

TEHNICI SIMBOLICE FUNDAMENTALE ALE LINGVISTICII COMPUTAȚIONALE

V. Cotelea, *dr.hab. conf.univ.*

Academia de Studii Economice din Moldova

Procesarea limbajului natural (PLN) și-a început dezvoltarea, în calitate de disciplină, imediat după cel de-al Doilea Război Mondial [9], ca un mecanism folosit pentru traducerea documentelor. Acesta a fost unul dintre primele obiective computaționale cele mai investigate. Dar eforturile premature depuse pentru analiza și modelarea limbajului uman s-au caracterizat printr-o aplicare a unei metode fără cunoștințe lingvistice și cu o performanță redusă a tehnicii de calcul din acele timpuri.

Potrivit lui Covington, "*Procesarea limbajului natural constă în utilizarea calculatoarelor pentru a înțelege limbile (naturale) umane, cum ar fi engleza, franceza sau japoneza. Prin înțelegere, nu se presupune că computerul are gânduri, sentimente sau cunoștințe asemenea omului, dar se subînțelege că computerul poate să recunoască și să utilizeze informații exprimate în limbajul uman*" [4].

2. NIVELURI DE CUNOSTINȚE ÎN PROCESAREA LIMBAJULUI NATURAL

Manaris și Slator definesc un sistem de PLN ca pe unul care încapsulează un model al limbajului natural în algoritmi adecvați și eficienți. În acest caz, tehnicile de modelare sunt larg asociate cu evenimente din multe alte domenii, printre care [10]:

- Informatica, care prevede metode de reprezentare a modelelor, proiectarea și implementarea algoritmilor pentru instrumentele software.
- Lingvistica, care contribuie cu noi modele și procese lingvistice.
- Matematica, care identifică modele formale și metodele.
- Neuroștiința, care explorează mecanismele mentale și alte activități creierului.

Dintre aceste domenii, lingvistica a oferit cunoștințele lingvistice despre limbile naturale. Această cunoaștere, în cadrul unui sistem de PLN,

poate fi împărțită în niveluri definite în termeni de caracteristici declarative (ce) și procedurale (cum), precum se arată în Tabelul 1. [10]. După cum se poate observa în acest tabel, cunoștințele lingvistice pot fi aranjate la diferite niveluri sau componente, deoarece structura oricărui limbaj uman se poate diviza natural între aceste niveluri [4].

Fonologic. Nivelul studiază sunetele limbii din punctul de vedere al valorii lor funcționale, stabilind inventarul de foneme ale unei limbi și caracterul diferitelor variante ale acestora. Fiecare limbă are un alfabet de sunete care se disting, acestea numindu-se foneme. Obiectul fonologiei îl constituie, așadar, sunetele ca realitate lingvistică, iar nu fizică sau fiziologică. Sunetele reprezintă materia sonoră. Astfel, nivelul fonologic tratează realizările acustice, de aceea, apar doar în sistemele de recunoaștere a vorbirii. Din punct de vedere tehnologic, prelucrarea vorbirii de către calculator este oarecum separată de restul prelucrării limbajului natural, dat fiind faptul că acest tratament al vorbirii este condiționat de analiza formei de undă a sunetului și de recunoaștere a formelor, în timp ce restul nivelurilor depind de o programare simbolică și un raționament automat.

Morfologic. Morfologia este ramura lingvisticii, care studiază regulile privind structura internă a cuvântului, adică regulile de combinare a morfemelor lexicale și gramaticale în cuvinte, stabilirea paradigmatelor lor în dependență de categoriile de gen, număr, caz etc. sau cuprinde regulile privitoare la modificările formale ale lor în diferite întrebuințări. Ideea generală se bazează pe faptul că morfemele individuale pot fi combinate pentru a forma cuvinte.

Morfemele sunt unitățile minimale de sens. Există două tipuri de morfeme: forma liberă, care poate să apară ca și cuvinte separate și forma legată, care nu poate apărea ca și cuvinte în sine. Acestea din urmă sunt, de obicei, numite afixe. De exemplu, cuvântul englezesc "*unselfish*" este compus din trei morfeme, "*un*", "*self*" și "*ish*". "*Self*" este o formă liberă, în timp ce "*un*" și "*ish*" morfeme legate. În special, "*un*" este aici un prefix, "*ish*" este un sufix și "*self*" este rădăcina.

