

UNIVERSITATEA „ALEXANDRU IOAN CUZA”
FACULTATEA DE INFORMATICĂ
DEPARTAMENTUL DE ÎNVĂȚĂMÂNT LA DISTANȚĂ

Gheorghe Grigoraș

PROIECTAREA COMPILATOARELOR

2006 – 2007

EDITURA UNIVERSITĂȚII „ALEXANDRU IOAN CUZA”
IAȘI

Adresa autorului: Universitatea Alexandru Ioan Cuza Iași
Facultatea de Informatică
str. Berthelot 16
700483, Iași
grigoras@infoiasi.ro, <http://www.infoiasi.ro/~grigoras>

Cuprins

1. LIMBAJE.....	2
1.1 Limbaje de programare.....	3
1.1.1 Descrierea sintaxei și semanticii unui limbaj.....	3
1.1.2 Implementarea limbajelor.....	7
1.1.3 Analiza lexicală și analiza sintactică.....	9
1.2 Limbaje formale.....	11
1.2.1 Expresii regulate.....	11
1.2.2 Automate finite.....	13
1.2.3 Gramatici independente de context.....	17
1.2.4 Arbori sintactici. Ambiguitate.....	18
2 ANALIZA LEXICALĂ.....	23
2.1 Un algoritm de trecere de la o expresie regulată la automatul echivalent.....	24
2.2 De la un automat cu ε-tranziții la automatul determinist echivalent.....	27
2.3 Proiectarea unui analizor lexical.....	28
2.4 Generatorul Lex.....	32
3 ANALIZA SINTACTICĂ DESCENDENTĂ.....	35
3.1 Rolul analizorului sintactic.....	36
3.2 Problema recunoașterii. Problema parsării.....	36
3.3 Analiza sintactică descendentă.....	37
3.3.1 Parser descendent general.....	38
3.3.2 Gramatici LL(k).....	40
3.3.3 O caracterizare a gramaticilor LL(1).....	41
3.3.4 Determinarea mulțimilor FIRST și FOLLOW.....	42
3.3.5 Tabela de analiză sintactică LL(1).....	45
3.3.6 Analizorul sintactic LL(1).....	45
3.3.7 Eliminarea recursiei stângi.....	47
4 ANALIZA SINTACTICĂ ÎN GRAMATICI LR.....	51
4.1 Gramatici LR(k).....	51
4.2 O caracterizare a gramaticilor LR(0).....	52
4.3 Algoritm de analiză sintactică LR(0).....	60
4.4 Gramatici și analizoare SLR(1).....	65
4.5 O caracterizare a gramaticilor LR(1).....	71
4.6 Analiză sintactică LR(1) și LALR(1).....	73
5 GENERATOARE DE ANALIZOARE SINTACTICE.....	79

5.1	Utilizarea generatorului de parsere YACC	82
5.2	Aplicații cu LEX și YACC	86
6	ANALIZA SEMANTICĂ	91
6.1	Gramatici cu atribute.....	91
6.1.1	Numere raționale în baza 2 (Knuth)	91
6.1.2	Declararea variabilelor într-un program	93
6.1.3	Definiția gramaticilor cu atribute	96
6.2	Grafuri atașate unei gramatici cu atribute	98
6.2.1	Graful DG_p	98
6.2.2	Graful BG_p	99
6.2.3	Graful DG_t	99
6.3	Metode de evaluare a atributelor	101
6.3.1	Evaluatorul rezultat din definiția gramaticii	101
6.4	Traducere bazată pe sintaxă.....	102
6.4.1	Reprezentări intermediare	103
6.4.2	Traducerea codului în notație postfix	105
6.4.3	Traducerea codului în arbore sintactic	106
6.4.4	Traducerea codului în secvențe de quaduple.....	107
6.5	Proiectarea unui interpreter cu flex și yacc.....	113
	BIBLIOGRAFIE	121

1. Limbaje

Acest capitol trece în revistă chestiuni legate de limbajele de programare de nivel înalt, precum și elemente de limbaje formale necesare abordării problemelor de analiză structurală în cadrul compilatoarelor. Este știut faptul că există și sunt folosite un număr impresionant de limbaje de programare. Vom examina caracteristicile câtorva din aceste limbaje, lucru util pentru compararea lor și, prin urmare, pentru alegerea limbajului potrivit pentru o aplicație anume. De asemenea, caracteristicile limbajelor de programare sunt importante pentru proiectarea și implementarea compilatoarelor.

Limbajul este instrumentul principal al comunicării. Gradul de succes sau de eșec al comunicării este influențat de caracteristicile limbajului folosit. Comunicarea între subiecți de o anumită specialitate se poate face prin intermediul unui limbaj ce încorporează o terminologie specială care face ca multe părți ale comunicării să fie neînțelese pentru un neinițiat. Tot așa, comunicarea cu computerul presupune utilizarea unui vocabular și a unei gramatici specializate în care regulile – chiar dacă sunt mai ușor de specificat decât pentru limbajele naturale – trebuie cunoscute de ambele părți. Oamenii comunică cu computerul prin limbaje specializate; există numeroase astfel de limbaje și se proiectează mereu altele noi. Aceste limbaje sunt numite, unele, limbaje de programare, altele limbaje de comandă sau limbaje de interogare a bazelor de date etc.

1.1 Limbaje de programare

Limbajele de programare sunt instrumente folosite pentru a construi descrieri formale ale algoritmilor. Un limbaj de programare trebuie să conțină, pe lângă operațiile de bază ce transformă o *stare inițială* dată într-o *stare finală*, și modul în care pașii ce conțin aceste operații sunt înlanțuiți pentru a rezolva o anumită problemă. Pentru a răspunde la aceste cerințe, un limbaj de programare trebuie să conțină trei componente:

- *tipuri de date, obiecte și valori* împreună cu *operațiile* ce se pot aplica acestora;
- reguli pentru stabilirea *relațiilor cronologice* între operațiile specificate;
- reguli ce stabilesc structura (statică) a unui *program*.

Aceste componente constituie nivelul de *abstracție* la care putem reprezenta algoritmi în limbajul respectiv. În raport cu acest nivel de abstracție, limbajele de programare se pot clasifica – simplificând destul de mult lucrurile – în două grupe mari: *limbaje de nivel scăzut* și *limbaje de nivel înalt*. Limbajele de nivel scăzut sunt mai apropiate de limbajele mașină: există o corespondență destul de puternică între operațiile implementate de limbaj și operațiile implementate de hardware-ul corespunzător. Limbajele de nivel înalt, pe de altă parte, sunt mai apropiate de limbajele folosite de oameni pentru a exprima probleme și algoritmi. Fiecare propoziție într-un limbaj de nivel înalt poate să fie echivalentă cu o succesiune de operații (mai multe propoziții) dintr-un limbaj de nivel scăzut. În realitate, există o mulțime de limbaje între cele două clase, de nivel scăzut și de nivel înalt, și nu se poate face o împărțire strictă a limbajelor doar în două clase.

Abstractizările într-un limbaj de programare permit programatorului să extragă proprietățile esențiale necesare pentru soluția problemei de rezolvat, ascunzând detaliile de implementare a soluției. Cu cât nivelul de abstractizare este mai înalt, programatorul se gândește mai puțin la mașina (hardware-ul) pe care va fi implementată soluția problemei. Cu alte cuvinte, ideile de programare pot fi separate în întregime de arhitectura sistemului pe care va fi executat programul. De pildă, în limbajele de programare de nivel înalt, programatorul lucrează cu nume de variabile simbolice și nu cu adrese numerice de memorie. Abstracțiile de date permit ca locațiile de memorie să fie privite drept celule de memorare pentru tipuri de date de nivel înalt; programatorul nu trebuie să-și pună problema reprezentării acestora. Abstracțiile de control permit programatorului să exprime algoritmi cu ajutorul structurilor de tipul *if*, *while* sau al procedurilor. Cum sunt implementate acestea în mașina respectivă, este un lucru care nu-l interesează, în general, pe programator.

Un program scris într-un limbaj de nivel înalt este tradus în cod mașină echivalent, cu ajutorul unui program special numit *compilator*. În această carte vom prezenta unele din tehnicile cunoscute pentru proiectarea unui compilator.

1.1.1 Descrierea sintaxei și semanticii unui limbaj

Limbajele, formal vorbind, sunt mulțimi de șiruri de caractere dintr-un alfabet fixat. Acest lucru este valabil atât pentru limbajele naturale, cât și pentru cele artificiale. Șirurile de caractere ale limbajului se numesc *propoziții* (*sentences*) sau *instrucțiuni*, *afirmații* (*statements*). Orice limbaj are un număr de reguli sintactice care stabilesc care din șirurile de caractere ce se pot forma cu simboluri din alfabet fac parte din limbaj. Dacă pentru limbajele naturale regulile sintactice sunt în general complicate, limbajele de programare sunt relativ simple din punct de vedere sintactic. Există în fiecare limbaj de programare unitățile sintactice de nivelul cel mai mic – numite *lexeme* – care nu sunt cuprinse în descrierea formală a sintaxei limbajului. Aceste unități sunt cuprinse într-o specificare lexicală a limbajului care precede descrierea sintaxei. Lexemele unui limbaj de programare cuprind identificatorii, literalii, operatorii, cuvintele rezervate, semnele de punctuație. Un program poate fi privit ca un șir de lexeme și nu un șir de caractere;

obținerea șirului de lexeme din șirul de caractere (programul obiect) este sarcina *analizorului lexical*. Un *token* al unui limbaj este o categorie de lexeme. De exemplu, un identificator este un token care are ca instanțe lexeme ca: *suma*, *total*, *i*, *par5*, etc. Tokenul PLUS pentru operatorul aritmetic “+” are o singură instanță. Instrucțiunea:

```
delta = b2 - ac
```

este compusă din lexemii și tokenurile specificate în tabela următoare:

Lexem	Token
delta	IDENTIFICATOR
=	ASIGNARE
b2	IDENTIFICATOR
-	MINUS
ac	IDENTIFICATOR

Un limbaj de programare poate fi descris, formal, prin așa-zisele mecanisme de generare a șirurilor din limbaj. Aceste mecanisme se numesc gramatici și au fost descrise de Chomsky în anii 1956 – 1959. Nu mult după ce Chomsky a descoperit cele patru clase de limbaje formale, grupul ACM – GAMM a proiectat limbajul Algol 58 iar la o conferință internațională John Backus a prezentat din partea grupului acest limbaj folosind o metodă formală nouă de descriere a sintaxei. Această nouă notație a fost modificată de Peter Naur la descrierea limbajului Algol 60. Metoda de descriere a sintaxei limbajelor de programare așa cum a fost folosită de Naur poartă numele de **forma Backus - Naur** sau simplu **BNF**. De remarcat faptul că BNF este aproape identică cu mecanismul generativ al lui Chomsky – gramatica independentă de context. BNF este un metalimbaj. Acesta folosește abstracțiile pentru a descrie structurile sintactice. O abstracție este scrisă între paranteze unghiulare, ca de exemplu: <asignare>, <variabilă>, <expresie>, <if_stm>, <while_stm> etc. Definiția unei abstracții este dată printr-o *regulă* sau *producție* care este formată din:

- partea stângă a regulii ce conține abstracția care se definește.
- o săgeată \rightarrow (sau caracterele ::=) care desparte partea stângă a regulii de partea dreaptă.
- partea dreaptă a regulii care este un șir format din abstracții, tokenuri, lexeme.

De exemplu, regula:

```
<asignare>  $\rightarrow$  <variabilă> = <expresie>
```

definește sintaxa unei asignări: abstracția <asignare> este o instanță a abstracției <variabilă> urmată de lexemul = urmat de o instanță a abstracției <expresie>. O propoziție care are structura sintactică descrisă de această regulă este:

```
total = suma1 + suma2
```

Într-o regulă BNF abstracțiile se numesc *simboluri neterminale* sau simplu *neterminali*, lexemii și token-urile se numesc *simboluri terminale* sau simplu *terminali*. O descriere BNF sau o gramatică (independentă de context) este o colecție de reguli.

Dacă se întâmplă ca un neterminal să aibă mai multe definiții acestea se pot insera într-o singură regulă în care partea dreaptă conține părțile drepte ale definițiilor sale despărțite prin simbolul |. De exemplu, definițiile:

```
<if_stm>  $\rightarrow$  if <expr_logică> then <stm>
```

```
<if_stm>  $\rightarrow$  if <expr_logică> then <stm> else <stm>
```

se scriu:

```
<if_stm>  $\rightarrow$  if <expr_logică> then <stm>
```

```
| if <expr_logică> then <stm> else <stm>
```

Iată un exemplu de gramatică (BNF) ce descrie un minilimbaj de programare:

```
<program>  $\rightarrow$  begin <lista_de_instructiuni> end
```

```
<lista_de_instructiuni>  $\rightarrow$  <instructiune>
```

	<lista_de_instructiuni>; <instructiune>
<instructiune>	→ <variabila> = <expresie>
<variabila>	→ LITERA
<expresie>	→ <expresie> + <expresie>
	<expresie> - <expresie>
	<variabila>
	NUMAR

Aici apar două tokenuri: LITERA și NUMAR. Acestea pot la rândul lor să fie descrise prin același mecanism:

LITERA → a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z

NUMAR → <cifra> | NUMAR <cifra>

<cifra> → 0|1|2|3|4|5|6|7|8|9

Se observă că de data aceasta gramatica este de tip 3 (sau gramatică regulată) ceea ce înseamnă că tokenurile se pot descrie folosind expresiile regulate iar recunoașterea lor (*analiza lexicală*) se poate face cu automate finite deterministe.

BNF a fost extinsă în diverse moduri mai ales din rațiuni de implementare a unui generator de analizor sintactic. Orice extensie a BNF este numită EBNF (Extended BNF). Extensiile BNF nu măresc puterea descriptivă a mecanismului ci reprezintă facilități privind citirea sau scrierea unităților sintactice. Extensiile ce apar frecvent sunt:

- includerea, în partea dreaptă a unei reguli, a unei părți opționale ce se scrie între paranteze drepte. De exemplu, instrucțiunea if ar putea fi descrisă astfel:

$$\langle \text{if_stm} \rangle \rightarrow \text{if } \langle \text{expr_logică} \rangle \text{ then } \langle \text{stm} \rangle [\text{else } \langle \text{stm} \rangle]$$
- includerea în partea dreaptă a unei reguli a unei părți ce se repetă de zero sau mai multe ori. Această parte se scrie între acolade și poate fi folosită pentru a descrie recursia:

$$\langle \text{lista_de_instructiuni} \rangle \rightarrow \langle \text{instructiune} \rangle \{ ; \langle \text{instructiune} \rangle \}$$
- includerea în partea dreaptă a unei reguli a opțiunii de a alege o variantă dintr-un grup. Acest lucru se face punând variante în paranteze rotunde, despărțite prin |. De exemplu, o instrucțiune for se poate defini astfel:

$$\langle \text{for_stm} \rangle \rightarrow \text{for } \langle \text{var} \rangle = \langle \text{expr} \rangle (\text{to} \mid \text{downto}) \langle \text{expr} \rangle \text{ do } \langle \text{stm} \rangle$$

Trebuie precizat că parantezele de orice fel din descrierea EBNF fac parte din metalimbaj, ele nu sunt terminali ci doar instrumente de a nota o anume convenție. În cazul în care unele din aceste simboluri sunt simboluri terminale în limbajul descris, acestea vor fi puse între apostrof: ‘[’ sau ‘{’ etc.

EBNF nu face altceva decât să simplifice într-un anume fel definițiile sintactice. Spre exemplificare iată cum se descrie expresia aritmetică folosind cele două mecanisme:

BNF:

$$\begin{aligned} \langle \text{expresie} \rangle &\rightarrow \langle \text{expresie} \rangle + \langle \text{termen} \rangle \\ &\quad | \langle \text{expresie} \rangle - \langle \text{termen} \rangle \\ &\quad | \langle \text{termen} \rangle \\ \langle \text{termen} \rangle &\rightarrow \langle \text{termen} \rangle * \langle \text{factor} \rangle \\ &\quad | \langle \text{termen} \rangle / \langle \text{factor} \rangle \\ &\quad | \langle \text{factor} \rangle \\ \langle \text{factor} \rangle &\rightarrow \text{'('} \langle \text{expresie} \rangle \text{' } \\ &\quad | \langle \text{identificator} \rangle \end{aligned}$$

EBNF:

$$\begin{aligned} \langle \text{expresie} \rangle &\rightarrow \langle \text{termen} \rangle \{ (+ \mid -) \langle \text{termen} \rangle \} \\ \langle \text{termen} \rangle &\rightarrow \langle \text{factor} \rangle \{ (* \mid /) \langle \text{factor} \rangle \} \\ \langle \text{factor} \rangle &\rightarrow \text{'('} \langle \text{expresie} \rangle \text{' } \mid \langle \text{identificator} \rangle \end{aligned}$$

Un exemplu de program în minilimbajul de programare descris mai sus este următorul:

```
begin a = b + c ; d = 5 - a end
```

O *derivare* a acestui program în gramatica dată este:

```
<program>  => begin <lista_de_instructiuni> end
=> begin <lista_de_instructiuni>; <instructiune> end
=> begin <instructiune>; <instructiune> end
=> begin <variabila>=<expresie>; <instructiune> end
=> begin LITERA=<expresie>; <instructiune> end
=> begin LITERA=<expresie>+<expresie>; <instructiune> end
=> begin LITERA=LITERA+<expresie>; <instructiune> end
=> begin LITERA=LITERA+LITERA; <instructiune> end
=> begin LITERA=LITERA+LITERA; <variabila>=<expresie>end
=> begin LITERA=LITERA+LITERA; LITERA=<expresie>end
=> begin LITERA=LITERA+LITERA; LITERA=<expresie>-<expresie> end
=> begin LITERA=LITERA+LITERA; LITERA=NUMAR-<expresie> end
=> begin LITERA=LITERA+LITERA; LITERA=NUMAR-LITERA end
```

Arborele de derivare pentru acest program este dat în figura 1.1. Se observă o strânsă legătură între derivările dintr-o gramatică și arborii de parsare.

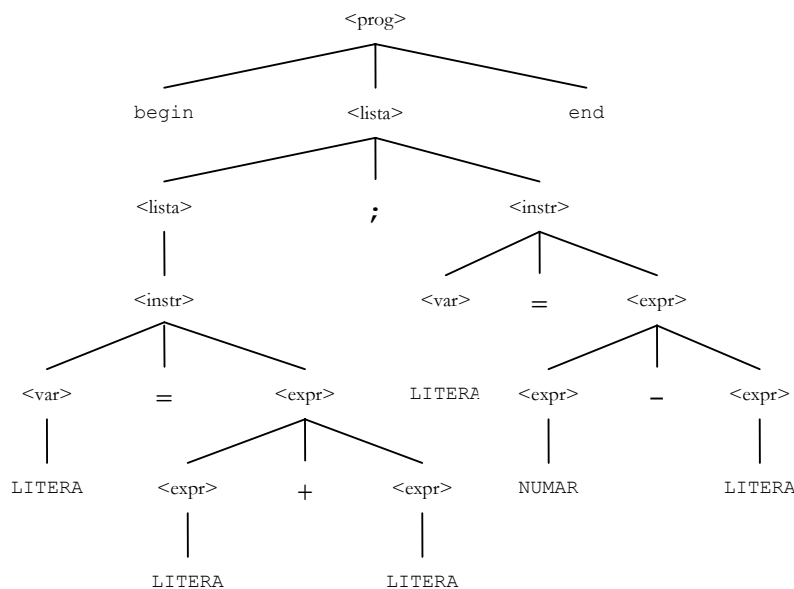


Figura 1.1

Există reguli ale limbajelor de programare care nu pot fi descrise utilizând BNF. De pildă faptul că o variabilă nu poate fi utilizată înainte de a fi declarată este o caracteristică ce nu poate fi descrisă prin BNF. Reguli de acest fel fac parte din *semantica statică* a limbajului: reguli ce țin indirect de înțelesul unui program relativ la execuție. Semantica statică este denumită așa pentru că poate fi verificată înainte de execuție, la momentul compilării. Un mecanism care descrie semantica statică a fost introdus de Knuth în 1968 și este denumit gramatică cu atribute. O *gramatică cu atribute* este o gramatică independentă de context, care descrie sintaxa, la care s-au adăugat următoarele:

- fiecărui simbol al gramaticii i se atașează o mulțime de *atribute*, unele dintre ele sunt *sintetizate* (valorile lor depind de valorile atributelor descendenților), altele *moștenite* (valorile lor depind de valorile atributelor părinților și fraților);
- fiecărei reguli din gramatică i se asociază o mulțime de *funcții semantice* care determină valorile atributelor sintetizate ale părții stânga și valorile atributelor moștenite ale simbolurilor din partea dreaptă.

Un arbore de derivare într-o gramatică cu atribute este arborele de derivare din gramatica ce descrie sintaxa împreună cu mulțimea valorilor atributelor atașate fiecărui nod al arborelui. Pentru a evita *circularitatea* în obținerea valorilor atributelor, deseori se pune restricția ca un atribut moștenit să depindă doar de valorile atributelor părții stânga a regulii și de atributele simbolurilor din stânga sa în șirul din partea dreaptă a regulii. Această restricție permite proiectarea unui evaluator care parcurge arborele de derivare în preordine. Nodurile frunze din arborele de derivare au atribute sintetizate ale căror valori se determină în afara arborelui de derivare; acestea se numesc *atribute intrinseci*. Spre exemplu, tipul unei instanțe a unei variabile într-un program se obține din tabela simbolurilor care este construită la analiza lexicală și conține numele variabilelor, tipul lor etc.. În capitolul *Analiza semantică* vor fi date exemple de gramatici cu atribute și modul de evaluare a atributelor.

Descrierea înțelesului unui program înseamnă *semantica dinamică*. Este vorba despre înțelesul expresiilor, a instrucțiunilor, a unităților de program. Dacă pentru descrierea sintaxei formalismul BNF este universal acceptat, nu există încă o notație universal acceptată pentru semantică (în continuare semantica dinamică o să fie numită semantică). Și aceasta în condițiile în care fiecare programator trebuie să cunoască atunci când folosește un limbaj, ce transformări face fiecare din instrucțiunile sale și cum sunt obținute aceste transformări. De obicei ei “învață” semantica citind text în limbaj natural care descrie fiecare instrucțiune. Or se știe că, din multe motive, aceste explicații sunt imprecise și / sau incomplete. Pe de altă parte, cei care implementează limbajele deseori folosesc explicațiile în limbaj natural atunci când proiectează compilatoarele. Utilizarea semanticii în limbaj natural atât de către programatori, cât și de către implementatori se poate explica prin aceea că mecanismele formale de descriere a semanticii sunt în general complexe și greu de utilizat practic. De aceea domeniul acesta este încă în atenția multor cercetători.

1.1.2 Implementarea limbajelor

Două din componentele primare ale unui computer sunt *memoria internă* și *procesorul*. Memoria internă este folosită pentru păstrarea programelor și a datelor. Procesorul este o colecție de circuite care implementează operații sau instrucțiuni mașină. În absența altor componente software un computer înțelege doar propriul limbaj mașină. Pentru a fi folosit, computerul are nevoie de un sistem de operare - care reprezintă interfața pentru alte sisteme software - și de un număr de limbaje de programare care să faciliteze rezolvarea de probleme diverse.

Limbajele de programare pot fi implementate prin trei metode generale: compilare, interpretare și sisteme de implementare hibride.

Un *compilator* este un sistem software care traduce programele scrise într-un limbaj de programare (*programe sursă*) în programe scrise în limbaj mașină (*programe obiect*) ce pot fi executate direct de către computer. Avantajul acestei metode de implementare este acela că execuția programelor este foarte rapidă. Procesul de compilare, în cele mai multe cazuri, se desfășoară după schema din figura 1.2. *Analizorul lexical* identifică în textul sursă *unitățile lexicale* (tokenurile) și transformă acest text într-un șir de unități lexicale,

ignorând spațiile albe și comentariile. Unități lexicale într-un limbaj de programare sunt identificatorii, operatorii, cuvintele rezervate, semnele speciale.

Analizorul sintactic preia șirul de unități lexicale de la analizorul lexical și construiește arborele de derivare (de parsare) al acestui șir în gramatica ce descrie sintaxa limbajului de programare. De fapt, de cele mai multe ori, se obține o reprezentare a acestui arbore de derivare, un set de informații din care poate fi construit acest arbore.

Tabela de simboluri, construită de cele două analizoare amintite mai sus, cuprinde informații privind numele, tipul, atributele tuturor identificatorilor definiți de utilizator în program. Aceste informații sunt necesare analizorului semantic și generatorului de cod.

Analizorul semantic este de regulă un generator de cod intermediar care are rolul de a verifica erorile ce nu pot fi detectate de analizorul sintactic (de exemplu cele referitoare la tipul variabilelor) și de a obține din arborele de derivare un program echivalent cu cel inițial, dar într-un limbaj intermediar între limbajul sursă și limbajul mașină. O parte opțională a procesului de compilare este optimizarea. Aceasta are rolul de a îmbunătăți programul obținut, în ceea ce privește mărimea sa, prin eliminarea unor instrucțiuni sau în ceea ce privește viteza de calcul. Optimizarea, atunci când se face, este mai ușor de aplicat codului intermediar decât codului mașină.

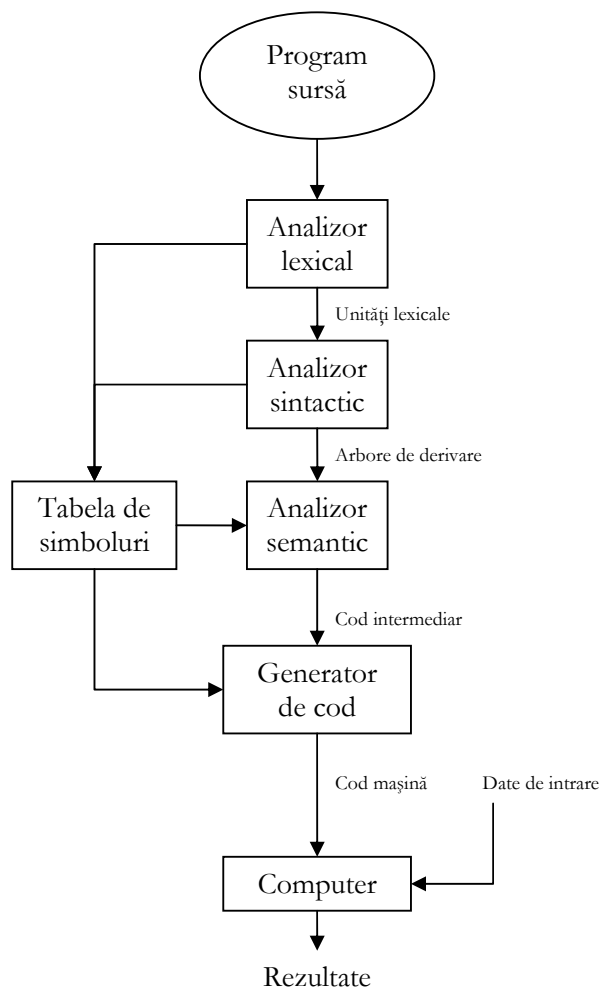


Figura 1.2

Generatorul de cod traduce codul intermediar (optimizat) într-un program echivalent scris în cod mașină. Acest program cod mașină are nevoie, de cele mai multe ori, pentru a fi executat, de module program oferite de sistemul de operare (cum ar fi programe de introducere sau extragere a datelor în/din memorie). Procesul de colectare a programelor sistem necesare execuției programului utilizator compilat și de legare a lor la acesta se face de către un program special numit *linker*.

O tehnică destul de folosită în implementarea limbajelor de programare este *interpretarea* sau *interpretarea pură* cum mai este cunoscută în literatura de specialitate. Programele scrise în limbajul sursă sunt interpretate de către un alt program, numit *interpreter*, care acționează ca un simulator software al unei mașini ce execută instrucțiune cu instrucțiune programul scris în limbajul sursă fără ca acesta să fie tradus în limbaj mașină. Simulatorul software oferă o mașină virtuală pentru limbajul sursă în cauză.

Avantajul interpretării este acela că permite implementarea ușoară a operațiilor de depanare a programelor, pentru că erorile de execuție se referă la unități ale limbajului sursă și nu țin de limbajul mașină. Din păcate, dezavantajele sunt mai numeroase: timpul de execuție este serios diminuat în cazul interpretării față de compilare (un program interpretat se execută de 10 până la 100 ori mai lent decât unul compilat); spațiul folosit de un program interpretat este mai mare; în cazul interpretării, semantica unei construcții trebuie determinată din programul sursă la momentul execuției. Aceste dezavantaje fac ca doar limbajele cu structuri simple să fie implementate prin interpretare (comenzi sistem în UNIX, APL, LISP). Există și excepții: cu toate că nu este un limbaj simplu, JavaScript este interpretat.

Sistemele de *implementare hibridă* folosesc atât compilarea, cât și interpretarea pentru implementarea limbajelor. Această tehnică se rezumă la obținerea unui cod intermediar pentru programul scris în limbaj de nivel înalt (compilare) și apoi interpretarea acestui cod intermediar. Desigur, se alege limbajul intermediar așa fel ca el să poată fi ușor interpretat. Avantajul implementărilor hibride rezultă din aceea că obținerea codului intermediar se poate face automat iar execuția este mai rapidă decât în cazul interpretării pure. Ca exemple de implementări hibride consemnăm limbajul Perl și implementările de început ale limbajului Java care este tradus în cod intermediar numit “byte code” și poate fi portat pe orice mașină care dispune de un interpreter de “byte code” asociat cu un sistem de execuție, în fapt ceea ce se numește *mașină virtuală Java*.

1.1.3 Analiza lexicală și analiza sintactică

Sintaxa unui limbaj de programare se descrie folosind formalismul gramaticilor independente de context sau așa zisa BNF – Backus Naur Form. Utilizarea BNF în descrierea sintaxei este avantajoasă pentru că este un formalism clar și concis și poate fi ușor transcris în sistemele software care implementează limbajul respectiv. În aproape toate compilatoarele, verificarea corectitudinii sintactice se face în două etape distincte: *analiza lexicală* și *analiza sintactică*. *Analiza lexicală* se ocupă de construcțiile simple ale limbajului precum: identificatori, literalii numerici, operatori, cuvinte rezervate, semne speciale. *Analiza sintactică* se ocupă de construcțiile evaluate din limbaj precum expresiile, instrucțiunile, unitățile de program. Tehnicile pentru analiza lexicală sunt mult mai puțin complexe față de cele ale analizei sintactice. Asta face ca implementarea separată a analizei lexicale să ducă la obținerea unor module relativ mici și ușor de întreținut pentru fiecare din cele două faze. În esență un analizor lexical este un program pentru recunoașterea subșirurilor de caractere dint-un șir dat, subșiruri ce se potrivesc cu un anume șablon (problema este cunoscută sub numele de pattern matching și se întâlnește frecvent în editoarele de text). Analizorul lexical citește textul sursă și colectează caracterele în grupări logice care se numesc *lexeme*. Acestor grupări logice li se asociază coduri interne care se numesc *tokenuri*. De pildă, textul următor:

alpha = beta + 734;

în urma analizei lexicale se prezintă astfel:

Lexem	Token
alpha	IDENT
=	ASIGN
beta	IDENT
+	PLUS
734	INT_LIT
;	PV

IDENT, ASIGN etc. sunt niște nume pentru niște coduri numerice care vor fi returnate de către analizor atunci când grupul de caractere citit îndeplinește condiția de a fi identificator, operator de asignare etc. Pentru că spațiile albe și comentariile nu sunt relevante pentru execuția unui program, acestea sunt ignorate de către analizorul lexical. Rolul analizorului lexical, în contextul compilării unui program, este de a citi textul sursă și de a oferi analizorului sintactic tokenurile.

Partea de analiză a textului sursă care se referă la unitățile sintactice precum expresii, instrucțiuni, blocuri etc. poartă numele de *analiză sintactică* sau *parsare*. Rolul unui analizor sintactic este dublu. În primul rând, este acela de a verifica dacă programul este corect sintactic. Dacă se întâlnește o eroare atunci analizorul produce un diagnostic și analizează mai departe textul astfel încât să fie semnalate cât mai multe din erorile de sintaxă existente în text. Al doilea rol este acela de a construi arborele de derivare al programului respectiv în gramatica ce descrie limbajul. De fapt, sunt obținute informațiile din care se poate construi arborele de derivare: derivarea extrem stângă sau derivarea extrem dreaptă corespunzătoare, care însemna moduri de traversare a arborilor de derivare. Sunt cunoscute două tipuri de parsare, în funcție de direcția în care este construit arborele de derivare: *parsare top – down* (*descendente*) care construiesc arborele de la rădăcină către frunze și *parsare bottom – up* (*ascendente*) care construiesc arborele de la frunze către rădăcină. Un parser *top – down* trasează derivarea extrem stângă a cuvântului de analizat. Arborele de derivare este construit în preordine: se începe cu rădăcina și apoi fiecare nod este vizitat înainte ca descendenții săi să fie generați. Aceștia sunt vizitați în ordinea de la stânga la dreapta. În termenii derivării, acest lucru se poate sintetiza astfel: dată o formă propozițională (un cuvânt dintr-o derivare extrem stângă), sarcina parserului este de a determina următoarea formă propozițională în derivare. Dacă forma propozițională curentă este $uA\alpha$ unde u este un cuvânt format din terminali (tokenuri), A este neterminal iar α este un cuvânt format din terminali și neterminali, atunci sarcina parserului este de a înlocui pe A cu partea dreaptă a unei reguli ce are pe A în partea stângă. Dacă regulile corespunzătoare lui A au în partea dreaptă respectiv β_1 , β_2, \dots, β_n și se alege β_k , atunci noua formă propozițională este $u\beta_k\alpha$. În general, alegerea corectă a variantei de rescriere a lui A se face prin analiza simbolului următor lui u în cuvântul de analizat uv . Dacă acest lucru este posibil, atunci gramatica se numește gramatică LL(1) și metoda de parsare se numește parsare LL(1). Implementarea se face fie prin transformarea fiecărei reguli din gramatică (forma BNF) într-un număr de linii de cod (parsare recursiv descendentă) fie prin utilizarea unei table de parsare pentru implementarea regulilor BNF. În capitolul al treilea vom trata pe larg parsarea LL(1).

Un parser *bottom-up* construiește arborele de derivare pornind de la frunze și înaintând spre rădăcină. Acesta produce reversul unei derivări extrem drepte a cuvântului de analizat. Parsarea bottom-up se poate descrie succint astfel: dată forma propozițională dreaptă (un cuvânt din derivarea extrem dreaptă) γ parserul trebuie să găsească în α partea dreaptă a unei reguli, iar aceasta se reduce la partea stângă și se obține forma propozițională precedentă. Dacă se găsește $\gamma = \alpha\beta u$ și β este partea dreaptă a regulii ce

are în stânga A, atunci precedenta formă propozițională este $\alpha A u$. Algoritmii de analiză bottom-up sunt cei bazați pe relații de precedență sau, cei mai utilizați, cei din familia LR.

1.2 Limbaje formale

În acest paragraf trecem în revistă, succint, noțiuni de limbaje formale necesare în abordarea analizei lexicale, sintactice și semantice în cadrul procesului de compilare. Un limbaj (formal) este o mulțime de cuvinte peste un alfabet Σ . Notăm prin Σ^* limbajul total: mulțimea tuturor cuvintelor peste Σ . Dacă L_1 și L_2 sunt limbaje, atunci la fel sunt și reuniunea lor, intersecția, complementara lui L_1 față de Σ^* . De asemenea, dacă vom considera produsul uv a două cuvinte u și v ca fiind un nou cuvânt ce se obține prin concatenarea celor două, notăm prin $L_1 L_2$ produsul celor două limbaje care înseamnă mulțimea cuvintelor uv unde u este din L_1 iar v este din L_2 . Vom nota prin ϵ cuvântul nul, cel fără de litere. Atunci putem vorbi de puterea L^n a unui limbaj L folosind produsul și stabilind că L^0 este limbajul ce conține doar ϵ . Considerăm și iterația unui limbaj L , notată L^* ca fiind reuniunea limbajelor L^n pentru $n = 0, 1, 2, \dots$. Este cunoscută ierarhia lui Chomsky în care limbajele se clasifică în *limbaje regulate*, *limbaje independente de context*, *limbaje dependente de context* și *limbaje de tip 0*. Limbajele regulate joacă un rol important în analiza lexicală. Unitățile lexicale pot fi descrise cu ajutorul *expresiilor regulate* (algebric) sau al gramaticilor regulate (generativ) și pot fi recunoscute cu ajutorul *automatelor finite*. În esență, un *analizor lexical* este implementarea unui automat finit. Vom trece în revistă pentru început chestiunile legate de expresii regulate automate finite, fără a intra în detalii. Apoi, vom aminti câte ceva despre gramatici independente de context, instrumentul generativ pentru descrierea sintaxei limbajelor de programare. Cititorul este îndemnat a parcurge lucrări precum [Juc99], [Gri86] pentru a afla amănunte, mai ales de ordin teoretic. Noi suntem interesați aici de utilizarea acestor modele de calcul pentru dezvoltarea de algoritmi de analiză structurală a textelor sursă încât, vom prezenta atât cât este necesar pentru acest lucru.

1.2.1 Expresii regulate

Definiția 1.2.1 Fie Σ un alfabet, simbolurile $\epsilon, \Phi, |, \bullet, *, (,)$ care nu aparțin lui Σ și E un cuvânt peste alfabetul $\Sigma \cup \{ \epsilon, \Phi, |, \bullet, *, (,) \}$. O *expresie regulată* peste Σ se definește inductiv astfel:

1. E este un *atom regulat* peste Σ dacă E este un simbol din $\Sigma \cup \{ \epsilon, \Phi \}$ sau este de forma (E_1) unde E_1 este o expresie regulată peste Σ ;
2. E este un *factor regulat* peste Σ dacă E este un atom regulat peste Σ sau este de forma E_1^* unde E_1 este un factor regulat peste Σ ;
3. E este un *termen regulat* peste Σ dacă E este un factor regulat peste Σ sau este de forma $E_1 \bullet E_2$ unde E_1 este un termen regulat, iar E_2 este un factor regulat peste Σ ;
4. E este o *expresie regulată* peste Σ dacă E este un termen regulat peste Σ sau este de forma $E_1 | E_2$, unde E_1 este o expresie regulată, iar E_2 este un termen regulat peste Σ .

Aici ϵ, Φ sunt privite ca simple simboluri fără vreo semnificație. Mai jos, interpretarea acestor expresii va fi desigur limbajul $\{\epsilon\}$ respectiv limbajul vid.

În definiția de mai sus, pe lângă noțiunea de expresie regulată se dau cele de termen regulat, factor regulat, atom regulat care reflectă precedența și asociativitatea operatorilor $|, \bullet, *$. În cele ce urmează vom omite scrierea operatorului \bullet așa cum se obișnuiește la scrierea operatorului de înmulțire.

Definiția 1.2.2 Limbajul $L(E)$ descris (sau notat, sau denotat) de expresia regulată E peste alfabetul Σ este definit inductiv astfel:

1. $L(\Phi) = \Phi$, $L(\varepsilon) = \{\varepsilon\}$, $L(a) = \{a\}$, pentru orice a din Σ ;
2. $L((E)) = L(E)$, pentru orice expresie regulată E peste Σ ;
3. $L(E^*) = (L(E))^*$, pentru orice factor regulat E peste Σ ;
4. $L(E_1 E_2) = L(E_1) L(E_2)$, pentru orice E_1 , termen regulat, și pentru orice E_2 factor regulat peste Σ ;
5. $L(E_1 \mid E_2) = L(E_1) \cup L(E_2)$, pentru orice E_1 expresie regulată, și pentru orice E_2 termen regulat peste Σ .

Proprietățile 1 și 2 de mai sus asigură faptul că $L(E_1 E_2)$ și $L(E_1 \mid E_2)$ sunt bine definite.

Exemplul 1.2.1 Fie $E_1 = \varepsilon \mid (a \mid b)^*(ab \mid baa)$ o expresie regulată.

Atunci limbajul descris de E_1 este:

$$\begin{aligned}
 L(E_1) &= L(\varepsilon \mid (a \mid b)^*(ab \mid baa)) = L(\varepsilon) \cup L((a \mid b)^*(ab \mid baa)) = \\
 &= \{\varepsilon\} \cup L((a \mid b)^*) L(ab \mid baa) = \\
 &= \{\varepsilon\} \cup (L(a \mid b))^* L(ab \mid baa) = \\
 &= \{\varepsilon\} \cup (L(a \mid b))^* (L(ab) \cup L(baa)) = \\
 &= \{\varepsilon\} \cup (L(a) \cup L(b))^* (L(a)L(b) \cup L(b)L(a)L(a)) = \\
 &= \{\varepsilon\} \cup (\{a\} \cup \{b\})^* (\{a\}\{b\} \cup \{b\}\{a\}\{a\}) = \\
 &= \{\varepsilon\} \cup \{a, b\}^* \{ab, baa\}.
 \end{aligned}$$

Analog pentru $E_2 = (a \mid b)(a \mid b)$ obținem $L(E_2) = \{aa, ab, ba, bb\}$

Este binecunoscut faptul că mulțimea limbajelor descrise de expresiile regulate coincide cu mulțimea limbajelor regulate [Juc99]. Asta înseamnă că pentru fiecare expresie regulată E poate fi construit (efectiv) un automat finit care să recunoască exact limbajul descris de expresia E . De asemenea, pentru orice automat finit, există o expresie regulată (care poate fi obținută efectiv) ce descrie limbajul recunoscut de el. Cum suntem interesați în proiectarea unui analizor lexical, parte a unui compilator, vom descrie algoritmic în paragrafele următoare trecerea de la o expresie regulată la un automat finit; algoritmi diferă substanțial de cei utilizați în demonstrațiile teoretice.

Definiția 1.2.3 Două expresii regulate E_1, E_2 peste alfabetul Σ , sunt echivalente dacă $L(E_1) = L(E_2)$. Notăm acest lucru prin $E_1 = E_2$.

Exemplul 1.2.2 Se verifică ușor că au loc următoarele :

$$a \mid b = b \mid a \quad \text{și} \quad \varepsilon \mid (a \mid b)^* = (a \mid b)^*.$$

Lema 1.2.1 (Proprietăți algebrice ale expresiilor regulate). Fie E, E_1, E_2, E_3 expresii regulate peste alfabetul Σ . Sunt adevărate următoarele:

1. $E_1 \mid E_2 = E_2 \mid E_1$
2. $(E_1 \mid E_2) \mid E_3 = E_1 \mid (E_2 \mid E_3)$
3. $(E_1 E_2) E_3 = E_1 (E_2 E_3)$
4. $E_1 (E_2 \mid E_3) = E_1 E_2 \mid E_1 E_3$
5. $(E_1 \mid E_2) E_3 = E_1 E_3 \mid E_2 E_3$
6. $\varepsilon E = E$, $E \varepsilon = E$, $E \Phi = \Phi$, $\Phi E = \Phi$
7. $E^* = (E \mid \varepsilon)^*$, $\Phi^* = \varepsilon$
8. $(E^*)^* = E^*$
9. $(E_1^* E_2^*)^* = (E_1 \mid E_2)^*$

Demonstrație Rezultă din definiția anterioară și proprietățile operațiilor cu mulțimi.

□

În procesul de descriere a limbajelor, un rol important îl joacă noțiunea de ambiguitate. Dacă \mathcal{D} este o descriere a unui limbaj iar $L(\mathcal{D})$ este limbajul descris de \mathcal{D} , spunem că \mathcal{D} este ambiguă dacă există măcar un cuvânt din $L(\mathcal{D})$ care este descris în două moduri diferite.

De exemplu expresia regulată $E = ab \mid (a \mid b)^*$ este ambiguă deoarece cuvântul ab în $L(E)$ are două descrieri: una în expresia ab și alta în $(a \mid b)^*$.

Formal, neambiguitatea unei expresii regulate este definită inductiv astfel:

Definiția 1.2.4 (Expresii regulate neambigue)

1. Expresiile regulate Φ , ϵ , a și (E) sunt neambigue, pentru orice a din Σ și oricare ar fi expresia neambiguă E peste Σ ;
2. Pentru orice factor E peste Σ , E^* este neambiguă dacă E este neambiguă și pentru orice $u \in L(E^*)$ există exact un număr $n \geq 0$ și o secvență unică (u_1, u_2, \dots, u_n) de cuvinte din $L(E)$ astfel încât $u_1 u_2 \dots u_n = u$ (pentru cazul $n = 0$, $u_1 u_2 \dots u_n = \epsilon$);
3. Pentru orice termen E_1 și pentru orice factor E_2 peste Σ , $E_1 E_2$ este neambiguă dacă:
 - a. $L(E_1 E_2) = \Phi$, sau
 - b. E_1 și E_2 sunt neambigue și pentru orice $u \in L(E_1 E_2)$ există o pereche unică (u_1, u_2) în $L(E_1) \times L(E_2)$ astfel încât $u_1 u_2 = u$.
4. Pentru orice expresie E_1 și orice termen E_2 , expresia $E_1 \mid E_2$ este neambiguă dacă E_1 și E_2 sunt neambigue și în plus $L(E_1) \cap L(E_2) = \Phi$.

1.2.2 Automate finite

Definiția 1.2.5 Un automat finit este sistemul $A = (Q, \Sigma, \delta, q_0, F)$ unde Q și Σ sunt mulțimi finite, nevide, numite mulțimea stărilor respectiv *alfabetul de intrare*, $q_0 \in Q$ este *starea inițială*, $F \subseteq Q$ este *mulțimea stărilor finale* iar δ este o funcție

$$\delta: Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q,$$

numită *funcția de tranziție* (unde prin 2^Q am notat mulțimea părților lui Q).

Observația 1.2.1 Modelul definit mai sus este cel cunoscut în literatură și sub denumirea de *sistem tranzițional* (sau automat *nedeterminist cu ϵ -tranziții*). Prin particularizări ale lui δ vom obține modelul de *automat nedeterminist* (fără ϵ -tranziții) și cel de *automat determinist*.

Un automat finit poate fi reprezentat prin *tabela de tranziție* (funcția δ) sau prin *graful de tranziție*. În reprezentarea grafului de tranziție facem convenția ca stările care nu sunt finale să le reprezentăm prin cercuri iar cele finale prin pătrate. De asemenea ϵ -tranzițiile sunt reprezentate prin arce neetichetate.

Exemplul 1.2.3 Fie $Q = \{0, 1, 2\}$, $\Sigma = \{a, b, c\}$, $F = \{2\}$, $q_0 = 0$, iar δ este dată astfel:

Tabela de tranziție:

δ	a	b	c	ϵ
0	$\{0\}$	Φ	Φ	$\{1\}$
1	Φ	$\{1\}$	Φ	$\{2\}$
2	Φ	Φ	$\{2\}$	Φ

Graful de tranziție:

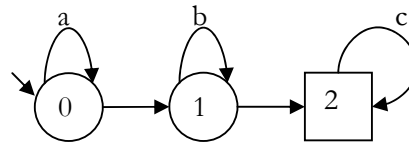


Figura 1.3

Exemplul 1.2.4 $Q = \{0, 1, 2\}$, $\Sigma = \{a, b\}$, $F = \{2\}$, $q_0 = 0$, iar δ este:

Tabela de tranziție:

δ	a	b	ϵ
0	{0, 1}	Φ	Φ
1	{1}	{2}	Φ
2	{2}	Φ	{0, 1}

Graful de tranziție:

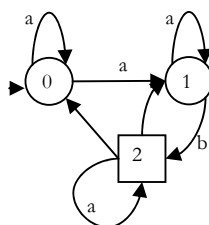


Figura 1.4

Exemplul 1.2.5 $Q = \{0, 1, 2, 3\}$, $\Sigma = \{a, b\}$, $F = \{0\}$, $q_0 = 0$, iar δ este dată în figura 1.5.

