow F$	$F + \overset{\circ}{<}(\overset{\circ}{T} \overset{\circ}{=}) * a * a$
$F \rightarrow (E), E \rightarrow T$	$F + \overset{\circ}{<}F * \overset{\circ}{<} \overset{\circ}{a} * a$
$F \rightarrow a$	$F + \overset{\circ}{<} \overset{\circ}{F} * \overset{\circ}{F} * a$
$T \rightarrow T * F, T \rightarrow F$	$F + \overset{\circ}{<} T \overset{\circ}{<} * \overset{\circ}{<} \overset{\circ}{a}$
$F \rightarrow a$	$F + \overset{\circ}{<} \overset{\circ}{T} * \overset{\circ}{F}$
$T \rightarrow T * F$	$\frac{F + T}{E}$
$E \rightarrow E * T, E \rightarrow T, T \rightarrow F$	

Figura 3.22: Reconstrucția derivării cuvântului $p = a + (a + a) * a * a$

stivei analizorul execută una din următoarele acțiuni:

- **Shift:** Mută primul simbol din sirul de intrare în stivă;
- **Reduce:** Alege o regulă a gramaticii, de exemplu $X -> ABC$, extrage de pe stivă simbolurile C, B, A și pune în vîrful stivei neterminalul X. Acțiunea corespunde reducerii frazei simple stângi din algoritmul general boottom up.

Initial, stiva este vidă și cuvântul de examinat constituie sirul de intrare al analizorului. Decizia asupra acțiunii ce trebuie efectuate (shift/reduce) se ia cu ajutorul unui automat finit determinist care examinează conținutul stivei. Acest automat încearcă recunoașterea frazei simple stângi (adică partea dreaptă a ultimei reguli aplicate în derivarea extrem dreaptă) atunci când aceasta apare în vîrful stivei. Să considerăm gramatica ce generează expresii aritmetice simple ce conțin doar adunarea (3.3) la care s-a adăugat simbolul \$ ca sfârșit de sir.

$$\left\{ \begin{array}{ll} S \rightarrow E\$ & (0) \\ E \rightarrow T + E & (1) \\ E \rightarrow T & (2) \\ T \rightarrow F * T & (3) \\ T \rightarrow T & (4) \\ F \rightarrow (E) & (5) \\ F \rightarrow x & (6) \end{array} \right.$$

Analiza cuvântului $p = (x + x) * x + x\$$ decurge ca în figura 3.23 și folosește automatul finit descris în tabelul 3.24.

Semnificația elementelor tabloului este următoarea:

- **sn** Shift și analizorul intră în starea n.
- **gn** Analizorul trece în starea n. (Goto n).
- **rk** Efectuează reducerea folosind regula cu numărul k. Dacă X este neterminalul nou ce trebuie pus pe stivă atunci se urmează acțiunea indicată de ultima stare la care s-a ajuns înainte de primul simbol al părții drepte reduse și neterminalul X .
- **a** Acceptă cuvântul de analizat.

Elementele necompletate din tabel corespund situației de găsire a unei erori sintactice.

În practică se folosesc gramatici ce admit o formă echivalentă de tip LR(1).

3.3.1 Construcția analizorului pentru gramatici LR(0)

Cazul $k = 0$ corespunde situației când decizia se poate lua prin examinarea stivei, mai exact când în vârful stivei se regăsește fraza simplă stângă. Deși există puține gramatici cu această proprietate, modul de construcție a tablei de analiză este ușor de adaptat pentru cazurile uzuale.

Definim noțiunea de **element** sau *item* printr-o regulă a gramaticii ce conține în partea dreaptă un marcat suplimentar, de obicei notat