Există trei procese morfologice principale utilizate în formarea cuvintelor:

• **Flexiunea.** Morfologia flexionară este preocupată de crearea cuvintelor noi, lăsându-le în aceeași categorie sintactică, dar schimbând relațiile gramaticale, cum ar fi, de exemplu, plural, timpul trecut și posesia. Cuvintele se formează cu ajutorul afixelor de inflexiune, care sunt conectate la morfemul liber. De exemplu, atât „pom”, ca și „pomi” (morfemul liber "pom" are afixul de plural "i") sunt substantive.

• **Derivarea.** Morfologia derivațională descrie modul în care sunt create cuvinte noi cu ajutorul unor afixe, trecându-le dintr-o categorie sintactică în alta. De exemplu, adjectivul „național” se derivă din substantivul „națiune”.

• **Compoziția.** Compoziția se ocupă cu construcția cuvintelor noi prin combinarea morfemelor libere.

Tabelul 1. Niveluri de cunoștințe în prelucrarea limbajului natural.

Nivel	Caracteristici	
	Declarativ (ce)	Procedural (cum)
Fonologic	Sunete vorbite	Formarea morfemelor
Morfologic	Unități de cuvinte, cuvinte	Formarea cuvintelor, Derivarea unităților de sens.
Sintactic	Rolul structural al cuvintelor (sau colecții de cuvinte)	Formarea frazelor
Semantică	Semnificația independentă de context	Derivarea semnificației frazelor
Discurs	Rolul structural al frazelor (sau colecții de fraze)	Formarea dialogurilor
Pragmatic	Semnificația dependentă de context	Derivarea semnificației frazelor ținând cont de discursul ambiant

Sintactic. Sintaxa, sau construcția propozițiilor, reprezintă nivelul cel mai de jos la care limbajul uman este, în mod constant, în proces de creare. Vorbitorii unei limbi creează mult mai rar unități fonice și lexicale. În schimb, sunt concepute în mod constant noi propoziții și fraze. Noam Chomsky (1957) a fost primul care a evidențiat acest aspect. El a introdus „gramatica generativă”, în care frazele sunt descrise de anumite reguli, și nu prin listarea lor și a structurilor în mod direct. Astfel de reguli au devenit standard nu numai în lingvistică, ci și în informatică, cu precădere în proiectarea compilatoarelor.

Cunoștințele sintactice reprezintă o componentă de bază a oricărui sistem de PLN, care este responsabilă pentru recunoașterea frazelor gramaticale și atribuirea unor structuri acestora. Procesul de recunoaștere a structurii unei fraze de către calculator se numește analiza sintactică computațională sau „parsare”.

Semantic. Semantica este o ramură a lingvisticii, al cărei obiect de studiu este *sensul*, unitate greu de abordat dintr-o perspectivă unică și unitară. La acest nivel, semantica computațională

Granița dintre flexiune (care furnizează diferitele forme ale unui cuvânt) și derivare (care produce cuvinte noi pornind de la cele existente) este, uneori, neclară. O diferență esențială o constituie faptul că numai derivarea poate introduce o schimbare de sens (prin introducerea de cuvinte noi). O altă deosebire constă în faptul că formele derivative ar putea să nu existe, în timp ce formele flexionare nu lipsesc aproape niciodată.

Indiferent dacă avem de-a face cu morfologia flexionară sau cu cea derivațională, putem spune că, din punct de vedere computațional, nivelul morfologic al limbii se ocupă de modul în care sunt alcătuite cuvintele pornindu-se de la unitățile de bază numite *morfeme*.

face un studiu al *sensului independent de context*. Cu alte cuvinte, interesează sensul pe care o propoziție îl are fără legătura cu contextul în care ea a fost utilizată. Semantica frazelor este o parte esențială a oricărui sistem, pentru că fără ea nu am putea atribui semnificație structurilor analizate.

Discursul tratează aspectele de interpretare afectate de frazele pronunțate anterior. La acest nivel, se acumulează cunoștințele care se referă la legarea sensului frazei izolate pentru a se integra în unități mai mari. În particular, această cunoaștere este folosită pentru a interpreta pronumele anaforic, a soluționa *elipsele* și a interpreta aspectele temporale.

Anafora, în limbajul natural, constă dintr-o expresie, care se referă la o expresie anterioară a unui discurs. În general, se folosește un prenume pentru a se referi la persoane, locuri sau lucruri menționate mai sus. De exemplu: "Pasărea a murit. Ea era foarte bătrână". În propoziția a doua „ea” se refera la pasăre. Pe de altă parte, elipsa: se referă la situații ale căror fraze sunt scurtate sau este eliminat vreun constituent, lăsând o parte din ele să fie înțelese din context. De exemplu, când cineva

este întrebat "*Care este numele tău?*" și se răspunde "*Ion Potcoavă*", ultima este o formă eliptică a propoziției "*Numele meu este Ion Potcoavă*".

