

# Instrumente de Procesare a Limbajului Natural Pentru Studiarea Limbilor Străine

Bobicev V., Carcea L.  
Catedra Informatica Aplicată  
Universitatea Tehnică a Moldovei  
Cișinău, Moldova  
victoria.bobicev@gmail.com, carcea@mail.utm.md

**Abstract** — This paper presents an online concordancer that was created at the Applied Informatics department for teaching and learning of French language. It also contains examples of possible activities for teachers and students that can be used in the educational process.

Besides tools specifically developed for learning and based on classical theories of language teaching, information technology offer other fundamentally new resources and tools that can be used in the educational process.

Corpora are the most important resources that are not yet fully utilized in language teaching.

The most difficult problem is the effective use of corpora. The contribution of TAL (Traitement Automatique des Langues / Natural Language Processing) is necessary for the development of tools and methods which simplify use of corpora. The concordancer is one of such tools. It can effectively process large amounts of text and reveal the necessary elements: keywords, phrases or grammatical constructions.

**Termeni cheie** — Concordansier, predarea limbii, studierea limbii, corpusuri de texte, procesarea limbajului natural.

## I. INTRODUCERE

Tehnologiile informaționale oferă noi posibilități, care nu sunt pedepplin aplicate în domeniul educațional. Vom menționa diferite aspecte de utilizare oferite de învățământul la distanță și tehnologiile informației și a comunicației:

- posibilitatea de a obține în orice loc sau orice moment cunoștințe actuale și moderne;
- studii în regim autonom cu materiale în format electronic la un calculator personal, sau un dispozitiv mobil;
- posibilitatea de a dezvolta resurse Web educaționale;
- crearea unei comunități de utilizatori distribuiți în rețele sociale, ceea ce conduce la activități educaționale globale;
- posibilitatea de a interacționa la distanță, de a obține ajutor sau a fi evaluat;

Furnizorul principal de informații și sfaturi sub toate aspectele de utilizare a tehnologiei de învățare a limbilor este EUROCALL [2]. Această asociație profesională a fost oficializată în 1993. Trei direcții de dezvoltare au fost identificate pentru EUROCALL:

- utilizarea calculatoarelor în sala de studii;
- pregătirea profesorilor;
- dezvoltarea și evaluarea produselor soft;

## II. PROCESAREA AUTOMATĂ A LIMBAJULUI NATURAL

Procesarea automată a limbajului natural (PALN) este o disciplină plasată la frontierele lingvisticii, informaticii și a inteligenței artificiale. Ea ține de aplicațiile produselor soft și tehnicilor informatice în toate aspectele limbajului natural (uman).

Să nu se confunde cu lingvistica informatică, care se bazează pe o modelare simbolică sau matematică a limbajului natural.

Primele cercetări au demarat la mijlocul secolului XX în scopul traducerii automate a textului. În continuare cercetătorii din inteligența artificială presupuneau că prelucrarea limbajului natural prin intermediul calculatoarelor este o problemă de scurt timp. Lipsa de rezultate a curmat acest punct de vedere și utilizarea algoritmilor euristici. Au fost necesare cercetări în domeniul formalizării gramaticelor limbajelor umane. Lucrările lui Chomsky [1] au condus la ipoteza de similitudine între limbajele naturale și cele formale informatice. Dar și gramaticile formale au întâmpinat numeroase obstacole în cazul limbajului uman. A urmat încercarea de a aplica modele bazate pe cunoștințe, care nu au permis obținerea rezultatelor rapide și universale. În consecință cercetătorii au acceptat dificultatea problemelor procesării limbajului natural și au devenit mai modești.

Totuși, creșterea vertiginoasă a circulației documentelor, analiza, prelucrarea și traducerea lor forțază căutarea soluțiilor viabile și imediate a problemei automatizării acestor procese. Apariția Internetului și utilizarea lui masivă provoacă urgent elaborarea instrumentelor ce ar putea prelucra informația de care au nevoie cercetătorii.

În aceste condiții a apărut pe scenă și încearcă să rezolve problemele existente, fără a intra în complexitatea funcționării limbajului natural procesarea automată a limbajului natural (PALN). În general se aplică metode statistice, care conduc la rezultate satisfăcătoare. De altfel aceste metode pot fi ușor realizate în formă de produs soft. Condiția principală a succesului este cantitatea majoră (aproape gigantică) necesară pentru crearea modelului statistic adecuat. Tocmai acest lucru a și devenit disponibil la momentul actual.