Stiva	Input	Actiune
1	$(x + x) * x + x\$$	shift
1(2	$x + x) * x + x\$$	shift
1(2x8	$+x) * x + x\$$	reduce $F \rightarrow x$
1(2F5	$+x) * x + x\$$	reduce $T \rightarrow F$
1(2T10	$+x) * x + x\$$	shift
1(2T10 + 17	$) * x + x\$$	shift
1(2T10 + 17x19	$) * x + x\$$	reduce $F \rightarrow x$
1(2T10 + 17F5	$) * x + x\$$	reduce $T \rightarrow F$
1(2T10 + 17T10	$) * x + x\$$	reduce $E \rightarrow T$
1(2T10 + 17E18	$) * x + x\$$	reduce $E \rightarrow E + T$
1(2E3	$) * x + x\$$	shift
1(2E3)6	$*x + x\$$	reduce $F \rightarrow (E)$
1F4	$*x + x\$$	shift
1F4 * 15	$x + x\$$	shift
1F4 * 15x8	$+x\$$	reduce $F \rightarrow x$
1F4 * 15F4	$+x\$$	reduce $T \rightarrow F$
1F4 * 15T16	$+x\$$	reduce $T \rightarrow F * T$
1T9	$+x\$$	shift
1T9 + 11	$x\$$	shift
1T9 + 11x8	$\$$	reduce $F \rightarrow x$
1T9 + 11F4	$\$$	reduce $T \rightarrow F$
1T9 + 11T9	$\$$	reduce $E \rightarrow T$
1T9 + 11E12	$\$$	reduce $E \rightarrow T + E$
1E7	$\$$	ACCEPT

Figura 3.23: Analiza sintactică *shift – reduce* pentru cuvântul $p = (x + x) * x + x\$$.

<i>Stare</i>	<i>x</i>	()	+	*	\$	<i>E</i>	<i>T</i>	<i>F</i>
1	<i>s8</i>	<i>s2</i>					<i>g7</i>	<i>g9</i>	<i>g4</i>
2	<i>s8</i>	<i>s2</i>					<i>g3</i>	<i>g10</i>	<i>g5</i>
3			<i>s6</i>						
4				<i>r4</i>	<i>s15</i>	<i>r4</i>			
5				<i>r4</i>	<i>r4</i>	<i>s13</i>			
6				<i>r6</i>	<i>r6</i>	<i>r6</i>			
7						<i>a</i>			
8				<i>r5</i>	<i>r5</i>	<i>r5</i>			
9					<i>s11</i>		<i>r2</i>		
10				<i>r2</i>	<i>s17</i>				
11	<i>s8</i>	<i>s2</i>					<i>g12</i>	<i>g9</i>	<i>g4</i>
12							<i>r1</i>		
13	<i>s8</i>	<i>s2</i>						<i>g14</i>	<i>g5</i>
14				<i>r3</i>	<i>r3</i>				
15	<i>s8</i>	<i>s2</i>						<i>g16</i>	<i>g4</i>
16					<i>r3</i>		<i>r3</i>		
17	<i>s19</i>	<i>s20</i>						<i>g18</i>	<i>g10</i>
18					<i>r1</i>				<i>g5</i>
19					<i>r5</i>	<i>r5</i>	<i>r5</i>		
20	<i>s19</i>	<i>s20</i>						<i>g3</i>	<i>g10</i>
									<i>g5</i>

Figura 3.24: Funcția de evoluție a automatului finit determinist pentru analiza shift-reduce.

cu punct ":" pentru a indica poziția curentă a analizorului în sirul de intrare reprezentat de forma propozițională curentă din derivarea extrem dreapta. Vom folosi în continuare denumirea de **input** pentru acest sirul ramas de citit din cuvântul de analizat, aflat în forma propozițională după continutul actual al stivei analizorului. O **stare** a analizorului este o mulțime de elemente (*itemi*).

Să considerăm gramatica din figura 3.3.1.

$$\left\{ \begin{array}{ll} S' \rightarrow S\$ & (0) \\ S \rightarrow (L) & (1) \\ S \rightarrow x & (2) \\ L \rightarrow S & (3) \\ L \rightarrow L, S & (4) \end{array} \right.$$

Înital, analizorul are stiva vidă iar input va fi un sir generat pornind de la S urmat de marcajul $\$$, adică un sir generat de S' . Notăm acest lucru prin itemul $S' \leftarrow .S\$$, unde punctul indică poziția curentă a analizorului. În această situație, când sirul input începe cu terminale derivate din S , înseamnă că rezultatul derivării va începe cu simboluri obținute prin aplicarea de reguli ce au S în stânga, adică obținute din părțile drepte ale regulilor lui S . Vom nota această situație prin

$$\boxed{\begin{array}{l} S' \rightarrow .S\$ \\ S \rightarrow .x \\ S \rightarrow .(L) \end{array}}^1$$

și o vom numi *starea 1*.