Această componentă este necesară sistemelor, pentru ca acestea să posedă cunoștințe din contextul comunicativ, în care sunt produse mesajele și să țină cont de aspectele pragmatice, precum intențiile expeditorului și destinatarului.

Pragmatic. Se referă la utilizarea limbii în context. În general, pragmatica include aspecte ale cunoștințelor conceptuale ale lumii, care merg mai departe de condițiile reale ale fiecărei propoziții. Această cunoaștere este considerată atunci când se comunică într-o singură limbă.

Sunt folosite, pentru a înțelege, o mulțime de informații, subînțelese, dar nu și exprimate în mod explicit în propoziții. În timp ce sintaxa și semantica studiază propozițiile, pragmatica studiază "*acțiunile discursului*" și situațiile în care limba este utilizată.

Multe cuvinte și fraze pot fi ambigue și să aibă mai mult de un sens, semnificația lor poate fi falsă sau să producă implicații false. Semnificația depinde de principiile pe care oamenii le folosesc atunci când vorbesc (de exemplu, să fie relevante și se pune accentul pe fraze adevărate [6]).

În acest sens, pragmatica are două concepte importante: implicarea și presupoziția propozițiilor. Implicarea unei propoziții conține informații care nu fac parte din sensul acesteia, dar trebuie să fie deduse de către un ascultător rezonabil. Presuposițiile unei propoziții sunt lucrurile care trebuie să fie adevărate pentru ca propoziția să fie adevărată sau falsă. Adică, în baza presuposițiilor (ipotezelor) (cunoașterea adevărată a unui domeniu), oamenii interpretează fraze și derivă cunoștințe (implicații), care pot fi sau nu adevărate.

2. PROCESAREA LIMBAJULUI NATURAL ȘI LINGVISTICA COMPUTAȚIONALĂ

În procesarea limbajului natural nivelul pragmatic tratează folosirea propozițiilor în diverse situații (contexte), precum și modul în care contextul influențează interpretarea unei propoziții.

Cunoștințele lingvistice sunt încorporate în sistemul de PLN începând cu anii șaizeci, și a devenit una dintre componentele sale majore. Pornind din acel moment a fost definită aria de cunoștințe numit Lingvistica Computațională.

2.1. Lingvistica computațională

Potrivit cercetătorilor în domeniu lingvistica computațională este o disciplină care se bazează pe două lucruri: limbile naturale și calculatoare. Multe

direcții de cercetare împărtășesc ambele obiective, dar din perspective diferite. Ca întotdeauna, orice obiect nou de studiu se confruntă cu definirea terminologiei științifice. Termenul lingvistica computațională este echivalent cu PLN și nu este echivalent cu Lingvistica informatică sau Ingineria lingvistică.

Lingvistica informatică: o disciplină care se referă la utilizarea calculatoarelor în legătură cu limbaje și limbi. Include toate tipurile de instrumente care asistă studierea limbilor străine și a lingvisticii. Lingvistica computațională este o parte a lingvisticii informatice.

Ingineria lingvistică se referă la potențialele aplicații comerciale care implică utilizarea noilor tehnologii. Include ediții electronice (dicționare, cărți), produse multimedia etc.

Conform lui Grishman, lingvistica computațională poate fi definită ca "*studiul sistemelor informatice utilizate pentru înțelegerea și generarea limbilor naturale*" [7]. Allen oferă o definiție echivalentă pentru procesarea limbajului natural "*Scopul acestei cercetări constă în crearea modelelor de calcul suficient de detaliate, care ar permite scrierea programelor care realizează diferite sarcini în ce privește limbajul natural*" [1]. Prin urmare, LC și PLN tratează același lucru: dezvoltarea programelor de calculator care simulează capacitatea lingvistică umană.

Inteligența Artificială (IA) este responsabilă de codificarea într-un program a facultăților cognitive cum ar fi inferența, luarea deciziilor, achiziția cunoștințelor etc. În acest sens, LC este parte integrantă a IA, în același mod, cum, pentru mulți lingviști, Lingvistica mai face parte din psihologie pentru tratarea unei dintre capacitățile cognitive, prin excelență, limba.

În continuare, termenii PLN și LC sunt interșanjabili. Cu toate acestea, termenul PLN apare mai des, deoarece este mai bine înțeles.