Tabela de tranziție:

δ	a	b
0	1	3
1	0	2
2	3	1
3	2	0

Graful de tranziție:

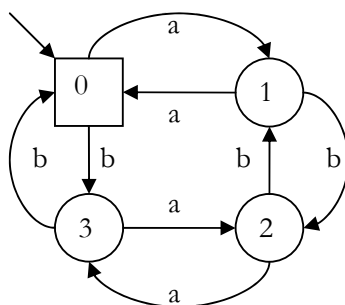


Figura 1.5

Exemplul 1.2.6 $Q = \{0, 1, 2\}$, $\Sigma = \{a, b\}$, $F = \{2\}$, $q_0 = 0$, iar δ este:

Tabela de tranziție:

δ	a	b	c
0	{0, 1, 2}	{1, 2}	{2}
1	Φ	{1, 2}	{2}
2	Φ	Φ	{2}

Graful de tranziție:

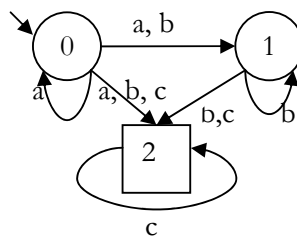


Figura 1.6

Definiția 1.2.6 Un automat $A = (Q, \Sigma, \delta, q_0, F)$ se numește:

1. *nedeterminist* (fără ε -tranzii), dacă $\delta(q, \varepsilon) = \Phi, \forall q \in Q$
2. *determinist*, dacă $\delta(q, \varepsilon) = \Phi, \forall q \in Q$ și $|\delta(q, a)| \leq 1, \forall q \in Q, \forall a \in \Sigma$

În literatura de specialitate, de obicei, un automat finit determinist are proprietatea că $|\delta(q, a)| = 1, \forall q \in Q, \forall a \in \Sigma$. Condiția $|\delta(q, a)| \leq 1$, dată în definiția precedentă, permite ca un automat determinist să fie parțial definit (eventual) dar nu este restrictivă relativ la puterea de calcul: putem obține automatul total definit dacă adăugăm la automat o stare specială Φ (numită *stare moartă* sau *stare de blocare*) și definim δ astfel:

$$\delta(q, a) = \Phi, \text{ dacă } |\delta(q, a)| < 1 \text{ și } \delta(\Phi, b) = \Phi, \forall b \in \Sigma.$$

Așadar, putem subînțelege, atunci când este nevoie, că este îndeplinită condiția $|\delta(q, a)| = 1$ (de exemplu când un automat trebuie să „citească” în întregime orice cuvânt, chiar dacă nu-l acceptă).

Pentru definirea *limbajului acceptat (recunoscut)* de un automat finit, să extindem funcția de tranziție la cuvinte. Vom nota mai întâi, pentru $q \in Q$:

$Cl(q) = \{q' \mid q' \in Q, \text{ în graful automatului } A \text{ există un drum de la } q \text{ la } q' \text{ de lungime } k \geq 0 \text{ ale cărui arce sunt etichetate cu } \varepsilon\}.$

Facem precizarea că $q \in Cl(q)$ pentru că q este legat de q printr-un drum de lungime 0. Dacă $S \subseteq Q$, atunci notăm:

$$Cl(S) = \bigcup_{q \in S} Cl(q) \text{ iar } \delta(S, a) = \bigcup_{q \in S} \delta(q, a).$$

Definiția 1.2.7 Dacă $A = (Q, \Sigma, \delta, q_0, F)$ este un automat, atunci extensia lui δ la cuvinte este funcția $\hat{\delta}: Q \times \Sigma^* \rightarrow 2^Q$ care se definește astfel:

1. $\hat{\delta}(q, \varepsilon) = Cl(q), \forall q \in Q$;
2. $\hat{\delta}(q, ua) = Cl(\delta(\hat{\delta}(q, u), a)), \forall q \in Q, \forall u \in \Sigma^*, \forall a \in \Sigma.$

Lema 1.2.2 Fie $A = (Q, \Sigma, \delta, q_0, F)$ un automat finit, o stare $q \in Q$ și un cuvânt $w \in \Sigma^*$, unde $w = a_1 a_2 \dots a_n$. Atunci $q' \in \hat{\delta}(q, w)$ dacă și numai dacă în graful automatului, există un drum de la q la q' de lungime $m \geq n$ cu arcele etichetate respectiv y_1, y_2, \dots, y_m (în această ordine) astfel încât $y_1 y_2 \dots y_m = w$.

Demonstrație.

1. \Leftarrow Se demonstrează prin inducție după lungimea lui w .
2. \Rightarrow Se demonstrează prin inducție după lungimea drumului de la q la q' .

□

Observația 1.2.2 Pentru automatele fără ε -tranzii (în care $\delta(q, \varepsilon) = \Phi, \forall q \in Q$), deoarece $Cl(q) = \{q\}$ are loc $\hat{\delta}(q, a) = \delta(q, a), \forall q \in Q, \forall a \in \Sigma$. Din acest motiv, pentru astfel de automate $\hat{\delta}$ va fi notată de asemenea cu δ .

În cazul automatelor cu ε -tranzii vom păstra notația $\hat{\delta}$ pentru extensie pentru că, în general, $\hat{\delta}(q, \varepsilon) \neq \delta(q, \varepsilon)$ și, de asemenea, $\hat{\delta}(q, a) \neq \delta(q, a), a \in \Sigma$. De pildă, în Exemplul 1.2.4 avem:

$$\hat{\delta}(0, a) = \{0, 1, 2\} \text{ iar } \delta(0, a) = \{0\} \text{ și } \hat{\delta}(0, \varepsilon) = \{0, 1, 2\} \text{ iar } \delta(0, \varepsilon) = \{1\}.$$

Lema 1.2.3 Fie $A = (Q, \Sigma, \delta, q_0, F)$ un automat finit și cuvintele $u, v \in \Sigma^*$. Atunci are loc:

$$\hat{\delta}(q, uv) = \hat{\delta}(\hat{\delta}(q, u), v).$$

Demonstrație. Inducție după $|v|$.

□

Definiția 1.2.8 *Limbajul acceptat (recunoscut)* de automatul $A = (Q, \Sigma, \delta, q_0, F)$ este mulțimea notată $L(A)$ dată de expresia:

$$L(A) = \{w \mid w \in \Sigma^*, \hat{\delta}(q_0, w) \cap F \neq \Phi\}.$$

Așadar, un cuvânt w este recunoscut de un automat A dacă, după *citirea* în întregime a cuvântului w , automatul (pornind din starea inițială q_0) poate să ajungă într-o stare finală. În cazul automatelor finite deterministe, putem scrie:

$$L(A) = \{w \mid w \in \Sigma^*, \delta(q_0, w) \in F\}.$$

deoarece $\delta(q_0, w)$ este o mulțime formată dintr-un singur element și identificăm această mulțime cu elementul respectiv.

Teorema 1.2.1 Pentru orice expresie regulată E peste Σ există un automat finit (cu ε -tranziții) A , astfel încât $L(E) = L(A)$.

Demonstrație Procedăm prin inducție asupra complexității expresiei E .

Dacă $E \in \{\Phi, \varepsilon, a\}$ ($a \in \Sigma$) atunci automatul corespunzător este respectiv:

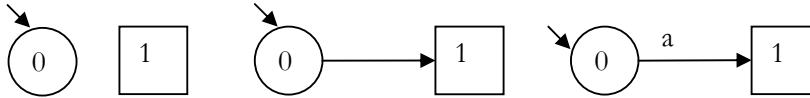


Figura 1.7

Dacă $E = E_1 \mid E_2$, $E = E_1 E_2$ sau $E = E_1^*$ atunci presupunem că există automatele finite $A_1 = (Q_1, \Sigma, q_{01}, \delta_1, \{f_1\})$, $A_2 = (Q_2, \Sigma, q_{02}, \delta_2, \{f_2\})$, cu Q_1, Q_2 disjuncte și $\delta_1(f_1, a) = \delta_2(f_2, a) = \Phi$, $\forall a \in \Sigma$ (ipoteza inductivă care are loc pentru $E = \Phi, \varepsilon$ sau a), care recunosc limbajele $L(E_1)$ respectiv $L(E_2)$. Să considerăm două stări noi q_0, f (ce nu sunt din $Q_1 \cup Q_2$). Atunci, automatul $A = (Q, \Sigma, \delta, q_0, F)$ se construiește astfel:

1. $E = E_1 \mid E_2$: $Q = Q_1 \cup Q_2 \cup \{q_0, f\}$, $F = \{f\}$, $\delta(q_0, \varepsilon) = \{q_{01}, q_{02}\}$,
 $\delta(f_1, \varepsilon) = \delta(f_2, \varepsilon) = \{f\}$, $\delta(q, a) = \delta_i(q, a) \forall q \in Q_i, \forall a \in \Sigma, i = 1, 2$.
2. $E = E_1 E_2$: $Q = Q_1 \cup Q_2$, $q_0 = q_{01}$, $F = \{f_2\}$, $\delta(f_1, \varepsilon) = q_{02}$,
 $\delta(q, a) = \delta_i(q, a) \forall q \in Q_i, \forall a \in \Sigma, i = 1, 2$.
3. $E = E_1^*$: $Q = Q_1 \cup \{q_0, f\}$, $\delta(q_0, \varepsilon) = \{q_{01}, f\}$, $\delta(f_1, \varepsilon) = \{q_{01}, f\}$,
 $\delta(q, a) = \delta_1(q, a) \forall q \in Q_1, \forall a \in \Sigma$.

Schematic, acest lucru se exprimă ca în figura 1.8.

Se dovedește ușor că automatele construite mai sus recunosc respectiv limbajele $L(E_1 \mid E_2)$, $L(E_1 E_2)$, $L(E_1^*)$.

□

Clasa limbajelor recunoscute de automatele finite coincide cu clasa limbajelor descrise de expresii regulate și această clasă este cunoscută sub numele de clasa limbajelor regulate (sau de tip 3). Această clasă conține numai limbaje neambigue: un automat finit determinist este un mecanism neambiguu. De aici rezultă că pentru orice expresie regulată E există o expresie regulată E' neambiguă, echivalentă cu E . Teoretic vorbind, trecerea de la o expresie ambiguă E la echivalenta sa neambiguă E' se face astfel: se construiește A , automatul finit determinist echivalent cu E (trecând prin automatul cu ε -tranziții și apoi prin cel nedeterminist fără ε -tranziții) pentru ca apoi să se determine expresia regulată E' care descrie limbajul $L(A)$; aceasta este neambiguă. Practic însă, descrierile ambigue pot să fie mai concise și mai ușor de înțeles decât echivalentele lor neambigue (care costă și timp substanțial pentru a fi obținute).

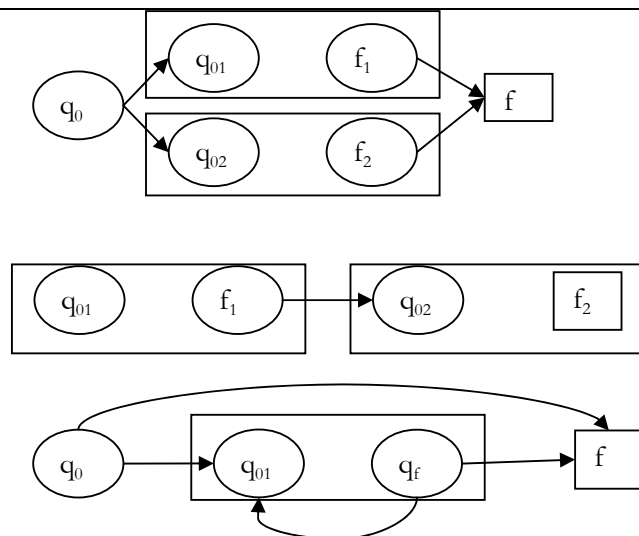


Figura 1.8

Exemplul 1.2.7 Limbajul $L = \{w \in \{0,1\}^* \mid w \text{ conține } 000 \text{ sau } 111\}$ este descris de expresia ambiguă $E = (0 \mid 1)^*(000 \mid 111)(0 \mid 1)^*$. Expresia E' neambiguă, echivalentă cu E , în afară de faptul că se obține greoi, nu mai exprimă în mod evident forma cuvintelor din $L(E')$. Cititorul se poate convinge de acest lucru printr-un exercițiu.

1.2.3 Gramatici independente de context

Definiția 1.2.9 O gramatică independentă de context este sistemul $G = (V, T, S, P)$ unde V și T sunt mulțimi nevide, finite, disjuncte de *neterminali* (*variabile*), respectiv *terminali*, $\Sigma = V \cup T$ se numește *vocabularul* gramaticii, $S \in V$ este o variabilă specială numită *simbol de start* sau *axioma gramaticii*, iar $P \subseteq V \times (V \cup T)^*$ este *mulțimea regulilor de producție* (*derivare*) sau simplu *mulțimea producțiilor*.

Convenții de notație. În scopul simplificării expunerii, vom face câteva convenții care vor fi păstrate de-a lungul acestei lucrări, cu excepția situațiilor în care se fac precizări suplimentare.

1. Următoarele simboluri notează *neterminali*:
 - S , folosit pentru a nota *axioma* gramaticii;
 - $A, B, C \dots$, (literele majuscule de la începutul alfabetului);
 - numele italicizate scrise cu minuscule: *expresie*, *instrucțiune*,...
2. Următoarele simboluri notează *terminali*:
 - literele mici de la începutul alfabetului: a, b, c, \dots
 - operatori: $+, -, *, /$, etc.
 - simboluri de punctuație, paranteze.
 - cifrele: $0, 1, \dots, 9$
 - numele scrise cu fontul courier (unitățile lexicale): *id*, *if*, *while*, *begin*, etc..
3. Literele mari de la sfârșitul alfabetului, X, Y, Z, \dots notează simboluri gramaticale (neterminali sau terminali, adică elementele din Σ);
4. Literele mici de la sfârșitul alfabetului, u, v, x, y, z, w notează cuvinte din T^* (formate din terminali);

5. Literele grecești $\alpha, \beta, \gamma, \dots$ reprezintă cuvinte peste Σ . Așadar, o producție a gramaticii se va nota prin $A \rightarrow \alpha$ (în loc de $(A, \alpha) \in P$, convenim să scriem $A \rightarrow \alpha$).
6. Dacă $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_n$ sunt toate producțiile care au A în partea stângă (numite A -producții), le vom nota prin:

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$
7. O gramatică va fi dată prin producțiile sale. De aici se subînțeleg mulțimile V și T , iar simbolul de start este partea stângă a primei producții.

Exemplul 1.2.8 Gramatica:

$$E \rightarrow EOE \mid (E) \mid -E \mid id$$

$$O \rightarrow + \mid - \mid * \mid /$$

este constituită din elementele $V = \{E, O\}$, $T = \{id, +, -, *, /, (,)\}$, E este simbolul de start, iar producțiile sunt cele indicate mai sus.

Definiția 1.2.10 Fie $G = (V, T, S, P)$ o gramatică independentă de context. Relația de *derivare* în gramatica G este o relație din $\Sigma^*V\Sigma^* \times \Sigma^*$, notată \Rightarrow_G , și este definită astfel:

$$\delta_1 \Rightarrow_G \delta_2 \text{ dacă și numai dacă } \delta_1 = \alpha A \gamma, \delta_2 = \alpha \beta \gamma \text{ și } A \rightarrow \beta.$$

Vom nota prin \Rightarrow_G^* (\Rightarrow_G^+) închiderea reflexivă și tranzitivă (închiderea tranzitivă) a relației \Rightarrow_G . Atunci când nu există confuzii, vom renunța la a indica G în aceste relații (scriem simplu $\Rightarrow, \Rightarrow^*, \Rightarrow^+$). Dacă $\alpha \Rightarrow \beta$ spunem că α *derivatează* într-un pas β ; dacă $\alpha \Rightarrow_G^* \beta$ ($\alpha \Rightarrow_G^+ \beta$) spunem că α *derivatează* β (α *derivatează* propriu β):

$$\alpha \Rightarrow_G^* \beta \text{ dacă și numai dacă există } \delta_0, \delta_1, \dots, \delta_n, n \geq 0 \text{ astfel ca:}$$

$$\alpha = \delta_0 \Rightarrow \delta_1 \Rightarrow \dots \Rightarrow \delta_n = \beta$$

$$\alpha \Rightarrow_G^+ \beta \text{ dacă și numai dacă există } \delta_0, \delta_1, \dots, \delta_n, n \geq 1 \text{ astfel ca:}$$

$$\alpha = \delta_0 \Rightarrow \delta_1 \Rightarrow \dots \Rightarrow \delta_n = \beta$$

Dacă $S \Rightarrow_G^* \gamma$ spunem că γ este *formă propozițională* în gramatica G , iar dacă γ este format din terminali, $\gamma = w$, spunem că w este o frază în gramatica G .

Definiția 1.2.11 *Limbajul generat* de gramatica G (notat $L(G)$) este mulțimea frazelor gramaticii G , $L(G) = \{w \mid w \in T^*, S \Rightarrow_G^+ w\}$.

Pentru că în această lucrare ne vom ocupa numai de gramatici independente de context, le vom numi pe acestea simplu gramatici. Pentru studiul ierarhiei lui Chomsky, cititorul este invitat a consulta lucrările [Gri86], [Juc99]. De asemenea, vom considera că gramaticile cu care lucrăm în continuare sunt gramatici reduse, adică orice simbol $X \in \Sigma$ este *simbol util*:

- este *accesibil*: $S \Rightarrow^+ \alpha X \beta$, pentru anumiți $\alpha, \beta \in \Sigma^*$;
- este *productiv*: $X \Rightarrow^* w$, pentru un anume $w \in T^*$.

Este cunoscut faptul că pentru orice gramatică G există o gramatică G' , cu $L(G) = L(G')$ (G și G' se zic echivalente în acest caz), astfel încât orice simbol al vocabularului gramaticii G' este util (vezi [Gri86], [Juc99], [JuA02]). O gramatică fără simboluri inutile (simboluri care nu sunt utile) se numește *gramatică redusă*. În continuare vom considera numai gramatici reduse.

1.2.4 Arbori sintactici. Ambiguitate

Un *arbore sintactic* ilustrează într-o gramatică modul în care axioma generează o (se rescrie într-o) frază a gramaticii. Ideea care stă la baza atașării unui arbore unei derivări este aceea că, dacă într-o derivare se folosește o producție de forma $A \rightarrow X_1 X_2 \dots X_n$, atunci în arbore există un nod interior etichetat A care are n descendenți etichetați respectiv cu X_1, X_2, \dots, X_n de la stânga la dreapta.

Definiția 1.2.12 Fie $G = (V, T, S, P)$ o gramatică. Un *arbore sintactic* (*arbore de derivare*, *arbore de parsare*) în gramatica G este un arbore ordonat, etichetat, cu următoarele proprietăți:

1. rădăcina arborelui este etichetată cu S ;
2. fiecare frunză este etichetată cu un simbol din T sau cu ε ;
3. fiecare nod interior este etichetat cu un neterminat;
4. dacă A etichetează un nod interior care are n succesori etichetați de la stânga la dreapta respectiv cu X_1, X_2, \dots, X_n , atunci $A \rightarrow X_1X_2 \dots X_n$ este o producție. Cazul în care producția este $A \rightarrow \varepsilon$ este un caz special: nodul etichetat A are un singur descendent etichetat ε .

Frontiera unui arbore de derivare este cuvântul $w = a_1a_2\dots a_n$ unde $a_i, 1 \leq i \leq n$ sunt etichetele nodurilor frunză în ordinea de la stânga la dreapta.

Exemplul 1.2.10 Să considerăm gramatica care descrie expresiile aritmetice ce se construiesc cu $+$, $*$, id și paranteze:

$$E \rightarrow E+E \mid E*E \mid (E) \mid id.$$

Un arbore de derivare în această gramatică este prezentat în figura 1.9.

Să observăm că acestui arbore de derivare îi putem atașa derivările:

$$E \Rightarrow E+E \Rightarrow id+E \Rightarrow id+E*E \Rightarrow id+id*E \Rightarrow id+id*id$$

$$E \Rightarrow E+E \Rightarrow E+E*E \Rightarrow E+E*id \Rightarrow E+id*id \Rightarrow id+id*id$$

$$E \Rightarrow E+E \Rightarrow E+E*E \Rightarrow E+id*E \Rightarrow id+id*E \Rightarrow id+id*id$$

Este clar că cele trei derivări (mai pot fi scrise și altele) sunt distincte. Ele diferă prin ordinea de aplicare a regulilor. În prima derivare, la fiecare pas s-a rescris cea mai din stânga variabilă a formei propoziționale (curente); astfel de derivări se numesc *derivări*

extrem stângi. Vom nota o derivare extrem stângă prin \Rightarrow_{st} :

$$uA\alpha \Rightarrow_{st} u\beta\alpha \text{ dacă } u \in T^*, A \rightarrow \beta \in P, \alpha \in \Sigma^*$$

Așadar, prima din derivările de mai sus se scrie: $E \xRightarrow[st]{+} id+id*id$.

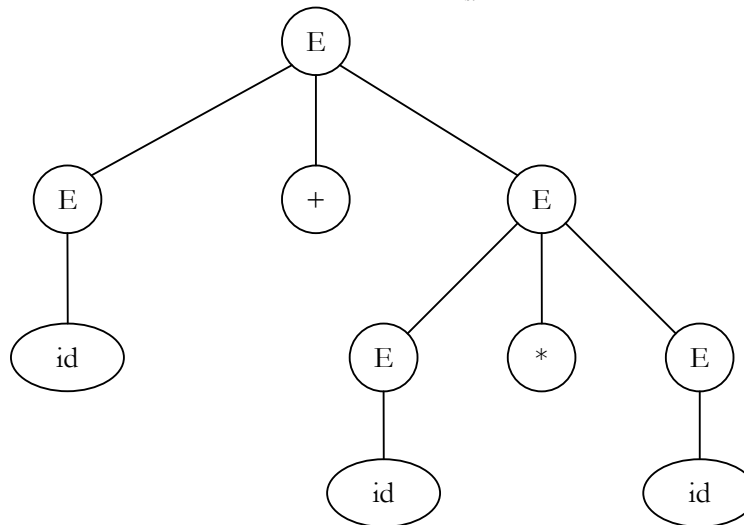


Figura 1.9

Cea de-a doua derivare are proprietatea că, la fiecare pas, s-a rescris cea mai din dreapta variabilă din forma propozițională. Vom numi astfel de derivări *derivări extrem drepte* și le vom nota prin \Rightarrow_{dr} :

$$\alpha Au \Rightarrow_{dr} \alpha \beta u \text{ dacă } u \in T^*, A \rightarrow \beta \in P, \alpha \in \Sigma^*$$

Prin urmare cea de-a doua derivare se scrie: $E \Rightarrow_{dr}^+ id+id*id$.

Așadar, arborele sintactic este o reprezentare grafică a unei derivări în care a dispărut alegerea ordinii în care se aplică regulile de producție. Unui arbore sintactic îi corespunde o singură derivare extrem stângă (respectiv o singură derivare extrem dreaptă). Se demonstrează relativ ușor următoarea leamnă (vezi [Gri86] și [Juc99]):

Lema 1.2.4. Fie G o gramatică și $w \in T^*$. Atunci, există o derivare $S \Rightarrow^+ w$ (sau echivalent $w \in L(G)$) dacă și numai dacă există un arbore de derivare cu frontiera w . Mai mult, dacă $A \in V$, $\gamma \in \Sigma^*$ are loc: $A \Rightarrow^+ \gamma$ dacă și numai dacă există un arbore (construit după aceleași reguli ca arborele de derivare) în care rădăcina este etichetată cu A , iar frontiera este γ .

□

Să observăm acum că, pentru cuvântul $w = id+id*id$ (gramatica din exemplul precedent), se mai poate construi un arbore de derivare distinct de cel de mai sus, anume cel din figura 1.10. Acestui arbore îi corespunde derivarea extrem stângă:

$$S \Rightarrow E * E \Rightarrow E + E * E \Rightarrow id + E * E \Rightarrow id + id * E \Rightarrow id + id * id$$

care diferă de derivarea extrem stângă precedentă. Cele două derivări extrem stângi (respectiv cei doi arbori de derivare) corespund aceluiași cuvânt din $L(G)$. Fenomenul este cunoscut sub numele de *ambiguitate*.

Definiția 1.2.13 O gramatică $G = (V, T, S, P)$ se zice că este *ambiguă* dacă există $w \in L(G)$ pentru care se pot construi măcar două derivări extrem stângi distincte (respectiv doi arbori de derivare distincți).

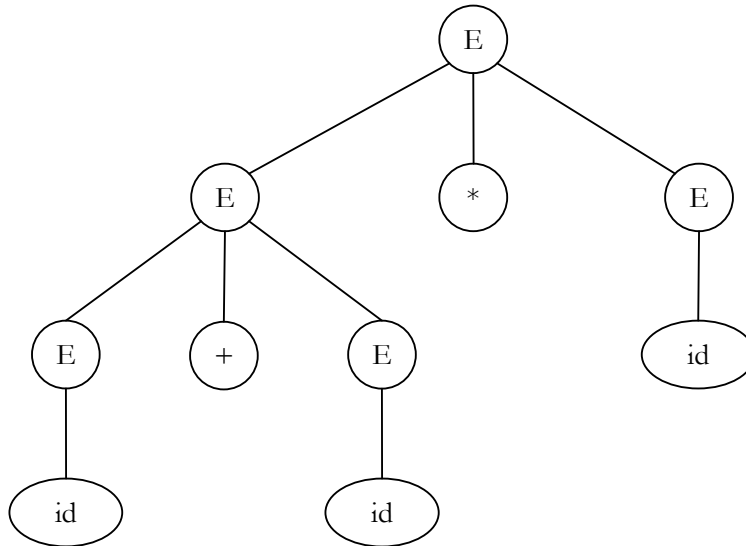


Figura 1.10

Exemplul 1.2.10 Gramatica următoare ce descrie construcțiile if-then și if-then-else este ambiguă:

$$I \rightarrow \text{if } e \text{ then } I \mid \text{if } e \text{ then } I \text{ else } I \mid a$$

În adevăr, pentru fraza:

$$w = \text{if } e \text{ then if } e \text{ then } a \text{ else } a$$

putem construi două derivări extrem stângi distincte:

$$\begin{aligned} I \Rightarrow \text{if } e \text{ then } I &\Rightarrow \text{if } e \text{ then if } e \text{ then } I \text{ else } I \Rightarrow \\ &\Rightarrow \text{if } e \text{ then if } e \text{ then } a \text{ else } I \Rightarrow \\ &\Rightarrow \text{if } e \text{ then if } e \text{ then } a \text{ else } a \end{aligned}$$

$$\begin{aligned} I \Rightarrow \text{if } e \text{ then } I \text{ else } I &\Rightarrow \text{if } e \text{ then if } e \text{ then } I \text{ else } I \Rightarrow \\ &\Rightarrow \text{if } e \text{ then if } e \text{ then } a \text{ else } I \Rightarrow \\ &\Rightarrow \text{if } e \text{ then if } e \text{ then } a \text{ else } a \end{aligned}$$

care corespund, evident, la arbori de derivare distincți.

Este de dorit ca, atunci când abordăm analiza sintactică, să o facem pentru o specificare neambiguă. Uneori, o gramatică ambiguă se poate transforma într-o gramatică echivalentă neambiguă. Această transformare este legată de anumite reguli care se impun. De pildă, în gramatica ambiguă:

$$E \rightarrow E+E \mid E * E \mid (E) \mid \text{id}$$

nu este specificată nici o ordine de prioritate între operatorii + și *. Dacă impunem ca * să aibă prioritate asupra lui +, atunci o nouă gramatică se poate scrie pentru descrierea expresiilor aritmetice:

$$E \rightarrow E+T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid \text{id}$$

În această nouă gramatică, noii neterminali au semnificația naturală: T notează ceea ce se cheamă termen, F este notația pentru factor, încât putem spune că de data aceasta gramatica exprimă mai fidel ceea ce se numește expresie. Dacă încercăm să impunem reguli similare pentru cazul if-then-else, putem ajunge la o exprimare de forma:

$$I \rightarrow I_1 \mid I_2$$

$$I_1 \rightarrow \text{if } e \text{ then } I_1 \text{ else } I_1 \mid a$$

$$I_2 \rightarrow \text{if } e \text{ then } I \mid \text{if } e \text{ then } I_1 \text{ else } I_2$$

Prin aceasta s-a exprimat regula: se asociază else celui then precedent, cel mai apropiat care nu are corespondent.

În capitolele următoare vom construi o serie de algoritmi. Pentru descrierea acestora vom folosi un limbaj(algoritm) bazat pe limbajul C: vom utiliza operatori specifici limbajului C ca ++(incrementare), ==(egalitate logică), !=(neegalitate logică), !(negație) blocurile vor fi specificate prin acolade, instrucțiunile while, for, if care vor fi folosite funcționează ca în C.

2 Analiza lexicală

Rolul pe care îl are un analizor lexical este de a citi textul sursă, caracter cu caracter, și a-l transforma într-o secvență de unități primitive (elementare) numite unități lexicale (tokens în limba engleză). Fiecare unitate lexicală descrie o secvență de caractere cu o anumită semnificație și este tratată ca o entitate logică. Pentru fiecare limbaj de programare se stabilesc (atunci când se proiectează limbajul) unitățile lexicale. În majoritatea limbajelor se disting următoarele unități lexicale:

Nr. Crt.	Unitatea Lexicală	Exemple
1	Constanta	573 -5.89 2e+3
2	Identificator	alpha un_ident
3	Operator	+ * / < >
4	Cuvânt rezervat	begin while
5	Semn special	; . :

În acest capitol vom discuta tehnicile de analiză lexicală, proiectarea și implementarea unui analizor lexical. Problema analizei lexicale prezintă cel puțin două aspecte. Mai întâi, trebuie găsită o modalitate de descriere a unităților lexicale. Se constată că expresiile regulate – modelul algebric de descriere a limbajelor regulate - sunt mecanismele care pot descrie orice unitate lexicală. Un al doilea aspect este cel al recunoașterii acestor unități lexicale (analiza lexicală propriu-zisă). Dacă descrierea se face cu expresii regulate atunci mecanismul de recunoaștere este automatul finit determinist. Se știe din teoria limbajelor formale că mulțimea limbajelor descrise de expresii regulate coincide cu cea a limbajelor recunoscute de automate finite (clasa limbajelor regulate). Vom descrie în acest capitol algoritmi (eficienți) de construcție a unui automat finit (determinist) echivalent cu o expresie regulată și, apoi, implementarea unui automat finit (analizorul lexical).

Unitățile lexicale sunt de două categorii:

- unități care descriu un șir de caractere anume (de exemplu `if`, `while`, `++`, `:=`, `;`);
- unități care descriu o clasă de șiruri: (identificatori, constante etc.).

În cel din urmă caz, vom trata o unitate lexicală ca fiind o pereche formată din *tipul* și *valoarea unității lexicale*. Pentru unitățile lexicale care descriu un șir anume, convenim că tipul este acel șir, iar valoarea coincide cu tipul. De exemplu caracterul „(” este de tip paranteză stângă iar `alpha` este unitatea lexicală de tip *identificator* care are valoarea `alpha`. În literatura de specialitate `alpha` se mai numește instanță a tokenului identificator sau *lexem*.

Vom descrie unitățile lexicale prin expresii regulate. De pildă, o constantă întreagă fără semn este descrisă de expresia:

$$(1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)(0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9)^*$$

Pentru comoditatea notațiilor, se pot da nume unor expresii regulate iar cu ajutorul acestora pot fi definite alte expresii regulate. Fie Σ un alfabet, d_1, d_2, \dots, d_n nume distincte date unor expresii regulate iar E_1, E_2, \dots, E_n expresii regulate astfel încât E_1 este expresie regulată peste Σ iar E_i ($1 < i \leq n$) este expresie regulată peste alfabetul $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$. Atunci o definiție regulată este un șir de definiții de forma:

$$d_1 \rightarrow E_1$$

$$d_2 \rightarrow E_2$$

...

$$d_n \rightarrow E_n$$

De exemplu, mulțimea identificatorilor PASCAL poate fi dată prin definiția regulată următoare (din comoditate s-a scris ... dar trebuie să înțelegem caracterele corespunzătoare):

$$\langle \text{literă} \rangle \rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z$$

$$\langle \text{cifră} \rangle \rightarrow 0 \mid 1 \mid \dots \mid 9$$

$$\langle \text{identificator} \rangle \rightarrow \langle \text{literă} \rangle (\langle \text{literă} \rangle \mid \langle \text{cifră} \rangle)^*$$

Constantele fără semn în PASCAL se definesc astfel:

$$\langle \text{cifră} \rangle \rightarrow 0 \mid 1 \mid \dots \mid 9$$

$$\langle \text{cifre} \rangle \rightarrow \langle \text{cifră} \rangle \langle \text{cifră} \rangle^*$$

$$\langle \text{fracție} \rangle \rightarrow .(\langle \text{cifre} \rangle \mid \epsilon)$$

$$\langle \text{exponent} \rangle \rightarrow (E \mid e)(+ \mid - \mid \epsilon) \langle \text{cifre} \rangle$$

$$\langle \text{număr} \rangle \rightarrow \langle \text{cifre} \rangle \mid \langle \text{cifre} \rangle \langle \text{fracție} \rangle \langle \text{exponent} \rangle$$

Operatorul unar $*$ semnifică repetarea de zero sau mai multe ori a unei instanțe. Dacă folosim operatorul unar $+$ pentru a desemna operația „cel puțin o instanță din”, adică $E^+ = EE^*$, atunci definiția numelui $\langle \text{cifre} \rangle$ se poate scrie:

$$\langle \text{cifre} \rangle \rightarrow \langle \text{cifră} \rangle^+$$

2.1 Un algoritm de trecere de la o expresie regulată la automatul echivalent

Să transpunem construcțiile din teorema 1.2.1 într-un algoritm pentru transformarea unei expresii regulate în automat finit. Mai întâi să observăm că fiecare din aparițiile operatorilor $|$ și $*$ dintr-o expresie regulată E introduce două noi stări în automatul construit, pe când operatorul \bullet nu introduce alte stări (figura 1.8). De asemenea, pentru orice apariție a unui simbol din Σ , cât și pentru ϵ , dacă acesta apare explicit în E , este nevoie de 2 stări în automatul construit. Așadar, dacă n este numărul de simboluri din E iar m este numărul de paranteze împreună cu aparițiile simbolului \bullet , atunci numărul stărilor automatului echivalent cu E este $2(n - m)$.

Să mai observăm că, din orice stare a automatului, se fac cel mult două tranziții: fie o tranziție cu un simbol din Σ , fie una sau două ϵ - tranziții, fie zero tranziții. Atunci, reprezentarea automatului echivalent cu o expresie regulată se poate face cu un tablou de dimensiune $p \times 3$ unde p este numărul stărilor (acestea sunt numerotate de la 1 la p), prima coloană conține simbolurile cu care se fac tranzițiile iar următoarele două coloane conțin stările rezultate în urma tranzițiilor. Pentru descrierea algoritmului o să identificăm cele trei coloane prin vectorii simbol, next1, next2.

Algoritmul pe care-l descriem mai jos este datorat lui Rytter (Ryt91).

Algoritmul 2.1.1

Intrare: Expresia regulată E cu n simboluri dintre care m sunt paranteze și apariții ale operatorului produs;

Ieșire: Vectorii simbol, next1, next2 de dimensiune $p = 2(n - m)$ ce descriu automatul cu ϵ - tranziții echivalent cu E ;

Metoda:

1. Se construiește arborele atașat expresiei E (o metodă se dă mai jos);

2. Se parcurge arborele în preordine și se atașează nodurilor vizitate, exceptând pe cele etichetate cu \bullet , respectiv numerele 1, 2, ..., n-m;
3. Se parcurge arborele în postordine și se atașează fiecărui nod N o pereche de numere (i, f) care reprezintă starea inițială respectiv finală a automatului corespunzător subarborelui cu rădăcina N, astfel:
 - 3.1. Dacă nodul are numărul k (de la pasul 2) atunci $N.i = 2k-1$, $N.f = 2k$;
 - 3.2. Dacă nodul este etichetat \bullet atunci $N.i = S.i$ iar $N.f = D.f$ (S și D sunt fii lui N, stâng respectiv drept);
4. for(1 ≤ j ≤ 2(n - m)) {simbol[j] = ' ', next1[j] = next2[j] = 0}
5. Se parcurge din nou arborele obținut. Dacă N este nodul curent iar S și D sunt fii săi, atunci, în funcție de eticheta lui N, se execută următoarele:
 - 5.1. Dacă N este etichetat cu | :
 $next1[N.i] = S.i$, $next2[N.i] = D.i$, $next1[S.f] = N.f$, $next1[D.f] = N.f$
 - 5.2. Dacă N este etichetat cu \bullet atunci $next1[S.f] = D.i$
 - 5.3. Dacă N este etichetat cu * (D nu există în acest caz):
 $next1[N.i] = S.i$, $next2[N.i] = N.f$, $next1[S.f] = S.i$, $next2[S.f] = N.f$
 - 5.4. Dacă N este etichetat cu a (deci este frunză):
 $simbol[N.i] = 'a'$, $next1[N.i] = N.f$

Construcția arborelui:

Intrare: Expresia regulată $E = e_0e_1\dots e_{n-1}$
 Precedența operatorilor: $prec(|) = 1$, $prec(\bullet) = 2$, $prec(*) = 3$.

Ieșire: Arborele asociat t.

Metoda: Se consideră două stive: STIVA1 – stiva operatorilor, STIVA2 – stiva operanzilor (care va conține arborii parțiali construiți).

```

1. i = 0;
2. while(i < n) {
3.   c = ei;
4.   switch(c) {
5.     case '(' : {STIVA1.push(c); break;}
6.     case operand : {STIVA2.push(c); break;}
7.     case ')' : {
8.       do{build_tree();}while(STIVA1.top() != '(');
9.       STIVA1.pop(); break;
10.    } //endcase
11.    case operator: {
12.      while(prec(STIVA1.top()) >= prec(c)) build_tree();
13.      STIVA1.push(c); break
14.    } //endcase
15.  } //endswitch
16. } //endwhile
17. while(STIVA1 !=  $\Phi$ ) build_tree();
18. t = STIVA2.pop();
19. build_tree() {
20.   op = STIVA1.pop();
21.   D = STIVA2.pop();
22.   switch(op) {
23.     case '*': {
24.       t = tree(op, D, NULL);
25.       STIVA2.push(t); break;
26.     }
27.     case '|': case '•': {
28.       S = STIVA2.pop();
29.       t = tree(op, S, D);
30.       STIVA2.push(t); break;
31.     }
32.   }
33. }

```

Exemplul 2.1.1 Să considerăm expresia $E = a*b \mid bb(a \mid c)^*$.

1. Arborele expresiei este prezentat în figura 2.1.
2. Nodurilor arborelui le sunt atașate respectiv numerele 1, 2, ..., 10 prin parcurgerea în preordine și exceptând nodurile produs.
3. Numărul stărilor automatului va fi $p = 20$. După parcurgerea în postordine a arborelui, fiecare nod are atașat o pereche (i, j) (figura 2.2).
4. Se inițializează vectorii simbol, next1 și next2.
5. Tabela de tranziție a automatului, în urma aplicării procedurilor descrise la 5 în fiecare nod al arborelui, este dată în continuare.

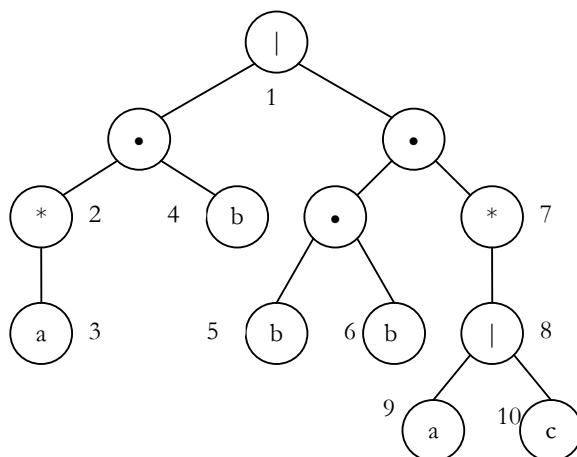


Figura 2.1

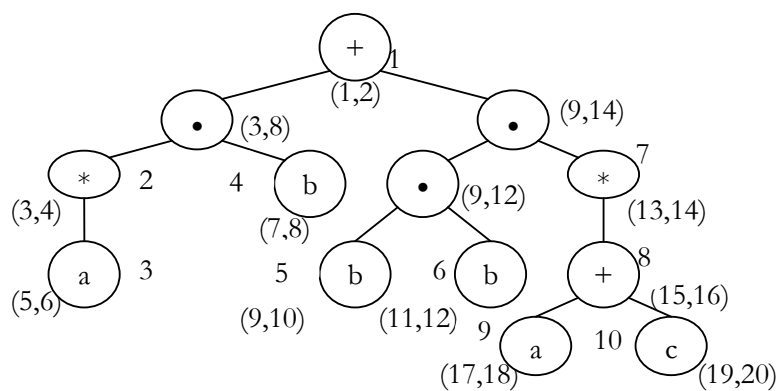


Figura 2.2

Tabela de tranziție a automatului:

p	simbol[p]	next1[p]	next2[p]
1		3	9
2		0	0
3		5	4
4		7	0
5	a	6	0
6		5	4
7	b	8	0
8		2	0
9	b	10	0

10		11	0
11	b	12	0
12		13	0
13		15	14
14		2	0
15		17	19
16		15	14
17	a	18	0
18		16	0
19	c	20	0
20		16	0

Teorema 2.1.1 Algoritmul 2.1.1 este corect: automatul cu ε - tranziții obținut este echivalent cu expresia regulată E.

Demonstrație Să observăm că modul în care au fost alese perechile (i, f) de stări pentru fiecare nod al arborelui construit corespunde construcțiilor din teorema 1.2.1. De asemenea tranzițiile care se definesc în pasul 5 al algoritmului urmăresc construcția din teorema amintită. Așadar, automatul obținut este echivalent cu expresia dată la intrare.

□

2.2 De la un automat cu ε -tranziții la automatul determinist echivalent

În paragraful precedent am descris procedeul de trecere de la o expresie regulată la automatul finit echivalent. Programul care descrie activitatea acestui automat finit pentru un cuvânt de intrare este de fapt analizorul lexical. Pentru ca acest analizor să fie eficient este de dorit ca automatul obținut să fie determinist. Cum prin procedeul descris în teorema precedentă se obține un automat cu ε - tranziții (nedeterminist deci), să indicăm o modalitate de trecere de la un automat cu ε - tranziții la automatul determinist echivalent. Mai întâi să descriem un algoritm pentru calculul mulțimii $Cl(S) = \hat{\delta}(S, \varepsilon)$.

Algoritmul 2.2.1

Intrare: Automatul (cu ε - tranziții) $A = (Q, \Sigma, \delta, q_0, F)$, $S \subseteq Q$.

Ieșire: $Cl(S) = \hat{\delta}(S, \varepsilon)$.

Metoda: Stiva STIVA se inițializează cu S și pentru fiecare q din top, stările din $\delta(q, \varepsilon)$ ce nu au fost puse în $Cl(S)$ se adaugă în $Cl(S)$ și în stivă.

```

1. R = S; STIVA = S;
2. while (STIVA  $\neq \Phi$ ) {
3.   q = STIVA.pop(); // Se extrage q din stivă
4.   T =  $\delta(q, \varepsilon)$  ;
5.   if (T  $\neq \Phi$ ) {
6.     for ( p  $\in$  T ) {
7.       if ( p  $\notin$  R ) {
8.         R = R  $\cup$  {p} ;
9.         STIVA.push(p); //Se adaugă p in stiva
        } //endif
      } //endfor
    } //endif
  } //endwhile
10. Cl(S) = R;
```

Algoritmul 2.2.2 Transformarea unui automat cu ε - tranziții în automat finit determinist.

Intrare: Automatul (cu ε - tranziții) $A = (Q, \Sigma, \delta, q_0, F)$.

Procedura de calcul pentru $Cl(S)$ (Algoritmul 2.2.1).

Ieșire: Automatul determinist $A' = (Q', \Sigma, \delta', q_0', F')$, echivalent cu A .

Metoda: Se construiesc stările lui A' începând cu q_0' și continuând cu cele accesibile prin tranziții cu simboluri din alfabet.

```

1.  $q_0' = Cl(q_0)$ ;  $Q' = \{q_0'\}$  ;
2.  $marcat(q_0') = false$ ;  $F' = \Phi$  ;
3. if (  $q_0 \cap F \neq \Phi$  ) then  $F' = F' \cup \{q_0'\}$  ;
4. while (  $\exists q' \in Q' \ \&\& \ !marcat(q')$  ) { //  $q'$  este nemarcat
5.   for (  $a \in \Sigma$  ) {
6.      $p' = Cl(\delta(q', a))$ ; //  $= \hat{\delta}(q', a)$ 
7.     if (  $p' \neq \Phi$  ) {
8.       if (  $p' \notin Q'$  ) {
9.          $Q' = Q' \cup \{p'\}$ ;
10.         $marcat(p') = false$ ;
11.         $\delta'(q', a) = p'$  ;
12.        if (  $p' \cap F \neq \Phi$  ) then  $F' = F' \cup \{p'\}$ ;
13.      } //endif
14.    } //endif
15.  } //endfor
16.   $marcat(q') = true$ ;
17. } //endwhile

```

2.3 Proiectarea unui analizor lexical

Definiția 2.4.1 Fie Σ un alfabet (al unui limbaj de programare). O *descriere lexicală* peste Σ este o expresie regulată $E = (E_1 \mid E_2 \mid \dots \mid E_n)^+$, unde n este numărul unităților lexicale, iar E_i descrie o unitate lexicală, $1 \leq i \leq n$.

Exemplul 2.4.1 Să considerăm $\Sigma = \{a, b, \dots, y, z, 0, 1, \dots, 9, :, =\}$ și următoarele expresii regulate:

Litera $\rightarrow a \mid b \mid c \mid \dots \mid z$
 Cifra $\rightarrow 0 \mid 1 \mid \dots \mid 9$
 Id \rightarrow Litera (Litera \mid Cifra)^{*}
 Intreg \rightarrow Cifra⁺
 Asignare $\rightarrow :=$
 Egal $\rightarrow =$
 Douapuncte $\rightarrow :$

O descriere lexicală care cuprinde aceste unități lexicale este:

$E = (Id \mid Intreg \mid Asignare \mid Egal \mid Douapuncte)^+$

Definiția 2.4.2 Fie E o descriere lexicală peste Σ ce conține n unități lexicale și $w \in \Sigma^+$. Cuvântul w este *corect relativ la descrierea* E dacă $w \in L(E)$. O *interpretare* a cuvântului $w \in L(E)$ este o secvență de perechi $(u_1, k_1), (u_2, k_2), \dots, (u_m, k_m)$, unde $w = u_1 u_2 \dots u_m$, $u_i \in L(E_{k_i})$, $1 \leq i \leq m$, $1 \leq k_i \leq n$.

Uneori, în loc de k_i în (u_i, k_i) , punem chiar E_{k_i} , asta însemnând faptul că u_i este o unitate de tipul E_{k_i} .

Exemplul 2.4.2 Fie E din Exemplul 2.4.1 și $w = \text{alpha} := \text{beta} = 542$. O interpretare a cuvântului w este:

$(\text{alpha}, \text{Id}), (:=, \text{Asignare}), (\text{beta}, \text{Id}), (=, \text{Egal}), (542, \text{Intreg})$

Sigur că aceasta nu este singura interpretare. Următoarea secvență este de asemenea o interpretare a aceluiași cuvânt: $(\text{alpha}, \text{Id}), (:=, \text{Douapuncte}), (=, \text{Egal}), (\text{b}, \text{Id}), (\text{eta}, \text{Id}), (=, \text{Egal}), (5, \text{Intreg}), (42, \text{Intreg})$.

Faptul că am obținut interpretări diferite pentru același cuvânt este consecință a ambiguității expresiei regulate E : unitatea lexicală Asignare are două descrieri distincte în E : Asignare și produsul DouapuncteEgal . La fel se întâmplă cu Id și cu Intreg . Din teorie se știe că, pentru orice expresie regulată E , există o expresie regulată E' neambiguă, echivalentă cu E . Soluția de a construi E' pentru a obține interpretări unice pentru fiecare cuvânt nu este convenabilă pentru că, pe de o parte, transformarea E în E' consumă timp iar, pe de altă parte, E' nu descrie în mod natural limbajul în cauză. În toate cazurile, se preferă să se dea niște reguli în plus pentru a putea alege o singură interpretare a unui cuvânt atunci când descrierea lexicală este ambiguă.

Definiția 2.4.3 Fie E o descriere lexicală peste Σ și $w \in L(E)$. O interpretare a cuvântului $w, (u_1, k_1)(u_2, k_2), \dots, (u_m, k_m)$, este *interpretare drept-orientată* dacă $(\forall i) 1 \leq i \leq m$, are loc:

$$|u_i| = \max \{ |v| \mid v \in L(E_1 \mid E_2 \mid \dots \mid E_n) \cap \text{Pref}(u_1 u_{i+1} \dots u_m) \}.$$

(unde prin $\text{Pref}(w)$ am notat mulțimea prefixelor cuvântului w).

Observația 2.4.1 Există descrieri lexicale E în care nu orice cuvânt din $L(E)$ admite o interpretare drept-orientată. Fie descrierea $E = (a \mid ab \mid bc)^+$ peste alfabetul $\Sigma = \{a, b\}$ și $w = abc$. Singura interpretare a cuvântului w este $(a, 1), (bc, 3)$ dar aceasta nu este drept orientată pentru că a nu este cel mai lung cuvânt din $L(a \mid ab \mid bc) \cap \text{Pref}(abc)$, acesta din urmă fiind ab .