Acțiuni SHIFT. În starea 1, analizăm efectul introducerii simbolului x în stiva analizorului, adică *shift* x . Aceste tip de acțiune va fi indicată prin mutarea punctului peste x în regula $S \rightarrow x$. Celelalte două elemente ale stării nu sunt relevante pentru această acțiune deoarece punctul nu se află în fața lui x , astfel că sunt ignorate. Se obține o nouă stare după shiftare. $\boxed{S \rightarrow x.}^2$.

Dacă din starea 1 se shiftează o paranteză stângă, în vîrful stivei apare paranteza stângă și poziția analizorului va fi la începutul unui sir derivat din neterminantul L urmat de o paranteză dreapta. Un astfel de sir începe cu terminale ce se pot deriva din L , ceea ce conduce la introducerea în noua stare a elemelor provenite din regulile lui L , deci în consecință și elemente ce provin din regulile lui S . Efectul este

o nouă stare,

$$\boxed{\begin{array}{l} S \rightarrow (.L) \\ L \rightarrow .L, S \\ L \rightarrow .S \\ S \rightarrow .x \\ S \rightarrow .(L) \end{array}}^3.$$

Acțiuni Goto. În starea 1 vom analiza efectul obținut prin găsirea în input a unui sir de terminale derivat din neterminálul S (Acest lucru se poate întâmpla dacă s-a shiftat un x , eventual urmat de sir derivat din L apoi). Se vor efectua reduceri succesive ale regulilor, pana la neterminálul S , ceea ce echivaleaza cu o secvență de extrageri de simboluri din stivă, anume cele ce corespund cu partea dreaptă a regulilor, iar în final S ajunge în vîrful stivei iar în input $\$$. Din starea 1 acțiunea **goto** pentru neterminálul S va fi simulată prin mutarea punctului ;in item peste simbolul S , ceea ce conduce la o nouă stare 4: $\boxed{S' \rightarrow S.\$}$ ⁴.

Acțiuni Reduce. În starea 2 avem punctul la sfârșitul unui element. Asta înseamnă că în vîrful stivei se află partea dreaptă a regulii $S \rightarrow x$ pregătită pentru reducere, deci efectul va fi extragerea simbolurilor părții drepte și depunerea pe stivă a neterminálului S aflat în stânga regulii.

Operațiile de bază efectuate asupra stărilor analizorului sunt **Closure(I)** și **Goto(I, X)**, unde I este o mulțime de elemente, iar X este un simbol al gramaticii. **Closure** adaugă elemente noi la o mulțime când punctul se află în fața unui neterminál, iar **Goto** mută punctul peste simbolul X în toate elementele unei mulțimi. Formal, cele două operațiile sunt descrise astfel:

$\text{Closure}(I) =$ repeat for orice item $(A \rightarrow \alpha.X\beta) \in I$ for orice item $A \rightarrow \alpha.X\beta$ din I for orice regulă $X \rightarrow \gamma$ $I \leftarrow I \cup \{X \rightarrow .\gamma\}$ until I nu se modifică. return $I \square$	$\text{Goto}(I, X) =$ $J \leftarrow \{\}$ adaug $A \rightarrow \alpha.X.\beta$ la J return $\text{Closure}(J) \square$
---	--

Cu aceste două proceduri simple se poate descrie algoritmul de construcție a tabelei de analiză sintactică LR(0). Se consideră gramat-

ica augmentată cu regula suplimentară $S' \leftarrow S\$$, mulțimea stărilor T și mulțimea acțiunilor (arce automatului finit deterministic) E .

```

Inițializează  $T$  cu  $\{ \text{Closure}(\{S' \leftarrow .S\$}\})\}$ 
Inițializează  $E$  cu mulțimea vidă
repeat
    for pentru fiecare stare  $I$  din  $T$ 
        for pentru fiecare item  $A \rightarrow \alpha.X\beta$  din  $I$ 
            let  $J \leftarrow \text{Goto}(I, X)$ 
             $T \leftarrow T \cup \{J\}$ 
             $E \leftarrow E \cup \{I, XJ\}$ 
until  $E$  și  $T$  nu se modifică în iterare  $\square$ 

```

Pentru simbolul nou introdus $\$$ nu se calculează $\text{Goto}(I, \$)$, ci se schimbă starea în *Accept*.