## III. CORPUSURI DE TEXTE

Corpusurile sunt în același timp materie primă pentru studiul limbilor și a textelor cât și mijloace pentru a testa modelele propuse. Un corpus este reprezentat printr-un număr

mare de materiale lingvistice: cărți, reviste, rapoarte la lucrări de laborator ale studenților, înregistrări audio și video, discuții, etc. Au apărut mai multe publicații ce țin de lingvistica corpusului înainte de era calculatorului, dar este problematic să găsești dovezi pornind de la corpus după ce a fost imprimat pe hârtie. Chomski a criticat corpusul ca mod de studiere a limbii. Mai mult ca atât această critică nu era unicul reproș. Argumentul principal era caracterul infinit al limbii umane. În toate lucrările de la debutul lingvisticii corpusului au fost susținute două ipoteze fundamentale, pentru moment imperfecte: frazele unui limbaj natural sunt limitate și pot fi colectate și numărate.

Însă numărul de fraze al unui limbaj natural este nu numai mare, el tinde spre infinit. Chomski a estimat că singura modalitate de elaborare a unei gramatici a unei limbi este posibilă prin descrierea regulilor sale și nu prin enumerarea frazelor sale.

Alți cercetători [8] afirmă corpusul este o metodă mult mai puternică din punct de vedere științific, fiindcă el este deschis unei verificări obiective a rezultatelor. Francis și Kucera [6] au început să lucreze asupra celebrului corpus Brown, care a durat circa două decenii. Pas cu pas corpusul a devenit lizibil și pentru calculator, ceea ce a simplificat considerabil manipularea și utilizarea sa. Lingvistica corpusului modern definește corpusul ca un ansamblu de texte care au patru caracteristici:

- el trebuie să fie lizibil pentru calculator;
- el trebuie să fie reprezentativ pentru domeniul elaborat;
- el trebuie să fie finit;
- toate textele trebuie să fie codificate și adnotate în același mod;

Majoritatea corpusurilor îndeplinesc aceste cerințe, dar unele sunt în creștere continuă. De exemplu, corpusurile ce conțin colecții de ediții periodice electronice cum ar fi reviste și ziare. Unele corpusuri sunt create în scopuri diverse. Corpusurile reprezentative au tendința de a reprezenta ansamblul limbii cu toată diversitatea sa particulară. Există corpusuri naționale pentru unele limbi, cum ar fi British National Corpus<sup>1</sup> or Russian National Corpus<sup>2</sup>.

Corpusurile speciale reprezintă fenomenele specifice ale unei limbi. De exemplu corpusul operelor lui Alexandre Dumas, sau corpusul corespondenței de serviciu.

Diverse metode de creare a corpusurilor au fost dezvoltate, dar mai este mult de lucru pentru a obține un corpus calitativ. Cele mai prețioase se consideră corpusurile marcate sau adnotate, dar marcajul calitativ necesită un lucru manual de profesioniști calificați. Adnotarea corpusului înseamnă adăugarea informațiilor lingvistice interpretative acestui corpus. De exemplu, o adnotare de tip comun este adăugarea etichetelor ce indică clasa morfologică a cuvintelor, altfel spus

marcarea categoriilor gramaticale și poate fi utilă pentru căutarea tuturor formelor unui cuvânt.

Alte tipuri de adnotare sunt: adnotarea lexicală, adnotarea fonetică, adnotarea semantică, adnotarea sentimentelor, etc.

Corpusurile sunt utilizate în lexicografie fiindcă exemplele pot fi ușor extrase și organizate pentru analiza ulterioară cu un produs soft. Ele sunt instrumente bune pentru cercetarea sintactică din motivul cuantificării reprezentând toată varietatea unei limbi. Avantajul principal este faptul că datele păstrate sunt în mare parte naturaliste. În acest mod corpusul constituie una din sursele fiabile de origine naturală ce poate fi examinată.

Totuși corpusurile de texte sunt resurse foarte importante care nu sunt pe deplin valorizate și nu sunt din plin utilizate în studierea și predarea limbilor.

