

# Instrumente informatice implicate în cercetarea terminologică

Cristina Varga  
(Universitatea Babeș-Bolyai, Cluj Napoca)

## Introducere

Există o mare varietate de instrumente informatice ce intervin în cercetarea lingvistică în general și în cercetarea terminologică în special. Acest din urmă domeniu este unul dintre domeniile predilecte ale lingvisticii aplicate, printre altele, și pentru faptul că rezultatele cercetării terminologice se materializează în instrumente de lucru pe teren lingvistic cu o largă utilizare<sup>1</sup>. Elementele fundamentale care stau la baza cercetării în domeniul terminologiei sunt: *crearea, manipularea, explorarea și gestionarea* de corpusuri de texte.

Analiza și explorarea unui corpus de texte se face, prin forța lucrurilor, utilizând instrumente electronice, deoarece acestea oferă avantajele vitezei de lucru, exactității și eficienței, rezultatele procesului de cercetare putând fi utilizate în varii domenii precum: **terminologia** (detectarea și extracția termenilor, selectarea unui context valid pentru un termen); **predarea limbilor străine, traducere și interpretare, lexicologie / lexicografie** (generală sau specializată), **lingvistică generală, lingvistică contrastivă, redactare de texte, analiza discursului, în procesul didactic** (în cadrul cursurilor de formare de traducători, interpreți și terminologi). În toate aceste domenii cercetarea științifică apelează la instrumente electronice capabile să recunoască, să extragă, să compare segmente lingvistice, care apoi sunt interpretate de către specialistul în domeniu în funcție de finalitatea cercetării. Alegerea instrumentelor de lucru, din multitudinea de alternative existente, trebuie să aibă în vedere deci finalitatea cercetării și gradul de eficacitate pe care îl are acesta într-un context clar determinat.

Acest articol își propune să prezinte câteva instrumente electronice complexe, capabile să gestioneze și să manipuleze corpusuri. Este vorba despre instrumente destinate activității de cercetare, cu distribuție gratuită în Internet, suficient de complexe ca să permită dezvoltarea unui proiect de investigație lingvistică sau terminologică. Sunt create în mediul de cercetare academic și fac obiectul de studiu al unor cursuri de specialitate în diverse universități europene. Cunoașterea lor și pe teren românesc, cunoașterea facilităților pe care le oferă precum și a limitelor lor putând constitui atât o ilustrare a direcțiilor de evoluție a instrumentelor electronice în cadrul cercetării lingvistice cât și un model ilustrativ pentru cei care, implicați fiind în proiecte de cercetare lingvistică pe teren românesc, ar dori să dezvolte instrumente de explorare și exploatare de corpusuri de texte dedicate limbii române.

## Instrumente electronice de investigație lingvistică

Ne propunem ca în acest articol să prezentăm trei instrumente electronice cu aplicații în domeniul cercetării terminologice, programe a căror utilitate, mai ales în ceea ce privește cercetarea terminologică, nu a fost încă suficient pusă în evidență pe teren românesc.

---

<sup>1</sup> De cele mai multe ori, analiza și exploatarea de corpusuri de texte au ca rezultat elaborarea a diverse tipuri de materiale lingvistice (*dicționare generale, dicționare-tezaur, glosare specializate, etc.*).

Acestea sunt: **SCP (Simple Concordance Program)**, **Lexico3 și Corpografo**. Sunt programe cu distribuție gratuită în scopul cercetării, două dintre acestea se instalează local pe stații de lucru, în timp ce al treilea, **Corpografo** este un program cu acces on-line<sup>2</sup>. Toate sunt caracterizate printr-o interfață grafică accesibilă, ușor de înțeles și de utilizat. De asemenea, prezintă un grad de complexitate<sup>3</sup> corespunzător necesităților din domeniul cercetării academice. Fiind de producție diferită<sup>4</sup>, fiecare dintre acestea prezintă elemente specifice ce ilustrează liniile de dezvoltare ale cercetării lingvistice în centrul academic de proveniență. Analiza lor contrastivă ne-a permis să facem o ierarhizare a acestor programe în funcție de gradul de complexitate și de facilitățile pe care le oferă în domeniul cercetării lingvistice, ierarhizare care va face ca prezentarea lor să debuteze cu cel mai simplu dintre ele și să continue treptat cu cele care au un grad mai mare de complexitate.

### **SCP (Simple Concordance Program)**

Program de analiză de texte și recuperare de informație dintr-un corpus textual, cu distribuție gratuită, creat de Alan Reed și ajuns actualmente la versiunea 4.09. Se poate obține de la URL: <http://www.textworld.com/>.