Diagrama de stări a automatului ce definește acțiunile analizorului este dată în figura 3.25

iar tabela de analiză este cea din figura 3.26.

Determinarea acțiunilor *reduce* din tabelă s-a făcut folosind algoritmul următor:

```

 $R \leftarrow \{ \}$  (inițializează multimea acțiunilor reduce)
for pentru fiecare stare  $I$  din  $T$ 
    for pentru fiecare item  $A \rightarrow \alpha$  din  $I$ 
         $R \leftarrow R \cup \{(I, A \rightarrow \alpha)\}$ 
     $\square$ 

```

3.3.2 Analizoare SLR

Deoarece foarte puține limbaje permit analiza $LR(0)$ datorită existenței mai multor acțiuni pentru un element al tablei de analiză, se folosește o variantă îmbunătățită numită **simple LR**. Aceasta elimină intrările multiple de tipul *shift/reduce*, ca în cazul gramaticii următoare 3.3.2:

$$\begin{cases} S \rightarrow E\$ & (0) \\ E \rightarrow T + E & (1) \\ E \rightarrow T & (2) \\ T \rightarrow x & (3) \end{cases}$$

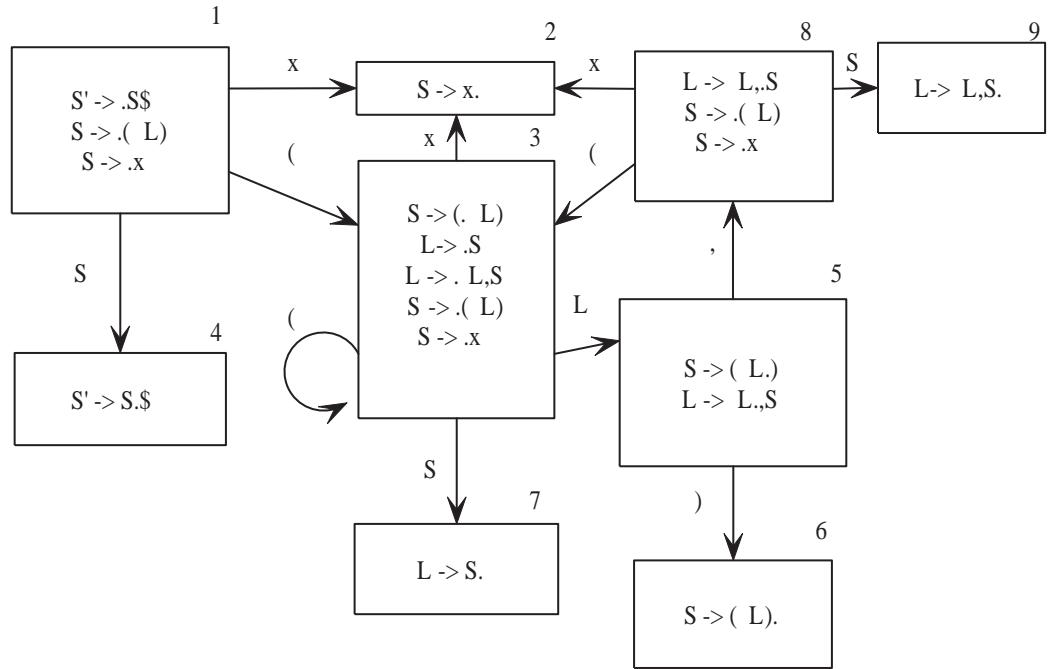


Figura 3.25: Starile automatului finit LR0

Stare	()	x	,	$\$$	S	L
1	s_3	s_2			g_4	
2	r_2	r_2	r_2	r_2	r_2	
3	s_3	s_2			g_7	g_5
4				a		
5		s_6	s_8			
6	r_1	r_1	r_1	r_1	r_1	
7	r_3	r_3	r_3	r_3	r_3	
8	s_3	s_2			g_9	
9	r_4	r_4	r_4	r_4	r_4	

Figura 3.26: Funcția de evoluție a automatului finit determinist pentru analiza shift-reduce.