#### IV. CONCORDANSIERE

Utilizarea eficientă a concordansierelor este dificilă motivul fiind volumul lor foarte mare. Pe de altă parte filologii nu sunt obișnuiți cu aceste surse de informații. Utilizarea eficientă a corpusurilor în lingvistică, filologie, lexicografie în studierea și predarea limbilor este posibilă cu instrumente și metode specifice. Concordansierul reprezintă unul din aceste instrumente. El permite scoaterea din joc a unor cantități mari de informație și vizualizarea elementelor necesare: cuvinte-cheie, fraze, sau construcții gramaticale.

Termenul 'concordanță' provine din engleză și semnifică o listă alfabetică de cuvinte dintr-o carte sau un set de cărți în care fiecare cuvânt poate fi găsit și deseori explicată utilizarea sa. Elaborarea concordansierelor în mod manual fără calculator a fost o problemă enormă și dificilă. Astăzi un calculator poate produce un concordansier în câteva minute.

Una din primele referințe la utilizarea concordansierelor electronice a fost [4], dar și alte exemple cu utilizare practică în sala de studii sunt descrise de [7]. În 1991 Tim Johns a ameliorat utilizarea concordansierului în clasa de studii cu conceptul "studierea pilotată de date" (Data Driven Learning - DDL) [5]. DDL încurajează elevii să lucreze în baza propriilor reguli vis-a-vis de sensul și utilizarea cuvintelor cu ajutorul unui concordansier pentru a găsi exemple într-un corpus de texte autentice. De asemenea este posibil ca profesorul poate să găsească locații tipice și să genereze exerciții bazate pe exemplele găsite.

În DDN procesul de studiere nu se bazează numai pe inițiativa profesorului de a veni cu subiecte, materiale metodice, însușirea explicită a regulilor, dar și pe descoperirea propriilor reguli, principii și modele de utilizare a limbii străine.

Principalul avantaj al concordansierului este rezumat prin fraza frecvent utilizată "cunoașteți cuvântul după vecinii săi" sau "spunemi cine sunt vecinii ca să-ți spun cine ești". Cu alte cuvinte vom găsi nu numai sensul unui cuvânt dar și utilizarea sa cu alte cuvinte în textul dat.

În acest mod o concordanță este un cuvânt numit cuvânt-cheie scos dintr-o lucrare în limbaj autentic (corpus) reprezentat cu părți contextuale în care el se întâlnește.

---

<sup>1</sup> <http://www.natcorp.ox.ac.uk/>

<sup>2</sup> <http://www.ruscorpora.ru/en/>

Cuvântul este recunoscut ca cuvânt-cheie în concordanță de conext sau KWIC-concordanță.

#### V. FUNCȚII TIPICE ALE UNUI CONCORDANSIER

Un concordansier are ca regulă mai multe funcții (instrumente) pentru predarea și studierea unei limbi străine:

- Lista de cuvinte. Generează lista cuvintelor din textul dat și le ordonează la afișare. Permite vizualizarea celor mai frecvente cuvinte în text și celor care se întâlnesc o singură dată.
- Lista cuvintelor-cheie. Acest instrument selectează cuvintele frecvente în textul dat cu algoritmi sofisticăți. Pentru început sunt eliminate cele mai frecvente cuvinte, care de obicei sunt prepoziții, articole, pronume, etc. Apoi se elimină cele mai rar întâlnite, care nu sunt specifice textului dat. În final se compară frecvența cuvintelor rămase cu frecvența lor în baza corpusului. Cuvintele care au frecvența mai mare în text de cât în corpusul de referință sunt considerate cuvinte-cheie.
- Clustere. Această funcție selectează clustere în baza unei condiții de căutare. De obicei se selectează expresii curente în textul dat.
- Concordanța cu un cuvânt și cu un fragment. Acest instrument va afișa rezultatele căutării în format KWIC (cuvinte-cheie în context). Concordanțele cu un fragment pot fi utilizate pentru căutarea cuvintelor care au aceiași rădăcină.
- Co-occurențe cu un cuvânt și cu un fragment.

Toate aceste instrumente sunt bazate pe statistică și produc un rezultat bun în cazul unui volum mare de texte. Există totuși și o problemă cu concordansierul on-line: Volumul textelor transmis prin Internet. Traficul insuficient poate provoca probleme pentru texte mari.