**Simple Concordance Program (SCP)** este cel mai „simplu” dintre cele trei instrumente electornice care fac obiectul acestui articol. Înainte de a pune în evidență utilitatea sa în contextul cercetării lingvistice, este de remarcat faptul că **SCP** s-a dovedit a fi un instrument foarte bun în procesul didactic. Foarte ușor de utilizat și cuprinzând principalele programe și funcții de analiză și explorare de corpus, în unele instituții academice a fost inclus în programul de studii pentru discipline ca: **Informatica pentru traducători** sau **Lingvistica de corpus**<sup>5</sup>. Acest program s-a transformat într-un instrument de lucru excelent pentru studenții<sup>6</sup> care se inițiază în cercetarea terminologică.

Funcțiile sale de bază se referă la recuperarea și extragerea de informații (*ocurențe, cuvinte cheie, termeni*, etc.) dintr-un corpus de texte propus de utilizator și care trebuie prezentat în format .TXT. Prezentarea unităților lexicale care fac obiectul cercetării se poate face în context (la nivel de frază – KWIC sau la nivel de linie - LINE). Sunt permise de asemenea afișarea listei unităților lexicale care compun textul în diverse moduri: *ordine alfabetică*, în funcție de *frecvența cuvintelor* sau în *ordinea apariției lor în text*. Programul beneficiază, de asemenea de funcții ce permit crearea unui *profil de frecvență* al unei unități lexicale sau prezentarea *proprietăților statistice* ale corpusului analizat.

---

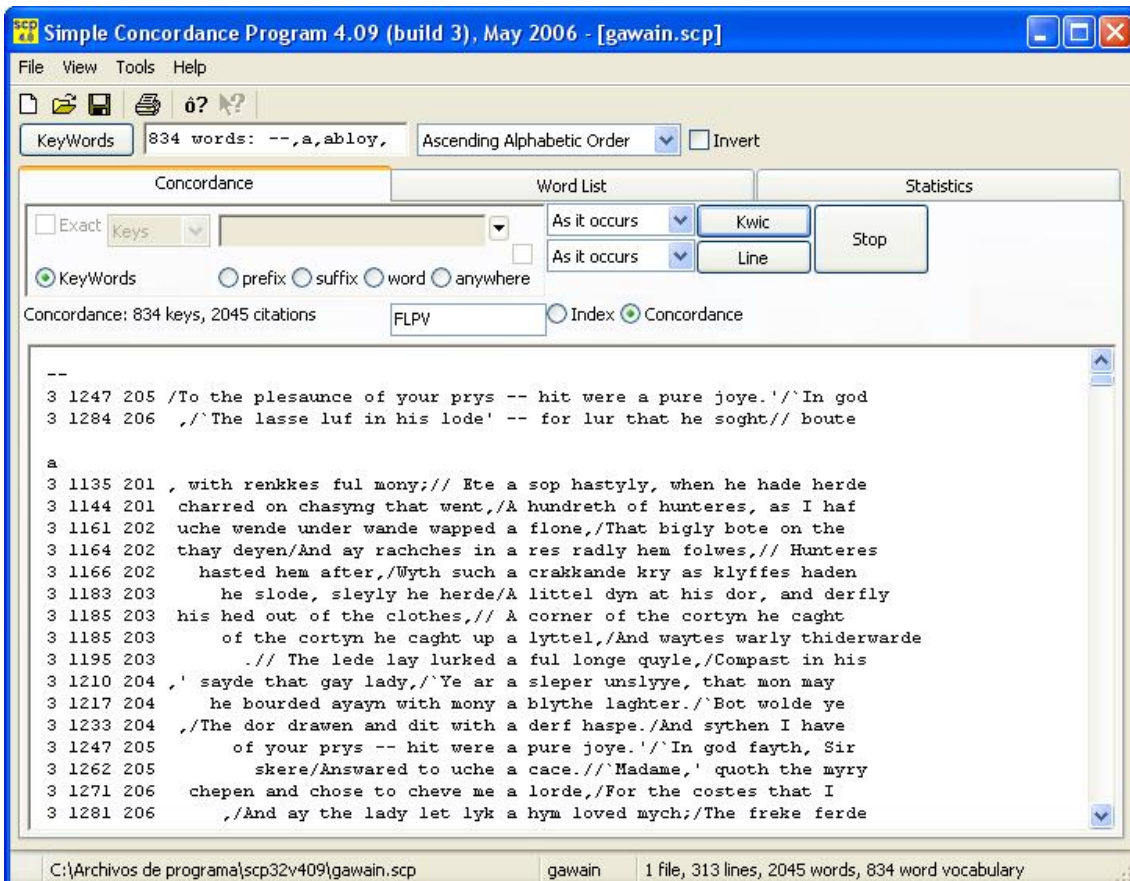
<sup>2</sup> Pentru a putea lucra cu **Corpografo** este necesară înscrierea în lista de utilizatori ai acestui program și obținerea unei chei de acces.