Definiția 2.4.4 O descriere lexicală E este *bine-formată* dacă orice cuvânt w din limbajul $L(E)$ are exact o interpretare drept-orientată.

Definiția 2.4.5 Fie $E = (E_1 \mid E_2 \mid \dots \mid E_n)^+$ o descriere lexicală bine formată peste Σ . Un *analizor lexical (scanner)* pentru E este un program ce recunoaște limbajul $L(E)$ și produce, pentru fiecare $w \in L(E)$, interpretarea sa drept-orientată.

Data o descriere lexicală E peste alfabetul Σ , proiectarea unui analizor lexical cuprinde următoarele proceduri:

1. Se construiește automatul finit (cu ε - tranziții) A , astfel ca $L(E) = L(A)$ (utilizând Algoritmul 2.1.1).
2. Se aplică Algoritmul 2.2.1 și se obține automatul determinist echivalent cu E , fie acesta A' .
3. (Opțional) Se aplică Algoritmul 2.3.1 pentru a obține automatul minimal echivalent cu A' .
4. Se scrie un program care implementează evoluția automatului obținut.

Ne vom ocupa în continuare de modificările ce trebuiesc aduse construcțiilor 1-3 pentru a putea obține o interpretare orientată dreapta pentru un cuvânt w din $L(E)$. Distingerea între clasele de unități lexicale E_1, E_2, \dots, E_n , poate fi făcută adăugând fiecărei expresii E_i o marcă de sfârșit $\#_i$, $1 \leq i \leq n$. Noua descriere lexicală, notată $E\#$ este:

$$E\# = (E_1\#_1 \mid E_2\#_2 \mid \dots \mid E_n\#_n)^+, \text{ peste alfabetul } \Sigma \cup \{\#_1, \dots, \#_n\}.$$

Automatul determinist ce recunoaște $L(E\#)$ va avea, desigur, tranziții pentru simbolurile $\#_i$. Pentru că aceste simboluri nu apar în șirul de intrare (se analizează cuvintele $w \in \Sigma^+$), vom trata aceste tranziții într-un mod special (în programul scris la etapa 4): când apare o tranziție de tipul $\#_i$, programul raportează o unitate lexicală din clasa E_i .

Exemplul 2.4.3 Să considerăm descrierea lexicală:

$\text{litera} \rightarrow a \mid b \mid \dots \mid z$ $\text{cifra} \rightarrow 0 \mid 1 \mid \dots \mid 9$
 $\text{identificator} \rightarrow \text{litera} (\text{litera} \mid \text{cifra})^*$
 $\text{semn} \rightarrow + \mid -$ $\text{numar} \rightarrow (\text{semn} \mid \varepsilon) \text{cifra}^+$
 $\text{operator} \rightarrow + \mid - \mid * \mid / \mid < \mid > \mid <= \mid >= \mid < >$
 $\text{asignare} \rightarrow :=$ $\text{doua_puncte} \rightarrow :$
 $\text{cuvinte_rezervate} \rightarrow \text{if} \mid \text{then} \mid \text{else}$
 $\text{paranteze} \rightarrow) \mid ($

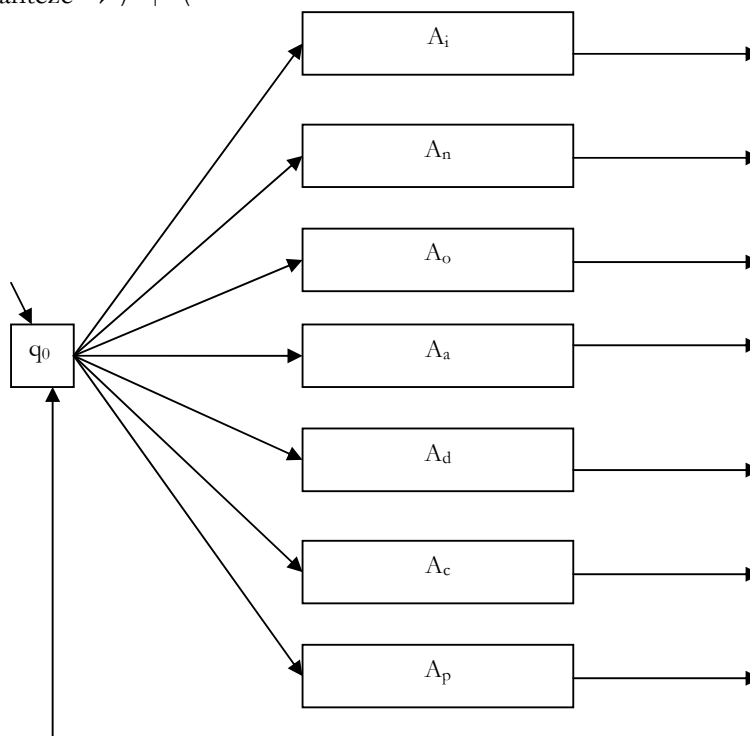


Figura 2.5

După cele discutate în paragraful precedent, există automatele (cu ε - tranziții) A_i , A_n , A_o , A_a , A_d , A_c , A_p care recunosc respectiv limbajele descrise de expresiile regulate din E. Atunci, un automat cu ε - tranziții ce recunoaște E se construiește după schema din figura 2.5.

Pentru acest automat cu ε - tranziții, există un automat determinist echivalent care recunoaște $L(E)$. Acest automat determinist poate fi încă “simplificat” prin găsirea automatului cu cele mai puține stări care recunoaște același limbaj. În sfârșit, se scrie un program - analizorul lexical - care simulează acest automat. Automatul determinist (minimal) echivalent cu cel descris în fig. 2.5 și cu tranzițiile corespunzătoare simbolurilor $\#_i$, $\#_n$, $\#_o$, $\#_a$, $\#_d$, $\#_c$, $\#_p$, care desemnează sfârșitul unei unități lexicale de tip *identificator*, *număr*, *operator*, *asignare*, *doua_puncte*, *cuvinte_rezervate*, *paranteze*, este dat în figura 2.6.

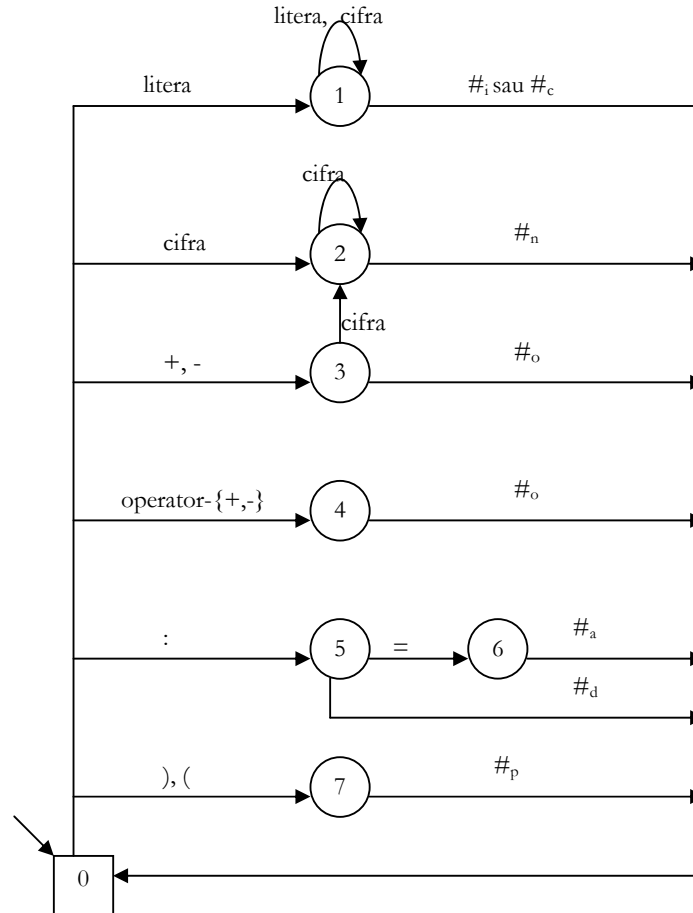


Figura 2.6

Așadar, 0 este stare inițială și finală, tranzițiile notate cu *literă*, *cifă* înseamnă că sunt tranziții pentru orice simbol din clasa *literă* respectiv *cifă*. Tranziția cu un delimitator # semnifică faptul că s-a identificat o unitate lexicală; aceasta se raportează și se reia activitatea automatului din starea inițială.

Programul (analizorul lexical) constă din câte un segment de program pentru fiecare stare a automatului. Dacă vom nota prin *buffer* zona de memorie unde se încarcă o unitate lexicală, *getnext(ch)*, *store(ch)* funcțiile ce citesc respectiv memorează *ch* în *buffer*, atunci pentru starea *q* care are tranziții respectiv în stările *q_i* pentru simbolurile *a_i*, $1 \leq i \leq m$ și tranziția cu *#_h* în *q₀*, segmentul de program este dat mai jos.

```

q : switch(ch) {
    case: a1
        store(ch);
        getnext(ch);
        goto q1;
    case: a2
        store(ch);
        getnext(ch);
        goto q2;
    .
    .
    case: am
        store(ch);
        getnext(ch);
        goto qm;
    default:
        write(buffer, i);
        empty(buffer);
        goto q0;
}

```

Așadar, în cazul în care caracterul citit este *a_i*, se memorează acesta, se citește alt caracter și se continuă programul pentru starea *q_i*, $1 \leq i \leq m$, iar dacă nu este nici unul dintre simbolurile *a₁*, *a₂*, ..., *a_m*, s-a depistat unitatea lexicală de tip *i*, se raportează aceasta (

`write(buffer,i)`), se golește zona `buffer` pentru a primi o nouă unitate lexicală după care se trece la programul stării inițiale.

Dacă din starea q nu există $\#$ - tranziție, grupul default se înlocuiește prin:

```
default:
    write(buffer,"eroare");
    empty(buffer);
    goto q0;
```

Pentru starea inițială q_0 (care este în același timp și stare finală, grupul default se înlocuiește prin:

```
default:
    if (EOF) {
        write("Sfarsitul analizei");
        exit(0);
    } else {
        write(buffer,"eroare");
        empty(buffer);
        getnext(ch);
        goto q0;
    }
```

2.4 Generatorul Lex

LEX este un instrument software deosebit de util celor care doresc să implementeze limbaje de programare, dar poate fi folosit și în multe alte aplicații. Acesta a fost dezvoltat la Bell Laboratories în anul 1975 de către M.E. Lesk și E. Schmidt. LEX citește un fișier de intrare (în general cu extensia `.l`) care conține expresii regulate ce definesc unitățile lexicale și generează un program C care realizează analiza lexicală a unui text, în conformitate cu specificațiile date.

LEX a devenit instrument standard UNIX începând cu versiunea a 7-a a acestui sistem de operare. Proiectul GNU al Fundației "Free Software" distribuie FLEX (**F**ast **LEX**ical Analyzer Generator) care este o îmbunătățire a generatorului LEX. De asemenea există versiuni pentru sistemele de operare DOS și WINDOWS; una dintre acestea este PCLEX lansat de Abraxax Software Inc..

Pentru a utiliza LEX, se parcurg trei pași:

1. Scrierea specificațiilor LEX într-un fișier cu extensia `.l` care reprezintă descrierea lexicală pentru care dorim să construim un analizor lexical.

2. Executarea programului LEX (sau PCLEX) cu intrarea fișierul construit:

```
lex [optiuni] <nume_fișier>
flex [optiuni] <nume_fișier>
pclex [optiuni] <nume_fișier>
```

Pentru a afla descrierea opțiunilor ce se pot folosi, executați programul `lex -h`.

Ca rezultat al execuției cu argumentul `<nume_fișier>` se obține un fișier cu același nume dar cu extensia `.c`, care este de fapt un program în limbaj C.

3. Se compilează acest program împreună, eventual, cu alte fișiere sursă.

Programul C obținut în urma executării LEX - ului este de fapt o funcție cu numele `yylex()` care, pentru a fi executată, trebuie inclusă într-o funcție `main()` (furnizată de utilizator) care conține apelul `yylex()` sau, se integrează această rutină cu altele (de exemplu într-o aplicație `yacc`).

O specificare LEX (fișierul de intrare) are structura:

```
Declarații
%%
Reguli
%%
Rutine auxiliare
```

Secțiunea *Declarații* precum și *Rutine auxiliare* pot lipsi încât o specificare minimă trebuie să conțină cele două linii `%%` între care sunt scrise liniile ce conțin reguli de specificare a unităților lexicale.

Secțiunea *Declarații* conține două părți. Prima parte constă din declarații C pentru variabilele sau constantele care se vor folosi ulterior iar a doua parte conține definiții ce se folosesc în specificațiile LEX în secțiunea *Reguli*. O definiție (macro pentru expresii regulate) are forma:

`<nume> <expresie>`

unde `<expresie>` este o expresie (regulată) ce poate folosi orice caracter împreună cu următorii operatori:

`" \ [] ^ - ? . * + | () $ / { } % <>`

Acești operatori au o anumită semnificație iar dacă se dorește a-i folosi drept caractere trebuie să precedă de `\` (backslash). Semnificația operatorilor este:

`"` Caracterele text sunt incluse între ghilimelele `"..."`. Construcția `"c"` este echivalentă cu `\c`.

`[]` Definește o clasă de caractere sub forma unei liste alternative. În interior se poate folosi caracterul `-` pentru a indica un domeniu. De exemplu, o cifră întreagă se poate defini prin `[0123456789]` sau prin `[0-9]`. Expresiile următoare: `[-a+]`, `[+a-]`, `[a\--]` sunt echivalente și desemnează unul din caracterele `a`, `+` sau `-`. Dacă `-` se folosește în `[]` atunci el trebuie pus primul, ultimul sau sub forma `\-`.

`^` Acest simbol este considerat operator dacă este primul într-o anumită descriere, altfel este caracter text. Se folosește pentru a exclude din definiție anumite caractere. De exemplu expresia `[^ \t\n]` desemnează un caracter diferit de spațiu, tab sau linie nouă.

`?` Elementul precedent acestui operator este opțional. Expresia `ab?c` desemnează `abc` sau `ac` iar `a[b-c]?d` notează `abd`, `acd` sau `ad`.

`*` Operatorul iterație de la expresii regulate.

`+` Operatorul iterație proprie (repetare de 1 sau mai multe ori) de la expresii regulate.

`|` Operatorul alternativă (reuniune) de la expresii regulate.

`()` Operatorul `()` (grupare) de la expresii regulate.

`{ }` În macro-definiții acest operator se poate folosi pentru a specifica numărul de repetiții ale unei expresii. De pildă `AAA` și `A{3}` sunt echivalente, iar dacă scriem `[a-z]{1,5}` aceasta desemnează una până la 5 litere mici ale alfabetului.

`.` Acest caracter (punct) desemnează orice caracter în afară de `\n`. De exemplu, descrierea `a.b` reprezintă `a0b`, `a1b`, `a\b`, `axb` etc. Într-o descriere lexicală ce urmează a fi inclusă într-un fișier intrare pentru lex, alternativa default se descrie printr-o linie ce are punctul în prima coloană, asta însemnând orice alt caracter ce nu a fost inclus până la această linie în vreo definiție sau descriere.

Următoarele două linii în secțiunea *declarații*:

`cifra [0-9]`
`litera [a-zA-Z]`

sunt definiții ale expresiilor regulate `cifra`, `litera` care pot fi folosite pentru definirea altor expresii în secțiunea *Reguli*. Numele definite aici trebuie să folosească acolade: `{litera}`, `{cifra}`.

Tot în această secțiune pot fi scrise segmente de cod C care urmează a fi incluse în procedura `yylex` (de pildă linii pentru preprocesor de tipul `#include` sau `#define`). Aceste elemente trebuie să fie încorporate într-un bloc de forma `%{...%}`. De exemplu:

`%{`
`#include <stdio.h>`
`%}`

Secțiunea *Reguli* definește funcționalitatea analizatorului lexical. Fiecare regulă este compusă dintr-o *expresie regulată* și o *acțiune*. Așadar, această secțiune are forma:

m_1	$\{Ac\breve{t}iune_1\}$
m_2	$\{Ac\breve{t}iune_2\}$
.	
.	
.	
m_n	$\{Ac\breve{t}iune_n\}$

unde m_i sunt expresii regulate iar $Ac\breve{t}iune_i$ este un segment de program C care se va executa atunci cînd analizorul întâlnește un exemplar din unitatea lexicală descrisă de m_i ($1 \leq i \leq n$). Pentru scrierea expresiilor regulate se folosesc operatorii pe care i-am descris mai sus, numele unor expresii deja definite, cuprinse între acolade precum și orice alte caractere text. Spre exemplu, expresia $\{litera\}(\{litera\}|\{cifra\})^*$ desemnează un identificator.

Secțiunea a treia, *Rutine auxiliare*, conține toate procedurile auxiliare care ar putea fi utile în procesul de analiză lexicală. Dacă analizorul lexical creat este independent, această secțiune trebuie să conțină o funcție `main()` care face apel la funcția `yylex()`.

Exemplul 2.5.1 Iată conținutul fișierului de intrare pentru descrierea lexicală din paragraful precedent:

```
%{
/* scanner1.1
// Exemplu de fisier de intrare pentru flex */
# include <stdio.h>
}%
litera [a-zA-Z]
cifra [0-9]
cifre ({cifra})+
semn [+ -]
operator [+*/<>= -]
spatiu [' '\t'\n']
%%
"if" | "then" | "else"      {
    printf("%s  cuvant rezervat\n", yytext);}
({litera})({litera}|{cifra})* {
    printf("%s  identificator\n", yytext);}
{cifre}|({semn})({cifre})    {
    printf("%s  numar intreg\n", yytext);}
{operator} {printf("%c  operator\n", yytext[0]);}
\:=       {printf("%s  asignare\n", yytext);}
\:        {printf("%c  doua puncte\n", yytext[0]);}
\(\)|\(\) {printf("%c  paranteza\n", yytext[0]);}
{spatiu}   {}
.          {printf("%c  caracter ilegal\n", yytext[0]);}
%%
int main( ){
    yylex( );
    return 0;
}
```

Exemplul 2.5.2 Următorul fișier constituie o descriere lexicală pentru un limbaj ce conține cuvintele rezervate `set`, `sqrt`, `quit`, identificatori, numere și operatori.

```
%{
/* scanner2.1
// Scanner expresii */
# include <stdio.h>
# include <string.h>
# include <stdlib.h>
}%
%%
```

```

quit {printf(" Cuvant cheie: %s\n", yytext);}
set {printf(" Cuvant cheie: %s\n", yytext);}
sqrt {printf(" Cuvant cheie: %s\n", yytext);}
\+ {printf(" Operator: %c\n", yytext[0]);}
\- {printf(" Operator: %c\n", yytext[0]);}
\* {printf(" Operator: %c\n", yytext[0]);}
\/ {printf(" Operator: %c\n", yytext[0]);}
\( {printf(" Paranteza stanga: %c\n", yytext[0]);}
\) {printf(" Paranteza dreapta: %c\n", yytext[0]);}
\:= {printf(" Asignare: %s\n", yytext);}
((([0-9]+(\.[0-9]*)?)|(\.[0-9]+))((e|E)(\+|\-)?[0-9]+)? {
    printf(" Numar: %f\n", atof(yytext));}
[_A-Za-z]+ {printf(" Identificator: %s\n", yytext);}
[\t]+ ; /* Eliminare spatii */
\n {return 0;}
. {printf(" Caracter ilegal: %c\n", yytext[0]);}
%%

int main() {
    printf("*****\n");
    printf("***          Exemplu de scanner          **\n");
    printf("***          Cuvinte cheie: set sqrt quit          **\n");
    printf("***          Operator: + - * / ( )          **\n");
    printf("***          Operanzi: numere, identificatori          **\n");
    printf("*****\n");
    yylex();
    return 0;
}

```

3 Analiza sintactică descendentă

La proiectarea limbajelor de programare se stabilesc reguli precise care descriu structura sintactică a programelor. În limbajul Pascal spre exemplu, un program este format din blocuri, un bloc conține instrucțiuni, o instrucțiune are în componența sa expresii, iar o expresie este o înșiruire de o anumită formă de unități lexicale.

Sintaxa unui limbaj de programare poate fi descrisă cu ajutorul gramaticilor independente de context sau a notației BNF (Backus - Naur Form). Descrierea sintaxei prin gramatici independente de context are avantaje atât din punct de vedere al proiectării, cât și din cel al implementării limbajelor. Aceste avantaje pot fi sintetizate astfel:

- o gramatică independentă de context oferă o specificare sintactică precisă și în același timp ușor de înțeles;
- pentru anumite clase de gramatici este posibilă construirea automată a unui analizor sintactic eficient. Analizorul sintactic are rolul de a determina dacă un program sursă este corect din punct de vedere sintactic. Mai mult, în procesul de construcție a analizorului sintactic pot fi depistate ambiguități sintactice sau alte construcții dificile, care pot rămâne nedetectate în timpul proiectării unui limbaj;
- descrierea sintaxei printr-o gramatică independentă de context impune limbajului o structură care se dovedește a fi utilă pentru traducerea programului pusă în cod obiect corect, precum și pentru detectarea erorilor.

3.1 Rolul analizorului sintactic

Un analizor sintactic primește la intrare un șir de unități lexice și produce arborele de derivare al acestui șir în gramatica ce descrie structura sintactică (mai degrabă o anume reprezentare a acestui arbore). În practică, odată cu verificarea corectitudinii sintactice, se produc și alte acțiuni în timpul analizei sintactice: trecerea în tabela de simboluri a unor informații legate de unitățile sintactice, controlul tipurilor sau producerea codului intermediar și alte chestiuni legate de analiza semantică (unele din acestea le vom discuta în capitolul ce tratează semantica). Să ne referim în câteva cuvinte la natura erorilor sintactice și la strategiile generale de tratare a acestora.

Dacă un compilator ar trebui să trateze doar programele corecte, atunci proiectarea și implementarea sa s-ar simplifica considerabil. Cum însă programatorii scriu adesea programe incorecte, un compilator bun trebuie să asiste programatorul în identificarea și localizarea erorilor. Cea mai mare parte a specificărilor limbajelor de programare nu descriu modul în care implementarea (compilatorul) trebuie să reacționeze la erori; acest lucru este lăsat pe seama celor care concep compilatoarele. Erorile într-un program pot fi:

- *lexicale*: scrierea eronată a unui identificator, a unui cuvânt cheie, a unui operator etc. ;
- *sintactice*: omiterea unor paranteze într-o expresie aritmetică, scrierea incorectă a unei instrucțiuni, etc. ;
- *semantice*: aplicarea unui operator aritmetic unor operanzi logici, nepotrivirea tipului la atribuire etc. ;
- *logice*: un apel recursiv infinit.

3.2 Problema recunoașterii. Problema parsării

Problema recunoașterii în gramatici independente de context este următoarea: Dată o gramatică $G = (V, T, S, P)$ și un cuvânt $w \in T^*$, care este răspunsul la întrebarea $w \in L(G)$? Se știe că problema este decidabilă; mai mult, există algoritmi care în timp $O(|w|^3)$ dau răspunsul la întrebare (Cooke-Younger-Kasami, Earley, vezi [Gri86]). Problema *parsării* (*analizei sintactice*) este problema recunoașterii la care se adaugă: dacă răspunsul la întrebarea $w \in L(G)$ este afirmativ, se cere arborele sintactic (o reprezentare a sa) pentru w .

Fie $G = (V, T, S, P)$ o gramatică și $w \in L(G)$. Să notăm prin $\xRightarrow{p}_{st} (\xRightarrow{q}_{dr})$ relația de derivare extrem stângă (dreaptă) în care, pentru rescriere s-a aplicat producția $p \in P$ ($q \in P$). Atunci, pentru $w \in L(G)$ există o derivare extrem stângă:

$$S = \alpha_0 \xRightarrow{p^1}_{st} \alpha_1 \xRightarrow{p^2}_{st} \alpha_2 \xRightarrow{p^3}_{st} \dots \xRightarrow{p^n}_{st} \alpha_n = w$$

și o derivare extrem dreaptă:

$$S = \beta_0 \xRightarrow{q^1}_{dr} \beta_1 \xRightarrow{q^2}_{dr} \beta_2 \xRightarrow{q^3}_{dr} \dots \xRightarrow{q_m}_{dr} \beta_m = w$$

Atunci, arborele sintactic corespunzător se poate reprezenta prin secvența de producții $p_1 p_2 \dots p_n \in P^+$ (respectiv $q_1 q_2 \dots q_m \in P^+$) care s-au aplicat, în această ordine, în derivarea extrem stângă (dreaptă) pentru obținerea lui w .

Definiția 3.2.1 Secvența de producții $\pi = p_1 p_2 \dots p_n \in P^+$ pentru care $S \xRightarrow{\pi}_{st} w$ se numește *analizare stângă* (*parsare stângă*) pentru cuvântul $w \in L(G)$.

Secvența de producții $\pi = q_1 q_2 \dots q_m \in P^+$ pentru care $S \xRightarrow[\text{dr}]{\tilde{\pi}} w$ se numește *analizare dreaptă* (*parsare dreaptă*) pentru cuvântul w (unde $\tilde{\pi} = q_m q_{m-1} \dots q_1$).

Exemplul 3.2.4 Fie gramatica

$$1. E \rightarrow E+T \quad 2. E \rightarrow T \quad 3. T \rightarrow T*F,$$

$$4. T \rightarrow F \quad 5. F \rightarrow (E) \quad 6. F \rightarrow \text{id}$$

unde 1, 2, ..., 6 identifică cele 6 producții ale acesteia.

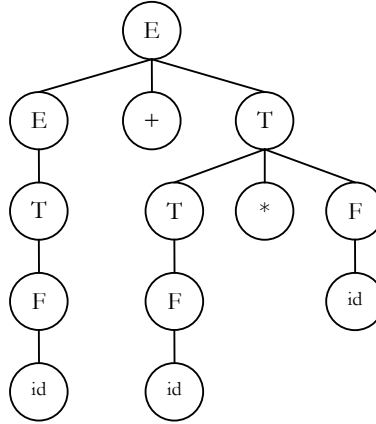


Figura 3.1

Pentru cuvântul $w = \text{id}+\text{id}*\text{id}$, arborele sintactic este dat în figura 3.1.

Analizarea stângă pentru acest cuvânt este: $\pi_1 = 12463466$ iar analizarea dreaptă este $\pi_2 = 64264631$. Aici am identificat fiecare producție prin numărul său.

3.3 Analiza sintactică descendentă

Analiza sintactică descendentă (*parsarea descendentă*) poate fi considerată ca o tentativă de determinare a unei derivări extrem stângi pentru un cuvânt de intrare. În termenii arborilor sintactici, acest lucru înseamnă tentativa de construire a unui arbore sintactic pentru cuvântul de intrare, pornind de la rădăcină și construind nodurile în manieră descendentă, în preordine (construirea rădăcinii, a subarborelui stâng apoi a celui drept). Pentru realizarea acestui fapt avem nevoie de următoarea structură (figura 3.2):

- o *bandă de intrare* în care se introduce cuvântul de analizat, care se parcurge de la stânga la dreapta, simbol cu simbol;
- o *memorie de tip stivă* (pushdown) în care se obțin formele propoziționale stângi (începând cu S). Prefixul formei propoziționale format din terminali se compară cu simbolurile curente din banda de intrare obținându-se astfel criteriul de înaintare în această bandă;
- o *bandă de ieșire* în care se înregistrează pe rând producțiile care s-au aplicat în derivarea extrem stângă care se construiește.

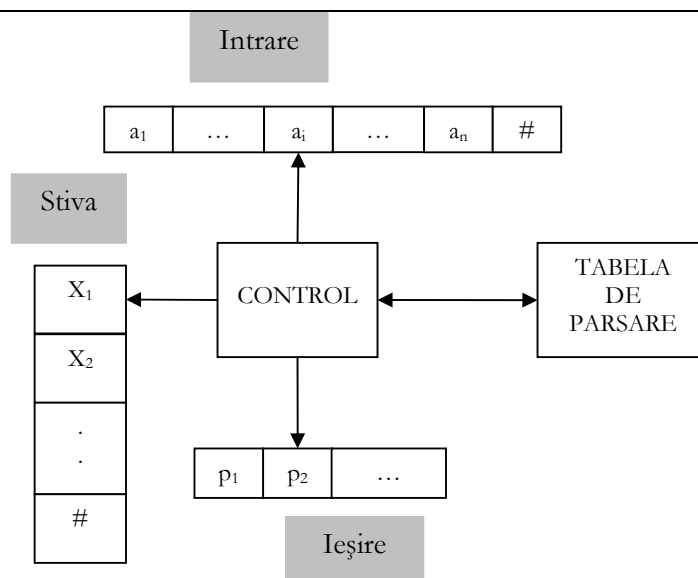


Figura 3.2

Criteriul de oprire cu succes este acela în care s-a parcurs întreaga bandă de intrare, iar memoria pushdown s-a golit. În acest caz în banda de ieșire s-a obținut *parsarea stângă* a cuvântului respectiv. Iată un exemplu de cum funcționează acest mecanism (considerăm gramatica din exemplul precedent și cuvântul $w = id+id*id$). Vom considera caracterul # pentru marcarea sfârșitului benzii de intrare (marca de sfârșit a cuvântului) precum și pentru a marca baza stivei.

Tabelul 1: Analiza sintactică descendentă pentru $w = id+id*id$

Pasul	Banda de intrare	Stiva de lucru	Banda de ieșire
1	$id+id*id\#$	$E\#$	
2	$id+id*id\#$	$E+T\#$	1
3	$id+id*id\#$	$T+T\#$	12
4	$id+id*id\#$	$F+T\#$	124
5	$id+id*id\#$	$id+T\#$	1246
6	$id*id\#$	$T\#$	1246
7	$id*id\#$	$T*F\#$	12463
8	$id*id\#$	$F*F\#$	124634
9	$id*id\#$	$id*F\#$	1246346
10	$id\#$	$F\#$	1246346
11	$id\#$	$id\#$	12463466
12	$\#$	$\#$	12463466

Se obține la ieșire $\pi = 12463466$.

3.3.1 Parser descendent general

Fie $G = (V, T, S, P)$ o gramatică care descrie sintaxa unui limbaj de programare. Să considerăm mulțimea $\mathcal{C} = T^*\# \times \Sigma^*\# \times P^*$ numită *mulțimea configurațiilor* (unui parser). Interpretarea pe care o dăm unei configurații este următoarea: dacă $c = (u\#, \gamma\#, \pi)$, atunci $u\#$ este cuvântul rămas de analizat (conținutul benzii de intrare la un moment dat), $\gamma\#$ este conținutul stivei, cu primul simbol din γ în top-ul stivei ($\gamma = X_1X_2\dots X_n$ ca în figura 3.2), iar $\#$ la baza acesteia, iar π este conținutul benzii de ieșire: dacă $\pi = p_1p_2\dots p_k$, până în acest moment s-au aplicat producțiile din π în această ordine.

Definiția 3.3.1 *Parserul (analizorul sintactic) descendent* atașat gramaticii $G = (V, T, S, P)$ este perechea (C_0, \vdash) unde $C_0 = \{(w\#, S\#, \varepsilon) \mid w \in T^*\} \subseteq C$ este mulțimea configurațiilor inițiale, iar $\vdash \subseteq C \times C$ este relația de tranziție între configurații dată de următoarea schemă:

1. $(u\#, A\gamma\#, \pi) \vdash (u\#, \beta\gamma\#, \pi r)$, unde $r = A \rightarrow \beta \in P$.
2. $(uv\#, u\gamma\#, \pi) \vdash (v\#, \gamma\#, \pi)$.
3. $(\#, \#, \pi)$ este *configurație de acceptare* dacă $\pi \neq \varepsilon$.
4. O configurație c pentru care nu există c' astfel ca $c \vdash c'$ (în sensul 1, 2) spunem că produce eroare.

Să considerăm relația $\vdash^+ (\vdash^*)$ închiderea tranzitivă (și reflexivă) a relației \vdash definită mai sus. Este clar că, dacă vom considera parserul ca un sistem de rescriere, acesta este nedeterminist așa cum este definită tranziția de tip 1: dacă în P există producțiile $A \rightarrow \beta_1$ și $A \rightarrow \beta_2$ cu $\beta_1 \neq \beta_2$, atunci au loc, în același timp: $(u\#, A\gamma\#, \pi) \vdash (u\#, \beta_1\gamma\#, \pi)$, $(u\#, A\gamma\#, \pi r_1) \vdash (u\#, \beta_2\gamma\#, \pi r_2)$ unde $r_1 = A \rightarrow \beta_1$, $r_2 = A \rightarrow \beta_2$. Pentru a fi corect, parserul definit mai sus va trebui să transforme (în mod nedeterminist, în general) configurația $(w\#, S\#, \varepsilon)$ în configurația $(\#, \#, \pi)$ pentru cuvintele $w \in L(G)$ și numai pentru acestea! Vom demonstra în continuare corectitudinea în acest sens.

Lema 3.3.1 Dacă în parserul descendent atașat gramaticii $G = (V, T, S, P)$ are loc calculul $(uv\#, \gamma\#, \varepsilon) \vdash^* (v\#, \psi\#, \pi)$, atunci în gramatica G are loc derivarea $\gamma \xRightarrow[st]{\pi} u\psi$, oricare ar fi $u, v \in T^*$, $\gamma, \psi \in \Sigma^*$, $\pi \in P^*$.

Demonstrație. Să demonstrăm afirmația din lema prin inducție asupra lungimii lui π (notată $|\pi|$):

- Dacă $|\pi| = 0$ înseamnă că în calculul considerat s-au folosit doar tranziții de tip 2. Atunci $\gamma = u\psi$ și are loc $\gamma \xRightarrow[st]{\pi} u\psi$, cu $\pi = \varepsilon$;
- Dacă $|\pi| > 0$ fie $\pi = \pi_1 r$, unde $r = A \rightarrow \beta$ este ultima producție care se folosește în calculul considerat și presupunem afirmația adevărată pentru calcule ce produc la ieșire π_1 . Așadar, putem scrie:

$$(uv\#, \gamma\#, \varepsilon) = (u_1 u_2 v\#, \gamma\#, \varepsilon) \vdash^* (u_2 v\#, A\gamma_1\#, \pi_1) \vdash (u_2 v\#, \beta\gamma_1\#, \pi_1 r) =$$

$$= (u_2 v\#, u_2 \psi\#, \pi_1 r) \vdash^* (v\#, \psi\#, \pi_1 r)$$

Atunci din $(u_1 u_2 v\#, \gamma\#, \varepsilon) \vdash^* (u_2 v\#, A\gamma_1\#, \pi_1)$, conform ipotezei inductive, are loc $\gamma \xRightarrow[st]{\pi} u_1 A \gamma_1$. Cum în $u_1 A \gamma_1$, $u_1 \in T^*$, iar $A \in V$ este cea mai din stânga variabilă, aplicând acesteia producția $A \rightarrow \beta$, obținem:

$$\gamma \xRightarrow[st]{\pi} u_1 A \gamma_1 \xRightarrow[st]{r} u_1 \beta \gamma_1 = u_1 u_2 \psi = u\psi \text{ și demonstrația este încheiată.}$$

□

Corolarul 3.3.1 Dacă în parserul descendent are loc $(w\#, S\#, \varepsilon) \vdash^+ (\#, \#, \pi)$ atunci în gramatica G are loc $S \xRightarrow[st]{\pi} w$.

Demonstrație Se aplică lema precedentă pentru $u = w$, $v = \varepsilon$, $\gamma = S$, $\psi = \varepsilon$.

□

Lema 3.3.2 Dacă în gramatica G are loc derivarea $\gamma \xRightarrow[st]{\pi} u\psi$ și $1:\psi \in V \cup \{\varepsilon\}$ (vezi definiția 3.3.2), atunci în parserul descendent are loc calculul:

$$(uv\#, \gamma\#, \varepsilon) \vdash^* (v\#, \psi\#, \pi), \forall v \in T^*.$$

Demonstrație Procedăm de asemenea prin inducție după lungimea lui π :

- Dacă $|\pi| = 0$, atunci $\gamma = u\psi$ și calculul are loc doar cu tranziții de tipul al doilea: $(uv\#, \gamma\#, \varepsilon) = (uv\#, u\psi\#, \varepsilon) \vdash^* (v\#, \pi\#, \varepsilon)$.
- Dacă $|\pi| > 0$ fie $\pi = \pi_1 r$ unde $r = A \rightarrow \beta$ este ultima producție aplicată în derivarea considerată: $\gamma \xRightarrow[st]{\pi_1} u_1 \psi_1 = u_1 A \psi_2 \xRightarrow[st]{r} u_1 \beta \psi_2 = u_1 u_2 \psi = u\psi$.

Am folosit aici efectiv ipoteza referitoare la primul simbol din ψ : primul simbol din ψ_1 este o variabilă (și nu ε) pentru că se aplică efectiv un pas cu producția r . Cum $|\pi_1| < |\pi|$, din $\gamma \xRightarrow[st]{\pi_1} u_1 \psi_1$, conform ipotezei inductive are loc $(u_1 v_1\#, \gamma\#, \varepsilon) \vdash^* (v_1\#, \psi_1\#, \pi_1)$

$\forall v_1 \in T^*$. Considerând $v_1 = u_2 v$ cu $v \in T^*$ oarecare, putem scrie:

$$\begin{aligned} (u_1 v_1\#, \gamma\#, \varepsilon) &= (u_1 u_2 v\#, \gamma\#, \varepsilon) && \vdash^* (u_2 v\#, \psi_1\#, \pi_1) &= \\ &= (u_2 v\#, A\psi_2\#, \pi_1) && \vdash (u_2 v\#, \beta\psi_2\#, \pi_1 r) &= \\ &= (u_2 v\#, u_2 \psi\#, \pi) && \vdash^* (v\#, \psi\#, \pi) \end{aligned}$$

□

Corolarul 3.3.2 Dacă în G are loc derivarea $S \xRightarrow[st]{\pi} w$ atunci în parserul descendent are loc calculul $(w\#, S\#, \varepsilon) \vdash^* (\#, \#, \pi)$.

Demonstrație În lema 3.3.2 se consideră $u = w$, $v = \varepsilon$, $\gamma = S$ și $\psi = \varepsilon$.

□

Teorema 3.3.1 (*Corectitudinea parserului descendent general*). Se consideră gramatica redusă $G = (V, T, S, P)$ și $w \in T^*$. Atunci, în parserul descendent general atașat acestei gramatici are loc $(w\#, S\#, \varepsilon) \vdash^* (\#, \#, \pi)$ (acceptare) dacă și numai dacă $w \in L(G)$ și π este o analizare sintactică stângă a frazei w .

Demonstrație Teorema sintetizează rezultatele din lemele 3.3.1 și 3.3.2.

□

3.3.2 Gramatici LL(k)

Parserul general introdus în paragraful precedent este un model nedeterminist iar implementarea sa pentru gramatici oarecare este ineficientă. Cea mai cunoscută clasă de gramatici pentru care acest model funcționează determinist este aceea a gramaticilor LL(k): *Parsing from Left to right using Leftmost derivation and k symbols lookahead*. Intuitiv, o gramatică este LL(k) dacă tranziția de tip 1 din Definiția 3.3.1 se face cu o unică regulă $A \rightarrow \beta$ determinată prin inspectarea a k simboluri care urmează a fi analizate în banda de intrare. Dar să definim în mod riguros această clasă de gramatici.

Definiția 3.3.2 Fie $G = (V, T, S, P)$ o gramatică, $\alpha \in \Sigma^*$ și $k \geq 1$ un număr natural. Atunci definim:

$$\begin{aligned} k:\alpha &= \text{if } |\alpha| \leq k \text{ then } \alpha \text{ else } \alpha_1, \text{ unde } \alpha = \alpha_1 \beta, |\alpha_1| = k, \\ \alpha:k &= \text{if } |\alpha| \leq k \text{ then } \alpha \text{ else } \alpha_2, \text{ unde } \alpha = \gamma \alpha_2, |\alpha_2| = k. \end{aligned}$$

Definiția 3.3.3 O gramatică independentă de context redusă este gramatică LL(k), $k \geq 1$, dacă pentru orice două derivări de forma:

$$\begin{aligned} S &\xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\beta_1\gamma \xRightarrow[st]{*} ux \\ S &\xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\beta_2\gamma \xRightarrow[st]{*} uy \end{aligned}$$

unde $u, x, y \in T^*$, pentru care $k:x = k:y$, are loc $\beta_1 = \beta_2$.

Definiția prezentată exprimă faptul că, dacă în procesul analizei cuvântului $w=ux$ a fost obținut (analizat) deja prefixul u și trebuie să aplicăm o producție neterminală A (având de ales pentru acesta măcar din două variante), această producție este determinată în mod unic de următoarele k simboluri din cuvântul care a rămas de analizat (cuvântul x).

□

3.3.3 O caracterizare a gramaticilor LL(1)

Pentru ca un parser SLL(k) să poată fi implementat, trebuie să indicăm o procedură pentru calculul mulțimilor $FIRST_k$ și $FOLLOW_k$. Pentru că în practică se folosește destul de rar (dacă nu chiar deloc) cazul $k \geq 2$, vom restrânge discuția pentru cazul $k = 1$. Vom nota în acest caz $FIRST_1$ și $FOLLOW_1$ prin $FIRST$ respectiv $FOLLOW$.

Așadar, dacă $\alpha \in \Sigma^+$, $A \in V$:

$$FIRST(\alpha) = \{a \mid a \in T, \alpha \xRightarrow[st]{*} au\} \cup \text{if } (\alpha \xRightarrow[st]{*} \varepsilon) \text{ then } \{\varepsilon\} \text{ else } \emptyset.$$

$$FOLLOW(A) = \{a \mid a \in T \cup \{\varepsilon\}, S \xRightarrow[st]{*} uA\gamma, a \in FIRST(\gamma)\}$$

Mai întâi să arătăm că gramaticile SLL(1) coincid cu gramaticile LL(1).

Teorema 3.3.10 O gramatică $G = (V, T, S, P)$ este gramatică LL(1) dacă și numai dacă pentru orice $A \in V$ și pentru orice producții $A \rightarrow \beta_1 \mid \beta_2$ are loc:

$$FIRST(\beta_1 FOLLOW(A)) \cap FIRST(\beta_2 FOLLOW(A)) = \emptyset$$

Demonstrație Să presupunem că G este LL(1). Acest lucru este echivalent cu (după Definiția 3.3.3):

$$\left. \begin{array}{l} S \xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\beta_1\gamma \xRightarrow[st]{*} uv_1 \\ S \xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\beta_2\gamma \xRightarrow[st]{*} uv_2 \\ 1:v_1 = 1:v_2 \end{array} \right\} \Rightarrow \beta_1 = \beta_2$$

Să presupunem că G nu este SLL(1): există producțiile $A \rightarrow \alpha_1, A \rightarrow \alpha_2$ și

$$a \in FIRST(\alpha_1 FOLLOW(A)) \cap FIRST(\alpha_2 FOLLOW(A))$$

Distingem următoarele cazuri:

- $a \in FIRST(\alpha_1) \cap FIRST(\alpha_2)$, adică $\alpha_1 \xRightarrow[st]{*} au_1, \alpha_2 \xRightarrow[st]{*} au_2$

Atunci, putem scrie derivările (G este redusă):

$$\begin{array}{l} S \xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\alpha_1\gamma \xRightarrow[st]{*} uau_1v \\ S \xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\alpha_2\gamma \xRightarrow[st]{*} uau_2v \end{array}$$

ceea ce contrazice faptul că G este LL(1).

- $a \in FIRST(\alpha_1) \cap FOLLOW(A)$ și $\alpha_2 \xRightarrow[st]{*} \varepsilon$

Atunci: $\alpha_1 \xRightarrow[st]{*} au_1$ și $S \xRightarrow[st]{*} uA\gamma, a \in FIRST(\gamma)$ adică $\gamma \xRightarrow[st]{*} au_2$, și putem scrie:

$$\begin{array}{l} S \xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\alpha_1\gamma \xRightarrow[st]{*} uau_1v \\ S \xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\alpha_2\gamma \xRightarrow[st]{*} u\gamma \xRightarrow[st]{*} uau_2 \end{array}$$

ceea ce contrazice de asemenea ipoteza.

- $a \in \text{FIRST}(\alpha_2) \cap \text{FOLLOW}(A)$ și $\alpha_1 \xRightarrow[st]{*} \varepsilon$, caz ce se tratează analog cu precedentul.
- $\alpha_1 \xRightarrow[st]{*} \varepsilon, \alpha_2 \xRightarrow[st]{*} \varepsilon$ și atunci $a \in \text{FOLLOW}(A)$.

În acest caz avem:

$$\begin{aligned} S &\xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\alpha_1\gamma \xRightarrow[st]{*} u\gamma \xRightarrow[st]{*} uau_2 \\ S &\xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\alpha_2\gamma \xRightarrow[st]{*} u\gamma \xRightarrow[st]{*} uau_2 \end{aligned}$$

Deci din nou se contrazice proprietatea LL(1) pentru gramatica G.

Invers, să presupunem că pentru orice două producții $A \rightarrow \beta_1, A \rightarrow \beta_2$ distincte are loc $\text{FIRST}(\beta_1\text{FOLLOW}(A)) \cap \text{FIRST}(\beta_2\text{FOLLOW}(A)) = \emptyset$ și, prin reducere la absurd, G nu este LL(1). Asta înseamnă că există două derivări:

$$\begin{aligned} S &\xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\alpha_1\gamma \xRightarrow[st]{*} uav \text{ și} \\ S &\xRightarrow[st]{*} uA\gamma \xRightarrow[st]{*} u\alpha_2\gamma \xRightarrow[st]{*} uav' \end{aligned}$$

cu $1:v = 1:v'$ dar $\alpha_1 \neq \alpha_2$.

Analizând cele două derivări, se constată ușor că:

$$a \in \text{FIRST}(\alpha_1\text{FOLLOW}(A)) \cap \text{FIRST}(\alpha_2\text{FOLLOW}(A))$$

ceea ce contrazice ipoteza.

□

3.3.4 Determinarea mulțimilor FIRST și FOLLOW

Vom indica în acest paragraf modalitatea de determinare a mulțimilor FIRST și FOLLOW pentru o gramatică G.

Un algoritm pentru determinarea mulțimilor FIRST(X) este descris mai jos. Cititorul este invitat să dovedească, pornind de la definiția lui FIRST, că acest algoritm este corect.

Algoritm 3.3.2

Intrare: Gramatica $G = (V, T, S, P)$ redusă;

Ieșire: FIRST(X), $X \in \Sigma$.

Metoda: Mulțimile FIRST sunt completate prin inspectarea regulilor gramaticii

```

1. for ( $X \in \Sigma$ )
2.   if ( $X \in T$ ) FIRST(X) = { X } else FIRST(X) =  $\Phi$ ;
3. for ( $A \rightarrow a\beta \in P$ )
4.   FIRST(A) = FIRST(A)  $\cup$  { a };
5. FLAG = true;
6. while(FLAG) { // FLAG marcheaza schimbarile in FIRST
7.   FLAG = false;
8.   for( $A \rightarrow X_1X_2...X_n \in P$ ) {
9.     i = 1;
10.    if ( (FIRST(X1)  $\not\subset$  FIRST(A)) {
11.      FIRST(A) = FIRST(A)  $\cup$  (FIRST(X1));
12.      FLAG = true;
13.    } //endif
14.    while(i < n &&  $X_i \xRightarrow[st]{+} \varepsilon$ )
15.      if ( (FIRST(Xi+1)  $\not\subset$  FIRST(A)) {
```

```

15.          FIRST(A) = FIRST(A) ∪ FIRST(Xi+1) ;
16.          FLAG = true; i++;
              } //endif
          } //endwhile
      } //endfor
  } //endwhile
17.  for (A ∈ V)
18.      if (A  $\xRightarrow[st]{+}$  ε) FIRST(A) = FIRST(A) ∪ { ε } ;

```

Pentru a verifica dacă o gramatică este LL(1), conform teoremei de caracterizare, este necesar să calculăm FIRST(α), unde $\alpha \in \Sigma^*$. Vom descrie în algoritmul următor acest lucru.

Algoritmul 3.3.3

Intrare: Gramatica $G=(V, T, S, P)$.
 Mulțimile FIRST(X), $X \in \Sigma$.
 $\alpha = X_1X_2 \dots X_n$, $X_i \in \Sigma$, $1 \leq i \leq n$.

Ieșire: FIRST(α).