Lingvistica computațională are mai multe aplicații practice, principalele fiind prezentate în clasificarea următoare:

1. Sisteme care încearcă să emuleze capacitatea omului de a procesa limbile naturale. În cadrul acestui grup, cele mai importante sunt: traducerea automată, recuperarea și extragerea informațiilor, interfețe om-mașină.

2. Sisteme care ajută la îndeplinirea sarcinilor lingvistice. Acest grup este format din instrumente care pot fi utilizate de către lingviști pentru facilitarea executării anumitor sarcini complexe. Unele aplicații, de acest tip, sunt: instrumentele de analiză textuală, bazele de date lexicografice, instrumentele de gestionare a corpusului. Corpusul

reprezintă o colecție de date lingvistice, de obicei, formată din mai multe texte. Această cantitate mare de texte, în limbaj natural, sunt utilizate pentru acumularea statisticilor necesare pentru analiza limbii.

3. Programe pentru ajutorarea la scrierea și compoziția textului. Aplicațiile incluse în acest grup au fost foarte dezvoltate și orice utilizator, obișnuit cu un procesor de text, este familiarizat cu ele: corectoare ortografice, corectoare sintactice și de stil.

4. Instruirea asistată de calculator. Acesta este un câmp de aplicație în continuă expansiune și are mai multe aspecte. Cel mai important este programul educațional pentru învățarea limbilor străine.

2.2. Limba din punct de vedere științific

Definiția limbii, din punct de vedere al științei, a determinat mulți lingviști să convină asupra diferitelor puncte. Covington le prezintă pe cele mai importante [4]:

- **Limba este formă, nu substanță.** Acest lucru înseamnă că limba nu este un set de pronunții sau de comportament, ci este un sistem de reguli care determină comportamentul. Un alt mod de a exprima acest lucru este de a distinge între competența vorbitorului (sistemul) și performanța (comportamentul observabil). Această distincție recunoaște că declarațiile accidentale, frazele întrerupte etc. nu sunt, realmente, instanțe de limba vorbită de o persoană, ci sunt derivații de limbă.

- **Limba este arbitrară.** O limbă constituie un set de simboluri, pe care oamenii sunt de acord să le utilizeze într-o manieră specifică.

- **Toate limbile umane utilizează modele de dualitate,** în care cuvintele sunt șiruri de sunete, iar propozițiile sunt șiruri de cuvinte. Cuvintele au o semnificație, sunetele, în sine, nu.

- **Toate limbile sunt aproape la fel de complicate, cu excepția dimensiunii vocabularului.** Limbile se schimbă în mod constant, dar fiecare schimbare este lentă. O limbă evoluează într-o anumită direcție sute de ani.

- **Toată lumea vorbește despre propria lor limbă.** Limba română vorbită de o persoană nu este complet identică cu limba română vorbită de tatăl ei. Acest lucru se datorează modului în care limba este învățată. În procesul de învățare a limbii, apar mici diferențe între indivizi, și deosebiri mari, inevitabile între grupuri sociale.

2.3. Probleme în utilizarea limbajului natural

Cunoașterea lumii este un factor important în sistemele de PLN. Prin urmare, un sistem de PLN ar trebui să impună limite cu privire la necesitatea de cunoaștere externă și experiența umană. Covington afirmă, în plus față de cele de mai sus, că PLN depinde de doi factori: primul se referă la puterea calculatoarelor. Apariția, în 1980, a microcalculatoarelor a schimbat situația. Anterior, PLN a fost atât de scumpă încât oamenii acceptau orice rezultat perfect, oricare ar fi fost atins. Această situație s-a schimbat și sistemele PLN, încă imperfecte, au devenit mai ieftine, iar utilizatorii găsesc aplicări bune pentru ele.

Al doilea factor, și poate cel mai important, constă în faptul că PLN depinde de cunoașterea exactă a modului în care limbajul uman funcționează, lucru care, acum, nu se cunoaște suficient. Până în ultimii ani, limba a fost studiată cvasi-exclusiv cu scopul de a o preda altei persoane. Principiul care stă la baza tuturor limbilor umane a era ignorat. Mai mult decât atât, știința lingvistică are doar câteva decenii vechime, și nici nu există încă un consens cu privire la unele aspecte de bază. Sistemele PLN trebuie să abordeze o varietate de probleme în ce privește limbajul natural [10]:

- **Inexactitatea,** inclusiv erorile de ortografie, punctuația incorectă, cuvintele transpuse și propozițiile negramaticale.

- **Incompletitudinea,** inclusiv în construcții eliptice, anafora etc.