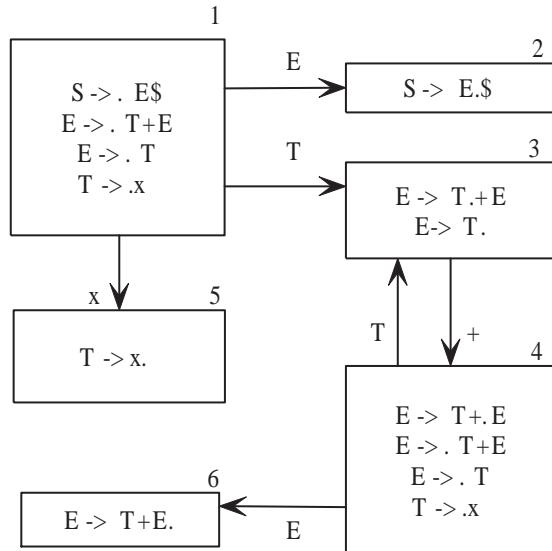
Figura 3.27: Starile analizorului **SLR** pentru gramatica 3.3.2.

Diagrama de stări a automatului finit determinist asociat analizei LR(0) este dată în figura 3.27 iar tabela de analiză din figura 3.28 conține două acțiuni pentru starea 3 și simbolul de intrare +. Eliminarea intrării multiple se face prin introducerea de acțiuni *reduce* în tabela de analiză LR(0) doar în coloanele corespunzătoare terminalelor ce fac parte din mulțimea **URM** (vezi algoritmul ?? de la analiza sintactică top down, secțiunea *Un pic de algebră*).

Algoritmul de calcul al mulțimii acțiunilor reduce devine:

```
R ← {} (initializează multimea acțiunilor reduce)
for pentru fiecare stare I din T
```

```
    for pentru fiecare item  $A \rightarrow \alpha$  din I
        for pentru fiecare terminal  $X \in URM(A)$ 
             $R \leftarrow R \cup \{(I, X, A \rightarrow \alpha)\}$ 
```

□

iar tabela SLR pentru gramatică considerată va fi cea din figura 3.29.

	x	$+$	$\$$	E	T
1	$s5$			$g2$	$g3$
2			a		
3		$s4, r2$	$r2$		
4	$s5$			$g6$	$g3$
5		$r3$	$r3$		
6			$r1$		

Figura 3.28: Tabela de analiza LR(0).

	x	$+$	$\$$	E	T
1	$s5$			$g2$	$g3$
2			a		
3		$s4$	$r2$		
4	$s5$			$g6$	$g3$
5		$r3$	$r3$		
6			$r1$		

Figura 3.29: Tabela de analiza SLR.

Notația $(I, X, A \rightarrow \alpha)$ indică faptul că din starea I , dacă simbolul X urmează în sirul input (*lookahead symbol*) atunci se face reducerea regulii $A \rightarrow \alpha$). Spunem că o gramatică este din clasa **SLR** dacă tabela de analiză SLR nu are cel mult o valoare pentru fiecare element al tabelului.

3.3.3 Gramatici LR(1) și analizoare LALR(1)

Majoritatea limbajelor de programare pentru care sintaxa este descrisă cu gramatici independente de context admit o formă a gramaticii de tip LR(1), adică se poate construi pentru acestea un analizor sintactic de tip LR(1). Algoritmul de construcție a tabelei de analiză este similar construcției tabelei SLR, dar noțiunea de *element* sau *item* este puțin mai complicată.

Un *item* LR(1) constă dintr-o regulă a gramaticii, o poziție a analizorului în partea dreaptă a regulii (marcată prin punct) și un simbol citit în avans din sirul input (*lookahead symbol*). Notația ușuală este $(A \rightarrow \alpha.\beta, x)$ și are semnificația obișnuită: sirul α este în vârful stivei, iar sirul de intrare poate fi derivat din βx . Desigur, o stare este o mulțime de elemente, iar primitivele de calcul **Closure** și **Goto** vor incorpora simbolul ce urmează la generare în sirul input. Pseudocodul asociat primitivelor va fi următorul:

Closure (I) =	Goto (I, X) =
repeat	$J \leftarrow \{\}$
for oricare $(A \rightarrow \alpha.X\beta, z) \in I$	for oricare $(A \rightarrow \alpha.X\beta, z) \in I$
for orice regula $X \rightarrow \gamma$	$J \leftarrow J \cup \{(A \rightarrow \alpha.X\beta, z)\}$
for $w \in PRIM(\beta z)$	return Closure (J) \square
$I \leftarrow I \cup \{(X \rightarrow .\gamma, w)\}$	
until I nu se modifică.	
return I \square	

jhsdghjsd Starea inițială a automatului finit va fi închiderea itemului $(S' \rightarrow .S\$, ?)$ unde $?$ înseamnă că simbolul citit în avans nu contează, deoarece terminatorul $\$$ nu va shiftat. Acțiunile **reduce** sunt selectate cu următorul algorithm:

$R \leftarrow \{\}$ (<i>initializează multimea acțiunilor reduce</i>)	
for pentru fiecare stare I din T	
for pentru fiecare item $(A \rightarrow \alpha., z)$ din I	
$R \leftarrow R \cup \{(I, z, A \rightarrow \alpha)\}$	
\square	

similar cazului gramaticilor SLR.

O gramatică ce nu este din categoria SLR, dar pentru care se poate construi o tabelă de analiză LR(1) are următoarele reguli (3.3.3):

$$\left\{ \begin{array}{ll} S' \rightarrow S\$ & (0) \\ S \rightarrow V = E & (1) \\ S \rightarrow E & (2) \\ E \rightarrow V & (3) \\ V \rightarrow x & (4) \\ V \rightarrow *E & (5) \end{array} \right.$$

	x	*	=	\$	S	E	V
1	s_8	s_6			g_2	g_5	g_3
2				a			
3			s_4	r_3			
4	s_{11}	s_{13}				g_9	g_7
5				r_2			
6	s_8	s_6				g_{10}	g_{12}
7				r_3			
8			r_4	r_4			
9				r_1			
10			r_5	r_5			
11				r_4			
12			r_3	r_3			
13	s_{11}	s_{13}				g_{14}	g_7
14				r_5			

Figura 3.30: Tabela de analiza LR(1).

Figura ?? conține diagrama de stări ale automatului finit deterministic calculat prin aplicarea primitivelor descrise în acest subcapitol. Atunci când apar mai multi itemi ce diferă doar prin simbolul citit în avans, aceștia au fost descriși pe scurt prin listarea simbolurilor în dreapta regulii. Este cazul stării 1 din diagramă.

Tabela de analiză LR(1) asociată (vezi figura 3.30) conține doar intrări simple, ceea ce va însemna că gramatica precedentă are tipul LR(1).

Deoarece tabela de analiză are de obicei dimensiuni mari se folosește o versiune prescurtată, numită tabelă **LALR(1)** obținută prin concatenarea în diagrama de stări a automatului finit LR(1) a celor stări care diferă doar prin simbolurile citite în avans. Este cazul stărilor 7,12 respectiv 8,11; 6,13; 10,14. Prin renumerotare vom obține următoarea tabelă (figura 3.31) iar analizorul va fi de tip LALR(1) (*lookahead LR(1)*).

	x	*	=	\$	S	E	V
1	$s8$	$s6$			$g2$	$g5$	$g3$
2				a			
3			$s4$	$r3$			
4	$s8$	$s6$				$g9$	$g7$
5				$r2$			
6	$s8$	$s6$				$g10$	$g7$
7			$r3$	$r3$			
8			$r4$	$r4$			
9				$r1$			
10			$r5$	$r5$			

Figura 3.31: Tabela de analiza LALR(1).

Cuprins

Bibliografie

1. Andrew W. Appel, *Modern compiler implementation in JAVA*, Cambridge University Press, 2002.
2. Octavian C. Dogaru, *Bazele informaticii. Limbaje formale*, Tipografia Universității din Timișoara, 1989.
3. Gheorghe Grigoraș, *Limbaje formale și tehnici de compilare*, Tipografia Universității "Alexandru Ioan Cuza", Iași, 1984.
4. J. E. Hopcroft , R. Motwani, J. D. Ullman, *Introduction to Automata Theory, Languages, and Computation. Second edition*, Addison Wesley, 2001.
5. Teodor Jucan, *Limbaje Formale*, Editura Academiei, București, 1993.
6. Ștefan Mărușter, *Curs de Limbaje formale și tehnici de compilare*, Tipografia Universității din Timișoara, 1980.
7. Luca Dan Șerbănați, *Limbaje de programare și compilatoare*, Editura Academiei, București, 1987.