Metoda:

```

1. FIRST( $\alpha$ ) = FIRST( $X_1$ ) - { ε }; i = 1;
2. while (i < n &&  $X_i \xRightarrow[st]{+}$  ε) {
3.   FIRST( $\alpha$ ) = FIRST( $\alpha$ ) ∪ (FIRST( $X_{i+1}$ ) - { ε }) ;
4.   i = i + 1 ;
   } //endwhile
5. if (i == n &&  $X_n \xRightarrow[st]{+}$  ε)
6.   FIRST( $\alpha$ ) = FIRST( $\alpha$ ) ∪ { ε } ;

```

Corectitudinea algoritmului descris rezultă din observația că, dacă pentru cuvântul $\alpha = X_1X_2 \dots X_n$ are loc $X_1X_2 \dots X_i \xRightarrow[st]{+}$ ε, $i < n$ (adică X_1, X_2, \dots, X_i se pot șterge), atunci

$$\left(\bigcup_{k=1}^{i+1} \text{FIRST}(X_k) - \{ \varepsilon \} \right) \subseteq \text{FIRST}(\alpha).$$

Dacă în plus $\alpha \xRightarrow[st]{+}$ ε, atunci FIRST(α) conține și ε.

Exemplul 3.3.1 Fie gramatica ce are următoarele reguli:

$E \rightarrow E+T \mid E-T \mid T, \quad T \rightarrow T^*F \mid T/F \mid F, \quad F \rightarrow (E) \mid a.$

Aplicând algoritmul pentru calculul mulțimilor FIRST, obținem:

FIRST(E) = FIRST(T) = FIRST(F) = { (, a },

FIRST(T^*F) = FIRST(T) = { (, a }.

Exemplul 3.3.2 Fie gramatica:

$S \rightarrow E \mid B \qquad E \rightarrow \varepsilon$
 $B \rightarrow a \mid \text{begin } SC \text{ end} \qquad C \rightarrow \varepsilon \mid ; SC$

FIRST(S) = { a, begin, ε }

FIRST(E) = { ε }

FIRST(B) = { a, begin }

FIRST(C) = { ;, ε }.

Aplicând Algoritmul 3.3.3 se găsește: FIRST(SEC) = { a, begin, ;, ε }, FIRST(SB) = { a, begin }, FIRST($;SC$) = { ; }.

Să trecem acum la determinarea mulțimilor FOLLOW(A), $A \in V$. Reamintind că $\text{FOLLOW}(A) = \{a \in T \cup \{\varepsilon\} \mid S \xRightarrow[st]{*} \alpha A \beta, a \in \text{FIRST}(\beta)\}$, să facem câteva observații care vor fi utile în determinarea lui FOLLOW:

- $\varepsilon \in \text{FOLLOW}(S)$ pentru că $S \xRightarrow[st]{*} S$.
- Dacă $A \rightarrow \alpha B \beta X \gamma \in P$ și $\beta \xRightarrow[st]{*} \varepsilon$, atunci $\text{FIRST}(X) - \{\varepsilon\} \subseteq \text{FOLLOW}(B)$. În adevăr, putem scrie $S \xRightarrow[st]{*} \alpha_1 A \beta_1 \xRightarrow[st]{*} \alpha_1 \alpha B \beta X \gamma \beta_1 \xRightarrow[st]{*} \alpha_1 \alpha B X \gamma \beta_1$, și atunci rezultă $\text{FIRST}(X) - \{\varepsilon\} \subseteq \text{FOLLOW}(B)$. Acest aspect, de fapt, se poate exprima mai simplu prin:
- Dacă $A \rightarrow \alpha B \beta \in P$ atunci $\text{FIRST}(\beta) - \{\varepsilon\} \subseteq \text{FOLLOW}(B)$.
- Dacă $A \rightarrow \alpha B \beta \in P$ și $\beta \xRightarrow[st]{*} \varepsilon$, atunci $\text{FOLLOW}(A) \subseteq \text{FOLLOW}(B)$. În adevăr, dacă $a \in \text{FOLLOW}(A)$ atunci există derivarea $S \xRightarrow[st]{*} \alpha_1 A \gamma$ și a este din $\text{FIRST}(\gamma)$. Din cele stipulate în ipoteză putem continua în această derivare:

$$S \xRightarrow[st]{*} \alpha_1 A \gamma \xRightarrow[st]{*} \alpha_1 \alpha B \beta \gamma \xRightarrow[st]{*} \alpha_1 \alpha B \gamma$$

și deci $a \in \text{FOLLOW}(B)$.

Ținând cont de aceste observații rezultă următorul algoritm:

Algoritm 3.3.4 (Determinarea mulțimilor FOLLOW)

Intrare: Gramatica $G = (V, T, S, P)$ redusă ; Procedura $\text{FIRST}(\alpha)$, $\alpha \in \Sigma^+$.

Ieșire: Mulțimile FOLLOW(A), $A \in V$.

Metoda:

```

1. for ( $A \in \Sigma$ ) FOLLOW(A) =  $\Phi$ ;
2. FOLLOW(S) = {  $\varepsilon$  } ;
3. for ( $A \rightarrow X_1 X_2 \dots X_n$ ) {
4.   i = 1;
5.   while (i < n) {
6.     while ( $X_i \notin V$ ) ++i;
7.     if (i < n) { //  $X_i$  este neterminal
8.       FOLLOW( $X_i$ ) = FOLLOW( $X_i$ )  $\cup$  ( $\text{FIRST}(X_{i+1} X_{i+2} \dots X_n) - \{\varepsilon\}$ );
9.       ++i;
10.    } //endif
11.  } //endwhile
12. } //endfor
13. FLAG = true;
14. while (FLAG) { // FLAG semnaleaza schimbarile in FOLLOW
15.   FLAG = false;
16.   for ( $A \rightarrow X_1 X_2 \dots X_n$ ) {
17.     i = n;
18.     while (i >= 1 &&  $X_i \in V$ ) {
19.       if (FOLLOW(A)  $\subsetneq$  FOLLOW( $X_i$ )) {
20.         FOLLOW( $X_i$ ) = FOLLOW( $X_i$ )  $\cup$  FOLLOW(A);
21.         FLAG = true;
22.       } //endif
23.       if ( $X_i \xRightarrow[st]{+} \varepsilon$ ) --i; // la  $X_{i-1}$  daca  $X_i$  se sterge
24.     } else continue; // la urmatoarea productie
25.   } //endwhile
26. } //endfor
27. } //endwhile

```

3.3.5 Tabela de analiză sintactică LL(1)

Pentru a implementa un analizor sintactic pentru gramatici LL(1), să considerăm o *tabelă de analiză*, sau *tabelă de parsare* LL(1):

$$M : V \times (T \cup \{ \# \}) \rightarrow \{ (\beta, p) \mid p = A \rightarrow \beta \in P \} \cup \{ \text{eroare} \}$$

construită după algoritmul următor:

Algoritmul 3.3.6 (Tabela de analiză sintactică)

Intrare: Gramatica $G = (V, T, S, P)$;

Mulțimile $\text{FIRST}(\beta)$, $\text{FOLLOW}(A)$, $A \rightarrow \beta \in P$.

Ieșire: Tabela de parsare M .

Metoda: Se parcurg regulile gramaticii și se pun în tabelă

```

1.  for(A ∈ V )
2.      for( a ∈ T ∪ { # } )
3.          M(A, a) = Φ ;
4.  for (p = A → β ∈ P) {
5.      for(a ∈ FIRST(β) - {ε})
6.          M(A, a) = M(A, a) ∪ { (β, p) } ;
7.      if(ε ∈ FIRST(β)) {
8.          for(b ∈ FOLLOW(A)) {
9.              if(b == ε) M(A, #) = M(A, #) ∪ { (β, p) } ;
10.             else M(A, b) = M(A, b) ∪ { (β, p) } ;
11.             } //endfor
12.         } //endif
13.     } //endfor
14. for(A ∈ V)
15.     for( a ∈ T ∪ { # } )
16.         if(M(A, a) = Φ) M(A, a) = {eroare} ;

```

Lema 3.3.6 (caracterizarea gramaticilor LL(1)) Gramatica redusă G este LL(1) dacă și numai dacă pentru orice $A \in V$ și orice $a \in T \cup \{ \# \}$ $M(A, a)$, dacă nu este eroare, conține cel mult o pereche (β, p) .

Demonstrație Dacă G este LL(1), atunci $\forall A \in V$, oricare ar fi două producții cu partea stângă A , $A \rightarrow \beta_1$, $A \rightarrow \beta_2$, are loc, conform teoremei 3.3.10, $\text{FIRST}(\beta_1 \text{FOLLOW}(A)) \cap \text{FIRST}(\beta_2 \text{FOLLOW}(A)) = \Phi$. Dacă ar exista în tabela construită $A \in V$ și $a \in V \cup \{ \# \}$ astfel încât $M(A, a) \supseteq \{ (\beta, p), (\beta', q) \}$, atunci, dacă $a \in T$ rezultă

$$a \in \text{FIRST}(\beta \text{FOLLOW}(A)) \cap \text{FIRST}(\beta' \text{FOLLOW}(A))$$

ceea ce contrazice ipoteza. Dacă $a = \#$, atunci ϵ ar fi comun în $\text{FIRST}(\beta \text{FOLLOW}(A))$ și $\text{FIRST}(\beta' \text{FOLLOW}(A))$, din nou absurd. Deci M îndeplinește condiția enunțată. În mod analog se dovedește implicația inversă.

□

3.3.6 Analizorul sintactic LL(1)

Parserul (analizorul sintactic) LL(1) este un *parser top-down* în care tranzițiile sunt dictate de tabela de parsare. Așadar, un astfel de parser are ca și configurație inițială tripleta $(w\#, S\#, \epsilon)$ unde w este cuvântul de analizat, iar tranzițiile se descriu astfel:

1. $(u\#, A\gamma\#, \pi) \vdash (u\#, \beta\gamma\#, \pi r)$ dacă $M(A, 1:u\#) = (\beta, r)$. (operația *expandare*)
2. $(uv\#, u\gamma\#, \pi) \vdash (v\#, \gamma\#, \pi)$. (operația *potrivire*)
3. $(\#, \#, \pi) \vdash \text{acceptare}$ dacă $\pi \neq \epsilon$.
4. $(au\#, b\gamma\#, \pi) \vdash \text{eroare}$ dacă $a \neq b$.
5. $(u\#, A\gamma\#, \pi) \vdash \text{eroare}$ dacă $M(A, 1:u\#) = \text{eroare}$.

Teorema de corectitudine a parserului descendent general împreună cu teorema de caracterizare LL(1), a modului în care s-a definit tabela M, dovedesc corectitudinea parserului LL(1). Suntem acum în măsură să indicăm modul de implementare al unui analizor sintactic LL(1). Presupunem că dispunem de o bandă de intrare (fișier de intrare) din care, cu o funcție `getnext()`, se obține caracterul (tokenul) următor. De asemenea dispunem de o stivă cu funcțiile specifice `pop()` și `push()`. Producțiile care se aplică pe parcursul analizei le scriem într-o bandă (fișier) de ieșire cu funcția `write()`. Atunci algoritmul de analiză sintactică se descrie astfel:

Algoritmul 3.3.6 (Parser LL(1))

Intrare: Gramatica $G = (V, T, S, P)$.

Tabela de analiză LL(1) notată M.

Cuvântul de intrare $w\#$.

Ieșire: Analiza sintactică stângă π a lui w dacă $w \in L(G)$,
eroare în caz contrar.

Metoda: Sunt implementate tranzițiile 1 – 5 folosind o stivă St.

```

1. St.push(#), St.push(S) // St = S#
2. a = getnext(),  $\pi = \epsilon$ ;
3. do {
4.   X = St.pop();
5.   if(X == a)
6.     if(X != '#') getnext();
7.     else{
8.       if ( $\pi \neq \epsilon$ ) {write("acceptare"); exit(0);}
9.       else {write("eroare"); exit(1);}
10.    } //endelse
11.   else{
12.     if(X ∈ T) {write("eroare"); exit(1);}
13.     else{
14.       if(M(X,a) == "eroare")
15.         {write("eroare"); exit(1);}
16.       else {
17.         // M(X,a)=( $\beta, r$ ),  $r=X \rightarrow \beta$ ,  $\beta=Y_1Y_2...Y_n$ 
18.         // $\beta$  inlocuiește pe X în stivă
19.         for(k = n; k>0; --k) push( $Y_k$ );
20.         write(r); //se adaugă r la  $\pi$ 
21.       } //endelse
22.     } //endelse
23.   } while(1);

```

Exemplul 3.3.4 Să reluăm gramatica din exemplele precedente:

- | | |
|-----------------------------|--|
| 1. $S \rightarrow E$ | 5. $B \rightarrow \text{begin SC end}$ |
| 2. $S \rightarrow B$ | 6. $C \rightarrow \epsilon$ |
| 3. $E \rightarrow \epsilon$ | 7. $C \rightarrow ;SC$ |
| 4. $B \rightarrow a$ | |

Mulțimile FIRST și FOLLOW sunt date în tabelul următor:

X	FIRST(X)	FOLLOW(X)
S	a begin ϵ	end ; ϵ
E	ϵ	end ; ϵ
B	a begin	end ; ϵ
C	; ϵ	end

Tabela de analiză LL(1) pentru această gramatică este dată mai jos:

M	a	begin	end	;	#
S	(B, 2)	(B, 2)	(E, 1)	(E, 1)	(E, 1)
E	eroare	eroare	(ε, 3)	(ε, 3)	(ε, 3)
B	(a, 4)	(begin SC end, 5)	eroare	eroare	eroare
C	eroare	eroare	(ε, 6)	(;SC, 7)	eroare

Se vede din tabelă că această gramatică este LL(1) (conform teoremei 3.3.13). Iată și două exemple de analiză, unul pentru cuvântul `begin a;;a end` care este din limbajul generat de gramatica dată, iar altul pentru `begin aa end` care nu este corect.

PAS	INTRARE	STIVA	OPERAȚIE	IEȘIRE
1	begin a;;a end#	S#	<i>expandare</i>	
2	begin a;;a end#	B#	<i>expandare</i>	2
3	begin a;;a end#	begin SC end#	<i>potrivire</i>	5
4	a;;a end#	SC end#	<i>expandare</i>	
5	a;;a end#	BC end#	<i>expandare</i>	2
6	a;;a end#	aC end#	<i>potrivire</i>	4
7	;;a end#	C end#	<i>expandare</i>	
8	;;a end#	;SC end#	<i>potrivire</i>	7
9	;a end#	SC end#	<i>expandare</i>	
10	;a end#	EC end#	<i>expandare</i>	1
11	;a end#	C end#	<i>expandare</i>	3
12	;a end#	;SC end#	<i>potrivire</i>	7
13	a end#	SC end#	<i>expandare</i>	
14	a end#	BC end#	<i>expandare</i>	2
15	a end#	aC end#	<i>potrivire</i>	4
16	end#	C end#	<i>expandare</i>	
17	end#	end#	<i>potrivire</i>	6
18	#	#	<i>acceptare</i>	

PAS	INTRARE	STIVA	OPERAȚIE	IEȘIRE
1	begin aa end#	S#	<i>expandare</i>	
2	begin aa end#	B#	<i>expandare</i>	2
3	begin aa end#	begin SC end#	<i>potrivire</i>	5
4	aa end#	SC end#	<i>expandare</i>	
5	aa end#	BC end#	<i>expandare</i>	2
6	aa end#	aC end#	<i>potrivire</i>	4
7	a end#	C end#	<i>eroare</i>	

3.3.7 Eliminarea recursiei stângi

O gramatică stâng recursivă nu este LL(k) pentru nici un număr k. Este cazul, spre exemplu, al gramaticii care generează expresiile aritmetice:

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

Prin eliminarea recursiei stângi, există posibilitatea obținerii unei gramatici din clasa LL. Acesta este un procedeu utilizat frecvent în proiectarea analizatoarelor sintactice. Să reamintim că o gramatică G este stâng recursivă dacă există un neterminal A pentru care $A \Rightarrow^+ A\alpha$. Spunem că un neterminal A este stâng recursiv imediat dacă există o

producție $A \rightarrow A\alpha$ în P . Să indicăm mai întâi un procedeu de eliminare a recursiei stângi imediate.

Lema 3.3.7 Fie $G = (V, T, S, P)$ o gramatică pentru care A este stâng recursiv imediat. Să considerăm toate A – producțiile gramaticii, fie ele $A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n$ cele cu recursie și $A \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_m$ cele fără recursie imediată (părțile drepte nu încep cu A). Există o gramatică $G' = (V', T, S, P')$ echivalentă cu G , în care A nu este stâng recursiv imediat.

Demonstrație Să considerăm un nou neterminal $A' \notin V$ și fie $V' = V \cup \{A'\}$ iar P' se obține din P astfel:

- Se elimină din P toate A -producțiile;
- Se adaugă la P următoarele producții:

$$A \rightarrow \beta_i A' \quad 1 \leq i \leq m.$$

$$A' \rightarrow \alpha_j A' \mid \varepsilon \quad 1 \leq j \leq n.$$

- Notăm cu P' mulțimea de producții obținută.

Să observăm că A nu mai este stâng recursiv în G' . Noul simbol neterminal introdus, notat cu A' , este de data aceasta *drept recursiv*. Acest lucru însă nu este un inconvenient pentru analiza sintactică descendentă.

Să observăm că, pentru orice derivare extrem stânga în gramatica G , de forma:

$$A \Rightarrow A\alpha_{k_n} \Rightarrow A\alpha_{k_{n-1}}\alpha_{k_n} \Rightarrow \dots \Rightarrow A\alpha_{k_1}\alpha_{k_2}\dots\alpha_{k_n} \Rightarrow \beta_h\alpha_{k_1}\alpha_{k_2}\dots\alpha_{k_n}$$

există în gramatica G' o derivare extrem dreapta de forma:

$$A \Rightarrow \beta_h A' \Rightarrow \beta_h\alpha_{k_1}A' \Rightarrow \dots \Rightarrow \beta_h\alpha_{k_1}\alpha_{k_2}\dots\alpha_{k_n}A' \Rightarrow \beta_h\alpha_{k_1}\alpha_{k_2}\dots\alpha_{k_n}.$$

De asemenea, afirmația reciprocă este valabilă.

Aceste observații conduc la concluzia că cele două gramatici sunt echivalente. □

Exemplul 3.3.5 În gramatica expresiilor aritmetice:

$$E \rightarrow E+T \mid E-T \mid -T \mid T$$

$$T \rightarrow T*F \mid T/F \mid F$$

$$F \rightarrow (E) \mid a$$

se fac transformările următoare:

- Producțiile $E \rightarrow E+T \mid E-T \mid -T \mid T$ se înlocuiesc cu:

$$E \rightarrow TE' \mid -TE', \quad E' \rightarrow +TE' \mid -TE' \mid \varepsilon$$

- Producțiile $T \rightarrow T*F \mid T/F \mid F$ se înlocuiesc cu:

$$T \rightarrow FT', \quad T' \rightarrow *FT' \mid /FT' \mid \varepsilon$$

Se obține astfel gramatica echivalentă:

$$E \rightarrow TE' \mid -TE'$$

$$E' \rightarrow +TE' \mid -TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid /FT' \mid \varepsilon$$

$$F \rightarrow (E) \mid a$$

Aplicând algoritmii pentru determinarea mulțimilor FIRST și FOLLOW se obține:

X	FIRST(X)	FOLLOW(X)
E	(a -	ϵ)
E'	+ - ϵ	ϵ)
T	(a	+ - ϵ)
T'	* / ϵ	+ - ϵ)
F	(a	* / + - ϵ)

Tabela de parsare pentru această gramatică este dată mai jos.

M	a	+	-	*	/	()	#
E	(TE', 1)	eroare	(-TE', 2)	eroare	eroare	(TE', 1)	eroare	eroare
E'	eroare	(+TE', 3)	(-TE', 4)	eroare	eroare	eroare	(ϵ , 5)	(ϵ , 5)
T	(FT', 6)	eroare	eroare	eroare	eroare	(FT', 6)	eroare	eroare
T'	eroare	(ϵ , 9)	(ϵ , 9)	(*FT', 7)	(/FT', 8)	eroare	(ϵ , 9)	(ϵ , 9)
F	(a, 11)	eroare	eroare	eroare	eroare	((E), 10)	eroare	eroare

Se observă de aici că gramatica este LL(1). Să urmărim analiza sintactică pentru expresia $-(a+a)*a$.

PASUL	INTRARE	STIVA	OPERAȚIE	IEȘIRE
1.	$-(a+a)*a\#$	E#	expandare	
2.	$-(a+a)*a\#$	-TE' #	potrivire	4
3.	$(a+a)*a\#$	TE' #	expandare	
4.	$(a+a)*a\#$	FT' E' #	expandare	6
5.	$(a+a)*a\#$	(E) T' E' #	potrivire	10
6.	$a+a)*a\#$	E) T' E' #	expandare	
7.	$a+a)*a\#$	TE') T' E' #	expandare	1
8.	$a+a)*a\#$	FT' E') T' E' #	expandare	6
9.	$a+a)*a\#$	aT' E') T' E' #	potrivire	11
10.	$+a)*a\#$	T' E') T' E' #	expandare	
11.	$+a)*a\#$	E') T' E' #	expandare	9
12.	$+a)*a\#$	+TE') T' E' #	potrivire	3
13.	$a)*a\#$	TE') T' E' #	expandare	
14.	$a)*a\#$	FT' E') T' E' #	expandare	6
15.	$a)*a\#$	aT' E') T' E' #	potrivire	11
16.	$) *a\#$	T' E') T' E' #	expandare	
17.	$) *a\#$	E') T' E' #	expandare	9
18.	$) *a\#$) T' E' #	potrivire	5
19.	$*a\#$	T' E' #	expandare	
20.	$*a\#$	*FT' E' #	potrivire	7
21.	$a\#$	FT' E' #	expandare	
22.	$a\#$	aT' E' #	potrivire	11
23.	$\#$	T' E' #	expandare	
24.	$\#$	E' #	expandare	9
25.	$\#$	#	acceptare	5

Să arătăm cum eliminăm în general recursia stângă într-o gramatică. Vom face ipoteza că gramatica G nu are ε - producții și nu are cicluri de forma $A \Rightarrow^+ A$. Aceasta nu este de fapt o restricție: este știut faptul că, pentru orice gramatică G există G' fără ε -producții și fără cicluri astfel încât $L(G') = L(G) - \{ \varepsilon \}$. (vezi [Gri86], [Juc99]). Așadar are loc următoarea teoremă:

Teorema 3.3.13 Pentru orice gramatică G fără ε - producții și fără cicluri, există o gramatică echivalentă G' care nu este stâng recursivă.

Demonstrație Să considerăm următorul algoritm:

Intrare: Gramatica $G = (V, T, S, P)$ fără ε -producții, fără cicluri.

Ieșire: Gramatica G' fără recursie stângă astfel ca $L(G)=L(G')$.

Metoda: După ordonarea neterminalilor se fac transformări ce nu schimbă limbajul gramaticii, inclusiv prin eliminarea recursiei imediate

```

1. Se consideră  $V = \{A_1, A_2, \dots, A_n\}$ 
2. for (i=1; i ≤ n; ++i) {
3.     for (j=1; j ≤ i - 1; ++j) {
4.         for ( $A_i \rightarrow A_j\gamma \in P$ ) {
5.              $P = P - \{ A_i \rightarrow A_j\gamma \};$ 
6.             for ( $A_j \rightarrow \beta \in P$ )  $P = P \cup \{ A_i \rightarrow \beta\gamma \};$ 
7.         } //endfor
8.     } //endfor
9. Se elimina recursia imediata pentru  $A_i$ ;
10. } //endfor
11.  $V' = V$ ;  $P' = P$ ; // V s-a modificat la 7 iar P la 4 - 6
```

Să observăm că transformările din 4 - 7 nu modifică limbajul generat de gramatică încât $L(G') = L(G)$. Să arătăm acum că G' nu are recursie stângă. Acest lucru rezultă din faptul că, după iterația $i - 1$ a buclei 2, fiecare producție de forma $A_k \rightarrow A_l\alpha$ pentru $k < i$ are proprietatea că $l > k$. Aceasta datorită transformărilor 4 - 6 și eliminării recursiei stângi imediate în linia 7: după execuția buclei 3 se obțin producții $A_i \rightarrow A_m\alpha$ cu $m \geq i$, iar dacă $m = i$, producția $A_i \rightarrow A_i\alpha$ se elimină în linia 7. Așadar, gramatica obținută este fără recursie stângă.

□

Exemplul 3.3.6 Fie gramatica:

$$S \rightarrow Aa \mid b \quad A \rightarrow Ac \mid Sd \mid a$$

Considerăm neterminalii în ordonați (S, A) . Aplicând bucla 2 pentru $i = 1$ nu se face nici o modificare deoarece în producția $S \rightarrow Aa$ ordinea neterminalilor este corectă. Pentru $i = 2$ trebuie modificată producția $A \rightarrow Sd$ (pentru că S este *mai mic* decât A). A - producțiile se modifică astfel:

$$A \rightarrow Ac \mid Aad \mid bd \mid a$$

Gramatica, în acest moment, este

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Aad \mid bd \mid a$$

și conține doar recursivitate imediată (A). Aplicând algoritmul de eliminare a recursiei stângi imediate (vezi Lema 3.3.7), obținem:

$$S \rightarrow Aa \mid b \quad A \rightarrow bdA' \mid aA' \quad A' \rightarrow cA' \mid adA' \mid \varepsilon$$

Se poate observa că rezultatul nu este o gramatică LL(1) deoarece pentru primele două producții $\text{FIRST}(Aa) \cap \text{FIRST}(a) = \{b\}$. Așadar eliminarea recursiei stângi nu conduce neapărat la proprietatea LL(1).

4 Analiza sintactică în gramatici LR

Vom prezenta în acest capitol o tehnică eficientă de analiză sintactică ascendentă care este utilizată pentru o clasă largă de gramatici: gramaticile LR(k). Denumirea LR(k) vine de la: *Left to right scanning of the input, constructing a Rightmost derivation in reverse, using k symbols lookahead*. Această metodă este cu siguranță cea mai des utilizată metodă de analiză sintactică, din următoarele motive:

- se pot construi analizoare sintactice LR pentru recunoașterea tuturor construcțiilor din limbajele de programare care se pot descrie printr-o gramatică independentă de context;
- clasa limbajelor ce pot fi analizate sintactic cu analizoare LR(1) coincide cu clasa limbajelor de tip 2 deterministe;
- metoda de analiză LR este o metodă de tip *deplasare-reducere* relativ ușor de implementat și eficientă în același timp;
- un analizor LR poate detecta o eroare de sintaxă cel mai rapid posibil parcurgând șirul de intrare de la stânga la dreapta.

Dezavantajul principal al metodei este acela că determinarea tabelului de analiză necesită un volum mare de muncă; există însă generatoare de analizoare de tip LR, precum *yacc* sau *bison*, care produc un astfel de analizor.

Vom considera în continuare o gramatică $G = (V, T, S, P)$ redusă și gramatica augmentată $G' = (V', T', S', P')$ unde $P' = P \cup \{S' \rightarrow S\}$, $V' = V \cup \{S'\}$, S' fiind un simbol nou. Gramatica G' este echivalentă cu G și are proprietatea că simbolul de start nu apare în nici o parte dreaptă a producțiilor din P' , condiție esențială pentru studiul gramaticilor LR(k). Să mai facem observația că, pentru gramaticile care au proprietatea amintită (simbolul de start nu apare în nici o parte dreaptă a producțiilor), nu este necesară augmentarea.

4.1 Gramatici LR(k)

Definiția 5.1.1 O gramatică G se numește *gramatică LR(k)*, $k \geq 0$, dacă pentru orice două derivări de forma:

$$\begin{aligned} S' &\Rightarrow S \xRightarrow{dr}^* \alpha A u \xRightarrow{dr} \alpha \beta u = \gamma u \\ S' &\Rightarrow S \xRightarrow{dr}^* \alpha' A' u' \xRightarrow{dr} \alpha' \beta' u' = \alpha \beta v = \gamma v \end{aligned}$$

pentru care $k:u = k':v$, are loc: $\alpha = \alpha'$, $A = A'$, $\beta = \beta'$.

Să vedem care este semnificația acestor condiții în contextul analizei sintactice. Să ne amintim că un pas în analiza sintactică ascendentă înseamnă:

- determinarea, în $\gamma = \alpha\beta u$, a obiectului derivării β ;
- reducerea lui β la A .

Cum în algoritmul general de *deplasare-reducere*, $\alpha\beta$ se află în stivă iar u se află în banda de intrare, rezultă că o gramatică G este LR(k) dacă, în orice pas al analizei sintactice ascendente, obiectul derivării β și producția $A \rightarrow \beta$ cu care se face reducerea sunt determinate în mod unic prin inspectarea primelor k simboluri care au rămas de analizat (pe banda de intrare). Definiția pe care am dat-o pentru gramatici LR(k) este dinamică: condiția ca o gramatică să fie LR(k) trebuie să fie îndeplinită de orice două perechi de derivări. Cum, în general, numărul acestora este infinit, nu se poate verifica dacă o gramatică este LR(k) folosind definiția. Vom încerca să găsim condiții echivalente cu

definiția pentru cazurile $k = 0$ și $k = 1$. Să stabilim în continuare câteva proprietăți ale gramaticilor LR(k).

Teorema 5.1.1 Dacă G este gramatică LR(k), $k \geq 0$, atunci G este neambiguă.

Demonstrație Fie $G = (V, T, S, P)$ o gramatică LR(k) și să presupunem, prin reducere la absurd, că ea este ambiguă. Există atunci un cuvânt $w \in L(G)$ care are două derivări extrem drepte distincte:

$$S' \Rightarrow S = \gamma_0 \xRightarrow{dr} \gamma_1 \xRightarrow{dr} \dots \xRightarrow{dr} \gamma_n = w$$

$$S' \Rightarrow S = \delta_0 \xRightarrow{dr} \delta_1 \xRightarrow{dr} \dots \xRightarrow{dr} \delta_m = w$$

Să observăm că $\gamma_n = \delta_m = w$. Cele două derivări sunt presupuse distincte, deci există un număr i , $1 \leq i \leq \min(m, n)$ astfel încât $\gamma_{n-i} \neq \delta_{m-i}$, iar $\gamma_{n-i} = \delta_{m-i}$. Avem atunci:

$$S' \Rightarrow S \xRightarrow{dr} \gamma_{n-i} = \alpha A u \xRightarrow{dr} \alpha \beta u = \gamma_{n-i}$$

$$S' \Rightarrow S \xRightarrow{dr} \delta_{m-i} = \alpha' A' u' \xRightarrow{dr} \alpha' \beta' u' = \alpha \beta u = \delta_{m-i}$$

Dar G este LR(k) și pentru că are loc $k:u = k':u$ rezultă $\alpha = \alpha'$, $A = A'$, $\beta = \beta'$ ceea ce este o contradicție. Așadar, presupunerea că există un număr i , $1 \leq i \leq \min(m, n)$ astfel încât $\gamma_{n-i} \neq \delta_{m-i}$ este falsă. Rezultă $m = n$ și $\gamma_i = \delta_i$, $1 \leq i \leq n$.

□

4.2 O caracterizare a gramaticilor LR(0)

Caracterizarea pe care o vom da pentru gramaticile LR(0) folosește un automat finit determinist care joacă un rol esențial în construcția analizatoarelor de tip LR.

Definiția 5.2.1 Fie $G = (V, T, S, P)$ o gramatică independentă de context redusă. Să presupunem că simbolul \bullet nu este în Σ . Un *articol* pentru gramatica G este o producție $A \rightarrow \gamma$ în care s-a adăugat simbolul \bullet într-o anumită poziție din γ . Notăm un articol prin $A \rightarrow \alpha \bullet \beta$ dacă $\gamma = \alpha \beta$. Un articol în care \bullet este pe ultima poziție se numește *articol complet*.

Definiția 5.2.2 Un *prefix viabil* pentru gramatica $G = (V, T, S, P)$ este orice prefix al unui cuvânt $\alpha \beta$ dacă $S \xRightarrow{dr} \alpha A u \xRightarrow{dr} \alpha \beta u$. Dacă $\beta = \beta_1 \beta_2$ și $\phi = \alpha \beta_1$ spunem că articolul $A \rightarrow \beta_1 \bullet \beta_2$ este *valid* pentru *prefixul viabil* ϕ .

Exemplul 5.2.1 Fie gramatica $S \rightarrow A$, $A \rightarrow aAa \mid bAb \mid c$. Articolele acestei gramatici sunt:

$S \rightarrow \bullet A$, $S \rightarrow A \bullet$, $A \rightarrow \bullet aAa$, $A \rightarrow a \bullet Aa$, $A \rightarrow aA \bullet a$, $A \rightarrow aAa \bullet$,
 $A \rightarrow \bullet bAb$, $A \rightarrow b \bullet Ab$, $A \rightarrow bA \bullet b$, $A \rightarrow bAb \bullet$, $A \rightarrow \bullet c$, $A \rightarrow c \bullet$.

Dacă, în plus, avem și regula $A \rightarrow \epsilon$, atunci singurul articol corespunzător acesteia este articolul (complet) $A \rightarrow \bullet$. În următorul tabel sunt date și câteva exemple de articole valide pentru prefixe viabile:

Prefixul viabil	Articole valide	Derivarea corespunzătoare
ab	$A \rightarrow b \bullet Ab$	$S \Rightarrow A \Rightarrow aAa \Rightarrow abAba$
	$A \rightarrow \bullet aAa$	$S \Rightarrow A \Rightarrow aAa \Rightarrow abAba \Rightarrow abaAaba$
	$A \rightarrow \bullet bAb$	$S \Rightarrow A \Rightarrow aAa \Rightarrow abAba \Rightarrow abbAbba$

ε	$S \rightarrow \bullet A$	$S \Rightarrow A$
	$A \rightarrow \bullet bAb$	$S \Rightarrow A \Rightarrow bAb$
	$A \rightarrow \bullet c$	$S \Rightarrow A \Rightarrow c$

Lema 5.2.1 Fie $G = (V, T, S, P)$ o gramatică și $A \rightarrow \beta_1 \bullet B \beta_2$ un articol valid pentru prefixul viabil γ . Atunci, oricare ar fi producția $B \rightarrow \beta$, articolul $B \rightarrow \bullet \beta$ este valid pentru γ .

Demonstrație Dacă $A \rightarrow \beta_1 \bullet B \beta_2$ este valid pentru γ există derivarea:

$$S \xRightarrow{dr}^* \alpha A u \xRightarrow{dr} \alpha \beta_1 B \beta_2 u = \gamma B \beta_2 u$$

Dar G este redusă, deci există derivarea $\beta_2 \xRightarrow{dr}^* v$ și atunci putem scrie:

$$S \xRightarrow{dr}^* \alpha A u \xRightarrow{dr} \alpha \beta_1 B \beta_2 u \xRightarrow{dr}^* \alpha \beta_1 B v u \xRightarrow{dr} \alpha \beta_1 \beta v u = \gamma \beta v u$$

ceea ce înseamnă că $B \rightarrow \bullet \beta$ este valid pentru γ . □

Teorema 5.2.1 Gramatica $G = (V, T, S, P)$ este gramatică LR(0) dacă și numai dacă, oricare ar fi prefixul viabil γ , sunt îndeplinite condițiile:

1. nu există două articole complete valide pentru γ .
2. dacă articolul $A \rightarrow \beta \bullet$ este valid pentru γ , nu există nici un articol $B \rightarrow \beta_1 \bullet a \beta_2$, $a \in T$, valid pentru γ .

Demonstrație Să presupunem că G este gramatică LR(0). Atunci, conform definiției, oricare ar fi două derivări:

$$S' \Rightarrow S \xRightarrow{dr}^* \alpha A u \xRightarrow{dr} \alpha \beta u = \delta u \quad (1)$$

$$S' \Rightarrow S \xRightarrow{dr}^* \alpha' A' u' \xRightarrow{dr} \alpha' \beta' u' = \delta v \quad (2)$$

are loc: $\alpha = \alpha'$, $A = A'$, $\beta = \beta'$.

Să presupunem, prin reducere la absurd, că există un prefix viabil γ pentru care nu sunt îndeplinite simultan 1 și 2. Dacă 1 nu are loc, există articolele $B \rightarrow \beta \bullet$ și $B' \rightarrow \beta' \bullet$, diferite, valide pentru γ . Asta înseamnă că:

$$S' \Rightarrow S \xRightarrow{dr}^* \varphi B x \xRightarrow{dr} \varphi \beta x = \gamma x \quad (3)$$

$$S' \Rightarrow S \xRightarrow{dr}^* \varphi' B' x' \xRightarrow{dr} \varphi' \beta' x' = \gamma x \quad (4)$$

Ori, (3) și (4) contrazic faptul că G este LR(0).

Dacă 2 nu are loc, există articolele $B \rightarrow \beta \bullet$, $C \rightarrow \beta_1 \bullet a \beta_2$ valide pentru γ , adică are loc derivarea (3) și derivarea:

$$S' \Rightarrow S \xRightarrow{dr}^* \psi C y \xRightarrow{dr} \psi \beta_1 a \beta_2 y = \gamma a \beta_2 y \quad (5)$$

Dacă în (5) $\beta_2 \in T^*$ atunci din (3), (5) și faptul că G este LR(0), rezultă că $\varphi = \psi$, $B = C$, $\beta = \beta_1 a \beta_2$. Dar $\gamma = \varphi \beta = \psi \beta_1$ ceea ce înseamnă că $\varphi \beta_1 a \beta_2 = \varphi \beta_1$ ceea ce este imposibil deoarece $a \in T$.

Dacă în (5) β_2 conține măcar o variabilă, cum G este redusă, există derivarea:

$$\beta_2 \xRightarrow{dr}^* u_1 D u_3 \xRightarrow{dr} u_1 u_2 u_3$$

Aplicând acest lucru în (5) se obține:

$$S' \Rightarrow S \xRightarrow{dr}^* \gamma a \beta_2 y \xRightarrow{dr}^* \gamma a u_1 D u_3 y \xRightarrow{dr}^* \gamma a u_1 u_2 u_3 y \quad (6)$$

Din (3) și (6), pentru că G este $LR(0)$, rezultă $\varphi = \gamma a u_1$, $B = D$, $\beta = u_2$. Dar din derivarea (3) are loc $\gamma = \varphi \beta$ și egalitatea $\varphi = \varphi \beta a u_1$ este imposibilă pentru că a este terminal (nu poate fi nul). Așadar, o implicație a teoremei este dovedită.

Să presupunem acum că oricare ar fi prefixul viabil γ au loc condițiile 1 și 2 din teoremă și, prin reducere la absurd, G nu este $LR(0)$. Asta înseamnă că există două derivări de forma (1) și (2) dar nu este adevărat că $\alpha = \alpha'$, $A = A'$ și $\beta = \beta'$. Fără a restrânge generalitatea, presupunem că în (1) și (2) avem îndeplinită condiția $|\delta| = |\alpha\beta| \leq |\alpha'\beta'|$. Distingem două cazuri:

Cazul 1: $|\alpha'| \leq |\delta|$. Schematic acest lucru arată astfel:

α	β	u
δ		
α'	β'	u'
δ		v

Rezultă de aici că $\beta' = \beta'_1 \beta'_2$, $v = \beta'_2 u'$, $\delta = \alpha' \beta'_1$ și $\beta'_2 \in T^*$. Atunci, din (1) rezultă că $A \rightarrow \beta \bullet$ este articol valid pentru δ iar din (2) rezultă că articolul $A' \rightarrow \beta'_1 \bullet \beta'_2$ este valid de asemenea pentru δ . Asta contrazice 1 dacă $\beta'_2 = \varepsilon$ respectiv 2 dacă $\beta'_2 \neq \varepsilon$. Așadar, în acest caz, presupunerea că G nu este $LR(0)$ este falsă.

Cazul 2: $|\alpha'| > |\delta|$. Din nou schematic condiția arată astfel:

α	β	u
δ		
α'	β'	u'
δ		v

Asta înseamnă că $\alpha' = \delta u_1$, $v = u_1 \beta' u' \in T^*$ ($|u_1| \geq 1$). În derivarea (2), punem în evidență prima formă propozițională care are prefixul δ ; aceasta există deoarece S nu are prefixul δ , iar $\alpha' \beta' u' = \delta v$ are prefixul δ . Așadar (2) devine:

$$S' \Rightarrow S \xRightarrow{dr}^* \alpha_1 A_1 u_1 \xRightarrow{dr}^* \alpha_1 \beta_1 u_1 = \alpha_1 \beta'_1 \beta''_1 u_1 = \delta \beta''_1 u_1 \xRightarrow{dr}^* \delta v \quad (7)$$

Din (7) rezultă că $A_1 \rightarrow \beta'_1 \bullet \beta''_1$ este articol valid pentru δ (iar din (1) rezultă că $A \rightarrow \beta \bullet$ este valid de asemenea pentru δ).

Distingem și aici trei subcazuri:

Cazul 2.1: $\beta''_1 = \varepsilon$, contradicție: condiția 1 din teoremă.

Cazul 2.2: $\beta''_1 = a \beta_2$, $a \in T$. Atunci $A_1 \rightarrow \beta'_1 \bullet a \beta_2$, $A \rightarrow \beta \bullet$ sunt valide pentru δ , ceea ce contrazice 2 din teoremă.

Cazul 2.3: $\beta''_1 = B \beta_2$, $B \in V$. Să observăm atunci că există un articol $C \rightarrow \bullet a \varphi$ valid pentru δ (eventual $C = B$) pentru că gramatica este redusă și există măcar o producție a cărei parte dreaptă începe cu un terminal (se folosește și Lema 5.2.1). Și în acest caz se contrazice 2 din teoremă. În concluzie, G este $LR(0)$. □

Teorema precedentă ne îndreptățește să cercetăm mai îndeaproape mulțimea prefixelor viabile pentru o gramatică G . Are loc:

Teorema 5.2.2 Fie $G = (V, T, S, P)$ o gramatică independentă de context. Mulțimea prefixelor viabile pentru gramatica G este limbaj regulat.

Demonstrație Să considerăm gramatica G' îmbogățită cu $S' \rightarrow S$. Vom construi un automat (nedeterminist cu ε -tranziții) care recunoaște mulțimea prefixelor viabile ale lui G . Considerăm automatul $M = (Q, \Sigma, \delta, q_0, Q)$, unde Q este mulțimea articolelor gramaticii G' , $\Sigma = V \cup T$, $q_0 = S' \rightarrow \bullet S$ iar funcția de tranziție $\delta : Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow 2^Q$ este definită astfel:

$$\delta(A \rightarrow \alpha \bullet B \beta, \varepsilon) = \{B \rightarrow \bullet \gamma \mid B \rightarrow \gamma \in P\}.$$

$$\delta(A \rightarrow \alpha \bullet X \beta, X) = \{A \rightarrow \alpha X \bullet \beta\}, X \in \Sigma.$$

$$\delta(A \rightarrow \alpha \bullet a \beta, \varepsilon) = \Phi, \forall a \in T.$$

$$\delta(A \rightarrow \alpha \bullet X \beta, Y) = \Phi, \forall X, Y \in \Sigma \text{ cu } X \neq Y.$$

Așadar, automatul are ε -tranziții numai din articolele (stările) care au un neterminal după punct iar X - tranziții, $X \in \Sigma$, din articolele care au X (același simbol) după punct. Să dovedim că are loc:

$A \rightarrow \alpha \bullet \beta \in \hat{\delta}(q_0, \gamma) \Leftrightarrow \gamma$ este prefix viabil și $A \rightarrow \alpha \bullet \beta$ este valid pentru γ .

Să presupunem mai întâi că γ este prefix viabil și $A \rightarrow \alpha \bullet \beta$ este valid pentru γ . Asta înseamnă că are loc o derivare de forma:

$$S \xRightarrow[\text{dr}]{*} \varphi A u \xRightarrow[\text{dr}]{*} \varphi \alpha \beta u = \gamma \beta u$$

Vom dovedi, prin inducție după k , lungimea derivării, că $A \rightarrow \bullet \alpha \beta \in \hat{\delta}(q_0, \varphi)$. Se deduce apoi că $\hat{\delta}(q_0, \varphi \alpha)$ conține articolul $A \rightarrow \alpha \bullet \beta$ (din definiția automatului și a extensiei lui δ).

$k = 1$: În acest caz derivarea este $S \xRightarrow[\text{dr}]{*} \alpha \beta$, deci $S \rightarrow \alpha \beta \in P$, $\varphi = u = \varepsilon$ și, după definiția lui δ are loc $S \rightarrow \bullet \alpha \beta \in \hat{\delta}(q_0, \varepsilon)$.

$k > 1$: În derivarea de lungime k punem în evidență pasul la care s-a introdus neterminalul A :

$$S \xRightarrow[\text{dr}]{*} \varphi_1 B u_2 \xRightarrow[\text{dr}]{*} \varphi_1 \alpha_1 A \beta_1 u_2 \xRightarrow[\text{dr}]{*} \varphi_1 \alpha_1 A u_1 u_2 = \varphi A u \xRightarrow[\text{dr}]{*} \varphi \alpha \beta u = \gamma \beta u$$

unde $B \rightarrow \alpha_1 A u_2 \in P$, $\beta_1 \xRightarrow[\text{dr}]{*} u_1$, $u_1 u_2 = u$, $\varphi_1 \alpha_1 = \varphi$. E clar că derivarea

$$S \xRightarrow[\text{dr}]{*} \varphi_1 B u_2 \xRightarrow[\text{dr}]{*} \varphi_1 \alpha_1 A \beta_1 u_2$$

are lungimea cel mult $k - 1$ și, după ipoteza inductivă, are loc $B \rightarrow \bullet \alpha_1 A u_2 \in \hat{\delta}(q_0, \varphi_1)$ iar din definiția extensiei lui δ avem și $B \rightarrow \alpha_1 \bullet A u_2 \in \hat{\delta}(q_0, \varphi_1 \alpha_1)$. Atunci $A \rightarrow \bullet \alpha \beta \in \delta(B \rightarrow \alpha_1 \bullet A u_2, \varepsilon)$ (definiția lui δ) și $A \rightarrow \bullet \alpha \beta \in \hat{\delta}(q_0, \varphi_1 \alpha_1) = \hat{\delta}(q_0, \varphi)$ (definiția lui $\hat{\delta}$). Să presupunem acum că $A \rightarrow \alpha \bullet \beta \in \hat{\delta}(q_0, \gamma)$. Atunci, în graful automatului M există măcar un drum de la q_0 la starea $A \rightarrow \alpha \bullet \beta$ etichetat cu $\gamma = \varphi \alpha$ (drum ce conține eventual și arce etichetate ε). Considerăm drumul de lungime minimă (unul din ele dacă sunt mai multe) cu această proprietate. Fie k lungimea acestuia. Demonstrăm prin inducție după k faptul că γ este prefix viabil pentru G iar articolul $A \rightarrow \alpha \bullet \beta$ este valid pentru γ .

$k = 1$: Drumul este de la $q_0 = S' \rightarrow \bullet S$ la o stare $S \rightarrow \bullet \beta$ sau la $S' \rightarrow S \bullet$. În primul caz $\gamma = \varepsilon$ care este prefix viabil iar $A \rightarrow \alpha \bullet \beta = S \rightarrow \bullet \beta$ este valid pentru ε pentru că $S' \xRightarrow[\text{dr}]{*} S$

$\Rightarrow_{dr} \beta$. În al doilea caz $\gamma = S$ este prefix viabil iar $A \rightarrow \alpha \bullet \beta = S' \rightarrow S \bullet$ este valid pentru S

pentru că $S' \Rightarrow_{dr} S$.

$k > 1$: În funcție de ultimul arc din drum distingem două cazuri:

Cazul 1: Ultimul arc este etichetat cu $X \in \Sigma$. Atunci el are forma din figura 5.2.1. Așadar $\gamma = \varphi \alpha_1 X = \gamma_1 X$.

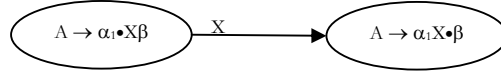


Figura 5.2.1

După ipoteza inductivă, γ_1 este prefix viabil iar articolul $A \rightarrow \alpha_1 \bullet X \beta$ este valid pentru el. De aici rezultă existența unei derivări

$$S \xRightarrow{dr} \varphi A u \xRightarrow{dr} \varphi \alpha_1 X \beta u = \gamma_1 X \beta u = \gamma \beta u$$

de unde rezultă faptul că γ este prefix viabil iar $A \rightarrow \alpha_1 X \bullet \beta$ este valid pentru γ .

Cazul 2: Ultimul arc al drumului este etichetat cu ε . Atunci $\alpha = \varepsilon$ iar arcul este de forma celui din figura 5.2.2.