- **Imprecizii,** inclusiv utilizarea de termeni relativi, fără un anumit punct de referință și cu utilizarea de termeni calitativi.

- **Ambiguitatea,** deoarece multe interpretări pot apărea la orice nivel de cunoștințe lingvistice (a se vedea tabelul 1). Ambiguitatea poate fi rezolvată folosind cunoștințe de un nivel mai înalt.

2.4. Modele simbolice de procesare a limbajului natural

Modelele și metodele PLN pot fi clasificate în: metode simbolice, empirice sau statistice, conexiuniste și abordări hibride. Primele două sunt numite modele matematice ale limbilor. Abordarea simbolică se bazează pe cunoaștere, utilizează reguli și algoritmi care funcționează pe structuri de date simbolice și reprezintă cunoașterea limbajului natural. Abordarea empirică sau statistică implică colecții de mostre selective de limbă (corpus), care sunt etichetate și folosite pentru a crea modele statistice pentru PLN. Tehnica conexiunistă

utilizează rețelele neuronale pentru reprezentarea cunoștințelor lingvistice. Pe de altă parte, tehnicile hibride combină una sau mai multe din modelele anterioare, pentru a completa beneficiile fiecăruia și a rezolva problemele domeniilor și aplicațiilor specifice.

Sistemele simbolice se bazează pe manipularea de simboluri. Ele au fost concepute de matematicieni pentru a capta într-o demonstrație riguroasă și sistematică a teoremelor matematice și logice. În lingvistică, Chomsky a fost primul care a introdus sistematic paradigma logicii formale.

De obicei, regulile de inferență, într-un sistem formal, permit să se concentreze pe sintaxa modelului, indiferent de interpretarea lui. Mulți lingviști cred că limba are o natură bazată pe reguli sau logică și că este ceea ce se încearcă să se reflecte în gramaticile formale. În general, aceste gramatici s-au dovedit eficiente în descrierea și explicarea fenomenelor legate de competență. Competența se referă la cunoștințele pe care fiecare vorbitor le are despre limba lui maternă.

3. GRAMATICI FORMALE

O gramatică formală este o specificație riguroasă și explicită a structurii unei limbi. Acest lucru este redat cu ajutorul unui formalism gramatical, adică cu o limbă artificială creată pentru descrierea limbilor naturale. Utilizarea lor se datorează faptului că un limbaj bine definit, riguros, facilitează evaluarea ipotezelor și permite elaborarea previziunilor.

Există diferite tipuri de gramatici, care sunt bine formalizate.

Acestea includ: gramatici generative, gramatici categoriale, gramaticilor de dependență,

gramatici de lanțuri lingvistice Harris și gramatici cu arbori adiacenți [7,14]. Cu toate acestea, cele mai răspândite sunt gramaticile generative, de asemenea, cunoscute sub numele de gramatici de structură a frazei sau sintagmatice propuse de Chomsky [3].

Gramaticile generative sunt constituite dintr-un set de reguli generative, care atribuie, în mod explicit, structura internă a propozițiilor. Aceste reguli, numite reguli de rescriere, operează pe mulțimi de elemente neterminale și terminale. Atât gramaticile transformazionale, cât și gramaticile de unificare sunt gramatici generative. Bach, în 1974, spunea că orice gramatică, care definește, în mod explicit și precis, frazele unei limbi este o gramatică generativă [2]. Ele sunt cele mai răspândite în lingvistică computațională.

Chomsky a propus o clasificare a tipurilor de gramatici, aplicată la gramaticile generative sau sintagmatice, care a devenit celebră cu numele creatorului său, Ierarhia Chomsky. Această ierarhie este organizată în conformitate cu "puterea generativă slabă". Conceptul de putere generativă sau formală se referă la capacitatea de predicție a unei gramatici. În special, preocupările puterii generative se referă la tipul de fraze ale gramaticii, care pot fi recunoscute drept gramaticale.

Există patru tipuri de gramatici generative, fiecare definit de un anumită clasă de reguli, pe care le conține.

Tabelul 2 rezumă fiecare din cele patru tipuri de gramatici ale lui Chomsky, clasele de limbaje pe care le generează, tipul de automat, care le recunoaște și forma regulilor de producție. În tabel, A reprezintă un simbol neterminal și α , β și χ - șiruri de terminali și neterminali, iar t - un simbol terminal. Șirurile α și β pot fi vide, dar χ - nu.