Concordansierul poate fi utilizat pentru predarea și studierea (învățarea) limbilor străine în diferite maniere:

- Profesorul poate utiliza concordansierul pentru a căuta exemple autentice, colocații tipice etc;
- Profesorul poate genera exerciții bazate pe exemple din diferite corpusuri;
- Elevii pot lucra cu reguli gramaticale de căutare a cuvintelor-cheie și a caracteristicilor lexicale elaborate de ei înșiși. În funcție de nivelul lor ei pot să formuleze întrebări ce țin de reguli în baza observațiilor modelelor limbii studiate.
- Elevii pot fi mult mai activi în învățarea vocabularului: ei pot fi invitați să descopere semnificații noi, să observe colocații obișnuite sau relațiile cuvintelor cu sintaxa, ba chiar să critice propunerile vocabularului;
- Elevii pot fi invitați să se pronunțe asupra utilizării limbii în general în baza propriilor analize a unui corpus de texte.

Există mai multe activități posibile în clasă, dar care trebuie pregătite de profesor cu concordansierul și textele respective.

Site-ul Gobuild Web sugerează generarea unui cuvânt-cheie apoi eliminarea lui. Elevii trebuie să găsească acest cuvânt, ceea ce îi va încuraja să determine cuvântul din context și nu din dicționar. Se generează o listă KWIC pentru un cuvânt particular și se va întreba care este prepoziția care precede sau cre urmează după acest cuvânt-cheie.

Îată câteva exemple de activități posibile cu concordansierul în clasă:

Activitatea 1: Determinați cuvântul misterios.

Scopul activității: Familiarizarea studenților cu concordanța KWIC și cu importanța conextelor din stânga din dreapta când se operează cu cuvinte cheie.

Foaie de lucru. Citiți grila din tabelul 1 unde cuvântul fără sens "gloop" a fost înscris în locul unui cuvânt corect. Sarcina este să decideți care este cuvântul corect.

Activitatea 2: Donc on peut dire que...

Scopul activității: conștientizarea elevilor cu deosebirile stilistice ale cuvântului "Donc" în limba franceză în texte scise formal față de cele neformale.

Foaie de lucru. În franceză puteți plasa un cuvânt ca "Donc" la începutul unei fraze (de exemplu "Donc on peut dire que..."). Dar puteți să-l întâlniți des între verbul principal și cel care urmează după verbul principal (de exemplu "On peut donc dire que ..."). Este unul mai bun ca altul? Unde trebuie inserat "Donc"? Pentru a cunoaște acest lucru analizați listele din tabelul 2.

Activitatea 3: Le roman s'agit d'un amour malheureux.

Scopul activității: a implica elevii în activitatea în care ei deduc regula utilizării expresiei «s'agit» care nu este utilizată nici o dată cu alt subiect decât cu «il» impersonal.

Foaie de lucru "Le roman (ou "le poème" ou "la pièce") s'agit d'un amour Malheureux ... " Profesorii de franceză vor corecta întotdeauna primele două cuvinte în astfel de fraze în eseurile studenților. Pentru a explica reacția cadrelor didactice analizați textele din tabelul 3 și decideți cum trebuie utilizată fraza «s'agit».

#### VI. CONCLUZII

Procesul de predare și învățare a limbilor străine asistate de calculator ne sugerează concluzia că calculatorul poate servi ca o varietate de utilizări pentru predarea limbilor. Calculatorul poate fi un tutore care propune exerciții ce țin de limba studiată; un stimulant pentru discuții și interacțiune; un instrument pentru a scrie și pentru cercetare. Cu apariția Internetului calculatorul poate deveni și un mijloc de comunicare mondială și o sursă de materiale autentice nelimitate.

Totuși după cum a fost subliniat în [3], utilizarea calculatorului nu este o metodă. Mai degrabă putem vorbi de un mediu în care o varietate de metode, de abordări de filosofii pedagogice pot fi lansate. Eficacitatea predării asistată de calculator nu poate să rezide în mediul său, dar numai în modul în care aduce profit.

TABEL I. EXEMPLE PENTRU EXERCİTIU 1.

1. pport critique sur certaines utilisations abusives de la	<b>gloop</b>	est devenu un geste banal plus qu’une décision.
2. que pour beaucoup d’entre nous le fait d’allumer une	<b>gloop</b>	.
3. laquelle on est pris pour gens qui “s’abrutissent à la	<b>gloop</b>	”, dans une proportion croissante depuis 1896
4. Tous les grands moments de	<b>gloop</b>	superposent un message recherché et un messa
5. sieurs postes et l’augmentation du temps de diffusion (	<b>gloop</b>	du matin et de la nuit).