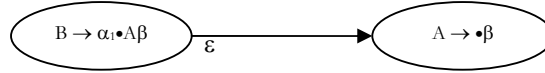


Figura 5.2.2

După ipoteza inductivă $B \rightarrow \alpha_1 \bullet A \beta_1$ este valid pentru γ , adică există derivarea:

$$S \xRightarrow{dr} \varphi B u \xRightarrow{dr} \varphi \alpha_1 A \beta_1 u = \gamma A \beta u$$

Dacă aplicăm $\beta \Rightarrow_{dr} v$ (gramatica este redusă) în derivarea de mai sus obținem:

$$S \xRightarrow{dr} \varphi B u \xRightarrow{dr} \varphi \alpha_1 A \beta_1 u = \gamma A \beta_1 u \xRightarrow{dr} \gamma A v u \xRightarrow{dr} \gamma \beta v u$$

ceea ce dovedește că și $A \rightarrow \bullet \beta$ este valid pentru γ . Cu aceasta demonstrația teoremei este încheiată. \square

Exemplul 5.2.2 Fie gramatica $S' \rightarrow S$, $S \rightarrow aSa \mid bSb \mid c$. Automatul care recunoaște prefixele viabile este reprezentat în figura 5.2.3. Arcele neetichetate constituie ε -tranziții. Condiția necesară și suficientă ca o gramatică G să fie $LR(0)$, dată în Teorema 5.2.1, se “traduce” relativ la automatul construit în Teorema 5.2.2 astfel:

1. Nu există în graful automatului două noduri cu etichete articole complete, respectiv $A \rightarrow \alpha \bullet$ și $B \rightarrow \beta \bullet$, și câte un drum de la q_0 la fiecare din acestea etichetate cu același cuvânt γ .
2. Nu există în graful automatului două noduri ce au etichete $A \rightarrow \alpha \bullet$ respectiv $B \rightarrow \beta \bullet a \beta'$, și câte un drum de la q_0 la fiecare din acestea etichetate cu același cuvânt γ ;

Aceste două condiții sunt greoi de verificat pe automatul (nedeterminist) cu ε -tranziții construit. Pentru acest automat, există un automat determinist echivalent, fie acesta $\mathcal{M} = (\mathcal{T}, \Sigma, g, t_0, \mathcal{T})$ unde $\mathcal{T} \subseteq 2^Q$ (o stare pentru automatul determinist este o submulțime de articole), t_0 conține toate articolele accesibile prin ε -tranziții din $q_0 = S' \rightarrow \bullet S$, iar $g: \mathcal{T} \times \Sigma \rightarrow \mathcal{T}$ este parțial definită (vezi 2.2.2). Cum M recunoaște mulțimea prefixelor viabile și articolul $A \rightarrow \alpha \bullet \beta$ este în $\hat{\delta}(q_0, \gamma)$ dacă și numai dacă $A \rightarrow \alpha \bullet \beta$ este valid pentru

prefixul viabil γ , înseamnă că, în automatul determinist, o stare $t \in \mathcal{T}$ conține toate articolele valide pentru o anumită mulțime de prefixe viabile și numai pe acestea. Așadar, dacă am construit \mathcal{M} , automatul determinist echivalent cu M , condițiile necesare și suficiente pentru ca G să fie gramatică LR(0) sunt:

1. Orice stare $t \in \mathcal{T}$ conține cel mult un articol complet.
2. Pentru orice $t \in \mathcal{T}$, dacă t conține un articol complet, atunci t nu conține nici un articol cu terminal după punct.

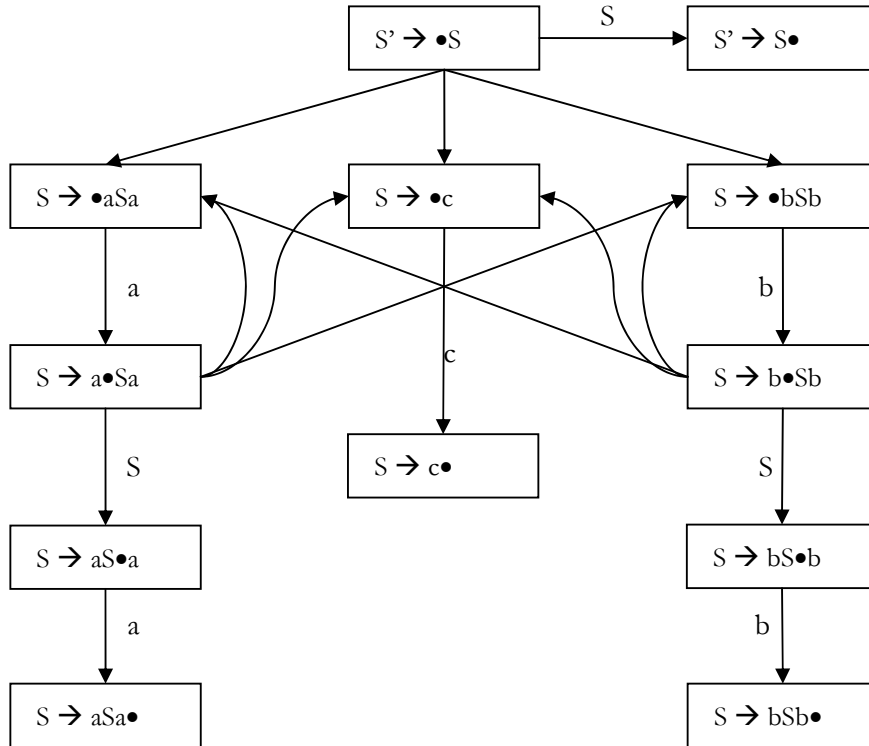


Figura 5.2.3

Aceste două condiții sunt ușor de verificat dacă \mathcal{M} este construit. \mathcal{M} se construiește aplicând algoritmul de trecere de la un automat (nedeterminist) cu ε -tranziții la automatul determinist echivalent (vezi Teorema 2.2.2). Pentru că această problemă este deosebit de importantă pentru analiza sintactică de tip LR, vom indica în continuare un algoritm pentru obținerea automatului \mathcal{M} , numit *automatul LR(0)* al gramaticii G . Mai întâi vom descrie o procedură, numită *închidere(t)* care are la intrare o mulțime t de articole (stări ale lui M) și obține la ieșire toate stările lui M accesibile din t prin ε .

Algoritmul 5.2.1 (procedura închidere(t))

Intrare: Gramatica $G = (V, T, S, P)$;
Mulțimea t de articole din gramatica G ;

Ieșire: $t' = \text{închidere}(t) = \{q \in Q \mid \exists p \in t, q \in \hat{\delta}(p, \varepsilon)\} = \hat{\delta}(t, \varepsilon)$

Metoda:

1. $t' = t$; flag = true;
2. while(flag) {
3. flag = false;
4. for ($A \rightarrow \alpha \bullet B \beta \in t'$) {

```

5.      for ( B → γ ∈ P )
6.          if ( B → •γ ∉ t' ) {
7.              t' = t' ∪ {B → •γ};
8.              flag = true;
              }//endif
          }//endforB
      }//endforA
  }//endwhile
9. return t';

```

Să observăm că algoritmul precedent nu folosește automatul M . Legătura cu acest automat este dată de:

Lema 5.2.2 Dacă $M = (Q, \Sigma, \delta, q_0, Q)$ este automatul cu ε -tranziții al gramaticii G și $t \subseteq Q$, atunci, în urma aplicării algoritmului precedent, se obține $t' = \hat{\delta}(t, \varepsilon)$.

Demonstrație Fie $t' = \text{închidere}(t)$. Să arătăm mai întâi că $t' \subseteq \hat{\delta}(t, \varepsilon)$. Fie articolul $A \rightarrow \alpha \bullet \beta \in t'$. Procedăm prin inducție asupra momentului când a fost adăugat acest articol la t' . Dacă acest lucru a fost făcut la linia 1, rezultă că $A \rightarrow \alpha \bullet \beta \in t$ evident $t \subseteq \hat{\delta}(t, \varepsilon)$ încât, $A \rightarrow \alpha \bullet \beta \in \hat{\delta}(t, \varepsilon)$. Dacă adăugarea s-a făcut la linia 7, înseamnă că în t' , în acel moment, există $B \rightarrow \gamma \bullet A \varphi$ care a fost adăugat anterior iar $\alpha = \varepsilon$. După ipoteza inductivă, $B \rightarrow \gamma \bullet A \varphi \in \hat{\delta}(t, \varepsilon)$ iar din definiția lui δ (Teorema 5.2.2) rezultă că $A \rightarrow \bullet \beta \in \delta(B \rightarrow \gamma \bullet A \varphi, \varepsilon)$ încât $A \rightarrow \bullet \beta$ este în $\hat{\delta}(t, \varepsilon)$.

Invers, fie $A \rightarrow \alpha \bullet \beta \in \hat{\delta}(t, \varepsilon)$. Prin inducție asupra lungimii drumului de la o stare din t la $A \rightarrow \alpha \bullet \beta$ se arată ușor că $A \rightarrow \alpha \bullet \beta$ se adaugă la t' (în linia 1 sau în linia 7).

□

Să descriem acum algoritmul pentru construcția automatului LR(0) pentru gramatica G .

Algoritmul 5.2.2 Automatul LR(0)

Intrare: Gramatica $G = (V, T, S, P)$ la care s-a adăugat $S' \rightarrow S$;

Ieșire: Automatul determinist $\mathcal{M} = (T, \Sigma, g, t_0, T)$ echivalent cu M .

Metoda:

```

1. t0=închidere(S'→•S); T={t0}; marcat(t0)=false;
2. while(∃t∈T && !marcat(t)) { // marcat(t) = false
3.     for( X ∈ Σ ) {
4.         t' = Φ;
5.         for( A → α•Xβ ∈ t )
6.             t' = t' ∪ {B → αX•β | B → α•Xβ ∈ t};
7.         if( t' ≠ Φ ){
8.             t' = închidere( t' );
9.             if( t' ∉ T ) {
10.                 T = T ∪ { t' };
11.                 marcat( t' ) = false;
            }//endif
12.             g(t, X) = t';
        }//endif
    }//endfor
    marcat( t ) = true;
  }// endwhile

```

Lema 5.2.3 Algoritmul 5.2.2 descris mai sus este corect: automatul \mathcal{M} obținut este determinist și este echivalent cu M .

Demonstrație Faptul că \mathcal{M} este determinist rezultă din linia 12. Condiția din linia 7 face ca g să fie parțial definită: dacă $\hat{\delta}(t, X) = \Phi$ în automatul M , atunci $g(t, X)$ nu este definită. Să mai observăm că, dacă $t' \neq \Phi$ în linia 7, atunci faptul că $g(t, X) = \hat{\delta}(t, X)$ ne îndreptățește să afirmăm că \mathcal{M} este echivalent cu M . \square

Exemplul 5.2.3 Pentru gramatica $S' \rightarrow S, S \rightarrow aSa \mid bSb \mid c$, automatul LR(0) este dat în figura 5.2.4. Se constată cu ușurință că automatul construit îndeplinește condițiile LR(0) deci gramatica este LR(0). Dacă înlocuim producția $A \rightarrow c$ cu $A \rightarrow \epsilon$, în starea inițială, în loc de $A \rightarrow \bullet c$ o să avem articolul $A \rightarrow \bullet$, care este un articol complet. Cum în aceeași stare există articolul $A \rightarrow \bullet aAa$, aceasta contrazice condiția 2 din caracterizarea LR(0); gramatica $S' \rightarrow S, S \rightarrow aSa \mid bSb \mid \epsilon$, nu este LR(0).

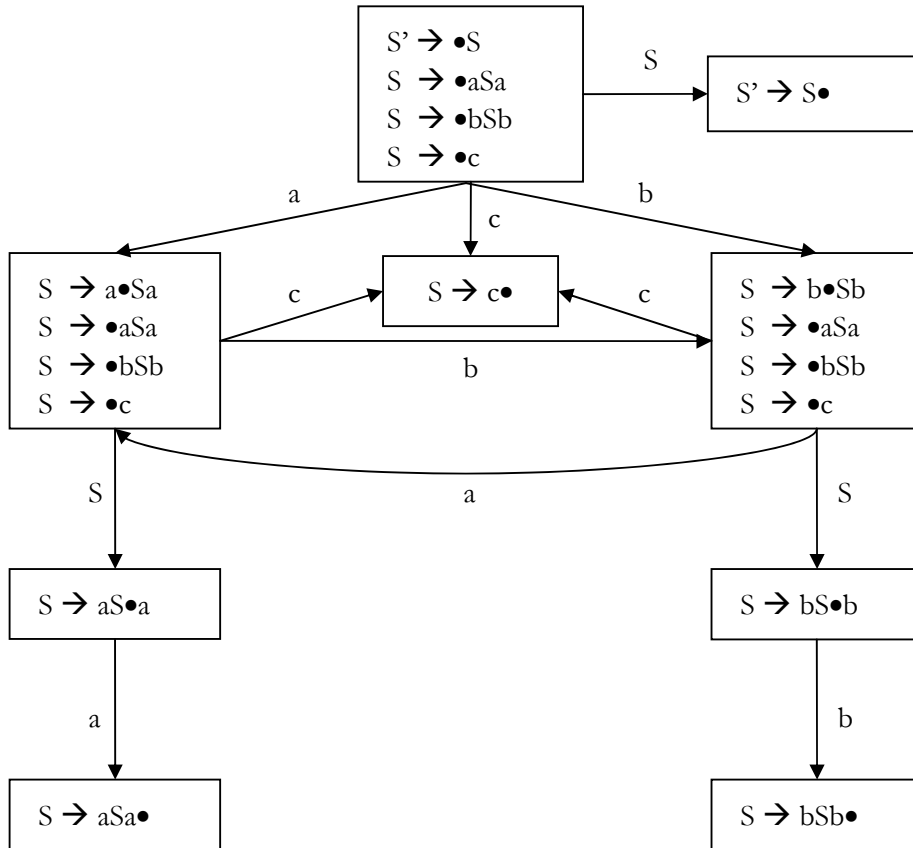
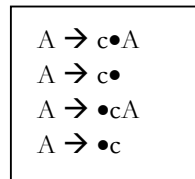


Figura 5.2.4

Exemplul 5.2.4 Fie gramatica: $S \rightarrow aAd \mid bAB$ $A \rightarrow cA \mid c$ $B \rightarrow d$. Automatul LR(0) corespunzător este dat în figura 5.2.5. Se observă că starea:



conține articolul complet $A \rightarrow c \bullet$ și articolele $A \rightarrow \bullet cA$, $A \rightarrow \bullet c$ care au terminal după punct; gramatica nu este LR(0). \square

□

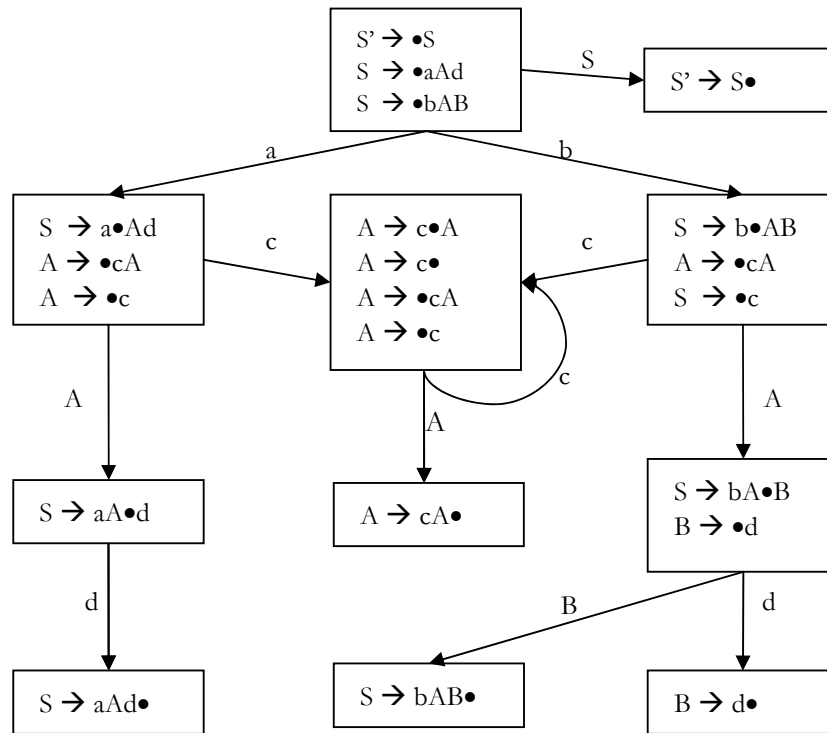


Figura 5.2.5

Clasa limbajelor LR(0) are mai multe caracterizări. Următoarea teoremă are o importanță (mai ales teoretică) deosebită; demonstrația ei se poate găsi în [Har78] pag 515.

Teorema 5.2.3 Fie $L \subseteq \Sigma^*$. Următoarele propoziții sunt echivalente:

1. L este limbaj LR(0) (există G , o gramatică LR(0), cu $L = L(G)$).
2. L este un limbaj independent de context determinist cu proprietatea:
 $\forall x \in \Sigma^+, \forall w, y \in \Sigma^*$, dacă $w \in L$ și $wx \in L$ iar $y \in L$, atunci $yx \in L$.
3. Există un automat pushdown determinist $A = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$, unde $F = \{q_f\}$ și există $Z_f \in \Gamma$ astfel ca:
 $L = T(A, Z_f) = T(A, \Gamma) = \{w \in \Sigma^* \mid (q_0, w, Z_0) \vdash^* (q_f, \varepsilon, Z_f)\}$.
4. Există limbajele strict deterministe L_0 și L_1 astfel încât $L = L_0 L_1^*$.

4.3 Algoritm de analiză sintactică LR(0)

Modelul de analizor sintactic LR(k) este dat în figura 5.3.1.

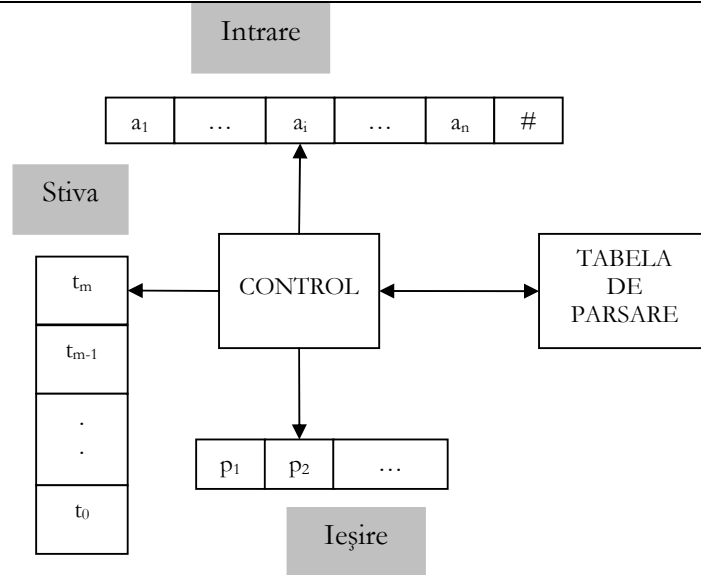


Figura 5.3.1

În cazul $k = 0$, tabela de parsare coincide cu automatul LR(0), \mathcal{M} . O configurație a analizorului LR(0) este o tripletă $(\sigma, u\#, \pi)$ unde $\sigma \in t_0T^*$, (σ este o secvență de stări care începe cu t_0), $u \in T^*$, $\pi \in P^*$. Configurația inițială este $(t_0, w\#, \varepsilon)$, unde w este cuvântul care trebuie analizat. Tranzițiile pe care le execută analizorul LR(0) sunt:

Deplasare: $(\sigma t, au\#, \pi) \vdash (\sigma tt', u\#, \pi)$ dacă $g(t, a) = t'$.

Reducere: $(\sigma t\sigma't', u\#, \pi) \vdash (\sigma tt'', u\#, \pi)$ dacă $A \rightarrow \beta\bullet \in t'$,

$r = A \rightarrow \beta$, $|\sigma't'| = |\beta|$ și $t'' = g(t, A)$.

Acceptare: $(t_0t_1, \#, \pi)$ este configurația de acceptare dacă $S' \rightarrow S\bullet \in t_1$. Analizorul se oprește cu acceptarea cuvântului de analizat, iar π este parsarea acestuia (șirul de reguli care s-a aplicat, în ordine inversă, în derivarea extrem dreaptă a lui w).

Eroare: o configurație căreia nu i se poate aplica nici o tranziție este configurație de eroare. Analizorul se oprește cu respingerea cuvântului de analizat.

Algoritm de analiză sintactică LR(0) funcționează astfel:

Se pornește cu starea inițială a automatului LR(0) în topul stivei și cuvântul de analizat în banda de intrare. Se execută o serie de operații de *deplasare - reducere*, până când, eventual, întreg cuvântul de intrare a fost parcurs și ultima reducere a fost făcută cu regula $S' \rightarrow S$. Dacă nu se obține acest lucru, atunci analiza sintactică eșuează: cuvântul nu este corect sintactic. Alegerea operației de reducere sau de deplasare este dictată de starea din top-ul stivei:

- dacă această stare nu conține un articol complet și în automat există tranziție din aceasta cu simbolul curent de intrare, se produce o deplasare: se înaintează cu o poziție în banda de intrare (în alte implementări simbolul curent din intrare este deplasat în stivă) și se adaugă în stivă starea obținută în automatul LR(0) prin tranziția cu simbolul în discuție. În acest mod se parcurg în automatul LR(0), începând cu t_0 , drumuri corespunzătoare prefixelor viabile din derivarea extrem dreaptă pentru cuvântul de analizat.
- dacă starea din top-ul stivei conține un articol complet $A \rightarrow \beta\bullet$, atunci se produce o reducere: în forma propozițională curentă, β se înlocuiește cu A . În termenii configurației de aici, asta înseamnă că vor fi eliminate din stivă un număr de stări egal cu lungimea lui β (stări care reprezintă în automatul LR(0)

drumul etichetat cu β) și se adăugă starea obținută prin tranziție din starea rămasă în topul stivei cu A. În acest caz se raportează la ieșire aplicarea producției $p = A \rightarrow \beta$. Să descriem algoritmul LR(0) și în pseudocod.

Algoritmul 5.3.1 (Analiză sintactică LR(0))

Intrare: Gramatica $G = (V, T, S, P)$ care este LR(0) augmentată cu $S' \rightarrow S$.

Automatul LR(0) al gramaticii G, notat $\mathcal{M} = (T, \Sigma, g, t_0, T)$.

Cuvântul de intrare $w \in T^*$.

Ieșire: Analiza sintactică (parsarea) ascendentă a lui w dacă $w \in L(G)$;

Eroare, în caz contrar.

Metoda: Fie STIVA o stivă, a simbolul curent,

```
char ps[] = "w#"; //ps este sirul de intrare w
i = 0; // pozitia in sirul de intrare
STIVA.push(t0); // se initializeaza stiva cu t0
while(true) { // se repeta pana la succes sau eroare
    t = STIVA.top();
    a = ps[i] // a este simbolul curent din intrare
    if( g(t, a) ≠ Φ { //deplasare
        STIVA.push(g(t, a));
        i++; //se inainteaza in intrare
    }
    else {
        if(A → X1X2...Xm• ∈ t) {
            if(A == 'S')
                if(a == '#')exit("acceptare");
            else exit("eroare");
            else // reducere
                for( i = 1; i ≤ m; i++) STIVA.pop();
                STIVA.push(g(top(STIVA), A));
        } //endif
        else exit("eroare");
    } //endelse
} //endwhile
```

Exemplul 5.3.1 Fie gramatica:

$$\begin{array}{lll} S' \rightarrow S & E \rightarrow E+T & T \rightarrow (E) \\ S \rightarrow E\$ & E \rightarrow T & T \rightarrow a \end{array}$$

Această gramatică generează expresiile aritmetice care au doar operatorul +, paranteze și operandul a. Producția $S \rightarrow E\$$ asigură că orice expresie se termină cu simbolul special \$ ceea ce face ca limbajul gramaticii să fie liber de prefixe: nici un prefix propriu al unui cuvânt din limbaj nu este în limbaj.

Automatul LR(0) pentru această gramatică este dat în figura 5.3.2. Stările automatului sunt notate 0, 1, 2,..., 10 iar tabela de tranziție a automatului, prezentată mai jos, este și tabelă de parsare LR(0).

G	a	+	()	\$	S	E	T
0	5		4			1	2	3
1								
2		7			6			
3								
4	5		4				8	3
5								
6								
7	5		4					9

8		7		10				
9								
10								

Marca \$ care este simbol terminal al gramaticii asigură că gramatica este LR(0), lucru care se poate verifica ușor inspectând cele 10 stări ale automatului și constatând că sunt îndeplinite condițiile LR(0).

În tabela de mai jos, în care se vizualizează atât configurațiile parserului, cât și acțiunea ce se execută (deplasare sau reducere), este detaliată parsarea expresiei $a+(a+a)\$$.

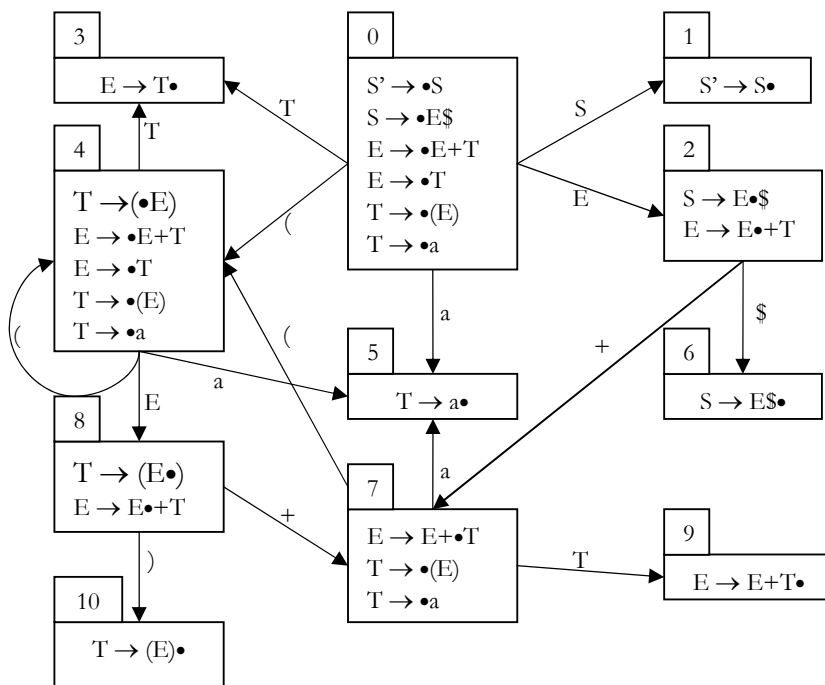


Figura 5.3.2

STIVA	INTRARE	ACȚIUNE	IEȘIRE
0	$a+(a+a)\$$	deplasare	
05	$+(a+a)\$$	reducere	$T \rightarrow a$
03	$+(a+a)\$$	reducere	$E \rightarrow T$
02	$+(a+a)\$$	deplasare	
027	$(a+a)\$$	deplasare	
0274	$a+a)\$$	deplasare	
02745	$+a)\$$	reducere	$T \rightarrow a$
02743	$+a)\$$	reducere	$E \rightarrow T$
02748	$+a)\$$	deplasare	
027487	$a)\$$	deplasare	
0274875	$)\$$	reducere	$T \rightarrow a$
0274879	$)\$$	reducere	$E \rightarrow E+T$
02748	$)\$$	deplasare	
02748(10)	$\$$	reducere	$T \rightarrow (E)$
0279	$\$$	reducere	$E \rightarrow E+T$
02	$\$$	deplasare	
026	$\#$	reducere	$S \rightarrow E\$$
01	$\#$	acceptare	

După acest exemplu putem să facem câteva considerații privind testul LR(0) pentru gramatici.

Definiția 5.3.1 Fie G o gramatică și \mathcal{M} automatul LR(0) atașat lui G . Spunem că o stare a lui \mathcal{M} are un *conflict reducere/reducere* dacă ea conține două articole complete distincte $A \rightarrow \alpha\bullet$, $B \rightarrow \beta\bullet$. Spunem că o stare a lui \mathcal{M} are un *conflict deplasare/reducere* dacă ea conține un articol complet $A \rightarrow \alpha\bullet$ și un articol cu terminal după punct de forma $B \rightarrow \beta\bullet\alpha\gamma$. Spunem că o stare este *consistentă* dacă ea nu conține conflicte și este *inconsistentă* în caz contrar. Cu aceste noțiuni, teorema de caracterizare a gramaticilor LR(0) se poate reformula astfel:

Teorema 5.3.1 Fie G o gramatică și \mathcal{M} automatul său LR(0). Gramatica G este LR(0) dacă și numai dacă automatul \mathcal{M} nu conține stări inconsistente. □

Dacă în gramatica din exemplul precedent în loc de producția $S \rightarrow E\$$ considerăm $S \rightarrow E$, starea 2 din automat va fi $\{S \rightarrow E\bullet, E \rightarrow E\bullet+T\}$ care este *inconsistentă* deoarece conține un conflict *deplasare/reducere*, deci gramatica nu este LR(0).

Pentru demonstrarea corectitudinii programului de analiză sintactică LR(0), să considerăm relația \vdash^+ care notează închiderea tranzitivă a relației \vdash definită pe mulțimea configurațiilor. Vom dovedi două rezultate mai generale, care prin particularizare conduc la corectitudinea algoritmului.

Lema 5.3.1 Fie $G = (V, T, S, P)$ o gramatică LR(0), $t_0\sigma$, $t_0\tau$ drumuri în automatul LR(0) etichetate cu \varnothing respectiv γ și $u, v \in T^*$. Atunci, dacă în parserul LR(0) are loc $(t_0\sigma, uv\#,$

$\varepsilon) \vdash^+ (t_0\tau, v\#, \pi)$, atunci în G are loc derivarea $\gamma \xRightarrow[\text{dr}]{\tilde{\pi}} \varnothing u$.

Demonstrație Procedăm prin inducție după lungimea lui π .

$|\pi| = 0$: În acest caz parserul produce doar deplasări și din definiția acestor acțiuni rezultă că $\gamma = \varnothing u$, și derivarea din concluzie are loc în mod trivial.

$|\pi| > 0$: Fie $\pi = r\pi'$ și $r = A \rightarrow \beta$. Calculul considerat se scrie (punând în evidență reducerea r) astfel:

$$\begin{aligned} (t_0\sigma, uv\#, \varepsilon) &= (t_0\sigma, u_1u_2v\#, \varepsilon) \vdash^+ (t_0\sigma\sigma_1, u_2v\#, \varepsilon) = \\ &= (t_0\sigma't'\sigma''t'', u_2v\#, \varepsilon) \vdash (t_0\sigma't't, u_2v\#, r) \vdash^+ (t_0\tau, v\#, r\pi') \end{aligned}$$

unde eticheta subdrumului $\sigma''t''$ este β , $\varnothing u_1 = \psi\beta$, ψ fiind eticheta drumului $t_0\sigma't'$, iar în automatul \mathcal{M} are loc $g(t', A) = t$. După felul cum au fost definite tranzițiile rezultă că are loc și:

$$(t_0\sigma't't, u_2v\#, \varepsilon) \vdash^+ (t_0\tau, v\#, \pi')$$

ceea ce înseamnă, conform ipotezei inductive:

$$\gamma \xRightarrow[\text{dr}]{\tilde{\pi'}} \psi A u_2 \quad (\psi A \text{ este eticheta drumului } t_0\sigma't't)$$

Aplicând în continuare regula $r = A \rightarrow \beta$ obținem:

$$\gamma \xRightarrow[\text{dr}]{\tilde{\pi'}} \psi A u_2 \xRightarrow[\text{dr}]{r} \psi \beta u_2 = \varnothing u_1 u_2 = \varnothing u,$$

ceea ce trebuia demonstrat. □

Corolarul 5.3.1 Dacă în parserul LR(0) are loc calculul $(t_0, w\#, \varepsilon) \vdash^+ (t_0t_1, \#, \pi)$ unde S'

$\rightarrow S\bullet \in t_1$, atunci în gramatica G are loc derivarea $S \xRightarrow[\text{dr}]{\tilde{\pi}} w$.

Demonstrație Lema 5.3.1 pentru cuvintele $\sigma = \varepsilon$, $u = w$, $v = \varepsilon$ și $t_0\tau = t_0t_1$. □

Lema 5.3.2 Fie $G = (V, T, S, P)$ o gramatică LR(0) și $\gamma \in \Sigma^+$ o formă propozițională a

lui G astfel încât $\gamma \xRightarrow[\text{dr}]{\tilde{\pi}} \varphi u$, cu $\varphi \in \Sigma^*$ și dacă este nenul atunci $\varphi:1$ este neterminal.

Atunci în parserul LR(0) are loc calculul

$$(t_0\sigma, uv\#, \varepsilon) \vdash^* (t_0\tau, v\#, \pi)$$

unde $t_0\sigma$ și $t_0\tau$ sunt drumuri în automatul LR(0) etichetate cu φ respectiv γ , pentru orice cuvânt $v \in T^*$.

Demonstrație Procedăm din nou prin inducție după lungimea lui π .

$|\pi| = 0$: În acest caz $\gamma = \varphi u$ și se verifică ușor că are loc calculul precizat cu tranziții de tip deplasare.

$|\pi| > 0$: Fie $\pi = r\pi'$, $r = A \rightarrow \beta$. Atunci $\tilde{\pi} = \tilde{\pi}'r$ și derivarea din G se scrie:

$$\gamma \xRightarrow[\text{dr}]{\tilde{\pi}'} \varphi_1 A u_2 \xRightarrow[\text{dr}]{r} \varphi_1 \beta u_2 = \varphi u_1 u_2 = \varphi u$$

Din $\varphi_1 A u_2 \xRightarrow[\text{dr}]{r} \varphi_1 \beta u_2 = \varphi u_1 u_2 = \varphi u$ și faptul că γ este o formă propozițională a lui G

(deci și $\varphi_1 A u_2$ este formă propozițională), rezultă că în automatul \mathcal{M} există un drum $t_0\sigma\sigma_1$ etichetat cu $\varphi_1\beta = \varphi u_1$ (pentru că acesta este un prefix viabil), și:

$$(t_0\sigma, u_1 u_2 v\#, \varepsilon) \vdash^* (t_0\sigma\sigma_1, u_2 v\#, \varepsilon) = (t_0\sigma't'\sigma''t'', u_2 v\#, \varepsilon) \vdash (t_0\sigma't't, u_2 v\#, r)$$

unde $A \rightarrow \beta \bullet \in t''$, $|\sigma''t''| = |\beta|$, iar drumul $t_0\sigma't't$ este etichetat cu φA . După ipoteza

inductivă, din $\gamma \xRightarrow[\text{dr}]{\tilde{\pi}'} \varphi_1 A u_2$ rezultă calculul:

$$(t_0\sigma't't, u_2 v\#, \varepsilon) \vdash^* (t_0\tau, v\#, \pi').$$

Asta înseamnă că putem scrie:

$$(t_0\sigma, u_1 u_2 v\#, \varepsilon) \vdash^* (t_0\sigma\sigma_1, u_2 v\#, \varepsilon) \vdash (t_0\sigma't't, v\#, r) \vdash^* (t_0\tau, v\#, r\pi')$$

ceea ce trebuia demonstrat. □

Corolarul 5.3.2 Dacă în gramatica G , care este LR(0), avem $S \xRightarrow[\text{dr}]{\tilde{\pi}'} w$, atunci în parserul

LR(0) are loc calculul $(t_0, w\#, \varepsilon) \vdash^+ (t_0 t_1, \#, \pi)$ unde $S' \rightarrow S \bullet \in t_1$.

Demonstrație Rezultă din lema precedentă luând $\gamma = S$, $\varphi = \varepsilon$, $u = w$, $v = \varepsilon$. □

Combinând cele două leme obținem teorema de corectitudine a algoritmului de analiză sintactică LR(0):

Teorema 5.3.2 Dacă G este gramatică LR(0) atunci, oricare ar fi cuvântul de intrare $w \in T^*$, parserul LR(0) ajunge la configurația de acceptare pentru w , adică $(t_0, w\#, \varepsilon) \vdash^+$

$(t_0 t_1, \#, \pi)$, dacă și numai dacă $S \xRightarrow[\text{dr}]{\tilde{\pi}} w$.

Demonstrație Rezultă din lemele 5.3.1, 5.3.2 și corolarele acestora. □

4.4 Gramatici și analizoare SLR(1)

Este mult prea optimist să credem că putem face analiza sintactică LR(0) pentru gramatici care descriu diverse construcții din limbajele de programare. În paragraful precedent am văzut că în cazul expresiilor aritmetice, pentru a obține o gramatică LR(0), trebuie să punem o marcă de sfârșit (\$) la aceste expresii. Dar expresiile aritmetice apar

în diverse contexte în limbajele de programare încât soluția aceasta nu este viabilă. Dacă renunțăm la marca \$ atunci, în automatul LR(0), se obține cel puțin o stare inconsistentă (cu conflict deplasare/reducere).

Ne propunem în acest paragraf să găsim o modalitate - dacă este posibil - de eliminare a conflictelor într-un automat LR(0). Dacă o anumită stare t conține un conflict *deplasare/reducere*, de pildă articolele $A \rightarrow \alpha\bullet$ și $B \rightarrow \beta_1\bullet a\beta_2$, înseamnă că există în gramatică un prefix viabil γ pentru care cele două articole sunt valide:

$$S \xRightarrow[\text{dr}]{*} \phi Au \xRightarrow[\text{dr}]{*} \phi \alpha u = \gamma u$$

$$S \xRightarrow[\text{dr}]{*} \phi' Bv \xRightarrow[\text{dr}]{*} \phi' \beta_1 a \beta_2 v = \gamma a \beta_2 v$$

Problema este să eliminăm conflictul, adică să găsim o modalitate de alegere în mod determinist a acțiunii care trebuie făcută, *deplasare* sau *reducere*. Aceasta ar putea fi făcută dacă simbolul 'a' nu este primul simbol din cuvântul u , ceea ce înseamnă că $a \notin \text{FOLLOW}(A)$. În acest caz am putea alege *deplasare* dacă simbolul din banda de intrare este 'a' și *reducere* dacă acesta ar fi unul din simbolurile din $\text{FOLLOW}(A)$.

În cazul unui conflict *reducere/reducere* într-o stare t , rezultă existența unui prefix viabil γ pentru care articolele diferite $A \rightarrow \alpha\bullet$, $B \rightarrow \beta\bullet$ din t sunt valide:

$$S \xRightarrow[\text{dr}]{*} \phi Au \xRightarrow[\text{dr}]{*} \phi \alpha u = \gamma u$$

$$S \xRightarrow[\text{dr}]{*} \phi' Bv \xRightarrow[\text{dr}]{*} \phi' \beta v = \gamma v$$

Eliminarea conflictului ar putea fi făcută dacă mulțimile $\text{FOLLOW}(A)$ și $\text{FOLLOW}(B)$ sunt disjuncte. Astfel, dacă în banda de intrare urmează un simbol din $\text{FOLLOW}(A)$, se face *reducere* cu regula $A \rightarrow \alpha$, dacă acesta este din $\text{FOLLOW}(B)$ se face *reducerea* cu $B \rightarrow \beta$, iar în celelalte cazuri se obține eroare. Din aceste considerente suntem îndreptățiți să introducem clasa de *gramatici SLR(1)* ("simplu LR(1)") prin:

Definiția 5.4.1 Fie G o gramatică pentru care automatul LR(0) conține stări inconsistente (decă G nu este LR(0)). Gramatica G este *gramatică SLR(1)* dacă oricare ar fi starea t a automatului LR(0) sunt îndeplinite condițiile:

1. Dacă $A \rightarrow \alpha\bullet$, $B \rightarrow \beta\bullet \in t$ atunci $\text{FOLLOW}(A) \cap \text{FOLLOW}(B) = \Phi$;
2. Dacă $A \rightarrow \alpha\bullet$, $B \rightarrow \beta\bullet a\gamma \in t$ atunci $a \notin \text{FOLLOW}(A)$.

Modelul de algoritm de analiză sintactică SLR(1) este același ca în cel descris în paragraful precedent cu deosebirea că tabela de analiză sintactică se construiește într-un anumit mod. Ea va conține două zone. Prima, numită *ACTIUNE*, determină dacă parserul va face deplasare respectiv *reducere*, în funcție de starea ce se află în topul stivei și de simbolul următor din intrare, iar cea de a doua, numită *GOTO*, determină starea ce se va pune în stivă în urma unei reduceri. Următorul algoritm descrie obținerea tabeli de parsare SLR(1).

Algoritmul 5.4.1 (construcția tabeli de parsare SLR(1))

Intrare: Gramatica $G = (V, T, S, P)$ augmentată cu $S' \rightarrow S$;

Automatul $\mathcal{M} = (\mathcal{T}, \Sigma, g, t_0, T)$;

Mulțimile $\text{FOLLOW}(A)$, $A \in V$ (ϵ se înlocuiește cu #).

Ieșire: Tabela de analiză SLR(1) compusă din două părți:

ACTIUNE(t, a), $t \in \mathcal{T}$, $a \in T \cup \{ \# \}$.

GOTO(t, A), $t \in \mathcal{T}$, $A \in V$.

Metoda:

1. for($t \in T$) //Se initializeaza tabela de parsare
2. for ($a \in T$) *ACTIUNE*(t, a) = "eroare";

```

3.      for (A ∈ V) GOTO(t, A) = "eroare";
4.  for(t ∈ T){
5.      for(A → α•aβ ∈ t)
6.          ACTIUNE(t,a)="D g(t, a)"; //deplasare in g(t, a)
7.      for(B → γ• ∈ t ) // acceptare sau reducere
8.          if(B == 'S') ACTIUNE(t, a) = "acceptare";
              else
9.              for(a∈FOLLOW(B)) ACTIUNE(t,a)="R B→γ";
              } // endfor
10.     for (A ∈ V) GOTO(t, A) = g(t, A);
    }

```

Toate intrările în tabela SLR(1) sunt inițial completate cu “eroare” (linia 1-3) astfel încât intrările nedefinite în liniile 6, 8, 9, 10 rămân cu această valoare. Să observăm că, dacă gramatica G este SLR(1), atunci tabela dată de algoritmul de mai sus este bine definită în sensul că o tentativă de adăugare a unei valori în partea *ACTIUNE* sau partea *GOTO* se face doar pentru intrările cu valoarea “eroare”. Algoritmul de analiză sintactică SLR(1) funcționează pe aceleași principii ca și cel LR(0). Deosebirea este că acțiunile sunt dictate de tabela de analiză SLR(1) încât tranzițiile pe care le execută acesta sunt:

Deplasare: $(\sigma t, au\#, \pi) \vdash (\sigma tt', u\#, \pi)$ dacă $ACTIUNE(t, a) = Dt'$;

Reducere: $(\sigma t\sigma't', u\#, \pi) \vdash (\sigma tt'', u\#, \pi p)$ dacă $ACTIUNE(t, a) = Rp$ unde $p = A \rightarrow \beta$, $|\sigma't'| = |\beta|$ și $t'' = GOTO(t, A)$;

Acceptare: $(tot, \#, \pi)$ dacă $ACTIUNE(t,a) = "acceptare"$; Analizorul se oprește cu acceptarea cuvântului de analizat iar π este parsarea acestuia (șirul de reguli care s-a aplicat, în ordine inversă, în derivarea extrem dreaptă a lui w).

Eroare: $(\sigma t, au\#, \pi) \vdash eroare$ dacă $ACTIUNE(t,a) = "eroare"$; Analizorul se oprește cu respingerea cuvântului de analizat.

Algoritmul 5.4.2 (Analiză sintactică SLR(1))

Intrare: Gramatica $G = (V, T, S, P)$ care este SLR(1);
Tabela de parsare SLR(1) (*ACTIUNE*, *GOTO*);
Cuvântul de intrare $w \in T^*$.

Ieșire: Analiza sintactică (parsarea) ascendentă a lui w dacă $w \in L(G)$;
eroare, în caz contrar.

Metoda: Se folosește stiva St pentru a implementa tranzițiile de mai sus.

```

char ps[] = "w#"; //ps este cuvântul de intrare w
int i = 0; // pozitia curenta in cuvântul de intrare
St.push(t0); // se initializeaza stiva cu t0
while(true) { // se repeta pana la succes sau eroare
    t = St.top();
    a = ps[i] // a este simbolul curent din intrare
    if(ACTIUNE(t,a) == "acceptare") exit("acceptare");
    if(ACTIUNE(t,a) == "Dt'"){
        St.push(t');
        i++; // se inainteaza in w
    } //endif
    else {
        if(ACTIUNE(t,a) == "R A → X1X2...Xm") {
            for( i = 1; i ≤ m; i++) St.pop();
            St.push(GOTO(top(STIVA), A));
        } //endif
        else exit("eroare");
    }
}

```

```

    }//endelse
  }//endwhile

```

Corectitudinea algoritmului SLR(1) se dovedește în același mod ca și cea a algoritmului LR(0). Așa cum am precizat, pentru o gramatică în care automatul LR(0) are stări inconsistente, prin construcția tabelului de analiză SLR(1) nu am făcut altceva decât să eliminăm conflictele de deplasare/reducere sau de reducere/reducere. Faptul că este corectă strategia propusă rezultă din analiza făcută înainte de a defini gramatica SLR(1). Lăsăm așadar cititorului, ca un exercițiu de acum simplu, demonstrarea corectitudinii algoritmului SLR(1).

Teorema 5.4.1 Dacă G este o gramatică SLR(1) atunci, oricare ar fi cuvântul de intrare w , algoritmul de analiză SLR(1) se oprește cu “acceptare” și obține la ieșire π dacă și

numai dacă $S \xRightarrow[\text{dr}]{\tilde{\pi}} w$.

□

Exemplul 5.4.1 Să aplicăm algoritmul de analiză sintactică SLR(1) pentru gramatica ce generează expresiile aritmetice:

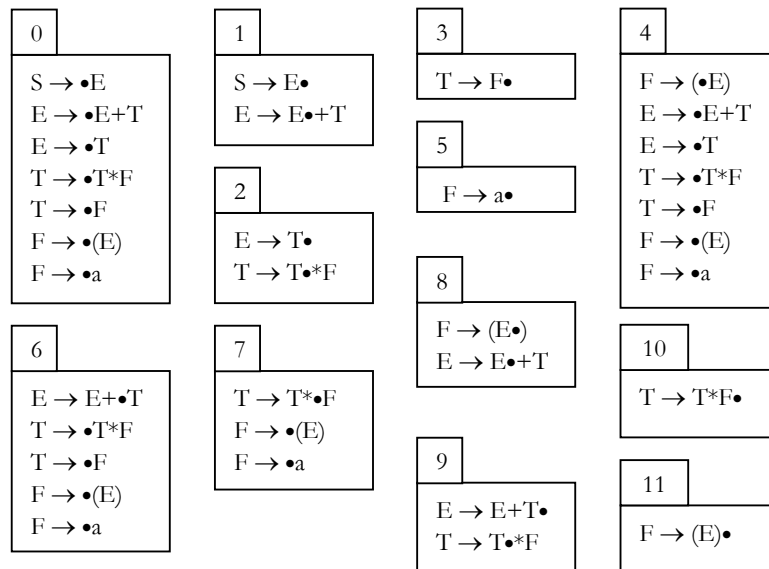
$$\begin{array}{ll}
 S \rightarrow E & T \rightarrow T * F \\
 E \rightarrow E + T & T \rightarrow F \\
 E \rightarrow T & F \rightarrow (E) \quad F \rightarrow a
 \end{array}$$


Figura 5.4.1

Tabela de tranziție a automatului LR(0) este:

g	a	+	*	()	E	T	F
0	5			4		1	2	3
1		6						
2			7					
3								
4	5			4		8	2	3
5								
6	5			4			9	3
7	5			4				10

8					11			
9			7					
10								
11								

Stările automatului LR(0) sunt date în figura 5.4.1. Din acest automat se deduce că G nu este LR(0): stările 1, 2, 9 conțin conflict de deplasare/reducere. Pentru a verifica dacă gramatica este SLR(1) avem nevoie de mulțimile FOLLOW(S) și FOLLOW(E). Se deduce ușor că: FOLLOW(S) = { ϵ } și FOLLOW(E) = { +,), ϵ }. Atunci gramatica este SLR(1) pentru că:

în starea 1: + \notin FOLLOW(S);

în starea 2: * \notin FOLLOW(E);

în starea 9: * \notin FOLLOW(E).

Pentru construcția tabelului de parsare avem nevoie de mulțimile FOLLOW pentru toți neterminalii (s-a înlocuit ϵ cu #):

NETERMINAL	FOLLOW
S	#
E	#, +,)
T	#, +, *,)
F	#, +, *,)

Tabela de analiză SLR(1) pentru această gramatică este dată mai jos.