Tabelul 2. Ierarhia Chomsky

Tip	Gramatici	Restricții asupra formei regulilor	Limbaje	Automate
0	Nerestricționate	Nicio restricție	Recursiv enumerabile	Mașinile Turing
1	Dependente de context	Partea dreaptă conține cel puțin simbolurile din partea stângă: $\alpha A \beta \rightarrow \alpha \chi \beta$	Dependente de context	Automatele linear limitate
2	Independente de context	Partea stângă poate avea doar un simbol neterminal: $A \rightarrow \chi$	Independente de context	Automatele cu stivă
3	Regulate sau cu stări finite	Regulile pot avea doar ultimele două forme: $A \rightarrow tB$ și $A \rightarrow t$ sau $A \rightarrow Bt$ și $A \rightarrow t$	Regulate	Automatele finite

Această clasificare este una teoretică, deoarece nu există tipuri pure de gramatici. În practică, gramaticile formale sunt modificate în funcție de nevoile specifice. În consecință, nu se poate ușor decide cărui tip

aparține gramatica. Aici indică distanța dintre realizările teoretice și practice, care este similar cu diferențele dintre lingviștii teoreticieni și lingviștii computaționali.

Lingvistică teoretică: se axează pe analiza competenței vorbitorilor, utilizează, în principal, introspecția pentru a obține date și are tendința de a ajunge la concluzii prin metode deductive. Obiectivele sale principale sunt de a obține o teorie gramaticală simplă, restrictivă, ținând cont de universalitățile lingvistice.

În lingvistica computațională, s-au obținut anumite rezultate și aplicații cu următoarele gramatici.

3.1. Gramatici regulate sau de stări finite

Gramaticile regulate sau, de asemenea, numite rețele de tranziție, sunt formate din noduri sau stări (reprezentate prin circumferințe) și arce (reprezentate prin săgeți) etichetate. Fiecare arc reprezintă o tranziție între două stări. Există două tipuri speciale de stări: Stările inițiale (marcate cu o mică săgeată de intrare) sunt singurele care nu primesc alte săgeți de arcuri și Stările finale (reprezentate cu o circumferință dublă), care sunt singurele, de la care nu pornesc tranziții la alte stări. Nu este posibilă tratarea recursivității cu o gramatică regulată, deoarece datele prelucrate sunt doar cele redade de starea în care se află.

Această gramatică de stări finite este aplicată la morfologie și recunoașterea vocabularului, lexicului, deoarece, în orice limbă, regulile de flexiune formează o mulțime aproape închisă și mult mai mică decât regulile de sintaxă. În definitiv, au fost folosite în mai multe scopuri limitate ale

limbajului, oferind o metodă foarte eficientă pentru calculator [11].

3.2. Gramatici independente de context

Termenul de Gramatică independentă de context (GIC) este un model specific propus pentru descrierea sintaxei unei limbi. Gramaticile pot fi folosite pentru descrierea oricărei limbi (naturale sau artificiale) și acestea trebuie să respecte niște constrângeri foarte simple. Constrângerile se referă la modul în care sunt generate clauzele. Strict vorbind, clauzele dependente trebuie să fie adiacente componentei de care depind. Majoritatea limbilor naturale par să urmeze acest comportament, cu o posibilă excepție idioma germană-elvețiană [5]. Mai multe teorii contemporane ale sintaxei limbajului natural sunt derivate ale schemei GIC.

Precum se vede în ierarhia Chomsky, GIC este o gramatică de tipul 2, care este formată dintr-un set de reguli (producții) și un set de intrări lexicale (sau lexicul). Cu aceste reguli, se poate descrie structura sintactică a mai multor fenomene ale limbilor naturale. Aceste gramatici sunt capabile să recunoască și să genereze propoziții.

Astfel, o GIC este o descriere formală a sintaxei unei limbi. În mod concret, descrierea este dată de setul de "reguli de producție", care definesc propozițiile bine formate ale limbii. Regulile, desigur, sunt ele însele, de asemenea, scrise într-un limbaj formal. Ca toate idiomele formale oficiale, acesta este definit de un vocabular și sintaxă. Tabelul 3 relevă definiția mai strictă a GIC.

Tabelul 3. Gramatici independente de context.