TABEL II. EXEMPLE PENTRU EXERCİTIU 2.

Il devient	donc	difficile de proposer un plan de gestion de ces troupeaux.
Elle permet	donc	la circulation des avions de grandes dimensions en cas de besoin.
Ce droit exclusif n’atténue	donc	pas du tout les droits de ces derniers puisque ils ont accès à toutes les
Les autochtones ont	donc	priorité quant à la récolte.
Il serait	donc	intéressant de comparer des données plus récentes afin de
Ces chiffres révèlent	donc	une tendance à la baisse entre 1976 et 1980 mais sûrement aussi
Les données furent	donc	suggéré de répartir les territoires de chasse selon des zones “

.Cher Mr Je viens	donc	vous écrire pour vous demander si je peux avoir de l’octroi pour
suite pourvu que je sois certain d’avoir ma prime	Donc	je me fie entièrement à vous pour régler cette affaire-là et je
mais cette année la Compagnie Fraser l’achète.	Donc	espérant recevoir une bonne réponse de vous le plus tôt
et ils demandent déjà un bon prix	donc	s’il vous plaît enseignez-moi les moyens à prendre pour le
Bien cordialement. Soyez	donc	assez bon de me dire s’il y a encore de bons lots à prendre
Je connais la terre étant fils de cultivateur.	Donc	Monsieur le Curé je sais que si vous le voulez je pourrais aller
bien si vous pouviez venir inspecter ce chemin	donc	je veux pas vous ennuyer avec cela

TABEL III. EXEMPLE PENTRU EXERCİTIU 3.

1. e de s’adapter au monde contemporain. Il	s’agit	de savoir si l’on table, oui ou non
2. es, de Beurs ni de Blacks (hélas). Il ne	s’agit	pas d’une bande dessinée mais
3. is gaulois. On en use à présent quand il	s’agit	d’évoquer les solutions apportées
4. oins comme ami, simplement parce qu’il	s’agit	de quelqu’un de différent
5. une femme. Et réciproquement. Quand il	s’agit	de cette différence-là il y a
6. olution partiellement dans le prototype. Il	s’agit	de définir les autorisations
7. uation de l’élève faite par le système. Il	s’agit	donc de recueillir ces informa

Una din principalele caracteristici ale învățării asistate de calculator este individualizarea procesului de învățare. Această funcție este utilizată în abordările mai recente care favorizează abordarea centrată pe profesor. Ultima este caracterizată prin utilizarea programelor de concordanță în clasa orelor de studiere a limbilor – abordare numită data Driven learning.

## BIBLIOGRAFIE

- [1] CHOMSKY, N., La Linguistique cartésienne suivi de La Nature formelle du langage, Éditions du Seuil, 1969.
- [2] DAVIES, G., The History of EUROCALL: an article produced to celebrate the dawning of the new millennium, 2000 - <http://www.camsoftpartners.co.uk/EuroHist.htm>
- [3] GARRETT, N., Technology in the service of language learning: Trends and issues. *Modern Language Journal*, 75(1), 74-101, 1991.
- [4] HIGGINS, J. & JOHNS, T. *Computers in language learning*, London: Collins, 1984.
- [5] Johns, T. & P. King (Eds.), (1991), *Classroom Concordancing*. *English Language. Research Journal*, 4. KETTEMANN, B. On the use of concordancing in ELT. *TELL & CALL*, 1995.
- [6] Henry Kučera; W. Nelson Francis *Computational Analysis of Present-Day American English* *International Journal of American Linguistics*, Vol. 35, No. 1 (Jan., 1969), pp. 71-75.
- [7] LAMY, M-N. & KLARSKOV MORTENSEN, H. J., Using concordance programs in the Modern Foreign Languages classroom. Module 2.4 in Davies G. (ed.) *Information and Communications Technology for Language Teachers (ICT4LT)*, Slough, Thames Valley University, 2011 [Online]: [http://www.ict4lt.org/en/en\\_mod2-4.htm](http://www.ict4lt.org/en/en_mod2-4.htm)
- [8] LEECH, G., *Corpora and theories of linguistic performance*. In Svartvik J. (ed.) *Directions in corpus linguistics*, Berlin: Mouton de Gruyter, 1992.