STARE	ACȚIUNE						GOTO		
	a	+	*	()	#	E	T	F
0	D5			D4			1	2	3
1		D6				acceptare			
2		R2	D7		R2	R2			
3		R4	R4		R4	R4			
4	D5			D4			8	2	3
5		R6	R6		R6	R6			
6	D5			D4				9	3
7	D5			D4					10
8					D11				
9		R1	D7		R1	R1			
10		R3	R3		R3	R3			
11		R5	R5		R5	R5			

Acțiunile analizorului SLR(1) pentru cuvântul de intrare $w = a^*(a+a)$ sunt date în tabela următoare:

STIVA	INTRARE	ACȚIUNE	IEȘIRE
0	$a^*(a+a)\#$	deplasare	
05	$*(a+a)\#$	reducere	$F \rightarrow a$
03	$*(a+a)\#$	reducere	$T \rightarrow F$
02	$*(a+a)\#$	deplasare	
027	$(a+a)\#$	deplasare	
0274	$a+a)\#$	deplasare	
02745	$+a)\#$	reducere	$F \rightarrow a$
02743	$+a)\#$	reducere	$T \rightarrow F$
02742	$+a)\#$	reducere	$E \rightarrow T$
02748	$+a)\#$	deplasare	
027486	$a)\#$	deplasare	
0274865	$)\#$	reducere	$F \rightarrow a$

0274863)#	reducere	$T \rightarrow F$
0274869)#	reducere	$E \rightarrow E+T$
02748)#	reducere	
02748(11)	#	reducere	$F \rightarrow (E)$
027(10)	#	reducere	$T \rightarrow T*F$
02	#	reducere	$E \rightarrow T$
01	#	acceptare	

□

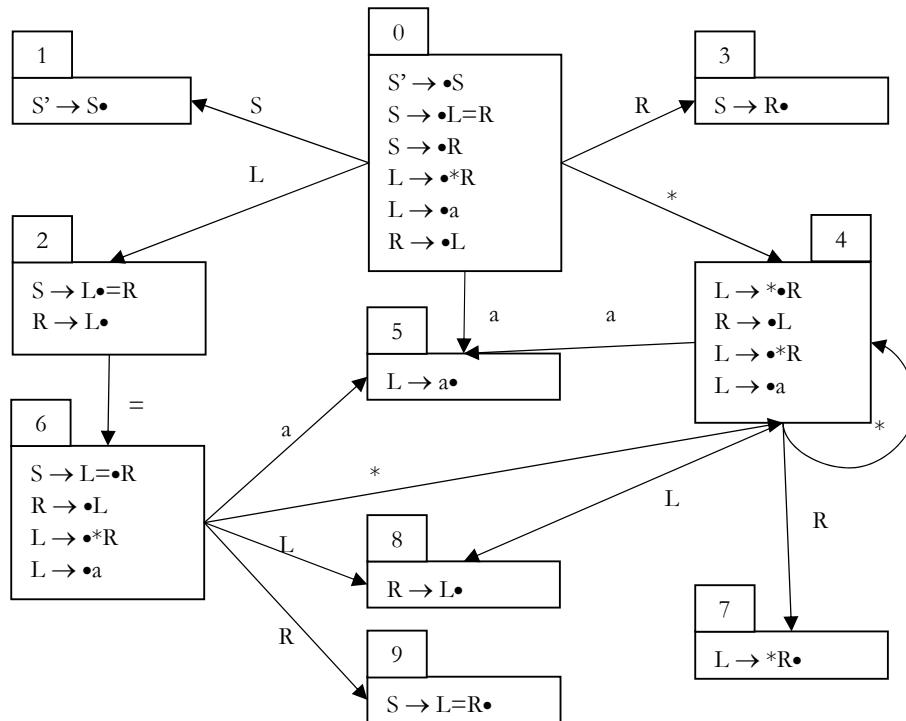


Figura 5.4.2

Vom vedea în paragraful următor că orice gramatică SLR(1) este gramatică LR(1), prin urmare este gramatică neambiguă. Există însă gramatici neambigue (chiar gramatici LR(1)) care nu sunt SLR(1). Vom da în continuare un exemplu de astfel de gramatică.

Exemplul 5.1.2 Fie gramatica dată de producțiile:

$$S \rightarrow L=R \mid R \quad L \rightarrow *R \mid a \quad R \rightarrow L$$

Se poate constata că această gramatică nu este ambiguă. Automatul LR(0) este dat în figura 5.4.2.

În starea 2 există un conflict deplasare/reducere: articolul $S \rightarrow L \bullet = R$ definește operația de deplasare (dacă următorul simbol este $=$) iar articolul $R \rightarrow L \bullet$ definește o reducere. Constatăm că $\text{FOLLOW}(R) = \{ \epsilon, = \}$ ceea ce înseamnă că gramatica nu este SLR(1) pentru că simbolul $'='$ este în această mulțime. Vom arăta în paragraful următor că această gramatică este LR(1).

□

4.5 O caracterizare a gramaticilor LR(1)

Definiția 5.5.1 Fie $G = (V, T, S, P)$ o gramatică independentă de context redusă. Un articol LR(1) pentru gramatica G este o pereche $(A \rightarrow \alpha \bullet \beta, a)$ unde $A \rightarrow \alpha \beta$ este un articol LR(0), iar $a \in \text{FOLLOW}(A)$ (se pune $\#$ în loc de ϵ).

Definiția 5.5.2 Articolul $(A \rightarrow \beta_1 \bullet \beta_2, a)$ este valid pentru prefixul viabil $\alpha \beta_1$ dacă are loc derivarea

$$S \xRightarrow[dr]{*} \alpha A u \xRightarrow[dr]{} \alpha \beta_1 \beta_2 u$$

iar $a = 1:u$ ($a = \#$ dacă $u = \epsilon$).

Teorema 5.5.1 O gramatică $G = (V, T, S, P)$ este gramatică LR(1) dacă și numai dacă oricare ar fi prefixul viabil φ , nu există două articole distincte, valide pentru φ , de forma $(A \rightarrow \alpha \bullet, a)$, $(B \rightarrow \beta \bullet \gamma, b)$ unde $a \in \text{FIRST}(\gamma b)$.

Demonstrație Să interpretăm mai întâi condiția (necesară și suficientă) pentru ca G să fie LR(1). Este de așteptat ca algoritmul de analiză LR(1) să funcționeze după aceleași principii ca și cel SLR(1). Asta înseamnă că un articol de forma $(A \rightarrow \alpha \bullet, a)$ invocă reducere cu regula $A \rightarrow \alpha$ dacă a este simbolul curent din intrare, iar un articol de forma $(B \rightarrow \beta \bullet \alpha \beta', b)$ invocă deplasare dacă a este la intrare. Așadar, condiția ca nici un prefix viabil să nu aibă două articole valide de forma indicată se traduce prin:

1. nu există conflict *deplasare/reducere*. Un astfel de conflict înseamnă două articole $(A \rightarrow \alpha \bullet, a)$ și $(B \rightarrow \beta \bullet \alpha \beta', b)$ valide pentru același prefix.
2. nu există conflict *reducere/reducere*. Un astfel de conflict înseamnă două articole complete $(A \rightarrow \alpha \bullet, a)$ și $(B \rightarrow \beta \bullet, a)$ valide pentru același prefix.

Demonstrația se face analog cu cea de la caracterizarea LR(0). □

Pentru a testa dacă o gramatică este LR(1) și apoi pentru obținerea tabelii de analiză sintactică LR(1) se construiește automatul LR(1) care recunoaște prefixele viabile ale gramaticii. Mai întâi, ca și în cazul LR(0), să definim o procedură închidere(I) care găsește toate articolele care sunt valide pentru aceleași prefixe pentru care sunt valide articolele din I .

Algoritmul 5.5.1 (procedura închidere(I))

Intrare: Gramatica $G = (V, T, S, P)$;

Mulțimea I de articole LR(1) din gramatica G ; Mulțimile FIRST.

Ieșire: I devine mulțimea tuturor articolelor valide pentru prefixele viabile date inițial.

Metoda: asemănătoare cu cea de la LR(0).

```
flag = true;
while( flag ) {
    flag = false;
    for ( (A → α•Bβ, a) ∈ I ) {
        for ( B → γ ∈ P )
            for( b ∈ FIRST(βa) ) {
                if ( (B → •γ, b) ∉ I ) {
                    I = I ∪ { (B → •γ, b) };
                    flag = true;
                } //endif
            } //endforb
        } //endforB
    } //endforA
} //endwhile
return I;
```

În continuare, folosind această procedură, descriem algoritmul pentru construirea automatului LR(1).

Algoritmul 5.5.2 (Automatul LR(1))

Intrare: Gramatica $G = (V, T, S, P)$ la care s-a adăugat $S' \rightarrow S$;

Ieșire: Automatul determinist LR(1) $\mathcal{M} = (T, \Sigma, g, t_0, T)$ ce recunoaște prefixele viabile și are ca stări mulțimi de articole LR(1).

Metoda:

```

1.  $t_0 = \text{închidere}((S' \rightarrow \bullet S, \#)); T = \{t_0\}; \text{marcat}(t_0) = \text{false};$ 
2. while( $\exists t \in T \ \&\& \ !\text{marcat}(t)$ ) { // marcat(t) = false
3.   for(  $X \in \Sigma$  ) {
4.      $t' = \Phi;$ 
5.     for(  $(A \rightarrow \alpha \bullet X \beta, a) \in t$  )
6.        $t' = t' \cup \{(B \rightarrow \alpha X \bullet \beta, a) \mid (B \rightarrow \alpha \bullet X \beta, a) \in t\};$ 
7.     if(  $t' \neq \Phi$  ) {
8.        $t' = \text{închidere}(t');$ 
9.       if(  $t' \notin T$  ) {
10.         $T = T \cup \{t'\};$ 
11.         $\text{marcat}(t') = \text{false};$ 
12.        } //endif
13.         $g(t, X) = t';$ 
14.      } //endif
15.    } //endfor
16.     $\text{marcat}(t) = \text{true};$ 
17.  } // endwhile

```

Să dovedim acum corectitudinea algoritmului descris mai sus.

Teorema 5.5.2 Automatul \mathcal{M} construit în algoritmul 5.5.2 este determinist și $L(\mathcal{M})$ coincide cu mulțimea prefixelor viabile ale lui G . Mai mult, pentru orice prefix viabil γ , $g(t_0, \gamma)$ reprezintă mulțimea articolelor LR(1) valide pentru γ . □

Automatul LR(1) pentru o gramatică G , se folosește pentru a verifica dacă G este LR(1). Conform teoremei de caracterizare LR(1) acest lucru decurge astfel:

- Conflict *reducere/reducere*: Dacă în \mathcal{T} există o stare ce conține articole de forma $(A \rightarrow \alpha \bullet, a), (B \rightarrow \beta \bullet, a)$ atunci gramatica nu este LR(1);
- Conflict *deplasare/reducere*: Dacă în \mathcal{T} există o stare ce conține articole de forma $(A \rightarrow \alpha \bullet, a)$ și $(B \rightarrow \beta_1 \bullet a \beta_2, b)$, atunci G nu este LR(1).

O gramatică este LR(1) dacă orice stare $t \in \mathcal{T}$ este liberă de conflicte.

Exemplul 5.5.1 Să considerăm din nou gramatica:

$$S \rightarrow L \mid R$$

$$L \rightarrow *R \mid a$$

$$R \rightarrow L$$

și să construim automatul LR(1) pentru aceasta.

În continuare, articolele ce au aceeași primă componentă le notăm condensat prin $(A \rightarrow \alpha \bullet \beta, \{a_1, a_2, \dots, a_n\})$. De exemplu, articolele $(A \rightarrow \alpha \bullet \beta, a), (A \rightarrow \alpha \bullet \beta, b)$, se scriu $(A \rightarrow \alpha \bullet \beta, \{a, b\})$.

Am văzut în paragraful precedent că această gramatică nu este SLR(1). Dacă analizăm stările automatului LR(1), prezentate în figura 5.5.1, constatăm că gramatica este LR(1): toate stările sunt libere de conflict.

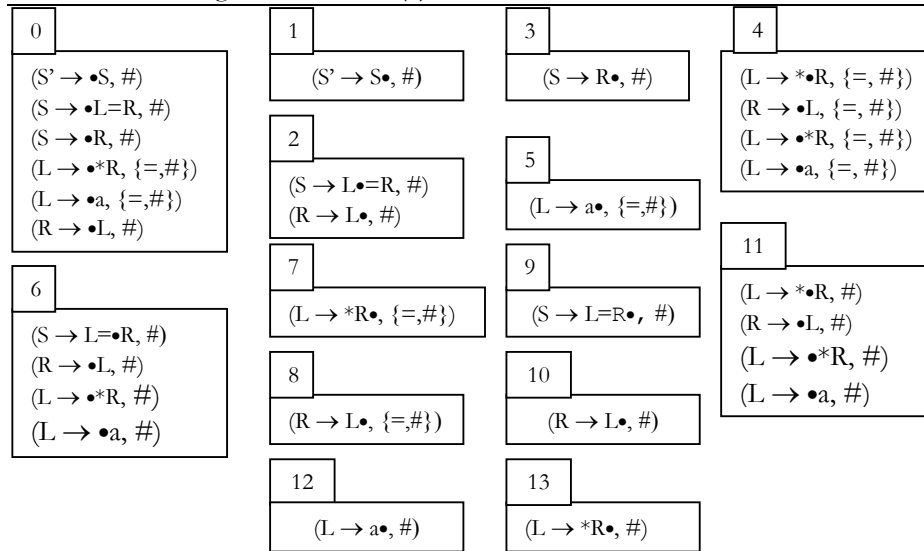


Figura 5.5.1

Tabela de tranziție a automatului este prezentată în continuare.

g	a	=	*	S	L	R
0	5		4	1	2	3
1						
2		6				
3						
4	5		4		8	7
5						
6	12		11		10	9
7						
8						
9						
10						
11	12		11		10	13
12						
13						

4.6 Analiză sintactică LR(1) și LALR(1)

Algoritmul de analiză sintactică LR(1) este, în esență, același cu algoritmul de analiză sintactică SLR(1). Ceea ce îl distinge este modul de construire a tabeli de analiză cu cele două părți, *ACȚIUNE* și *GOTO*. Vom indica în continuare construirea acestei tabeli, pornind de la automatul LR(1).

Algoritmul 5.6.1 (Construcția tabeli de analiză LR(1))

Intrare: Gramatica LR(1) $G = (V, S, T, P)$ augmentată cu $S' \rightarrow S$.

Automatul LR(1) $\mathcal{M} = (\mathcal{T}, \Sigma, g, t_0, \mathcal{T})$.

Ieșire: Tabela de analiză cu componentele:

$ACȚIUNE(t, a)$, $t \in \mathcal{T}$, $a \in T \cup \{\#\}$

$GOTO(t, A)$, $t \in \mathcal{T}$, $A \in V$.

Metoda:

```
for(t ∈ T){ //Se initializeaza tabela de parsare
```

```

for (a ∈ T) ACTIUNE(t, a) = "eroare";
for (A ∈ V) GOTO(t, A) = "eroare";
}
for(t ∈ T){
  for((A → α•aβ, L) ∈ t)
    ACTIUNE(t,a)="D g(t, a)";//deplasare in g(t, a)
  for((B → γ•, L) ∈ t ) {
    for(c ∈ L){
      if(B == 'S') ACTIUNE(t, c) = "acceptare";
      else
        ACTIUNE(t,c)="R B→γ";//reducere cu B→γ
    }// endforc
  }// endforB
} // endfor
for (A ∈ V) GOTO(t, A) = g(t, A);

```

Așadar, stările automatului LR(1) determină acțiunea analizorului sintactic :

- Dacă starea t din topul stivei conține un articol cu terminalul 'a' după punct și la intrare, se produce o deplasare, anume starea $g(t,a)$ se adaugă în stivă;
- Dacă starea din topul stivei conține un articol complet ($A \rightarrow \beta\bullet$, a) unde 'a' este terminalul curent din intrare, se produce reducere cu producția $A \rightarrow \beta$. În cazul în care $A \rightarrow \beta = S' \rightarrow S$, obținem *acceptare*;
- În celelalte cazuri se obține o *eroare*.

Exemplul 5.6.1 Să construim după algoritmul tocmai descris tabela de analiză pentru gramatica din exemplul 5.4.1:

0: $S' \rightarrow S$ 1: $S \rightarrow L=R$ 2: $S \rightarrow R$
 3: $L \rightarrow *R$ 4: $L \rightarrow a$ 5: $R \rightarrow L$

STARE	ACȚIUNE				GOTO		
	a	=	*	#	S	L	R
0	D5		D4		1	2	3
1				acceptare			
2		D6		R5			
3				R2			
4	D5		D4			8	7
5		R4		R4			
6	D12		D11			10	9
7		R3		R3			
8		R5		R5			
9				R1			
10				R5			
11	D12		D11			10	13
12				R4			
13				R3			

Aplicarea algoritmului (de deplasare/reducere) folosind această tabelă de analiză, asupra cuvintelor ****aa** respectiv ****a=a**, este ilustrată în tabelele următoare:

STIVA	INTRARE	ACȚIUNE	IEȘIRE
0	**aa#	deplasare D4	
04	*aa#	deplasare D4	
044	aa#	deplasare D5	
0445	a#	eroare	
STIVA	INTRARE	ACȚIUNE	IEȘIRE

0	**a=a#	deplasare D4	
04	*a=a#	deplasare D4	
044	a=a#	deplasare D5	
0445	=a#	reducere R4	$L \rightarrow a$
0448	=a#	reducere R5	$R \rightarrow L$
0447	=a#	reducere R3	$L \rightarrow *R$
048	=a#	reducere R5	$R \rightarrow L$
047	=a#	reducere R3	$L \rightarrow *R$
02	=a#	deplasare D6	
026	a#	deplasare D12	
026(12)	#	reducere R4	$L \rightarrow a$
026(10)	#	reducere R5	$R \rightarrow L$
0269	#	reducere R1	$S \rightarrow L=R$
01	#	acceptare	

Să observăm că automatul LR(1) construit în paragraful precedent, pentru gramatica care tocmai am discutat-o, are un număr mai mare de stări decât numărul stărilor automatului LR(0) pentru această gramatică. Unele din stările automatului LR(1) au ceva în comun:

starea 5:

$$(L \rightarrow a\bullet, \{=, \#\})$$

cu starea 12:

$$(L \rightarrow a\bullet, \#)$$

starea 4:

$$\begin{aligned} &(L \rightarrow *\bullet R, \{=, \#\}) \\ &(R \rightarrow \bullet L, \{=, \#\}) \\ &(L \rightarrow \bullet *R, \{=, \#\}) \\ &(L \rightarrow \bullet a, \{=, \#\}) \end{aligned}$$

cu starea 11:

$$\begin{aligned} &(L \rightarrow *\bullet R, \#) \\ &(R \rightarrow \bullet L, \#) \\ &(L \rightarrow \bullet *R, \#) \\ &(L \rightarrow \bullet a, \#) \end{aligned}$$

starea 7:

$$(L \rightarrow *R\bullet, \{=, \#\})$$

cu starea 13:

$$(L \rightarrow *R\bullet, \#)$$

starea 8:

$$(R \rightarrow L\bullet, \{=, \#\})$$

cu starea 10:

$$(R \rightarrow L\bullet, \#)$$

Aceste stări au respectiv prima componentă a articolelor aceeași.

Definiția 5.6.1 Fie t o stare în automatul LR(1) al unei gramatici G . *Nucleul* acestei stări este mulțimea articolelor LR(0) care apar ca prime componente în articolele LR(1) din t .

Definiția 5.6.2 Două stări t_1 și t_2 ale automatului LR(1) pentru gramatica G se numesc *echivalente* dacă ele au același nucleu.

Cum fiecare stare a automatului LR(1) este o mulțime de articole LR(1), având două stări t_1 și t_2 putem să vorbim de $t_1 \cup t_2$. De pildă, dacă

$$t_1 = \{(L \rightarrow *R\bullet, \{=, \#\})\}, t_2 = \{(L \rightarrow *R\bullet, \#)\}$$

atunci $t_1 \cup t_2 = t_1$ pentru că $t_2 \subset t_1$.

Definiția 5.6.3 Fie G gramatică LR(1) și $\mathcal{M} = (\mathcal{T}, \Sigma, g, t_0, \mathcal{T})$ automatul LR(1) corespunzător. Spunem că gramatica G este LALR(1) (**L**ook **A**head **LR**(1)) dacă oricare ar fi perechea de stări echivalente t_1, t_2 din automatul LR(1), starea $t_1 \cup t_2$ este liberă de conflicte.

Am introdus aceste noțiuni în scopul de a simplifica, atunci când este posibil, tabela de analiză LR(1), în sensul de a restrânge numărul liniilor tabelului. Acest lucru este posibil dacă prin reuniunea a două stări echivalente nu se introduc conflicte. În acest mod

putem obține o tabelă de analiză de tip LR(1) dar cu numărul de linii (stări ale automatului) egal cu cel al tabelii SLR(1), adică, cu cel al numărului de stări LR(0). În exemplul considerat mai sus avem:

$$4 \cup 11 = 4 \quad 5 \cup 12 = 5 \quad 7 \cup 13 = 7 \quad 8 \cup 10 = 8$$

Așadar, gramatica în discuție este LALR(1) pentru că nu s-au obținut stări cu conflicte. În continuare se reconstruiește automatul LR(1) prin reuniunea stărilor echivalente și se construiește tabela de analiză pornind de la acest nou automat LR(1). Mai facem observația că, dacă stările t_1 și t_2 sunt echivalente, atunci și stările $g(t_1, X)$, $g(t_2, X)$ sunt echivalente, pentru orice X din Σ . Să descriem în întregime acest algoritm.

Algoritmul 5.6.2 (Construcția tabelii de analiză LALR(1))

Intrare: Gramatica $G = (V, T, S, P)$ augmentată cu $S' \rightarrow S$;

Ieșire: Tabela de analiză LALR(1) (*ACȚIUNE* și *GOTO*).

Metoda:

1. Se construiește automatul LR(1), $\mathcal{M} = (\mathcal{T}, \Sigma, g, t_0, \mathcal{T})$ folosind Algoritmul 5.5.2. Fie $\mathcal{T} = \{t_0, t_1, \dots, t_n\}$. Dacă toate stările din \mathcal{T} sunt libere de conflict, urmează 2, altfel algoritmul se oprește: gramatica nu este LR(1).
2. Se determină stările echivalente din \mathcal{T} și, prin reuniunea acestora, se obține o nouă mulțime de stări $\mathcal{T}' = \{t'_0, t'_1, \dots, t'_m\}$
3. Dacă în \mathcal{T}' există stări ce conțin conflicte, algoritmul se oprește: gramatica G nu este LALR(1). Altfel se continuă cu 4.
4. Se construiește automatul $\mathcal{M}' = (\mathcal{T}', \Sigma, g', t'_0, \mathcal{T}')$, unde $\forall t' \in \mathcal{T}'$:
 Dacă $t' \in \mathcal{T}$ atunci $g'(t', X) = g(t, X)$, $\forall X \in \Sigma$;
 Dacă $t' = t_1 \cup t_2 \cup \dots, t_1, t_2, \dots \in \mathcal{T}$, atunci
 $g'(t', X) = g(t_1, X) \cup g(t_2, X) \cup \dots$
5. Se aplică algoritmul 5.6.1 pentru construirea tabelii LR(1) pornind de la automatul \mathcal{M}' . Tabela obținută se numește tabela LALR(1) pentru gramatica G .

Exemplul 5.6.2. Aplicând algoritmul de mai sus automatului \mathcal{M} din exemplul 5.5.1, obținem automatul \mathcal{M}' dat de tabela de tranziție de mai jos (stările au fost redenumite prin 0, 1, ..., 9).

g	a	=	*	S	L	R
0	5		4	1	2	3
1						
2		6				
3						
4	5		4		8	7
5						
6	5		4		8	9
7						
8						
9						

Pornind de la acest automat, aplicând algoritmul de construcție al abelei LR(1), obținem tabela de parsare LALR(1) pentru gramatica dată:

STARE	ACȚIUNE				GOTO		
	a	=	*	#	S	L	R
0	D5		D4		1	2	3

1				acceptare			
2		D6		R5			
3				R2			
4	D5		D4			8	7
5		R4		R4			
6	D5		D4			8	9
7		R3		R3			
8		R5		R5			
9				R1			

Analiza cuvântului $w = **a=a$ este ilustrată în tabela următoare:

STIVA	INTRARE	ACȚIUNE	IEȘIRE
0	$**a=a\#$	deplasare D4	
04	$*a=a\#$	deplasare D4	
044	$a=a\#$	deplasare D5	
0445	$=a\#$	reducere R4	$L \rightarrow a$
0448	$=a\#$	reducere R5	$R \rightarrow L$
0447	$=a\#$	reducere R3	$L \rightarrow *R$
048	$=a\#$	reducere R5	$R \rightarrow L$
047	$=a\#$	reducere R3	$L \rightarrow *R$
02	$=a\#$	deplasare D6	
026	$a\#$	deplasare D5	
0265	$\#$	reducere R4	$L \rightarrow a$
0268	$\#$	reducere R5	$R \rightarrow L$
0269	$\#$	reducere R1	$S \rightarrow L=R$
01	$\#$	acceptare	

Dacă vom compara tabela LR(1) cu tabela LALR(1) pentru gramatica din exemplul dat, vom observa că tabela LALR(1) coincide cu liniile 0-9 din tabela LR(1) cu deosebirea că în locul stărilor 10, 11, 12 și 13 apar stările 8, 4, 5, 7 respectiv. Acest fapt face ca erorile, în unele cuvinte, să fie semnalate mai târziu în analiza LALR(1) decât s-ar face în analiza LR(1). În exemplul precedent, pentru tabela LR(1) avem $ACȚIUNE(12, =) = eroare$ pe când $ACȚIUNE(5, =) = R4$ încât, pentru că stările 5 și 12 sunt echivalente, în tabela LALR(1) rămâne doar $ACȚIUNE(5, =) = R4$. Cititorul este invitat să verifice că în analiza LR(1), pentru cuvântul de intrare $*a=a\#$, se obține o eroare la configurația $(026(12), =\#, \pi)$, pe când algoritmul LALR(1) va ajunge, pentru același cuvânt, în configurația $(0265, =\#, \pi)$ care are tranziție în configurația $(0268, =\#, \pi4)$ apoi în $(0269, =\#, \pi45)$ și abia aceasta din urmă produce eroarea.

□

Avantajul metodei LALR(1) constă în aceea că dimensiunea tabelului LALR(1) este aceeași cu cea a tabelului SLR(1); tabela LR(1) pentru o gramatică care descrie părți ale sintaxei unui limbaj de programare poate fi foarte mare!

Să arătăm în continuare că există gramatici LR(1) care nu sunt LALR(1).

Exemplul 5.6.3 Fie gramatica:

$$S \rightarrow aAb \mid bAd \mid aBd \mid bBc$$

$$A \rightarrow e \quad B \rightarrow e$$

Stările automatului LR(1) pentru această gramatică sunt cele din figura 5.6.1.

0	1	2	3
$(S \rightarrow \bullet aAb, \#)$ $(S \rightarrow \bullet bAd, \#)$ $(S \rightarrow \bullet aBd, \#)$ $(S \rightarrow \bullet bBc, \#)$	$(S \rightarrow a\bullet Ab, \#)$ $(S \rightarrow a\bullet Bd, \#)$ $(A \rightarrow \bullet e, b)$ $(B \rightarrow \bullet e, d)$	$(S \rightarrow b\bullet Ad, \#)$ $(S \rightarrow b\bullet Bc, \#)$ $(A \rightarrow \bullet e, d)$ $(B \rightarrow \bullet e, b)$	$(S \rightarrow aA\bullet b, \#)$
5	6	7	8
$(A \rightarrow e\bullet, b)$ $(B \rightarrow e\bullet, d)$	$(S \rightarrow bA\bullet d, \#)$	$(S \rightarrow bB\bullet c, \#)$	$(A \rightarrow e\bullet, d)$ $(B \rightarrow e\bullet, b)$
9	10	11	12
$(S \rightarrow aAb\bullet, \#)$	$(S \rightarrow aBd\bullet, \#)$	$(S \rightarrow bAd\bullet, \#)$	$(S \rightarrow bBc\bullet, \#)$

Figura 5.6.1

Să observăm că toate stările automatului LR(1) sunt libere de conflicte; gramatica este LR(1). Stările 5 și 8 sunt stări echivalente. Constatăm că $5 \cup 8 = \{(A \rightarrow e\bullet, \{b, d\}), (B \rightarrow e\bullet, \{b, d\})\}$, iar această nouă mulțime conține conflicte de reducere/reducere, ceea ce înseamnă că gramatica considerată nu este gramatică LALR(1).

5 Generatoare de analizoare sintactice

Pentru a ușura munca implementatorilor de limbaje au fost concepute instrumente care generează automat analizoare (lexicale sau sintactice). Printre cele mai populare instrumente de acest fel se numără sistemul YACC (“**Y**et **A**nother **C**ompiler **C**ompiler”) conceput și realizat în 1975 de S. C. Johnson la Bell Laboratory în USA. Produsul a fost inclus în sistemul de operare UNIX și a fost folosit la scară industrială pentru implementarea multor limbaje. Un alt produs, distribuit în sistemele de operare Linux este Bison, compatibil în întregime cu YACC, scris de Robert Corbett și Richard Stallmann. În 1990 a apărut versiunea pentru MS-DOS a produsului YACC, numit PCYACC, dezvoltat de Abraxax Software Inc.

Există astăzi numeroase astfel de generatoare, construite de grupuri de cercetare din universități sau din societăți specializate pe software. Enumerăm câteva din acestea împreună cu resursele web corespunzătoare.

La adresa <http://dinosaur.compilertools.net/> se găsesc informații (inclusiv documentație) despre Lex, Yacc, Flex, și Bison.

La adresa <http://epaperpress.com/lexandyacc/index.html> se găsește un ghid compact care este util celor care vor să utilizeze lex și yacc.

Pagina oficială destinată proiectului GNU Flex (Fast lexical analyser generator) este <http://www.gnu.org/software/flex/>. Așa cum se preciza și în capitolul 3.5, Flex este un instrument pentru generare de programe care realizează pattern-matching pentru texte iar împreună cu Bison constituie perechea de instrumente (www.gnu.org/software/bison/) cel mai des utilizată pentru scrierea compilatoarelor. Bison, ca și Yacc, este un generator de parsere care convertește descrierea unei gramatici context – free LALR(1) într-un program C care realizează analiza sintactică în gramatica descrisă.

Iată și alte instrumente de acest fel:

- *Spirit C++ Parser Framework* - <http://spirit.sourceforge.net/>

Spirit este un instrument ce implementează parsarea recursiv descendentă pornind de la descrierea EBNF(Extended Backus Normal Form) a unei gramatici. Implementarea este orientată obiect în C++.

- *BtYacc: BackTracking Yacc* - <http://www.siber.com/btyacc/>

BtYacc este o versiune modificată a lui Yacc care suportă backtraking. Pentru corectarea bug-urilor vizitați <http://www.vendian.org/mncharity/dir3/btyacc/>

- *LEMON* - <http://www.hwaci.com/sw/lemon/>

LEMON este un generator de parsere LALR(1), este open-source, produce cod C și pretinde că produce un parser mai rapid decât yacc/bison.

- *SYNTAX*: <http://www-rocq.inria.fr/oscar/www/syntax/>

SYNTAX cuprinde un set de instrumente pentru proiectarea și implementarea părții front-end a unui compilator. Are toate capacitățile lex și yacc iar partea de procesare a erorilor este mult dezvoltată.

- *PCCTS* : <http://www.polhode.com/pccs.html>

PCCTS este un generator pentru parsare recursiv descendentă LL(k). Proiectul s-a transformat în ANTLR.

- *ANTLR*: <http://www.antlr.org>

ANTLR (ANother Tool for Language Recognition) este un instrument pentru construcția compilatoarelor care suportă descrierea acțiunilor în Java, C# sau C++. ANTLR conține module pentru construcția arborilor, traversarea arborilor, construcția translaatoarelor. ANTLR implementează parsarea LL(k) și este în domeniul public.

- *Visual Parse++ 5.0* - <http://www.sand-stone.com/>

Visual Parse permite proiectarea vizuală a analizoarelor lexicale și sintactice și utilizarea lor în aplicații C, C++, C#, Delphi, Java, sau Visual Basic sub sistemele de operare UNIX, Linux și Windows. Este un produs comercial.

- *AnaGram* - <http://www.parsifalsoft.com/>

AnaGram este un generator de parsere LALR(1) care produce C/C++ ce poate fi compilat pe orice platformă și poate fi executat sub Win9x/NT. Este un produs comercial dar firma oferă și o variantă gratuită. De asemenea, este disponibil XIDEK (Extensible Interpreter Development Kit) ce poate fi folosit pentru proiectarea și implementarea interpretelor.

- *Yacc++(R)* - <http://world.std.com/~compres/>

Yacc++(R) este o rescriere orientată obiect a instrumentelor Lex și Yacc pentru C++.

- *LLgen* - <http://www.cs.vu.nl/~ceriel/LLgen.html>

LLgen implementează eficient parsere recursiv descendente ale gramaticilor ELL(1). Dispune de facilități statice și dinamice pentru tratarea ambiguităților.

- *Elkhound* - <http://www.cs.berkeley.edu/~smcpeak/elkhound/>

Este un generator de parsere bazat pe tehnica GLR; fiind dezvoltat la Universitatea Berkeley este open source.

- *Grammatica* - <http://www.nongnu.org/grammatica/>

Grammatica este un generator de parsere LL(k) pentru C# și Java. Este distribuit sub licența GNU și suportă gramatici cu un număr nelimitat de token-uri look-ahead. De asemenea, are un modul puternic pentru recuperarea automată a erorilor.

- *jay* este o variantă Java pentru Yacc.

<http://www.informatik.uni-osnabrueck.de/alumni/bernd/jay/>

- *RDP* este un generator de parsere pentru gramatici LL(1) cu atribute.

<http://www.dcs.rhnc.ac.uk/research/languages/projects/rdp.shtml>

- *TPG* - <http://christophe.delord.free.fr/en/tpg/>

TPG (Toy Parser Generator) este un generator care, pentru o gramatică cu atribute produce un parser recursiv descendent în limbajul Python.

- *ProGrammar* - <http://www.programmar.com/main.shtml>

ProGrammar este un mediu vizual pentru proiectarea de parsere ce sunt independente de platforma și de limbajul de programare.

- *GOLD Parser* - <http://www.devincook.com/goldparser/>

GOLD (Grammar Oriented Language Developer) este un generator care implementează parsarea LALR(1) într-o nouă manieră (față de yacc de exemplu): descrierea gramaticii este analizată de un modul – builder – care creează tabela de parsare ce este salvată într-un fișier separat, independent de parser. Acesta este încărcat de motorul de parsare care este disponibil pentru limbajele Java, C# (.NET), ActiveX.

- *Win32 Lex/YACC* - <http://www.monmouth.com/~wstreet/lex-yacc/lex-yacc.html> Flex și Bison pentru platforma Win32.

- *Depot4* - <http://www.math.tu-dresden.de/wir/depot4/>

Depot4 este un generator de translator dezvoltat la Universitatea din Dresda.

- *IParse* - <http://home.planet.nl/~faase009/MM.html>

Program ce implementează un parser greedy backtracking. Are la intrare un fișier ce conține descrierea unei gramatici și un altul ce conține fraza de parsat în conformitate cu gramatica dată. Produce arborele abstract al frazei.

- *Styx* - <http://www.speculate.de/styx/>

Styx este un generator de scanner și parser LALR(1). Este disponibil atât sub Lynx, cât și sub Windows.

- *YaYacc* (*Yet another Yacc*) acceptă gramatici yacc și produce cod C++. <http://www.gradsoft.com.ua/eng/Products/YaYacc/yayacc.html>
- *Rie* - <http://www.is.titech.ac.jp/~sassa/lab/rie-e.html>

Rie este un generator de compilator (front end) bazat pe o clasă de gramatici LR cu atribut numită ECLR-attributed grammar. Generatorul este scris în C, este open source și este cu siguranță o extensie a sistemelor Yacc/Bison.

- *Coco/R*, <http://www.ssw.uni-linz.ac.at/Research/Projects/Compiler.html>. *Coco/R* este un generator de scanner și parser LL(1). Sunt disponibile versiuni Java, C#, Oberon, C, Pascal, Modula2.

- *PRECC* - <http://www.afm.sbu.ac.uk/precc/>
PREC(PREttier Compiler-Compiler) este un generator de parser pentru gramatici dependente de context care generează cod C. Este disponibil sub Unix sau sub MS-Dos.

- *SGLR*: - <http://www.cwi.nl/projects/MetaEnv/sglr/>
SGLR (Scannerless Generalized LR) este un generator de parser LR care nu utilizează scanarea: tabela de parsare generată de la o definiție a sintaxei conține informații suficiente pentru a realiza analiza sintactică.

- *SLK* - <http://home.earthlink.net/~slkpg/>
Este un generator de parser LL(k) ce produce cod C, C++, C#, și Java. Disponibil gratuit pentru Unix și Windows.

- *YooLex* (*Yet another Object-Oriented Lex*)
<http://yoolex.sourceforge.net/>. Produce cod C++ compatibil cu STL.

- *Happy* - <http://haskell.cs.yale.edu/happy/> Generator pentru Haskell.
- *VLCC* (*Visual Languages Compiler-Compiler*)

www.dmi.unisa.it/people/costagliola/www/home/ricerca/vlcc/vlcc.htm

- *Hapy* - <http://hapy.sourceforge.net/> generează un parser scris în C++ folosind descrierea sintaxei cu o interfață similară cu EBNF.

- *CppCC* (*C++ Compiler Compiler*) - <http://cppcc.sourceforge.net/>
Este un generator de scanner și parser LL(k).

- *ClearParse* - <http://www.clearjump.com/products/ClearParse.html>
Parser comercial.

- *TPLex/Yacc for Delphi 3* - <http://www.17slon.com/gp/gp/tply.htm>
- *JB2CSharp* - <http://sourceforge.net/projects/jb2csharp/>

Este un proiect al Universității Boulder – Colorado de a porta instrumentele Java-Bison/Flex pentru platforma .NET: acțiunile pot fi scrise în C#.

- *oolex* <http://www.inf.uos.de/alumni/bernd/oolex/>
oolex (object-oriented lexer) este un generator pentru analiza lexicală construit în manieră orientată obiect. Există și un succesor al lui *oolex* – *lolo* care se găsește la <http://www.inf.uos.de/alumni/bernd/lolo/index.html>. De altfel, împreună cu *oops* (<http://www.inf.uos.de/alumni/bernd/oops/>) aceste sisteme fac parte dintr-un proiect dezvoltat la Universitatea din Osnabruck privind construcția compilatoarelor folosind programarea orientată obiect în Java.

- *EAG* - <http://www.cs.ru.nl/~kees/eag/>
Acest sistem folosește formalismul Extended Affix Grammar (EAG) care descrie sintaxa și semantica limbajelor de programare.

- *Bison++*, *Flex++*: <http://www.kohsuke.org/flex++bison++/>
Sunt variantele Bison, Flex care produc cod C++.

5.1 Utilizarea generatorului de parsere YACC

YACC poate genera un “traducător” în modul următor: Specificațiile sintactice precum și unele elemente semantice sunt incluse într-un fișier de intrare pentru YACC. Acest fișier are în general extensia “.y”. Linia de comandă (în UNIX):

```
yacc nume_fisier.y
```

produce, plecând de la specificațiile cuprinse în fișierul de intrare nume_fisier.y, un program scris în limbaj C care implementează metoda de analiză sintactică LALR(1) pentru gramatica respectivă – în speță funcția yyparse(). Fișierul de intrare (programul sursă) pentru YACC are structura următoare:

Declarații

%%

Reguli

%%

Rutine C

Partea *Declarații* a fișierului sursă este opțională și poate cuprinde două secțiuni. O primă secțiune cuprinde între delimitatorii % { și % } declarații în limbajul C pentru variabilele care se utilizează în regulile de traducere sau în procedurile din cea de-a treia parte a fișierului. Așadar, textul cuprins între % { și % } se copie nealterat în fișierul C produs de YACC. A doua secțiune a acestei prime părți conține declarații ale unităților lexicale ale gramaticii, declarații de asociativitate și precedență a operatorilor, declarații ale tipurilor de date pentru valorile semantice ale simbolurilor gramaticii.

Simbolurile terminale ale gramaticii – unitățile lexicale – cu excepția celor formate dintr-un singur caracter (precum +, * etc.), trebuie declarate. Aceste simboluri sunt reprezentate în YACC prin niște coduri numerice; funcția yylex transmite codul unității respective funcției yyparse. Programatorul nu trebuie să știe aceste coduri numerice, ele sunt generate automat folosind opțiunea -d la lansarea lui YACC și sunt trecute într-un fișier nume.tab.h (sau yytab.h sub DOS). Unitățile lexicale se definesc prin linii de forma:

```
% token <nume _unitate _lexicală>
```

și pot fi utilizate în celelalte două părți ale fișierului de intrare. Tot în această secțiune se poate introduce simbolul de start al gramaticii, în cazul în care acesta nu este partea stângă a primei reguli sintactice:

```
start <nume _simbol _de _start>
```

Alte linii care pot fi incluse în această secțiune sunt:

type - pentru definirea tipului;

left - pentru asociativitate stânga a operatorilor;

right - pentru asociativitate dreapta a operatorilor;

prec - pentru precizarea precedenței operatorilor;

nonassoc - pentru declarațiile de neasociativitate.

Simbolurile gramaticii pot avea valori semantice. În mod implicit, aceste valori sunt valori întregi și sunt specificate în YYSTYPE. Pentru a specifica alt tip, în partea de declarații C se adaugă o directivă define pentru YYSTYPE:

```
#define YYSTYPE <nume tip>
```

declară tipul valorilor semantice pentru simbolurile gramaticii ca fiind <nume tip>. Utilizarea de tipuri diferite pentru simboluri diferite se realizează prin specificarea acestor tipuri într-o declarație union (specifică YACC) și declararea tipului simbolurilor cu type. Iată un exemplu:

```
%union {
    int intval;
```

```

        double doubleval;
        symrec *tptr;
    }
    type <intval> INT
    type <doubleval> REAL
    type <tptr> ID

```

Prin aceste declarații se specifică tipul valorilor pentru token-urile INT, REAL, ID ca fiind respectiv int, double, pointer la symrec (pointer în tabela de simboluri).

Partea a doua a fișierului de intrare, *Reguli*, este partea obligatorie. Aici se specifică regulile gramaticii care descriu sintaxa pe care dorim să o verificăm cu acest analizor. Fiecare regulă conține:

- partea stângă;
- partea dreaptă;
- partea de acțiune.

Partea stângă trebuie să conțină un neterminal al gramaticii, iar partea dreaptă șirul format din terminali și neterminali corespunzător unei reguli. Cele două părți sunt despărțite prin “:” (două puncte). Partea de acțiune, cuprinsă între { și }, conține un text scris în limbajul C care va fi inclus în analizor și reprezintă operațiile care se execută atunci când analizorul realizează o reducere cu regula specificată. O regulă se termină prin “;” (punct virgulă). Dacă sunt mai multe reguli care au aceeași parte stângă, acestea se scriu o singură dată și părțile drepte se despart prin “|”. Iată, spre exemplu, cum se scriu regulile gramaticii care descrie expresiile aritmetice:

```

expr : NUMAR
    | expr '+' expr
    | expr '-' expr
    | expr '*' expr
    | expr '/' expr
    | '(' expr ')'
;

```

Aici, expr este numele neterminalului care definește expresia aritmetică, terminalii +, -, *, /, (,) se pun între apostrof, iar NUMAR reprezintă token-ul (unitatea lexicală) număr (ce trebuie declarat în secțiunea declarații). Dacă dorim ca analizorul pe care-l construim să realizeze și acțiuni semantice, de pildă să evalueze expresiile (pentru șirul 25+15, pe lângă faptul că verifică sintaxa, raportează și rezultatul adunării, valoarea 40), se adaugă la reguli partea de acțiune:

```

expr : NUMAR      {$$ = $1;}
    | expr '+' expr {$$ = $1+$3;}
    | expr '-' expr {$$ = $1-$3;}
    | expr '*' expr {$$ = $1*$3;}
    | expr '/' expr {$$ = $1/$3;}
    | '(' expr ')' {$$ = $2;}
;

```

În partea de acțiune se scriu instrucțiuni în limbajul C. Aici, \$\$ reprezintă valoarea unui atribut al neterminalului expr din stânga regulii sintactice iar \$i reprezintă valoarea atributului celui de-al i-lea simbol din partea stângă a regulii.

Ultima parte a fișierului de intrare conține rutine scrise în limbaj C care se includ nealterate în analizorul obținut de YACC. În această parte trebuie furnizat un analizor lexical, adică o funcție yylex() precum și apelul (în programul principal) la funcția yyparse() pe care o creează YACC-ul.

Exemplul 6.1.1

```

/* Exemplul 1 de fisier intrare YACC */
%{
int yylex(void);
void yyerror(char* s);

```

```

%}
%token ID
%token WHILE
%token BEGIN
%token END
%token DO
%token IF
%token THEN
%token ELSE
%token SEMI
%token ASSIGN
%start prog
%%
prog: stmlist
    ;
stm: ID ASSIGN ID
    | WHILE ID DO stm
    | BEGIN stmlist END
    | IF ID THEN stm
    | IF ID THEN stm ELSE stm
    ;
stmlist: stm
        | stmlist SEMI stm
        ;
%%

```

Dacă se compilează acest fișier cu YACC, acesta construiește automatul LALR(1) pentru gramatica descrisă în secțiunea Reguli. Această gramatică are 3 neterminali: prog (simbolul de start al gramaticii), stm și stmlist. Terminalii sunt declarați prin directiva %token. YACC raportează conflictele de deplasare – reducere și cele de reducere-reducere dacă există. Aceste conflicte sunt rezolvate de către YACC astfel:

- conflictul de deplasare-reducere este rezolvat în favoarea deplasării;
- conflictul de reducere-reducere este rezolvat în favoarea reducerii cu regula care apare prima în lista de reguli

Pentru descrierea precedentă este raportat un singur conflict de deplasare reducere (este obținut datorită celor două forme ale instrucțiunii IF). Automatul LALR(1) obținut este vizibil dacă în linia de comandă pentru lansarea YACC – ului se scrie opțiunea -v (pentru a afla ce opțiuni pot fi folosite lanșați YACC -h). Iată cum se descrie automatul în exemplul dat:

```

- * - - * - - * - - * - - * - LALR PARSING TABLE - * - - * - - * - - * - -
+-----+ STATE 0 -----+
+ CONFLICTS:
+ RULES:
    $accept : ^prog $end
+ ACTIONS AND GOTOS:
    ID : shift & new state 4
    WHILE : shift & new state 5
    BEGIN : shift & new state 6
    IF : shift & new state 7
        : error
    prog : goto state 1
    stmlist : goto state 2
    stm : goto state 3
+-----+ STATE 1 -----+
+ CONFLICTS:
+ RULES:
    $accept : prog^$end
+ ACTIONS AND GOTOS:

```

```

$end : accept
      : error
+----- STATE 2 -----+
+ CONFLICTS:
+ RULES:
    prog : stmlist^ (rule 1)
    stmlist : stmlist^SEMI stm
+ ACTIONS AND GOTOS:
    SEMI : shift & new state 8
          : reduce by rule 1
.....
+----- STATE 21 -----+
+ CONFLICTS:
+ RULES:
    stm : IF ID THEN stm ELSE stm^ (rule 6)
+ ACTIONS AND GOTOS:
    : reduce by rule 6
===== SUMMARY =====
.....
-*--*-*-*-*-*- END OF TABLE -*--*-*-*-*-

```

1

Exemplul 6.1.2

```

/* Intrare YACC pentru expresii aritmetice de forma a*(a+a) */
%{
#define QUIT      ((double) 101010)
%}

%start S

%%

S:      { prompt(); }
      | S '\n'      { prompt(); }
      | S E '\n'    { if ($2 == QUIT) {
                      return(0);
                    } else {
                      printf("Expresia este corecta.\n");
                      prompt();
                    }
      }

;

E : E '+' T      {printf(" E->E+T\n");}
  | E '-' T      {printf(" E->E-T\n");}
  | T            {printf(" E->T\n");}

;

T : T '*' F      {printf(" T->T*F\n");}
  | T '/' F      {printf(" T->T/F\n");}
  | F            {printf(" T->F\n");}

;

F : '(' E ')'    {printf(" F->(E)\n");}
  | 'a'          {printf(" F->a\n");}

;

%%

#include <stdio.h>
#include <ctype.h>

int main(){
printf("\n*****\n");
printf("***                               **\n");
printf("***      Parser pentru expresii aritmetice      **\n");
printf("***      cu operanzi a si operatori +,-,*,/, ( )  **\n");
printf("***      La promterul READY>                        **\n");

```

```

printf("***   Introduceți o expresie. Ex. a*(a+a)-a   **\n");
printf("***           Pentru ieșire tastati quit           **\n");
printf("***                                           **\n");
printf("\n*****\n");
yyparse();
}
yyerror(s){
    printf("Expresia este incorecta\n");
}
yylex(){
    int c;
    while ((c=getchar()) == ' ' || c == '\t' )
        ;
    if(c==EOF)
        return 0;
    if (c == 'Q' || c == 'q')
        if ((c=getchar()) == 'U' || c == 'u')
            if ((c=getchar()) == 'I' || c == 'i')
                if ((c=getchar()) == 'T' || c == 't') {
                    yylval = QUIT;
                    return( EOF );
                }
            else return '?';
    return c;
}
void prompt( void ){
    printf("READY> ");
}

```

YACC poate fi folosit pentru proiectarea de aplicații diverse. Dezvoltarea unui proiect cu ajutorul YACC - ului presupune parcurgerea următoarelor etape:

1. **Identificarea problemei.** Acest lucru presupune o viziune de ansamblu asupra proiectului stabilindu-se arhitectura sistemului ce urmează a fi proiectat, descompunerea sistemului în funcțiile componente, etc;
2. **Definirea limbajului sursă.** Obiectul acestui pas este specificarea limbajului care urmează a fi tradus. Cel mai adecvat mecanism pentru această specificare este gramatica independentă de context pentru că se poate descrie această gramatică direct în limbajul YACC;
3. **Scrierea programului** de descriere a gramaticii în fișierul de intrare pentru YACC;
4. **Scrierea codului auxiliar.** În acest pas se stabilesc și se scriu acțiunile atașate regulilor sintactice, funcțiile necesare lansării aplicației (main(), yylex(), yyerror()) precum și alte instrucțiuni C care urmează a fi încorporate în programul generat de YACC;
5. **Obținerea funcției yyparse().** Odată creat fișierul de intrare, se lansează YACC pentru acesta și se obține programul sursă C. Dacă sunt necesare și alte funcții C, acestea se pot incorpora în programul obținut de YACC;
6. **Compilarea textului** obținut și obținerea programului executabil.