Vocabularul	Regulile conțin trei tipuri de simboluri: <ul style="list-style-type: none"> • Neterminale, care corespund componentelor limbajului descris. Unul din simbolurile neterminale are o poziție specială. Acesta se numește simbol distinctiv. • Terminale, care corespund cuvintelor limbajului descris. • \rightarrow - simbolul săgeata, care delimitează partea stângă a unei reguli de partea ei dreaptă.
Sintaxa	Propozițiile în limbajul GIC sunt reguli de producție (Sintaxa se referă aici la formatul regulilor înseși, nu la limba pe care o descriu). O regulă de producție are următoarele proprietăți: <ol style="list-style-type: none"> 1. Acesta constă dintr-o parte stângă (LHS) și o parte dreaptă (RHS), separate de o săgeată: LHS \rightarrow RHS 2. LHS constă dintr-un singur simbol neterminal. 3. RHS constă din unul sau mai multe simboluri neterminale sau terminale.

Aceste gramatici oferă propozițiilor o structura ierarhică internă. Pot fi descrise exprimarea alternanței și a opționalității. În afară de aceasta, gramaticile independente de context au proprietăți formale care facilitează proiectarea algoritmilor de parsare. Cu toate

construcțiile recursive, care nu pot fi tratate cu gramaticile regulate. GIC permit, de asemenea, acestea, există probleme cu tratamentul anumitor fenomene lingvistice, precum ar fi constituenții discontinui, subcategorizarea și concordanța.

Aici trebuie menționat că componentele discontinue sunt componentele care pot fi găsite în mai mult de o poziție structurală, iar subcategorizarea, în principiu, fiind un fenomen în care o structură a frazei lexico-semantice este prezisă în funcție de semantica verbală. Ea are importanță în sintaxă, deoarece se specifică combinațiile posibile de cuvinte. Verbele și unele adjective admit o structură de complemente, iar subcategorizarea se referă la numărul și categoriile complementelor fiecărui verb. În același timp, concordanța este un fenomen care are loc în mai multe limbi, în care cuvintele iau anumite flexiuni, în funcție de modul în care acestea se referă la alte cuvinte într-o propoziție. Această listă se referă la gen, număr. Un exemplu simplu are loc cu verbele la persoana a treia singular și subiecții săi: "*ea dansează*", "*tu dansezi*".

Din aceste motive, în practică, niciun sistem PLN, de o anumită acoperire, nu folosește versiunea pură a acestui tip de gramatici.

3.3. Gramatici de unificare cu caracteristici

Cele mai frecvente restricții care apar în gramaticile independente de context sunt reprezentate de fenomenele de concordanță și subcategorizare. Gramaticile de unificare cu caracteristici tratează ambele cazuri [8], astfel, încât acestea sunt considerate ca un model de calcul mai cuprinzător și restrictiv, în același timp, din cele cunoscute până în prezent.

Aceste gramatici sunt redată prin descrieri formale complexe cu utilizarea de caracteristici și folosesc o operațiune generală de combinare și verificare a informațiilor gramaticale, cunoscută sub numele de unificare [12].

Structura caracteristicilor este mecanismul de bază de reprezentare a informației despre unitățile lingvistice. Astfel, fiecare element de informații (unitate lingvistică) este asociat cu o structură caracteristică, unde o caracteristică reprezintă o pereche formată dintr-un atribut și o valoare. Atributul poartă un nume care identifică caracteristica, de exemplu, "*numărul = plural*", unde număr este atribut, iar plural - valoare. Valorile caracteristicilor pot fi valori complexe, și care, la rândul lor, pot fi structuri caracteristice.

Informațiile conținute într-o structură caracteristică este combinată într-o nouă structură, prin operația de unificare. Pentru ca acest lucru să se întâmple, structurile de informații trebuie să fie compatibile, deoarece, în caz contrar, nu pot fi unificate. Compatibilitatea ține de natura caracteristicii și valorii acesteia. Caracteristicile, care apar numai într-una din structurile unificate, sunt incluse în structura rezultatului unificării, reușind să combine informații comune și diferite. Acest fapt permite ca diferitele structuri informaționale să poată fi combinate coerent.

Astfel, tendința lingviștilor de a utiliza gramatici mai restrictive, a condus la faptul că gramaticile independente de context s-au extins cu caracteristici, folosind mecanismul de unificare [12].

3.4. Structura unui sistem PLN simbolic

Se consideră că orice software de PLN are două tipuri mari de cunoștințe stocate:

1. **Cunoștințele lingvistice** în formă de gramatică, vocabular și model conceptual al lumii. Gramatica este, pur și simplu, o definiție abstractă a unui set de elemente structurate și bine formate. Ar fi echivalentă cu competența noastră lingvistică și pragmatică.