5.2 Aplicații cu LEX și YACC

Instrumentele Lex și Yacc pot fi folosite împreună pentru a realiza analiza lexicală și sintactică a unui text și chiar pentru a obține un interpretor sau un compilator prin mecanismul acțiunilor semantice. Schema de utilizare împreună a celor două generatoare este dată în figura 6.2.1. Dacă fișierul calc.l descrie unitățile lexicale ale limbajului pe care

dorim să-l implementăm iar `calc.y` descrie sintaxa limbajului, atunci comenzile (Linux) pentru a obține fișierul `calc.exe` din `calc.l` și `calc.y` sunt următoarele:

```
yacc -d calc.y
lex calc.l
cc lex.yy.c y.tab.c -o calc.exe
```

Comanda `yacc` produce fișierul `y.tab.c` care conține funcția de parsare `yyparse()`. Opțiunea `-d` produce fișierul header `y.tab.h` care conține directivele ce atribuie unităților lexicale coduri numerice. Acestea vor fi folosite de analizorul lexical produs de `lex`. Fișierul `lang.l` va trebui să conțină directiva `#include y.tab.h`.

Comanda `lex` va produce fișierul `lex.yy.c` care conține definiția funcției `yylex()` care este invocată de `yyparse()` pentru obținerea, pe rând, a unităților lexicale din textul sursă. Cele două fișiere `.c` vor fi compilate împreună, și se obține în urma editării legăturilor executabilul `lang.exe` ce poate fi lansat pentru a analiza un text sursă din limbajul dat.

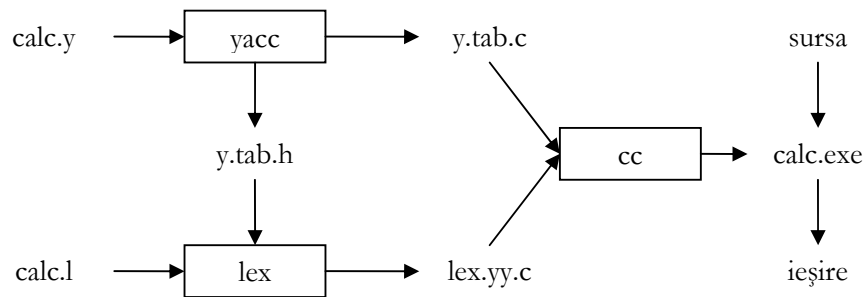


Figura 6.2.1

Exemplul 6.2.1 Implementarea unui calculator de buzunar

Vom construi fișierele pentru `lex` și `yacc` care descriu sintaxa expresiilor aritmetice ce se pot forma cu constante întregi (CINT), variabile (VAR), operațiile aritmetice uzuale și paranteze. Desigur, orice variabilă ce poate apărea într-o expresie trebuie să fie inițializată cu valoarea unei expresii. Astfel, un program în accepțiunea de aici este o succesiune de expresii aritmetice și/sau asignări. Gramatica ce descrie un astfel de limbaj este următoarea:

```

program      → program statement | ε
statement    → expression | VAR '=' expression
expression   → CINT | VAR
              | expression '+' expression
              | expression '-' expression
              | expression '*' expression
              | expression '/' expression
              | '(' expression ')'
```

Fișierul pentru `yacc` `calc.y` de mai jos este scris în așa fel încât acțiunile semantice au ca efect obținerea valorii numerice a unei expresii atunci când ea este transmisă la intrare programului `calc.exe` ce se obține. Tabloul `sym[26]` este folosit pentru a păstra valorile variabilelor ce apar în programe; o variabilă este aici o literă (va fi descrisă în fișierul pentru `lex`).

```
// Fișierul calc.y
%{
    #include <stdio.h>
    void yyerror(char *);
    int yylex(void);
    int sym[26];
%}
```

```

%token CINT VAR
%left '+' '-'
%left '*' '/'

%%

program:
    program statement '\n'
    | /* NULL */
    ;

statement:
    expression {printf("%d\n", $1);}
    | VAR '=' expression { sym[$1] = $3; }
    ;

expression:
    CINT
    | VAR { $$ = sym[$1]; }
    | expression '+' expression { $$ = $1 + $3; }
    | expression '-' expression { $$ = $1 - $3; }
    | expression '*' expression { $$ = $1 * $3; }
    | expression '/' expression { $$ = $1 / $3; }
    | '(' expression ')' { $$ = $2; }
    ;

%%

void yyerror(char *s) {
    fprintf(stderr, "%s\n", s);
}

int main(void) {
    yyparse();
}

```

Declararea token-urilor CINT și VAR prin %token INT VAR în fișierul calc.y produce obținerea fișierului y.tab.h (dacă la invocarea lui yacc se folosește parametrul -d) în care sunt definite valorile numerice pentru aceste token-uri. Acest fișier, care este listat mai jos, trebuie inclus în fișierul calc.l.

```

// Fisierul y.tab.h
#ifndef YYSTYPE
#define YYSTYPE int
#endif
#define CINT 258
#define VAR 259
extern YYSTYPE yylval;

```

Fișierul calc.l conține descrierea token-urilor VAR și CINT. Valorile asociate token-urilor sunt returnate de lex în variabila (standard) yylval. Pentru variabile yylval va însemna indexul în tabloul sym iar pentru constante yylval este numărul întreg corespunzător (scanat de analizorul lexical). Variabila yytext este pointer la unitatea lexicală tocmai depistată, de aceea pentru operatori (token-uri semn) definiția în fișierul lex este:

```
[-+()=/*\n] { return *yytext; }
```

```

// Fisierul calc.l

%{
    #include "y.tab.h"

```

```

#include <stdlib.h>
void yyerror(char *);

%}
%%

[a-z]      {
            yylval = *yytext - 'a';
            return VAR;
        }

[0-9]+     {
            yylval = atoi(yytext);
            return CINT;
        }

[+-()=/*\n] { return *yytext; }

[ \t]      ;      /* ignora spatiile */

.          yyerror("Caracter necunoscut");

%%

int yywrap(void) {
    return 1;
}

```

După obținerea fișierului calc.exe după schema din figura 6.2.1, o sesiune de lucru cu acest program ar putea fi:

```

utilizator:    x=4
calc:          y=12
utilizator:    (x-y) * (x-y)
calc:          64
utilizator:    a=3
utilizator:    b=7
utilizator:    u=y*y-x*x
utilizator:    v=b*b-a*a
utilizator:    u+v
calc:          168
utilizator:    (u+v) / 2 + (u-v) / 2
calc:          128
utilizator:    u
calc:          128

```


6 Analiza semantică

6.1 Gramatici cu atribute

În 1968 Donald Knuth [Knu68] introduce un formalism pentru specificarea semanticii limbajelor de programare: gramatici independente de context care exprimă *sintaxa* și se adaugă un sistem de atribute care exprimă *semantica*. Acest formalism s-a dovedit deosebit de util în *analiza semantică* în cadrul compilatoarelor. O *gramatică cu atribute* este o extindere a unei gramatici independente de context: fiecare simbol al acesteia are atașată o mulțime de *tribute* care pot lua valori într-o mulțime precizată. Valorile atributelor se calculează după niște reguli atașate producțiilor gramaticii, numite *reguli semantice*. O *regulă semantică* definește modul de calcul al unui atribut al părții stângi a producției – și atunci atributul se zice că este *sintetizat* – sau al unui atribut al unui simbol din partea dreaptă a producției – și atunci atributul se zice că e *moștenit*. Valoarea unui atribut este definită în funcție de valorile altor atribute ale simbolurilor din aceeași producție. Se justifică denumirea de *atribut sintetizat*: valorile sale depind de valorile atributelor fiilor într-un arbore de derivare (eventual și de unele atribute proprii) și se propagă în arbore de la frunze către rădăcină, precum și de *atribut moștenit*: valorile sale depind de valorile atributelor părinților și fraților, propagându-se de la rădăcină către frunze. Într-un arbore de derivare al gramaticii independente de context un anume atribut poate avea mai multe apariții, după cum simbolul la care aparține are mai multe apariții. Valorile acestora, desigur, pot fi diferite. Este posibil ca simboluri diferite ale gramaticii să aibă atribute cu același nume; acestea se consideră totuși atribute diferite. De pildă, dacă simbolurile X și Y au câte un atribut numite val , atunci acestea se exprimă prin $X.val$ respectiv $Y.val$. Simbolul de start al gramaticii are un atribut sintetizat care este desemnat a defini *semantica cuvântului* w din limbajul gramaticii. Este natural să considerăm că o gramatică cu atribute care are ca suport (sintactic) gramatica G este bine definită (lucru ce se referă la regulile semantice), dacă pentru orice cuvânt $w \in L(G)$ se poate determina semantica acestuia, respectiv, valoarea atributului desemnat pentru simbolul de start în arborele de derivare al cuvântului w . De asemenea, este natural să spunem că un anume atribut este *util* dacă el contribuie la calculul semanticii lui w , adică semantica lui w depinde, direct sau indirect, de valoarea acestui atribut. Astfel, o gramatică cu atribute se zice că este *bine definită (necirculară)* dacă, pentru orice arbore de derivare construit în gramatica suport, se pot calcula valorile tuturor aparițiilor atributelor utile în acest arbore (se mai spune că arborele poate fi decorat). Așadar, există două probleme importante legate de gramaticile cu atribute:

- **Problema circularității:** este gramatica cu atribute bine definită?
- **Problema evaluării:** care este cea mai bună strategie de evaluare a atributelor în gramatica dată (presupusă necirculară).

Abordarea acestor probleme conduce la identificarea unor clase de gramatici cu atribute pe care le vom trece în revistă în acest capitol. Înainte de a defini riguros noțiunea de gramatică cu atribute și a studia clasele importante, să prezentăm niște exemple.

6.1.1 Numere raționale în baza 2 (Knuth)

Se consideră gramatica independentă de context cu următoarele reguli:

$$\begin{array}{lll} N \rightarrow L & L \rightarrow LB & B \rightarrow 0 \\ N \rightarrow LL & L \rightarrow B & B \rightarrow 1 \end{array}$$

Se constată ușor că limbajul generat de această gramatică este:

$$L(G) = \{b_1b_2\dots b_n \mid b_i \in \{0, 1\}, 1 \leq i \leq n\} \cup \\ \cup \{b_1b_2\dots b_n \cdot b_{n+1}b_{n+2}\dots b_m \mid b_i \in \{0, 1\}, 1 \leq i \leq m, m > n \geq 1\}$$

Dorim să înzestram această gramatică cu atribute și reguli semantice astfel ca pentru fiecare cuvânt w generat de aceasta să putem determina numărul rațional (în baza 10) a cărui reprezentare în baza 2 este w . Pentru aceasta considerăm atributele:

- Pentru simbolul B (bit):
 - Un atribut sintetizat, numit **val**, care are ca domeniu al valorilor mulțimea numerelor raționale și care trebuie interpretat astfel: dacă B generează 0, atunci **B.val** are valoarea 0, iar dacă B generează 1 atunci **B.val** are valoarea 2^s unde s este locul lui 1 în reprezentarea binară în cauză (s se numește scală);
 - Un atribut moștenit numit **scala** care va defini exponentul lui 2 pentru calculul valorii lui B. Domeniul valorilor acestui atribut este mulțimea numerelor întregi;
- Pentru variabila L (listă de biți):
 - Două atribute sintetizate: **val** și **lung** care au ca domenii de valori mulțimea numerelor raționale, respectiv naturale iar ca interpretare, valoarea numerică zecimală, respectiv lungimea lui w dacă $L \Rightarrow w^+$;
 - Un atribut moștenit **scala** care are ca domeniu de valori mulțimea numerelor întregi și care va fi folosit pentru calculul scalei lui B.
- Pentru variabila N (număr):
 - Un atribut sintetizat **val** care reprezintă semantica cuvintelor generate în gramatică. Domeniul său de valori este mulțimea numerelor raționale.

Regulile semantice sunt atașate regulilor sintactice (producțiilor gramaticii) și sunt date în tabela de mai jos. Facem precizarea că, în producțiile în care L are apariții multiple ($L \rightarrow LB, N \rightarrow L.L$), aceste apariții au fost indexate pentru a evita confuzia.

Reguli sintactice	Reguli semantice
$B \rightarrow 0$	$B.val = 0$
$B \rightarrow 1$	$B.val = 2^{B.scala}$
$L \rightarrow B$	$L.val = B.val; L.lung = 1; B.scala = L.scala$
$L_1 \rightarrow L_2B$	$L_1.val = L_2.val + B.val; L_1.lung = L_2.lung + 1$ $L_2.scala = L_1.scala + 1; B.scala = L_1.scala$
$N \rightarrow L$	$N.val = L.val; L.scala = 0$
$N \rightarrow L_1.L_2$	$N.val = L_1.val + L_2.val;$ $L_1.scala = 0; L_2.scala = -L_2.lung$

Să determinăm, folosind regulile semantice date mai sus, semantica cuvântului $w = 101.11$. Arborele de derivare pentru acest cuvânt în gramatica suport este dat în figura 7.1.1

este de forma : type v_1, v_2, \dots, v_n iar o instrucțiune are forma use v_1, v_2, \dots, v_m . Iată exemple de programe corecte în această descriere sintactică:

- 1) begin type a, b; type c; use a, c; use b, c end
- 2) begin type a, b; type a, c; use a, c; use c end
- 3) begin type x, a; use x; use a, c end
- 4) begin type x, c; type x, z; use x; use z, x end

În descrierea (semantică) a limbajelor de programare există în general regulile:

- Orice variabilă din program trebuie declarată.
- Declararea multiplă a unei variabile nu este permisă.

După aceste reguli, programele 2, 3, 4 de mai sus nu sunt corecte. Ne propunem să descriem o gramatică cu atribute care să descrie prin regulile semantice cele precizate mai sus.

Variabilelor gramaticii de mai sus le atașăm atributele:

Variabila Z:

- atributul sintetizat **Răspuns** cu valori în mulțimea {true, false} care dă răspuns la întrebarea: este programul corect?

Variabila D:

- Atributul sintetizat **Dec**, care are ca valori submulțimi de variabile și eventual constanta literală *eroare*. **D.Dec** conține toate variabilele care sunt declarate în subarboarele lui D, și numai pe acelea, iar dacă o variabilă are declarații multiple, **D.Dec** va conține și *eroare*.

Variabila I:

- Atributul sintetizat **Răspuns** cu valori în {true, false}; **I.Răspuns** este true dacă și numai dacă variabilele folosite în instrucțiunile generate de I au fost declarate în același program.
- Atributul moștenit **Tabel** cu valori în 2^{Var} . **I.Tabel** conține toate variabilele declarate în programul curent.

Variabila V:

- Atributul sintetizat **Lista** cu valori în $2^{Var \cup \{eroare\}}$. **V.Lista** reprezintă lista variabilelor generate de V și eventual constanta *eroare*.
- Atributul moștenit **Loc** cu valori în mulțimea {dec, inst}. **V.Loc** este *dec* dacă V face parte din subarboarele cu rădăcina D și este *inst* dacă V face parte din subarboarele cu rădăcina I.

Regulile semantice pentru calculul valorilor atributelor sunt date în tabelul de mai jos.

SINTAXA	SEMANTICA
$Z \rightarrow \text{begin } D ; I \text{ end}$	$Z.\text{Răspuns} = \text{if } (eroare \in D.\text{Dec})$ then false else I.Răspuns $I.\text{Tabel} = D.\text{Dec} - \{eroare\}$
$D_1 \rightarrow D_2 ; \text{ type } V$	$D_1.\text{Dec} = D_2.\text{Dec} \cup V.\text{Lista} \cup$ $\cup \text{if } (V.\text{Lista} \cap D_2.\text{Dec})$ then {eroare} else Φ $V.\text{Loc} = dec$
$D \rightarrow \text{type } V$	$D.\text{Dec} = V.\text{Lista}; V.\text{Loc} = dec.$
$I_1 \rightarrow I_2 ; \text{ use } V$	$I_1.\text{Răspuns} = \text{if } (V.\text{Lista} \not\subseteq I_1.\text{Tabel})$ then false else I₂.Răspuns $I_2.\text{Tabel} = I_1.\text{Tabel}; V.\text{Loc} = inst.$
$I \rightarrow \text{use } V$	$I.\text{Răspuns} = \text{if } (V.\text{Lista} \subseteq I.\text{Tabel})$ then true else false $V.\text{Loc} = inst$
$V_1 \rightarrow V_2, v$	$V_1.\text{Lista} = V_2.\text{Lista} \cup \{v\} \cup$ $\cup \text{if } (V_1.\text{Loc} = dec \text{ and } v \in V_2.\text{Lista})$ then eroare else $\Phi;$

	$V_2.\mathbf{Loc} = V_1.\mathbf{Loc}$
$V \rightarrow v$	$V.\mathbf{Lista} = \{v\}$

Fie acum programul:

```
begin type x, a; use x ; use a, c end
```

Arborele de derivare pentru acesta este dat în figura 7.1.3. Din motive lesne de înțeles, etichetele nodurilor corespunzătoare terminalilor begin, end, type și use sunt reprezentate în arbore respectiv prin b, e, t și u.

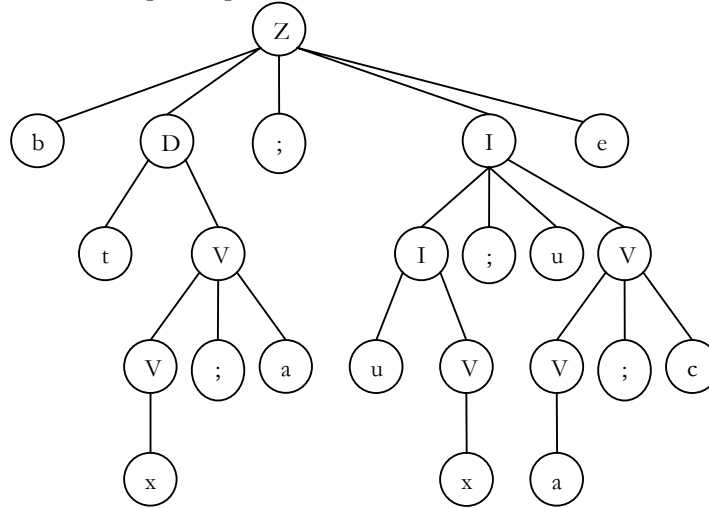


Figura 7.1.3

Dacă pentru I reprezentăm atributele în ordinea **I.Tabel**, **I.Răspuns** iar pentru V în ordinea **V.Lista**, **V.Loc**, atunci valorile calculate ale atributelor pentru acest arbore sunt date în graful din figura 7.1.4. Nodurile acestui graf corespund aparițiilor atributelor neterminabililor Z, D, I și V din arborele de mai sus, iar arcele reprezintă dependențele între valorile acestor atribute prin regulile semantice.

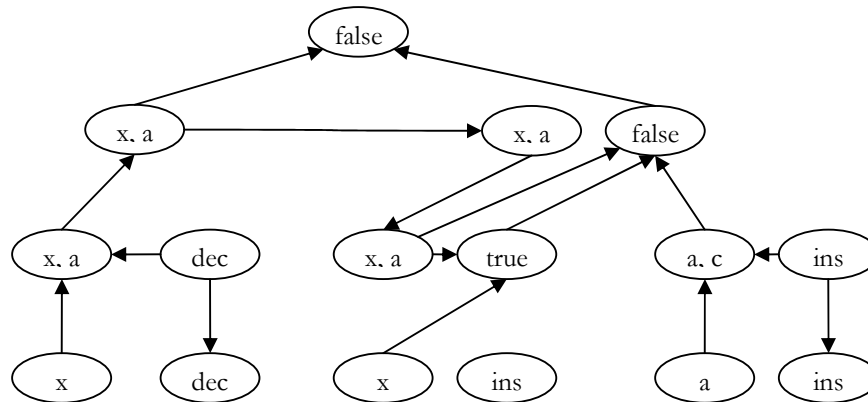


Figura 7.1.4

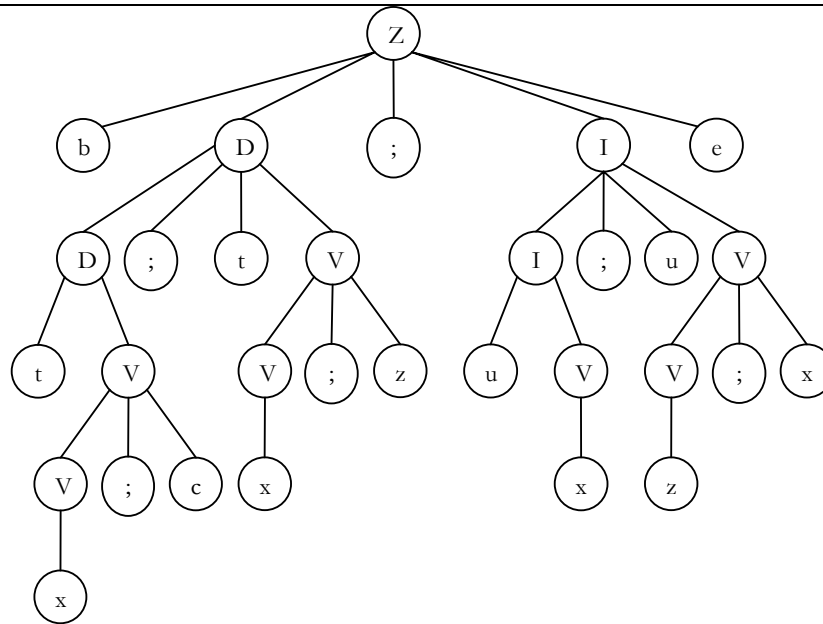


Figura 7.1.5

Pentru programul:

begin type x, c; type x, z; use x; use z, x end

arborele de derivare este dat în figura 7.1.5 iar valorile atributelor în figura 7.1.6.

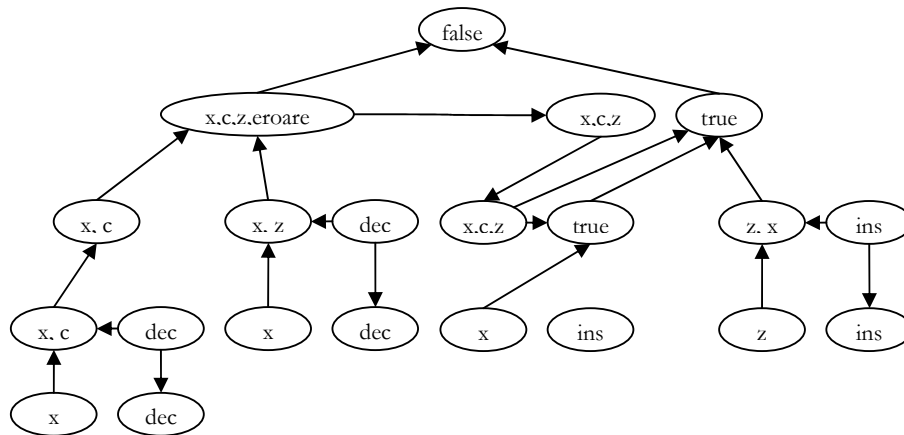


Figura 7.1.6

6.1.3 Definiția gramaticilor cu atribute

Definiția 7.1.1 O gramatică cu atribute este un 4-uplu $\mathcal{GA} = (G, \mathcal{D}, \mathcal{A}, \mathcal{R})$ unde:

1. $G = (V, T, Z, P)$ este o gramatică independentă de context redusă, numită *gramatica suport* pentru \mathcal{GA} . Notăm prin $\Sigma = V \cup T$ vocabularul gramaticii. Aici Z reprezintă simbolul de start al gramaticii.
2. $\mathcal{D} = (\mathcal{M}, \mathcal{F})$ este *domeniul semantic* format din:
 - 2.1. \mathcal{M} – o familie finită de mulțimi;
 - 2.2. \mathcal{F} – o familie finită de funcții total definite și calculabile de forma

$$f : M_1 \times M_2 \times \dots \times M_n \rightarrow M_0, \text{ unde } M_i \in \mathcal{M}, 0 \leq i \leq n.$$
3. $\mathcal{A} = (I, S, \text{type})$ este o *descriere a atributelor gramaticii*, unde;

- 3.1. $I = \bigcup_{x \in \Sigma} I(X)$ este *mulțimea atributelor moștenite*;
- 3.2. $S = \bigcup_{x \in \Sigma} I(X)$ este *mulțimea atributelor sintetizate*;
- 3.3. Pentru orice atribut $\alpha \in I \cup S$, $\text{type}(\alpha) \in \mathcal{M}$ este mulțimea din domeniul semantic în care α poate lua valori;
- 3.4. Atributele *moștenite* $I(X)$ și cele *sintetizate* $S(X)$ pentru simbolul X se consideră disjuncte, oricare ar fi $X \in \Sigma$. De asemenea, notând $\mathcal{A}(X) = I(X) \cup S(X)$, vom considera $\mathcal{A}(X) \cap \mathcal{A}(Y) = \Phi$, $\forall X, Y \in \Sigma$, $X \neq Y$, chiar dacă, uneori, atribute ale unor simboluri diferite au aceleași nume (ca în exemplele prezentate). Referirea la atributul α al simbolului X se face prin $X.\alpha$.
- 3.5. Simbolul de start Z al gramaticii are un atribut sintetizat desemnat a reprezenta *semantica* unei construcții $w \in L(G)$.
4. $\mathcal{R} = \bigcup_{p \in P} \mathcal{R}(P)$ este mulțimea (finită) a *regulilor semantice* (reguli de evaluare a atributelor). Dacă producția p este de forma: $X_{p0} \rightarrow X_{p1}X_{p2}\dots X_{pnp}$ atunci o regulă semantică este de forma $X_{pk}.a = f(X_{pk1}.a_1, X_{pk2}.a_2, \dots, X_{pkm}.a_m)$ unde :
 - 4.1. $f \in \mathcal{F}$, $f : \text{type}(a_1) \times \text{type}(a_2) \times \dots \times \text{type}(a_m) \rightarrow \text{type}(a)$;
 - 4.2. Pentru $k = 0$ atunci $a \in S(X_{p0})$ cu precizarea că $\mathcal{R}(P)$ conține câte o regulă și numai una pentru fiecare atribut din $S(X_{p0})$.
 - 4.3. Pentru $k > 0$ atunci $a \in I(X_{pk})$ iar $\mathcal{R}(p)$ conține câte o regulă și numai una pentru fiecare atribut din $\bigcup_{1 \leq k \leq np} I(X_{pk})$.
 - 4.4. $a_i \in I(X_{p0}) \cup (\bigcup_{1 \leq k \leq np} \mathcal{A}(X_{pk}))$, $1 \leq k \leq np$.

□

Exemplul 7.1.1 Să punem în evidență elementele din definiția gramaticii cu atribute pentru gramatica prezentată în paragraful 7.1.1.

1. $G = (\{N, L, B\}, \{0, 1, .\}, N, P)$ unde
 $P = \{N \rightarrow L, N \rightarrow L.L, L \rightarrow LB, L \rightarrow B, B \rightarrow 0, B \rightarrow 1\}$;
2. $\mathcal{M} = \{\mathbf{N}, \mathbf{Z}, \mathbf{Q}\}$, $\mathcal{F} = \{c_0, c_1, id, suc, minus, suma, exp\}$, unde c_0 și c_1 sunt funcțiile constante 0 respectiv 1, suc este funcția succesor, $minus(x) = -x$, $suma(x, y) = x + y$, $exp(x) = 2^x$. Domeniile acestor funcții sunt adecvate regulilor semantice descrise.
3. $\mathcal{A}(N) = S(N) = \{\mathbf{val}\}$, $\text{type}(\mathbf{val}) = \mathbf{Q}$
 $S(L) = \{\mathbf{val}, \mathbf{lung}\}$, $\text{type}(\mathbf{val}) = \mathbf{Q}$, $\text{type}(\mathbf{lung}) = \mathbf{N}$
 $I(L) = \{\mathbf{scala}\}$, $\text{type}(\mathbf{scala}) = \mathbf{Z}$
 $S(B) = \{\mathbf{val}\}$, $\text{type}(\mathbf{val}) = \mathbf{Q}$
 $I(B) = \{\mathbf{scala}\}$, $\text{type}(\mathbf{scala}) = \mathbf{Z}$

Observație. Aparent $I(L) = I(B) = \{\mathbf{scala}\}$. Trebuie precizat însă că cele două atribute cu același nume sunt diferite: $L.\mathbf{scala} \neq B.\mathbf{scala}$ și acest lucru se vede din definirea acestora prin regulile semantice. Același lucru este valabil pentru atributul \mathbf{val} care apare la S, L și B .

4. Regulile semantice sunt cele din tabela de la paragraful 7.1.1. Este ușor de verificat că au loc 6.1 – 6.4.

Exemplul 7.1.2 Să punem în evidență elementele din definiția gramaticii cu atribute pentru gramatica prezentată în 7.1.2.

1. $G = (\{Z, D, I, V\}, \{\text{begin, end, tip, use, }, a, b, \dots\}, S, P)$ unde
 $P = \{Z \rightarrow \text{begin } D; \text{end}, D \rightarrow D; \text{type } V, D \rightarrow \text{type } V, I \rightarrow I; \text{use } V,$
 $I \rightarrow \text{use } V, V \rightarrow v, V \rightarrow V, v\}$.
2. $\mathcal{M} = \{\text{Bool}, 2^{\text{Var}}, 2^{\text{Var} \cup \{\text{eroare}\}}, \{\text{dec, inst}\}\}$ unde $\text{Bool} = \{\text{true, false}\}$. \mathcal{F} conține toate funcțiile folosite în tabelul ce definește regulile semantice.

3. $\mathcal{A}(Z) = \mathcal{S}(Z) = \{\mathbf{Răspuns}\}$, $\text{type}(\mathbf{Răspuns}) = \text{Bool}$
 $\mathcal{A}(D) = \mathcal{S}(D) = \{\mathbf{Dec}\}$, $\text{type}(\mathbf{Dec}) = 2^{\text{Var} \cup \{\text{emare}\}}$
 $\mathcal{S}(I) = \{\mathbf{Răspuns}\}$, $\text{type}(\mathbf{Răspuns}) = \text{Bool}$
 $I(I) = \{\mathbf{Tabel}\}$, $\text{type}(\mathbf{Tabel}) = 2^{\text{Var}}$
 $\mathcal{S}(V) = \{\mathbf{Lista}\}$, $\text{type}(\mathbf{Lista}) = 2^{\text{Var} \cup \{\text{emare}\}}$
 $I(V) = \{\mathbf{Loc}\}$, $\text{type}(\mathbf{Loc}) = \{\text{dec}, \text{inst}\}$
4. Regulele semantice sunt cele din tabela de la 7.1.2.

6.2 Grafuri atașate unei gramatici cu atribute

În figurile 7.1.2, 7.1.4 și 7.1.6 sunt puse în evidență niște grafuri atașate arborilor de derivare pentru anumite cuvinte generate de gramatica suport a gramaticii cu atribute. Aceste grafuri (și altele pe care le vom defini în continuare) joacă un rol important în verificarea corectitudinii definirii regulilor semantice. În esență o regulă semantică definește o relație de dependență între atributele din gramatica respectivă, iar această relație definește un graf.

6.2.1 Graful \mathcal{DG}_p

Definiția 7.2.1 Fie $\mathcal{GA} = (G, \mathcal{D}, \mathcal{A}, \mathcal{R})$ o gramatică cu atribute și producția p dată de $X_0 \rightarrow X_1 X_2 \dots X_{np}$. Graful $\mathcal{DG}_p = (\mathcal{X}, \mathcal{E})$ asociat producției p , numit *graful de dependență* al producției p , este graful în care:

- nodurile sunt toate aparițiile atributelor simbolurilor din producția p , adică $\mathcal{X} = \bigcup_{1 \leq i \leq np} \mathcal{A}(X_i)$;
- arcele sunt de forma $(X_i.\alpha, X_j.\beta)$ dacă $X_j.\beta$ depinde printr-o regulă semantică de $X_i.\alpha$ (existentă în $\mathcal{R}(p)$), adică:
 $\mathcal{E} = \{(X_i.\alpha, X_j.\beta) \mid X_j.\beta = f(\dots, X_i.\alpha, \dots) \in \mathcal{R}(p)\}$

Exemplul 7.2.1 Fie gramatica din exemplul 7.1.1. Atunci, grafurile atașate producțiilor gramaticii sunt prezentate în figurile 7.2.1 – 7.2.6



Figura 7.2.1

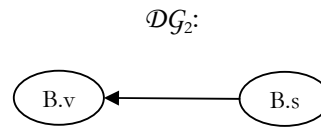


Figura 7.2.2

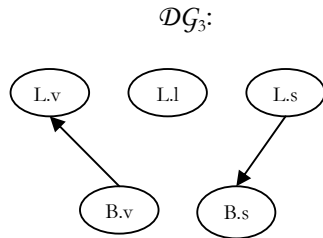


Figura 7.2.3

\mathcal{DG}_5 :

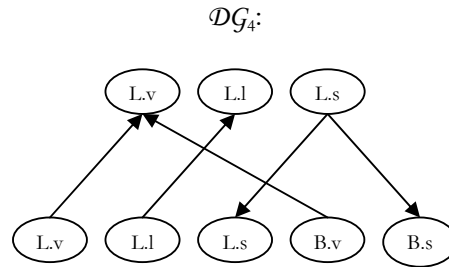


Figura 7.2.4

\mathcal{DG}_6 :

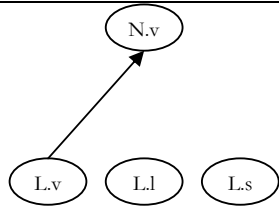


Figura 7.2.5

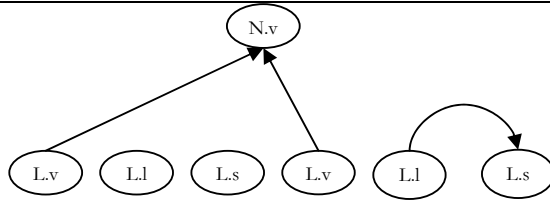


Figura 7.2.6

Cititorul este invitat să construiască grafurile \mathcal{DG}_p pentru gramatica din exemplul 7.1.2.

6.2.2 Graful \mathcal{BG}_p

Definiția 7.2.2 Fie $\mathcal{GA} = (G, \mathcal{D}, \mathcal{A}, \mathcal{R})$ o gramatică cu atribute și producția p dată de $X_0 \rightarrow X_1 X_2 \dots X_{n_p}$. Graful \mathcal{BG}_p asociat producției p este graful în care:

- nodurile sunt X_1, X_2, \dots, X_{n_p} , privite ca elemente distincte (dacă un simbol se repetă în p , atunci aparițiile lui se numerează);
- arcele sunt toate perechile (X_i, X_j) și numai acelea, pentru care există un atribut moștenit al lui X_j care depinde (printr-o regulă semantică) de un atribut sintetizat al lui X_i , $1 \leq i, j \leq n_p$.

Exemplul 7.2.2 Să considerăm gramatica din exemplul 7.1.2. Atunci graful \mathcal{BG}_1 , corespunzător producției $z \rightarrow \text{begin } D; I \text{ end}$, este dat în figura 7.2.7. Pentru gramatica din exemplul 7.1.1 graful \mathcal{BG}_6 este dat în figura 7.2.8.

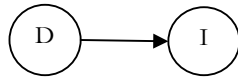


Figura 7.2.7

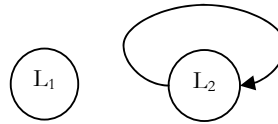


Figura 7.2.8

6.2.3 Graful \mathcal{DG}_t

Definiția 7.2.3 Fie $\mathcal{GA} = (G, \mathcal{D}, \mathcal{A}, \mathcal{R})$ o gramatică cu atribute și t un arbore de derivare în gramatica suport G . Graful \mathcal{DG}_t atașat arborelui t se definește recursiv astfel:

- Dacă $t = p \in P$ atunci $\mathcal{DG}_t = \mathcal{DG}_p$;
- Dacă t este un arbore de adâncime mai mare ca 1, cu rădăcina etichetată X_0 iar producția ce se aplică în rădăcină este $p = X_0 \rightarrow X_1 X_2 \dots X_{n_p}$, atunci \mathcal{DG}_t se obține din \mathcal{DG}_p și \mathcal{DG}_{t_i} $1 \leq i \leq n_p$, prin alipire: nodurile din \mathcal{DG}_p corespunzătoare atributelor lui X_i se identifică respectiv cu nodurile din \mathcal{DG}_{t_i} corespunzătoare acelorași atribute ale lui X_i . În cazul în care un anume X_i este terminal, graful \mathcal{DG}_{t_i} corespunzător are ca noduri atributele lui X_i și nu are nici un arc.

Exemplul 7.2.3 Pentru gramatica din exemplul 7.1.1 graful \mathcal{DG}_t pentru arborele din figura 7.2.9 este dat în figura 7.2.10. El se obține din grafurile \mathcal{DG}_4 (corespunzător producției din rădăcină) la care se alipesc grafurile \mathcal{DG}_3 și \mathcal{DG}_2 ca în figura 7.2.11.

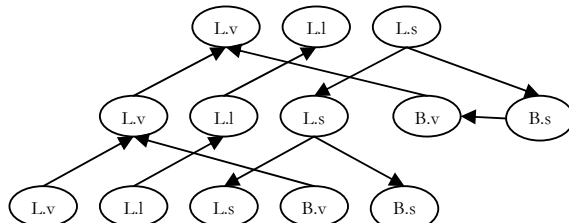
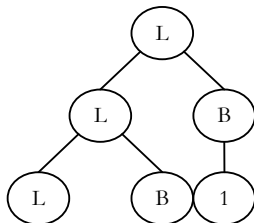


Figura 7.2.9

Figura 7.2.10

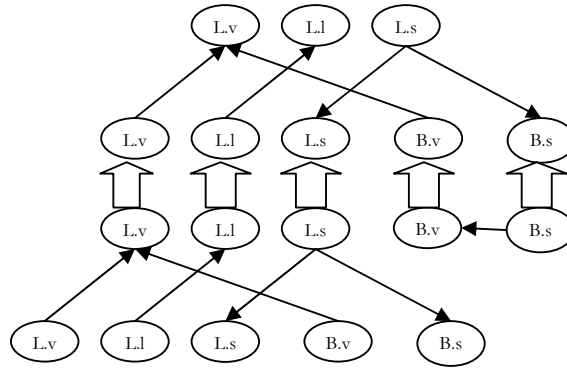


Figura 7.2.11

Pentru arborele din figura 7.1.1 graful \mathcal{DG}_t este obținut în același mod și este dat în figura 7.2.12.

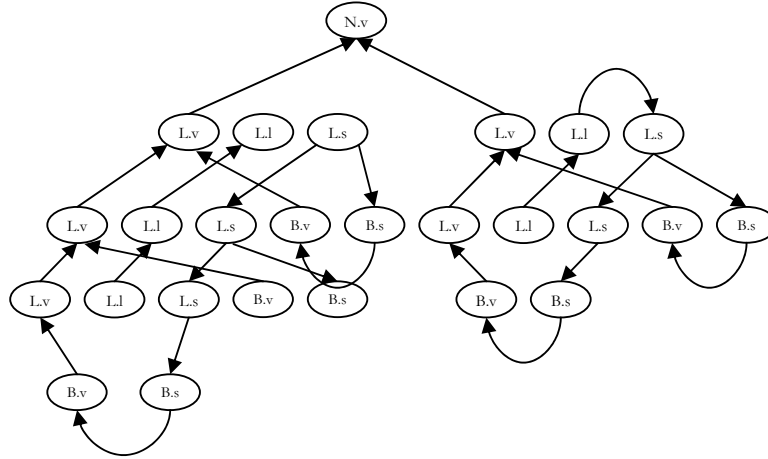


Figura 7.2.12

Graful \mathcal{DG}_t pentru un arbore de derivare t complet (în care rădăcina este etichetată cu simbolul de start Z iar frunzele sunt etichetate cu terminali) se obține, se înțelege, în același mod. Dacă acest graf nu are circuite, atunci el definește o ordine parțială în mulțimea aparițiilor atributelor în arborele t : dacă $(X.\alpha, Y.\beta)$ este un arc în \mathcal{DG}_t , atunci, pentru a putea calcula valoarea lui $Y.\beta$ avem nevoie de valoarea lui $X.\alpha$.

Definiția 7.2.4 O gramatică cu atribute \mathcal{GA} este *circulară* (sau nu este bine definită) dacă există măcar un arbore de derivare t pentru care graful \mathcal{DG}_t are circuite. O gramatică cu atribute \mathcal{GA} este *necirculară* (este *bine definită*) dacă pentru orice arbore de derivare t graful \mathcal{DG}_t nu are circuite.

Pentru o gramatică cu atribute necirculară, dat un arbore de derivare t , deci un cuvânt $w \in L(G)$, se pot calcula valorile tuturor aparițiilor (instanțelor) atributelor din arborele t , deci se poate determina semantica lui w . Așadar, două din problemele importante ce se pun referitor la o gramatică cu atribute sunt:

- *problema circularității*: dată \mathcal{GA} , este ea necirculară?
- *problema evaluării atributelor*: dată \mathcal{GA} necirculară, să se construiască o procedură care, atunci când este dat un arbore t , să determine valorile tuturor aparițiilor atributelor etichetelor arborelui.

Înainte de a aborda problema circularității, să mai definim niște grafuri ce vor fi utile rezolvării acesteia.

6.3 Metode de evaluare a atributelor

Un evaluator pentru o gramatică cu atribute este un program care, pentru un arbore de derivare t al gramaticii suport, determină valorile tuturor aparițiilor atributelor din arborele t . Un generator de evaluator este un program care, pentru o gramatică cu atribute ca intrare, o respinge (dacă ea este circulară) sau o acceptă și construiește pentru ea un evaluator.

În funcție de modul cum sunt evaluate atributele, un evaluator poate fi:

- recursiv sau iterativ;
- aplicativ sau imperativ;
- tree-walking sau tree-free;
- pass-oriented sau visit-oriented;

Un tip de evaluator corespunde unei clase de gramatici cu atribute pentru care se poate construi acel tip de evaluator. Vom trece în revistă câteva metode de evaluare împreună cu clasele de gramatici pentru care se aplică aceste metode.

6.3.1 Evaluatorul rezultat din definiția gramaticii

Pentru a putea calcula valorile tuturor aparițiilor atributelor într-un arbore de derivare t , graful \mathcal{DG}_t trebuie să fie liber de circuite, așa cum a fost arătat în 7.2.3. Să considerăm cazul general al relațiilor pentru a defini o ordine de evaluare.

Dacă V este o mulțime finită iar $R \subseteq V \times V$ este o relație peste V , atunci (V, R) este un graf orientat. Să notăm prin R^+ închiderea tranzitivă a relației R . Atunci R este *necirculară* (sau *aciclică*, sau *ireflexivă*) dacă $\forall v \in V, (v, v) \notin R^+$. Relația $T \subseteq V \times V$ este o *ordine totală* pe V dacă :

- $T^+ = T$ (T este tranzitivă);
- T este necirculară;
- T este totală: $\forall x \in V, \forall y \in V, x = y$ sau $(x, y) \in T$ sau $(y, x) \in T$.

Pentru orice ordine totală T pe mulțimea V , există o enumerare a elementelor lui V , să zicem v_1, v_2, \dots, v_n , astfel ca $(v_i, v_j) \in T$ dacă și numai dacă $i < j$. Această enumerare se numește *secvență de evaluare*.

Spunem că o ordine totală T este o *ordine de evaluare* a relației R pe V , dacă $R \subseteq T$. Dacă v_1, v_2, \dots, v_n este secvența de evaluare pentru T , atunci, $(v_i, v_j) \in R$ implică $i < j$. Un element $v \in V$ se numește *element minimal* pentru R dacă nici un element din V nu este în relație cu v ($\forall w \in V, (w, v) \notin R$). Analog, $v \in V$ este *maximal* pentru R dacă el nu este în relație cu nici un element ($\forall w \in V, (v, w) \notin R$).

Teorema 7.4.1 (de evaluare) Fie R o relație pe V . R este necirculară dacă și numai dacă R admite o ordine de evaluare.

Demonstrație. Dacă R admite o ordine de evaluare T , aceasta este necirculară (din definiția ordinii totale). Cum $R \subseteq T$ rezultă că și R este necirculară. Invers, să presupunem că R este o relație necirculară pe V . Următorul algoritm, numit de *sortare topologică*, determină o secvență de evaluare pe R , ceea ce este suficient pentru o ordine de evaluare pe R :

Algoritmul 7.4.1

Intrare: Mulțimea V și relația $R \subseteq V \times V$ necirculară;

Ieșire: O secvență de evaluare $eval = \{v_1, v_2, \dots, v_n\}$, unde $|V| = n$;

Metoda:

1. $eval = \Phi$; $k = 0$;

```

2. while (eval  $\neq$  V) {
3.   Se alege v un element minimal din V - eval;
4.   k = k + 1;
5.   vk = v;
6.   eval = eval  $\cup$  {vk};
}
```

Să observăm că în $V - \text{eval}$ există întotdeauna un element minimal pentru că R este necirculară pe V , deci și restricția sa la $V - \text{eval}$ este necirculară. Algoritmul este finit iar din modul cum s-au adăugat elementele în eval rezultă că aceasta reprezintă o secvență de evaluare pe R . □

Algoritmul 7.4.1 poate fi folosit pentru proiectarea unui evaluator pentru o gramatică cu atribute care este necirculară. Considerând un arbore de derivare t , aplicăm algoritmul precedent pentru relația ce definește graful \mathcal{DG}_t pentru a afla ordinea de evaluare a atributelor în t .

Evaluatorul P1

Intrare: Gramatica cu atribute $\mathcal{GA} = (G, \mathcal{D}, \mathcal{A}, \mathcal{R})$ necirculară;
 Graful $\mathcal{DG}_t = (\mathcal{N}, \mathcal{E})$ pentru un arbore t (\mathcal{N} este mulțimea tuturor aparițiilor atributelor din t);

Ieșire: Valorile tuturor aparițiilor atributelor în t ;

Metoda:

```

1. eval =  $\Phi$ ; k = 0;
2. while (eval  $\neq$   $\mathcal{N}$ ) {
3.   Se alege v un element minimal din  $\mathcal{N} - \text{eval}$ ;
4.   evaluare(v); // se aplica o regulă semantica
5.   k = k + 1; vk = v; eval = eval  $\cup$  {vk};
}
```

În linia 4 din algoritm se apelează o procedură care evaluează, conform cu regulile semantice, valorile atributului v . Acest lucru este posibil pentru că valorile atributului v depind de valorile unor atribute ce au fost calculate înainte, conform secvenței de evaluare.

Evaluatorul tocmai prezentat este important pentru faptul că reprezintă o caracterizare a gramaticilor cu atribute necirculare: o gramatică cu atribute este necirculară dacă și numai dacă P1 este un evaluator pentru aceasta. P1 este cunoscut în literatură ca și “*defining evaluator*” – el exprimă faptul că gramatica respectivă este bine definită.

6.4 Traducere bazată pe sintaxă

Una din cele mai frecvente tehnici folosite în analiza semantică este cea a evaluării atributelor în timpul analizei sintactice ascendente (LR). Am văzut în paragraful precedent cum pot fi evaluate atributele în cadrul unui parser SLR(1). Dacă toate atributele gramaticii sunt sintetizate atunci lucrurile se simplifică: la fiecare reducere cu o producție de forma $A \rightarrow \beta$ se aplică regulile semantice pentru evaluarea atributelor (sintetizate) ale lui A . Pentru analiza semantică a limbajelor de programare în acest mod, precum și pentru obținerea codului (intermediar) echivalent cu codul sursă, se folosesc *acțiunile semantice* atașate regulilor sintactice, care sunt aplicate atunci când analizorul sintactic realizează o reducere. Se obține astfel o structură care nu este altceva decât o gramatică cu atribute în care există un atribut care are ca valori secvențe de cod intermediar. Așa cum o să vedem în continuare, determinarea valorilor atributelor se poate face prin aplicarea unor proceduri specifice. Regulile semantice nu mai sunt chiar de forma pe care am specificat-o în definiția gramaticii cu atribute ci poate apărea, pe lângă definiții ale atributelor cu ajutorul unor funcții calculabile, și apelul unor proceduri

care generează cod intermediar. De aceea numim acțiuni semantice aceste proceduri atașate regulilor sintactice. O astfel de gramatică cu atribute este cunoscută în literatură sub numele de *schemă de traducere bazată (orientată) pe sintaxă*. Schemele de traducere bazate pe sintaxă sunt deosebit de utile celor ce implementează limbaje de programare; într-o astfel de schemă poate fi implementată generarea de cod intermediar pornind de la structura sintactică a limbajului sursă. Odată obținut codul intermediar, acesta urmează a fi transformat în cod obiect pentru o mașină anume. În acest paragraf vom prezenta mai întâi câteva tipuri de cod intermediar pentru ca apoi să vedem cum se abordează traducerea structurilor sintactice importante din limbajele de programare.

6.4.1 Reprezentări intermediare

Analiza sintactică oferă la ieșire reprezentarea programului ca arbore de derivare (parsare) în gramatica ce descrie sintaxa limbajului. Pentru a crea codul obiect echivalent cu cel sursă, în general se folosește o reprezentare intermediară a acestuia, un cod intermediar. Tipurile de cod intermediar utilizate în procesul de compilare sunt:

Notăția postfix (Notăția poloneză inversă). Această reprezentare se generează ușor de la arborele abstract printr-o traversare postordine, sau cu un parser bottom-up. Avantajul acestei reprezentări este acela că evaluarea expresiilor se face ușor și eficient folosind o stivă. Instrucțiunile pot fi privite ca și operatori (de o anumită aritate), încât reprezentarea unui program se face uniform. De pildă, dacă notăm prin ? operatorul ternar `if_then_else`, atunci expresia

`if a + b > c then x = a - c else x = a - b`

are reprezentarea în notația postfix astfel:

`a b + c > x a c - = x b c - = ?`

Forma postfix a unui program se poate obține folosind Yacc (vezi exemplul 6.1.3). Există și un dezavantaj major în această abordare: dacă ne referim la reprezentarea operatorului `if_then_else`, atunci evaluarea expresiei presupune evaluarea tuturor subexpresiilor, chiar dacă acest lucru nu este necesar. Pentru optimizare se introduc niște operatori de salt condiționat în cadrul reprezentării.

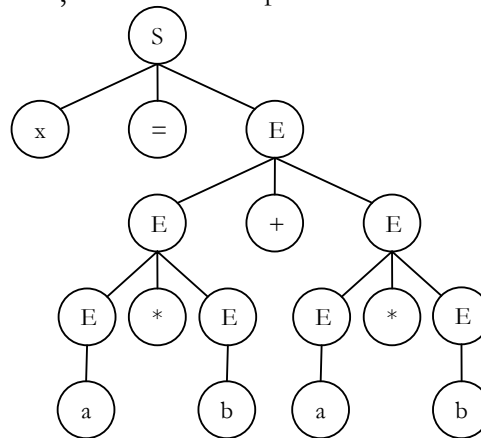


Figura 7.5.1

Arborele sintactic (abstract). Acesta se obține din arborele (concret) de parsare în care neterminalii sunt înlocuiți cu operatorii corespunzători (dacă ne referim la expresii). Să analizăm, spre exemplu, arborele de parsare (figura 7.5.1) pentru atribuirea `x = a * b + a * b` ce se obține în gramatica:

$S \rightarrow id = E \quad E \rightarrow E + E \mid E * E \mid (E) \mid id$

Arborele sintactic corespunzător este dat în figura 7.5.2.

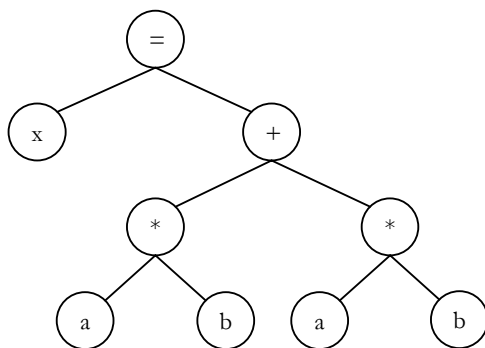


Figura 7.5.2

În general, arborele sintactic se poate defini astfel:

- Pentru o regulă sintactică de forma $A \rightarrow a$ (regulă ce produce terminali) arborele are un singur nod etichetat cu a .
- Pentru o regulă de forma $A \rightarrow B \text{ op } C$ (regulă corespunzătoare operatorilor binari), arborele are o rădăcină etichetată cu op , un subarbor stâng corespunzător regulii ce rescrie pe B și un subarbor drept corespunzător regulii ce rescrie pe C .
- Pentru o regulă de forma $A \rightarrow op \ B$ (regulă corespunzătoare operatorilor unari) arborele are o rădăcină etichetată cu op și un subarbor (stâng) corespunzător regulii ce rescrie pe B .
- Pentru o regulă corespunzătoare unui operator ternar, de exemplu cea care descrie instrucțiunea `if`, $S \rightarrow \text{if } B \text{ then } S \text{ else } S$, arborele sintactic are o rădăcină (etichetată cu numele operatorului, de exemplu `if`) și trei subarbori corespunzători, în ordine, celor trei operanzi.

În acest mod, oricărui program scris într-un limbaj de programare îi putem atașa un arbore sintactic. Acest arbore se poate obține în procesul analizei sintactice ascendente aplicând acțiuni semantice corespunzătoare fiecărei reduceri efectuate de analizorul sintactic.

Codul cu trei adrese Această reprezentare este avantajoasă pentru cazul în care instrucțiunile limbajului intermediar au cel mult 3 variabile și cel mult un operator. Acestea se traduc ușor în limbaj de asamblare. Reprezentarea se face prin quadruple în care, pe lângă adresele celor doi operanzi și a rezultatului, este păstrat și codul operației. De pildă, instrucțiunea

$\text{delta} := b * b - 4 * a * c$

se traduce prin codul:

$T_1 = b * b$
 $T_2 = 4 * a$
 $T_3 = T_2 * c$
 $T_4 = T_2 - T_3$
 $\text{delta} = T_4$

care este reprezentat de quadruplele:

	Operator	Argument 1	Argument 2	Rezultat
0	*	b	b	T_1
1	*	4	a	T_2
2	*	T_2	c	T_3
3	-	T_2	T_3	T_4

4	=	T ₄	-	delta
---	---	----------------	---	-------

Aici T_1, T_2, \dots sunt nume temporare care sunt adăugate la tabela simbolurilor pe măsură ce sunt create încât, în câmpurile relative la argumente și rezultat, apar pointeri la tabela de simboluri corespunzătoare numelor ce apar acolo.

6.4.2 Traducerea codului în notație postfix

O schemă de traducere pentru obținerea unui cod postfix este relativ ușor de obținut. O să exemplificăm acest lucru pentru atribuiri cu expresii aritmetice generate de gramatica:

$$S \rightarrow A := E \quad E \rightarrow E \text{ op } E \mid (E) \mid \text{id} \quad A \rightarrow \text{id}$$

Dacă vom nota prin id.Loc pointerul către id în tabela de parsare (numele identificatorului) și prin $E.\text{Cod}$ traducerea expresiei E în cod intermediar în notația postfix, atunci schema de traducere este dată în tabela de mai jos, unde prin $||$ am notat operația de concatenare a codurilor:

Regula sintactică	Acțiunea semantică
$S \rightarrow A := E$	$S.\text{Cod} = A.\text{Cod} E.\text{Cod} :=$
$E \rightarrow E_1 \text{ op } E_2$	$E.\text{Cod} = E_1.\text{Cod} E_2.\text{Cod} \text{op}$
$E \rightarrow (E_1)$	$E.\text{Cod} = E_1.\text{Cod}$
$E \rightarrow \text{id}$	$E.\text{Cod} = \text{id.Loc}$
$A \rightarrow \text{id}$	$A.\text{Cod} = \text{id.Loc}$

Această schemă de traducere poate fi implementată în cadrul unui parser bottom-up precizând care sunt acțiunile ce se aplică la fiecare reducere:

Regula sintactică	Acțiunea (segment de program)
$S \rightarrow A := E$	{ print " := " }
$E \rightarrow E_1 \text{ op } E_2$	{ print op }
$E \rightarrow (E_1)$	{ }
$E \rightarrow \text{id}$	{ print id.Loc }
$A \rightarrow \text{id}$	{ print id.Loc }

Traducerea textului $d := b + a * c$ cu un parser ascendent este dată mai jos:

Stiva	Intrare	Acțiune
#	$d := b + a * c \#$	Deplasare d
# d	$:= b + a * c \#$	Reduce și print d
# A	$:= b + a * c \#$	Deplasare $:=$
# $A :=$	$b + a * c \#$	Deplasare b
# $A := b$	$+ a * c \#$	Reduce și print b
# $A := E$	$+ a * c \#$	Deplasare $+$
# $A := E +$	$a * c \#$	Deplasare a
# $A := E + a$	$* c \#$	Reduce și print a
# $A := E + E$	$* c \#$	Deplasare $*$
# $A := E + E *$	$c \#$	Deplasare c
# $A := E + E * c$	#	Reduce și print c
# $A := E + E * E$	#	Reduce și print *
# $A := E + E$	#	Reduce și print +
# $A := E$	#	Reduce și print :=
# S	#	ACCEPTARE

Așadar, forma postfix a textului dat este $d \ b \ a \ c \ * \ + \ :=$.

6.4.3 Traducerea codului în arbore sintactic

Arborele de derivare (parsare) conține, în general informații redundante (relativ la generarea de cod obiect) care pot fi eliminate obținând astfel un arbore mai economic ce reprezintă codul sursă. Acest arbore se numește arbore sintactic (vezi figura 7.5.2).

Să arătăm cum se obține arborele abstract pentru aceeași gramatică ce generează instrucțiunile de atribuire:

$$S \rightarrow A := E \quad E \rightarrow E \text{ op } E \mid (E) \mid \text{id} \quad A \rightarrow \text{id}$$

Să notăm prin A.Val, E.Val respectiv S.Val acest arbore corepunzător lui A, expresiilor E respectiv asignărilor S. Să considerăm `arbinar(op, st, dr)` o funcție ce construiește un arbore cu rădăcina `op` și subarborii `st` respectiv `dr`, iar `frunza(id)` o funcție ce construiește arborele cu rădăcina `id` fără subarbori. Fiecare din aceste funcții returnează pointeri la arborii construiți. Atunci schema de traducere este următoarea:

Reguli sintactice	Acțiuni semantice
$S \rightarrow A := E$	S.Val = arbinar(:=, A.Val, E.Val)
$E \rightarrow E_1 \text{ op } E_2$	E.Val = arbinar(op, E ₁ .Val, E ₂ .Val)
$E \rightarrow (E_1)$	E.Val = E ₁ .Val
$E \rightarrow \text{id}$	E.Val = frunza(id.Loc)
$A \rightarrow \text{id}$	A.Val = frunza(id.Loc)

Traducerea aceluiași text `d := b + a*c`, ca în paragraful anterior, cu un parser ascendent este dată mai jos:

Stiva	Intrare	Acțiune
#	d := b+a*c#	Deplasare d
#d	:= b+a*c#	Reduce si t₁=frunza(d)
#A	:= b+a*c#	Deplasare :=
#A:=	b+a*c#	Deplasare b
#A:= b	+a*c#	Reduce si t₂=frunza(b)
#A:= E	+a*c#	Deplasare +
#A:= E+	a*c#	Deplasare a
#A:= E+a	*c#	Reduce si t₃=frunza(a)
#A:= E+E	*c#	Deplasare *
#A:= E+E*	c#	Deplasare c
#A:= E+E*c	#	Reduce si t₄=frunza(c)
#A:= E+E*c	#	Reduce si t₅=arbinar(*, t₃, t₄)
#A:= E+E	#	Reduce si t₆=arbinar(+, t₂, t₅)
#A:= E	#	Reduce si t₇=arbinar(:=, t₁, t₆)
#S	#	ACCEPTARE

Arborele construit este cel din figura 7.5.3. Traducerea codului obiect în arbore abstract este utilizată în proiectarea de interpretoare dar și de compilatoare folosind generatoare de parsare yacc sau bison. În paragraful 7.6 prezentăm un exemplu de interpreter construit în acest mod.

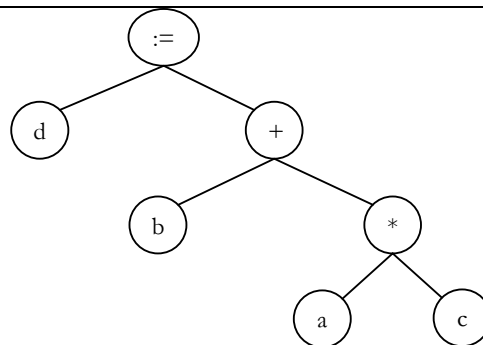


Figura 7.5.3

6.4.4 Traducerea codului în secvențe de quadrule

Codul cu trei adrese este o secvență de instrucțiuni simple care conțin cel mult trei referințe la o tabelă de simboluri. Odată transformat codul sursă în cod cu trei adrese, traducerea în cod mașină este relativ simplă deoarece, în general, există o corespondență între instrucțiunile codului intermediar cu trei adrese și cele ale mașinii. Avantajul este că, până la acest nivel, proiectarea compilatorului nu depinde de mașina pentru care se implementează. Instrucțiunile din codul cu trei adrese se implementează ca și structuri care conțin patru câmpuri: pe lângă cele trei adrese se reprezintă și operatorul ce intervine în instrucțiune. Iată cele mai uzuale instrucțiuni ale codului cu trei adrese:

```

A = B op C
A = op B
goto L
if A goto L
if A rel B goto L
A = B[I]
A[I] = B

```

În acest paragraf dorim să dezvoltăm o schemă de traducere care să transforme un text sursă în cod intermediar cu trei adrese. De exemplu, codul:

```

while (a != b)
    if(a < b) then b = b - a else a = a - b

```

va fi transformat în codul intermediar următor:

```

101 if a != b goto 103
102 go to 111
103 if a < b goto 105
104 goto 108
105 T1 = b - a
106 b = T1
107 goto 101
108 T2 = a - b
109 a = T2
110 goto 101
111 ...

```

Vom considera limbajul descris de următoarea gramatică:

```

S → A
   | begin L end
   | if B then S
   | if B then S else S
   | while B do S
L → L ; S

```

$$\begin{aligned}
 & \mid S \\
 A & \rightarrow id = E \\
 E & \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid - E \mid (E) \mid id \\
 B & \rightarrow B \text{ or } B \mid B \text{ and } B \mid \text{not } B \mid (B) \\
 & \mid id < id \mid id > id \mid id == id \mid id != id \mid id
 \end{aligned}$$

Acestei sintaxe îi vom adăuga o semantică (acțiuni semantice) care să permită traducerea programelor în cod intermediar.

Traducerea instrucțiunilor de atribuire

Vom începe descrierea schemei de traducere cu cea pentru instrucțiunile de atribuire. Intenția este de a traduce un text de forma $\text{delta} = b * b - 4 * a * c$ în secvența:

$$\begin{aligned}
 T_1 &= b * b \\
 T_2 &= 4 * a \\
 T_3 &= T_2 * c \\
 T_4 &= T_1 - T_3 \\
 \text{delta} &= T_4
 \end{aligned}$$

Pentru a obține această traducere să definim niște atribute pentru simbolurile E și id :

- $E.\text{Nume}$ este numele ce păstrează valoarea expresiei E . Pentru evaluarea expresiei vom folosi nume temporare care să păstreze valorile subexpresiilor. Aceste nume le vom nota cu T_1, T_2, \dots și ele sunt furnizate, în această ordine, de apelul succesiv al unei funcții pe care o notăm $\text{newtemp}()$ și care se poate implementa ușor.
- $id.\text{Nume}$ este numele instanței identificatorului id respectiv valoarea constantei atunci când id este instanțiat cu o constantă.

Pentru obținerea unei instrucțiuni din codul intermediar folosim o funcție pe care o notăm $\text{GEN}(I)$ unde I este o instrucțiune. De pildă, $\text{GEN}(E.\text{Nume} = E_1.\text{Nume} + E_2.\text{Nume})$ va genera instrucțiunea $T_1 = b * b$ dacă $E.\text{Nume}$ este T_1 , $E_1.\text{Nume}$ și $E_2.\text{Nume}$ sunt b . Cu acestea schema de traducere se descrie astfel:

Reguli sintactice	Acțiuni semantice
$A \rightarrow id = E$	$\text{GEN}(id.\text{Nume} = E.\text{Nume})$
$E \rightarrow E_1 + E_2$	$E.\text{Nume} = \text{newtemp}()$ $\text{GEN}(E.\text{Nume} = E_1.\text{Nume} + E_2.\text{Nume})$
$E \rightarrow E_1 - E_2$	$E.\text{Nume} = \text{newtemp}()$ $\text{GEN}(E.\text{Nume} = E_1.\text{Nume} - E_2.\text{Nume})$
$E \rightarrow E_1 * E_2$	$E.\text{Nume} = \text{newtemp}()$ $\text{GEN}(E.\text{Nume} = E_1.\text{Nume} * E_2.\text{Nume})$
$E \rightarrow E_1 / E_2$	$E.\text{Nume} = \text{newtemp}()$ $\text{GEN}(E.\text{Nume} = E_1.\text{Nume} / E_2.\text{Nume})$
$E \rightarrow -E_1$	$E.\text{Nume} = \text{newtemp}()$ $\text{GEN}(E.\text{Nume} = - E_1.\text{Nume})$
$E \rightarrow (E_1)$	$E.\text{Nume} = E_1.\text{Nume}$
$E \rightarrow id$	$E.\text{Nume} = id.\text{Nume}$

Codul intermediar se obține în timpul parsării folosind un parser bottom-up împreună cu această schemă de traducere: de fiecare dată când se aplică o reducere (regulă sintactică), se aplică și acțiunile semantice corespunzătoare. Vom exemplifica obținerea codului intermediar pentru instrucțiunea $\text{delta} = b * b - 4 * a * c$. Cuvântul ce trebuie

analizat este $id = id * id - id * id * id$ unde $id.Nume$ are respectiv valorile δ , b , b , 4 , a , c . Aceste valori sunt luate din tabela identificatorilor și a constantelor construite de analizorul lexical. De fapt, $id.Nume$ pot fi reprezentați ca și pointeri la tabela simbolurilor și pot fi păstrați în stiva de parsare împreună cu simbolurile corespunzătoare: o intrare în stivă este o pereche $(X, X.Nume)$. În tabela următoare este dată întreaga analiză sintactică și semantică:

Intrare	Simbol stivă	Simbol.Nume	Acțiune	Cod
$\delta = b * b - 4 * a * c \#$			<i>deplasare</i>	
$= b * b - 4 * a * c \#$	id	δ	<i>deplasare</i>	
$b * b - 4 * a * c \#$	id=	$\delta_$	<i>deplasare</i>	
$* b - 4 * a * c \#$	id=id	δ_b	reducere	
$* b - 4 * a * c \#$	id=E	δ_b	<i>deplasare</i>	
$b - 4 * a * c \#$	id=E*	$\delta_b_$	<i>deplasare</i>	
$- 4 * a * c \#$	id=E*id	δ_b_b	reducere	
$- 4 * a * c \#$	id=E*E	δ_b_b	reducere	$T_1 = b * b$
$- 4 * a * c \#$	id=E	δ_T_1	<i>deplasare</i>	
$4 * a * c \#$	id=E-	$\delta_T_1_$	<i>deplasare</i>	
$* a * c \#$	id=E-id	$\delta_T_1_4$	reducere	
$* a * c \#$	id=E-E	$\delta_T_1_4$	<i>deplasare</i>	
$a * c \#$	id=E-E*	$\delta_T_1_4_$	<i>deplasare</i>	
$* c \#$	id=E-E*id	$\delta_T_1_4_a$	reducere	
$* c \#$	id=E-E*E	$\delta_T_1_4_a$	reducere	$T_2 = 4 * a$
$* c \#$	id=E-E	$\delta_T_1_T_2$	<i>deplasare</i>	
$c \#$	id=E-E*	$\delta_T_1_T_2_$	<i>deplasare</i>	
$\#$	id=E-E*id	$\delta_T_1_T_2_c$	reducere	
$\#$	id=E-E*E	$\delta_T_1_T_2_c$	reducere	$T_3 = T_2 * c$
$\#$	id=E-E	$\delta_T_1_T_3$	reducere	$T_4 = T_1 - T_3$
$\#$	id=E	δ_T_4	reducere	$\delta = T_4$
$\#$	A			

Traducerea expresiilor logice

Evaluarea expresiilor logice poate fi făcută în diverse moduri. Dacă se codifică valorile de adevăr prin valori numerice (ca în limbajul C), atunci evaluarea se face analog cu evaluarea expresiilor aritmetice. În acest caz o schemă de traducere pentru expresii logice poate fi obținută ușor din schema de traducere a expresiilor aritmetice. O altă modalitate este aceea de a traduce expresia logică într-o succesiune de instrucțiuni (cod intermediar) de forma celor care conțin goto sau if și de a obține valorile *adevărat* respectiv *fals* prin atingerea unor poziții în acest cod intermediar. De exemplu, expresia:

$a < b$ or $b > c$ and not d

ar putea fi tradusă în codul:

```
101 if a < b goto ?
102 goto 103
103 if b > c goto 105
104 goto ?
105 if d goto ?
106 goto ?
```

cu precizarea că ieșirile (țintele lui goto) din acest cod la liniile 101 și 106 sunt corespunzătoare valorii *adevărat* pentru expresia dată iar cele de la 104 și 105 sunt corespunzătoare valorii *fals*. Dacă această expresie apare în contextul unui „program”, de pildă:

if($a < b$ or $b > c$ and not d) then $x = x + y$ else $x = x - y$

atunci traducerea acestuia ar putea fi:

```

101 if a < b goto 107
102 goto 103
103 if b > c goto 105
104 goto 110
105 if d goto 110
106 goto 107
107 T1 = x + y
108 x = T1
109 goto 113
110 T2 = x - y
111 x = T2
112 goto 113
113 ...

```

Vom indica în continuare o schemă de traducere a unei expresii logice într-o secvență de cod intermediar ce conține instrucțiuni de salt condiționat (if A goto L, if A rel B goto L) și salt necondiționat (goto L). În acest cod, unele din instrucțiuni vor avea țintele L către o locație dintr-o mulțime notată B.True în cazul în care expresia B are valoarea *adevărat* iar altele către o locație dintr-o mulțime B.False în cazul în care expresia B are valoarea *fals*. B.True și B.False vor fi așadar atribute ale simbolului B, care au ca valori liste de pointeri către quadruplele ce formează codul intermediar. Dacă, de exemplu, avem expresia $B = B_1 \text{ or } B_2$ atunci B.True este, deocamdată, B₁.True pentru că, odată ce am evaluat expresia B₁ și am constatat că are valoarea *adevărat*, atunci E are valoarea *adevărat*. Dacă B₁ are valoarea *fals*, atunci trebuie evaluată expresia B₂ pentru a obține valoarea lui B, încât B.False va fi B₁.False iar B.True este reuniunea celor două mulțimi (liste) B₁.True și B₂.True. Aici apare următoarea problemă: în quadruplele din B₁.False trebuie ca ținta L să aibă ca valoare pointerul la quadrupla de început al codului lui B₂. Ori acest pointer nu este cunoscut înainte de a termina generarea codului pentru B₁. Soluția este de a genera instrucțiuni goto care nu au completată ținta L; aceasta va fi stabilită atunci când se face o reducere cu regula $B \rightarrow B_1 \text{ or } B_2$ pentru că în acest punct se cunoaște codul generat atât pentru B₁, cât și pentru B₂. Asemănător se petrec lucrurile la expresiile de forma B₁ and B₂. sau la descrierea semanticii instrucțiunilor if sau while. Pentru a descrie o schemă de traducere în această idee, avem nevoie de niște funcții (vom păstra denumirea din AhU83):

- `makelist(i)`: funcția creează o nouă listă de quadruple care conține doar lista cu indexul i (pointer la tabloul quadruplelor ce se generează) și returnează pointer la lista creată;
- `merge(p1, p2)`: funcția concatenează listele pointate de p₁ și p₂ și returnează pointer la noua listă obținută;
- `backpatch(p, i)`: funcția completează cu i țintele instrucțiunilor goto din quadruplele listei pointată de p.

Vom nota de asemenea cu NEXTQUAD un pointer ce păstrează, la un moment dat, numărul primei quadruple ce trebuie generată, și care va permite identificarea punctului de intrare în codul unei subexpresii. Acest lucru se poate realiza prin introducerea unui neterminat M în sintaxa expresiilor, care să marcheze începutul codului unei subexpresii prin atributul M.Quad (=NEXTQUAD). Astfel, regulile sintactice $B \rightarrow B \text{ or } B$ și $B \rightarrow B \text{ and } B$ se transformă în $B \rightarrow B \text{ or } M B$, $B \rightarrow B \text{ and } M B$ și $M \rightarrow \epsilon$. Cu aceste pregătiri, schema de traducere pentru expresii logice este următoarea (am notat prin rel un operator relațional):

Reguli sintactice	Acțiuni semantice
$B \rightarrow B_1 \text{ or } M B_2$	backpatch($B_1.\text{False}$, $M.\text{Quad}$); $B.\text{True} = \text{merge}(B_1.\text{True}, B_2.\text{True})$; $B.\text{False} = B_2.\text{False}$;
$B \rightarrow B_1 \text{ or } M B_2$	backpatch($B_1.\text{True}$, $M.\text{Quad}$); $B.\text{True} = B_2.\text{True}$; $B.\text{False} = \text{merge}(B_1.\text{False}, B_2.\text{False})$;
$B \rightarrow \text{not } B_1$	$B.\text{True} = B_1.\text{False}$; $B.\text{False} = B_1.\text{True}$;
$B \rightarrow (B_1)$	$B.\text{True} = B_1.\text{True}$; $B.\text{False} = B_1.\text{False}$;
$B \rightarrow \text{id}_1 \text{ rel id}_2$	$B.\text{True} = \text{makelist}(\text{NEXTQUAD})$; $B.\text{False} = \text{makelist}(\text{NEXTQUAD} + 1)$; $\text{GEN}(\text{ if id}_1.\text{Nume rel id}_2.\text{Nume go to ...})$; $\text{GEN}(\text{go to ...})$;
$B \rightarrow \text{id}$	$B.\text{True} = \text{makelist}(\text{NEXTQUAD})$; $B.\text{False} = \text{makelist}(\text{NEXTQUAD} + 1)$; $\text{GEN}(\text{ if id.Nume go to ...})$; $\text{GEN}(\text{go to ...})$;
$M \rightarrow \varepsilon$	$M.\text{Quad} = \text{NEXTQUAD}$;

Să observăm că am tratat aici un caz particular, când expresiile relaționale sunt doar de forma id rel id . Cititorul este invitat să trateze cazul general al sintaxei $B \rightarrow E \text{ rel } E$ unde E este o expresie aritmetică.

Dacă realizăm analiza sintactică a expresiei

$$a < b \text{ or } b > c \text{ and not } d$$

și aplicăm acțiunile semantice la fiecare reducere, obținem rezultatele din tabela de mai jos. Pentru variabila B vom nota prin $(q_1, q_2, \dots; p_1, p_2, \dots)$ valorile atributelor $B.\text{True}$ respectiv $B.\text{False}$.

Intrare	Simbol stivă	Valorile atributelor	Cod
$a < b \text{ or } b > c \text{ and not } d$ #			
$< b \text{ or } b > c \text{ and not } d$ #	id	a	
$b \text{ or } b > c \text{ and not } d$ #	id <	a_	
$\text{or } b > c \text{ and not } d$ #	id < id	a_b	(1)
$\text{or } b > c \text{ and not } d$ #	B	(101; 102)	
$b > c \text{ and not } d$ #	B or	(101; 102)_	
$b > c \text{ and not } d$ #	B or M	(101; 102)_103	
$> c \text{ and not } d$ #	B or M id	(101; 102)_103b	
$c \text{ and not } d$ #	B or M id >	(101; 102)_103b_	
$\text{and not } d$ #	B or M id > id	(101; 102)_103b_c	(2)
$\text{and not } d$ #	B or M B	(101; 102) 103(103; 104)	
$\text{not } d$ #	B or M B and	(101; 102) 103(103; 104)_	
$\text{not } d$ #	B or M B and M	(101; 102) 103(103; 104)_105	
d #	B or M B and M not	(101; 102) 103(103; 104)_105_	
#	B or M B and M not id	(101; 102) 103(103; 104)_105_d	(3)
#	B or M B and M not B	(101; 102) 103(103; 104)_105_ (105;106)	
#	B or M B and M B	(101; 102) 103(103; 104)_105 (106;105)	(4)
#	B or M B	(101; 102) 103(106; 104,105)	(5)
#	B	(101, 106; 104,105)	

- (1) 101 if a<b goto...
102 goto...

-
- (2) 103 if b < c goto...
104 goto...
 - (3) 105 if d goto...
106 goto...
 - (4) În lista {103} se pune 105 ca țintă pentru goto
Se determină conform acțiunilor semantice B.True și B.False
 - (5) În lista {102} se pune 103 ca țintă pentru goto
Se determină conform acțiunilor semantice B.True și B.False

Urmare a acestor operații codul obținut este:

```

101 if a < b goto ...
102 goto 103
103 if b > c goto 105
104 goto ...
105 if d goto ...
106 goto ...

```

iar listele atașate lui B, corespunzătoare valorilor *adevărat* respectiv *fals*, sunt:

B.True = {101, 106}, B.False = {104, 105}.

Traducerea instrucțiunilor

Să dezvoltăm mai întâi o schemă de traducere pentru instrucțiunile de control. Dacă ar trebui să traducem o instrucțiune de forma **if B then S else S**, atunci codul obținut trebuie să conțină codul pentru B, urmat de codul pentru primul S apoi codul pentru al doilea S. În determinarea codului pentru B se obțin, așa cum am văzut mai sus, listele B.True și B.False. Țintele pentru goto din B.True trebuie să coincidă cu eticheta de început a codului pentru primul S iar cele din B.False trebuie să coincidă cu eticheta de început a codului celui de-al doilea S. De asemenea, în codul obținut trebuie să precizăm și faptul că, după execuția instrucțiunilor desemnate de S (atât primul, cât și al doilea) urmează instrucțiunea de după if. De aceea vom considera un atribut S.Next care reprezintă un pointer la o listă de quadruple ce conțin goto, condiționat sau nu, la instrucțiunea ce urmează după S în ordinea execuției. Mecanismul de completare a acestor ținte va fi cel de la traducerea expresiilor logice încât va trebui și aici să considerăm un neterminat ce marchează începutul codului lui S și un neterminat N care printr-un atribut N.Next permite saltul după execuția ramurii then a instrucțiunii if. Așadar, schema de traducere pentru instrucțiunile if se poate descrie astfel:

1. $S \rightarrow \text{if } B \text{ then } M_1 S_1 N \text{ else } M_2 S_2$
 $\text{backpatch}(B.\text{True}, M_1.\text{Quad});$
 $\text{backpatch}(B.\text{False}, M_2.\text{Quad});$
 $S.\text{Next} = \text{merge}(S_1.\text{Next}, N.\text{Next}, S_2.\text{Next});$
2. $N \rightarrow \epsilon$
 $N.\text{Next} = \text{makelist}(\text{NEXTQUAD});$
 $\text{GEN}(\text{goto...});$
3. $M \rightarrow \epsilon$
 $M.\text{Quad} = \text{NEXTQUAD};$
4. $S \rightarrow \text{if } B \text{ then } M S_1$
 $\text{backpatch}(B.\text{True}, M.\text{Quad});$
 $S.\text{Next} = \text{merge}(B.\text{False}, S_1.\text{Next});$

În mod analog tratăm cazul instrucțiunii `while`. Dacă sintaxa este **while B do S** atunci va trebui să stabilim în quadruplele din `B.True` din codul lui `B`, ca țintă, începutul codului lui `S`, iar în quadruplele din `B.False` ținta trebuie să fie începutul codului instrucțiunii ce urmează lui `while`. De asemenea, trebuie să ne asigurăm că, după execuția lui `S`, se reia evaluarea lui `B`. Atunci, traducerea acestei instrucțiuni se face după schema:

```
5. S → while M1 B do M2 S1
    backpatch(S1.Next, M1.Quad);
    backpatch(B.True, M2.Quad);
    S.Next = B.False;
    GEN(goto M1.Quad);
```

Să observăm că apare pentru prima dată generarea unei instrucțiuni `goto` care are ținta fixată: `goto M1.Quad`. Asta asigură faptul că, în toate cazurile, după execuția lui `S1` urmează evaluarea din nou a expresiei `B`.

Pentru traducerea celorlalte instrucțiuni din limbajul specificat la începutul paragrafului lucrurile sunt cât se poate de simple:

```
6. S → begin L end
    S.Next = L.Next;
7. S → A
    S.Next = makelist(); // lista vidă
8. L → L1; M S
    backpatch(L1.Next, M.Quad);
    L.Next = S.Next;
9. L → S
    L.Next = S.Next;
```

Să mai observăm că doar în cazurile $N \rightarrow \varepsilon$ și la instrucțiunea `while` sunt generate noi quadruple; codul este asociat în rest expresiilor aritmetice sau logice. Traducerea instrucțiunilor înseamnă în plus asocierea țințelor pentru `goto` care au fost generate fără a fi complete (apelul funcției `backpatch`).

Dacă aplicăm un parser bottom-up pentru textul:

```
while (a != b)
  if (a < b) then
    b = b - a
  else
    a = a - b
```

și realizăm inclusiv traducerea sa conform schemei de traducere prezentată, obținem un cod intermediar de forma celui anunțat la începutul acestui paragraf:

```
101 if a != b goto 103
102 go to 111
103 if a < b goto 105
104 goto 108
105 T1 = b - a
106 b = T1
107 goto 101
108 T2 = a - b
109 a = T2
110 goto 101
111 ...
```

6.5 Proiectarea unui interpretor cu flex și yacc

Vom descrie în acest paragraf modul cum se poate proiecta un interpretor folosind instrumentele `flex` și `bison`. Dorim să interpretăm programe de forma:

```

a = 36;
b = 42;
while (a != b)
    if(a > b) a = a - b;
    else b = b - a;
print a;

```

Pentru aceasta vom descrie fişierele de intrare pentru flex respectiv bison. Fişierul pentru flex descrie unităţile lexicale care vor fi transmise analizorului sintactic construit de bison. Fişierul flex este următorul:

```

/* Fisierul glang.l intrare pentru flex. Se foloseste comanda
   flex glang.l dupa comanda bison -d glang.y */
%{
#include <stdlib.h>
#include "glang.h"
#include "glang.tab.h"
void yyerror(char *);
}%

%%

[a-z]      { yylval.sIndex = *yytext - 'a'; return VAR; }

[0-9]+     { yylval.iVal = atoi(yytext); return INT; }
[-()<=>+*/;{}.] { return *yytext; }
">="      { return GE; }
"<="      { return LE; }
"=="      { return EQ; }
"!="      { return NE; }
"while"    { return WHILE; }
"if"       { return IF; }
"else"     { return ELSE; }
"print"    { return PRINT; }

[ \t\n]+   ; /* Spatiile sunt ignorate */

.          yyerror(" Caracter ilegal!\n");
%%
int yywrap(void) {
    return 1;
}

```

Fişierul glang.h descrie structurile de date folosite pentru reprezentarea constantelor, a identificatorilor precum şi a tipurilor de noduri din arborele abstract care va fi construit de analizorul sintactic. Identificatorii sunt definiţi aici ca fiind litere şi de aceea am folosit extern int sym[26] ca tabelă a identificatorilor. O să indicăm mai târziu cum se poate aborda cazul general.

```

/* Fisierul glang.h */
typedef enum { typeCon, typeId, typeOper } nodeEnum;

/* constante */
typedef struct {
    int val; /* valoarea constantei */
} conNodeType;

/* identificatori */

```

```

typedef struct {
    int i;          /* indice in tabela de simboluri */
} idNodeType;

/* operatori */
typedef struct {
    int oper;          /* operator */
    int nops;          /* nr. operanzi */
    struct nodeTypeTag *op[1]; /* operanzi */
} oprNodeType;

typedef struct nodeTypeTag {
    nodeEnum type;          /* tipul nodului */
    union {
        conNodeType con;    /* constante */
        idNodeType id;      /* identificatori */
        oprNodeType opr;    /* operatori */
    };
} nodeType;
extern int sym[26];

```

Fișierul `glang.tab.h` este obținut de către yacc (bison) atunci când se folosește opțiunea `-d`. Acest fișier definește codurile numerice pentru unitățile lexicale:

```

/* Fisierul glang.tab.h creat de bison cu optiunea -d */
#ifndef BISON_GLANG_TAB_H
#define BISON_GLANG_TAB_H
#ifndef YYSTYPE
typedef union {
    int iVal;          /* valoare intreaga */
    char sIndex;       /* index in tabela de simboluri */
    nodeType *nodPtr;  /* pointer la un nod in arbore */
} YYSTYPE;
#define YYSTYPE YYSTYPE
#define YYSTYPE_IS_TRIVIAL 1
#endif
#define INT 257
#define VAR 258
#define WHILE 259
#define IF 260
#define PRINT 261
#define IFX 262
#define ELSE 263
#define GE 264
#define LE 265
#define EQ 266
#define NE 267
#define UMINUS 268

extern YYSTYPE yylval;
#endif /* not BISON_GLANG_TAB_H */

```

Fișierul de intrare pentru bison descrie sintaxa limbajului de programare împreună cu acțiunile semantice ce trebuie aplicate în cazul reducerilor. Aici este vorba de apelul funcției de construire a arborelui abstract corespunzător. Iată mai jos fișierul de intrare pentru bison:

```

/* Fisierul glang.y intrare pentru bison. Se foloseste
   comanda : bison -d glang.y */
%{

```

```

#include <stdio.h>
#include <stdlib.h>
#include <stdarg.h>
#include "glang.h"

/* prototipurile functiilor */
nodeType *nodOper(int oper, int nops, ...);
        /* nod operator in arbore */
nodeType *nodId(int i);
        /* nod frunza identificator */
nodeType *nodCon(int value);
        /* nod frunza constanta */
void freeNode(nodeType *p);
        /* eliberare memorie */
int interpret(nodeType *p);
        /* functia de interpretare */
int yylex(void);
        /* functia creata de flex */

void yyerror(char *s);
int sym[26];          /* tabela de simboluri */
%}

%union {
    int iVal;          /* valoare intreaga */
    char sIndex;       /* index in tabela de simboluri */
    nodeType *nodPtr;  /* pointer la un nod in arbore */
};

%token <iVal> INT
%token <sIndex> VAR
%token WHILE IF PRINT
%nonassoc IFX          /* rezolvarea ambiguitatii if-then-else */
%nonassoc ELSE
%left GE LE EQ NE '>' '<'
%left '+' '-'
%left '*' '/'
%nonassoc UMINUS

%type <nodPtr> stmt expr stmt_list

%%
program: function          { exit(0); }
;

function: function stmt    { interpret($2); freeNode($2); }
| /* NULL */
;

stmt:      ';'             { $$ = nodOper(';', 2, NULL, NULL); }
| expr ';'               { $$ = $1; }
| PRINT expr ';'         { $$ = nodOper(PRINT, 1, $2); }
| VAR '=' expr ';'       {
    $$ = nodOper('=', 2, nodId($1), $3);
}
| WHILE '(' expr ')' stmt {
    $$ = nodOper(WHILE, 2, $3, $5);
}
| IF '(' expr ')' stmt %prec IFX
    { $$ = nodOper(IF, 2, $3, $5); }
| IF '(' expr ')' stmt ELSE stmt
    { $$ = nodOper(IF, 3, $3, $5, $7); }
| '{' stmt_list '}'      { $$ = $2; }
;

```

```

stmt_list:
    stmt                { $$ = $1; }
  | stmt_list stmt      { $$ = nodOper(';', 2, $1, $2); }
  ;

expr:
    INT                { $$ = nodCon($1); }
  | VAR                { $$ = nodId($1); }
  | '-' expr %prec UMINUS
    { $$ = nodOper(UMINUS, 1, $2); }
  | expr '+' expr
    { $$ = nodOper('+', 2, $1, $3); }
  | expr '-' expr
    { $$ = nodOper('-', 2, $1, $3); }
  | expr '*' expr
    { $$ = nodOper('*', 2, $1, $3); }
  | expr '/' expr
    { $$ = nodOper('/', 2, $1, $3); }
  | expr '<' expr
    { $$ = nodOper('<', 2, $1, $3); }
  | expr '>' expr
    { $$ = nodOper('>', 2, $1, $3); }
  | expr GE expr
    { $$ = nodOper(GE, 2, $1, $3); }
  | expr LE expr
    { $$ = nodOper(LE, 2, $1, $3); }
  | expr NE expr
    { $$ = nodOper(NE, 2, $1, $3); }
  | expr EQ expr
    { $$ = nodOper(EQ, 2, $1, $3); }
  | '(' expr ')'
    { $$ = $2; }
  ;

%%

#define SIZEOF_NODETYPE ((char *)&p->con - (char *)p)

nodeType *nodCon(int value) {
    nodeType *p;
    size_t nodeSize;

    /* alocare memorie pentru noul nod */
    nodeSize = SIZEOF_NODETYPE + sizeof(conNodeType);
    if ((p = malloc(nodeSize)) == NULL)
        yyerror("out of memory");

    /* copiere valoare constanta */
    p->type = typeCon;
    p->con.val = value;

    return p;
}

nodeType *nodId(int i) {
    nodeType *p;
    size_t nodeSize;

    /* alocare memorie pentru noul nod */
    nodeSize = SIZEOF_NODETYPE + sizeof(idNodeType);
    if ((p = malloc(nodeSize)) == NULL)
        yyerror("out of memory");

```

```

    /* copiere valoare indice */
    p->type = typeId;
    p->id.i = i;

    return p;
}

nodeType *nodOper(int oper, int nops, ...) {
    va_list ap;
    nodeType *p;
    size_t nodeSize;
    int i;
    /* alocare memorie pentru noul nod */
    nodeSize = SIZEOF_NODETYPE + sizeof(oprNodeType) +
        (nops - 1) * sizeof(nodeType*);
    if ((p = malloc(nodeSize)) == NULL)
        yyerror("out of memory");

    /* copiere informatii functie de nops */
    p->type = typeOper;
    p->opr.oper = oper;
    p->opr.nops = nops;
    va_start(ap, nops);
    for (i = 0; i < nops; i++)
        p->opr.op[i] = va_arg(ap, nodeType*);
    va_end(ap);
    return p;
}

void freeNode(nodeType *p) {
    int i;

    if (!p) return;
    if (p->type == typeOper) {
        for (i = 0; i < p->opr.nops; i++)
            freeNode(p->opr.op[i]);
    }
    free (p);
}

void yyerror(char *s) {
    fprintf(stdout, "%s\n", s);
}

int main(void) {
    yyparse();
    return 0;
}

```

Conform regulilor următoare (din fișierul tocmai prezentat):

```

function: function stmt
        { interpret($2); freeNode($2); }
| /* NULL */
;

```

un program în limbajul sursă descris este un șir de instrucțiuni (stmt). Fiecare din aceste instrucțiuni este interpretată: apelul funcției `interpret` pentru argumentul `$2` care este valoarea atributului lui `stm`, adică arborele construit pentru `stmt`, face ca această instrucțiune să fie executată. Prin `freeNode($2)` se eliberează memoria alocată la construirea arborelui amintit. Funcția `interpret()` este prezentată mai jos :


```

/* Functia interpret(nodeType *p) */
#include <stdio.h>
#include "glang.h"
#include "glang.tab.h"

int interpret(nodeType *p) {
    if (!p) return 0;
    switch(p->type) {
        case typeCon:      return p->con.val;
        case typeId:       return sym[p->id.i];
        case typeOper:
            switch(p->opr.oper) {
                case WHILE: while(interpret(p->opr.op[0]))
                             interpret(p->opr.op[1]); return 0;
                case IF:    if (interpret(p->opr.op[0]))
                             interpret(p->opr.op[1]);
                             else if (p->opr.nops > 2)
                                 interpret(p->opr.op[2]);
                             return 0;
                case PRINT: printf("%d\n", interpret(p->opr.op[0]));
                             return 0;
                case ';':   interpret(p->opr.op[0]);
                             return interpret(p->opr.op[1]);
                case '=':   return sym[p->opr.op[0]->id.i] =
                             interpret(p->opr.op[1]);
                case UMINUS: return -interpret(p->opr.op[0]);
                case '+':   return interpret(p->opr.op[0]) +
                             interpret(p->opr.op[1]);
                case '-':   return interpret(p->opr.op[0]) -
                             interpret(p->opr.op[1]);
                case '*':   return interpret(p->opr.op[0]) *
                             interpret(p->opr.op[1]);
                case '/':   return interpret(p->opr.op[0]) /
                             interpret(p->opr.op[1]);
                case '<':   return interpret(p->opr.op[0]) <
                             interpret(p->opr.op[1]);
                case '>':   return interpret(p->opr.op[0]) >
                             interpret(p->opr.op[1]);
                case GE:    return interpret(p->opr.op[0]) >=
                             interpret(p->opr.op[1]);
                case LE:    return interpret(p->opr.op[0]) <=
                             interpret(p->opr.op[1]);
                case NE:    return interpret(p->opr.op[0]) !=
                             interpret(p->opr.op[1]);
                case EQ:    return interpret(p->opr.op[0]) ==
                             interpret(p->opr.op[1]);
            }
    }
    return 0;
}

```

În urma lansării comenzilor `bison -d glang.y` și `flex glang.l` se obțin fișierele `lexyy.c` (analizorul lexical) și `glang.tab.c` (analizorul sintactic) care, împreună cu interpretul din fișierul `interpret.c`, vor fi compilate cu un compilator de limbaj C și se obține astfel un interpret pentru limbajul sursă descris. Dacă lansăm acest interpret pentru programul prezentat la începutul secțiunii, se obține rezultatul 6.

Bibliografie

- [AnG97] Andrei Șt., Grigoraș Gh. : *Tehnici de compilare. Lucrări de laborator*. Editura Universității „Al.I.Cuza”, Iași, 1995
- [App97] Appel, A., *Modern Compiler Implementation in C*, Cambridge University Press, 1997
- [ASU86] Aho, A.V., Sethi, R., Ullman, J.D.: *Compilers: Principles, Techniques, and Tools*. Addison-Wesley Publishing Company, U.S.A., 1986
- [Gri86] Grigoraș, Gh.: *Limbaje formale și tehnici de compilare*. Editura Universității „Al.I.Cuza”, Iași, 1986
- [Gri05] Grigoraș, Gh.: *Construcția compilatoarelor - Algoritmi fundamentali*, Editura Universității "Al. I. Cuza" Iasi, ISBN 973-703-084-2, 274 pg., 2005.
- [JuA97] Jucan, T., Andrei, Șt.: *Limbaje formale și teoria automatelor. Culegere de probleme*. Editura Universității „Al.I.Cuza”, Iași, 1997
- [Juc99] Jucan, T., *Limbaje formale și automate*, Editura MatrixRom, 1999.
- [Șer97] Șerbănați, L.D., *Limbaje de programare și compilatoare*, Editura Academiei, București, 1987.