2. **Programul sau parserul**, care conține instrucțiunile pentru procesarea datelor lingvistice. Parser-ul este un algoritm sau un set de instrucțiuni, care leagă secvențele de simboluri cu cunoștințele lingvistice stocate. Parserul este un mecanism de calcul, care a deduce structura secvențelor de cuvinte pornind de la cunoștințele stocate în gramatică și dicționar și stabilește dacă șirurile de caractere sunt, din punct de vedere gramatical, corecte sau incorecte. Problema pe care trebuie să o rezolve parserul este una pur sintactică: recunoașterea frazelor gramaticale și atribuirea unei structuri. Alte componente sunt responsabile pentru interpretare. Figura 1 prezintă structura generală a procesului de analiză [13].

Algoritmii de parsare sunt responsabili pentru determinarea regulilor ce trebuie aplicate și în ce ordine. Fiecare algoritm îmbină, de obicei, diferiți parametri și diferite structuri de lucru. Există mulți algoritmi, dar toți se bazează pe o

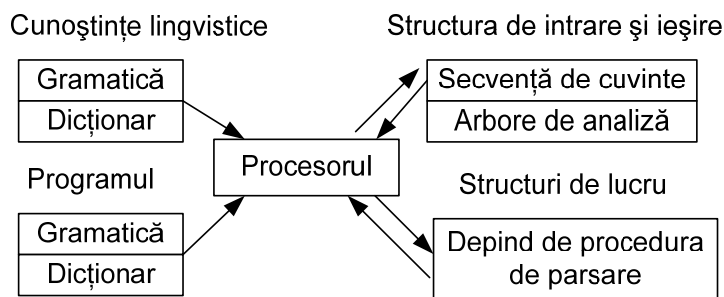


Figura 1. Structura generală a procesului de parsing

combinație de trei parametri esențiali care sunt luați în considerare: analiza descendentă (top-down) / analiza ascendentă (bottom-up), procesarea secvențială / procesarea paralelă, procesarea deterministă / procesarea nedeterministă. Având în vedere acești parametri, pot fi menționași câțiva algoritmi de parsing: (1) Algoritmul descendent în serie cu backtracking [7] și (2) Algoritmii cu Chart [1,14]

Modelele simbolice reprezintă paradigma predominantă în LC, repertoriul lor de concepte și metode este mai amplu și au fost aplicate în multe probleme și limbi. Dintre acestea, cele mai utilizate sunt automatele finite (pentru simplitate și eficiență de procesare) și gramaticile independente de context, completate cu gramatici de unificare cu caracteristici (pentru puterea expresivă de a ține cont de fenomenele lingvistice).

Bibliografie

1. **Allen, J.** *Natural Language Understanding*. Addison-Wesley, 2 edition 654 pages, 1994.
2. **Bach, E.** *Syntactic Theory*. Univ Prof Amer, 310 pages, February 1982.
3. **Chomsky, N.** *Syntactic structures*. De Gruyter; 2nd edition, 117 pages, 2002.
4. **Covington, Michael A.** *Natural Language Processing for Prolog Programmers*. Englewood Cliffs, NJ: Prentice Hall, 348 pages, 1994.
5. **Gazdar, G. Mellish, C.** *Natural Language Processing in PROLOG. An Introduction to Computational Linguistics*, Addison Wesley, 504 pages, 1989.
6. **Grice, H., P.** *Logic and conversation*. In Cole, P. and J.L. Morgan, eds. *Syntax and semantics*. vol. 3, *Speech acts*. NY: Academic Press, p. 41-58, 1975.
7. **Grishman, R.** *Computational Linguistics: an introduction*. Cambridge, Cambridge University Press, 193 pages, 1986.

8. **Kay, M.** *Parsing in functional unification grammar*. In Dowty, D., Karttunen, L. and Zwicky, A. pp 251-278, 1985.

9. **Locke W.N. and Booth A.D.** *Machine Translation of Languages*. Technology Press of MIT and Wiley, Cambridge, Mass., p.15-23, 1955.

10. **Manaris, Bill Z. and Sator, Brian M.** *Interactive Natural Language Processing: Building on Success*. IEEE Computer, vol.29, Nr.7, p.28-32, 1996.

11. **Roche, E. and Schabes, Y.** *Finite-State Language Processing*. Cambridge, The M.I.T. Press, 482 pages, 1997.

12. **Shieber, S.** *An Introduction to unification-based approaches to grammar*. Chicago, Chicago University Press, 120 pages, 2001.

13. **Winograd, T.** *Language as a Cognitive Process: Syntax*. Reading, Addison-Wesley, 654 pages, 1983.

14. **Winograd, T.** *Understanding Natural Language*. New York: Academic Press, 191 pages, 1972.

Recomandat spre publicare: 16.05.2013.