<sup>3</sup> În cazul fiecărui instrument vorbim despre un complex de aplicații unificate într-un pachet de programe, ceea ce individualizează profilul fiecărui instrument.

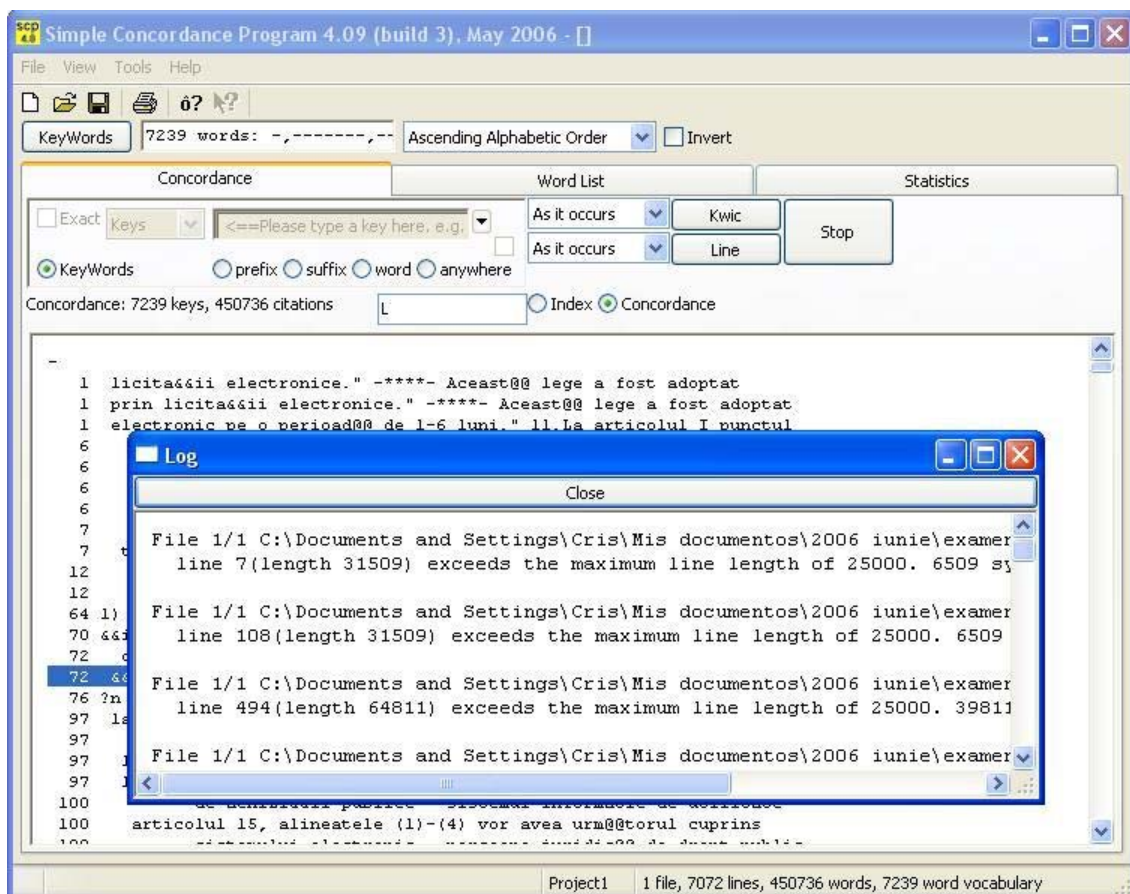
<sup>4</sup> Fiecare dintre aceste programe reprezintă orientările cercetărilor lingvistice din altă țară. Astfel, **SCP** este produs de un grup de cercetători din SUA, **Corpografo** reprezintă un proiect în desfășurare a unui grup de cercetare din Portugalia, iar **Lexico3** este dezvoltat de către cercetători de la Université Paris 3 – Sorbonne Nouvelle, Franța.

<sup>5</sup> Este cazul **Universității Pompeu Fabra**, Barcelona, unde acest program se studiază în cadrul a două specializări.

<sup>6</sup> De remarcat faptul că este utilizat și în centre de cercetare terminologică de prestigiu cum ar fi de exemplu IULA (**Institut Universitari de Lingüística Aplicada**), Barcelona.



Fiecare sesiune de lucru debutează cu crearea unui *proiect SCP* sau cu modificarea unuia existent, *proiect* în care se poate încărca un corpus de texte. Există situații în care pot apare mesaje de eroare, este cazul în care lungimea unei linii a corpusului de texte depășește limita de 25.000 de caractere. Semnele care depășesc această limită vor fi ignorate. Programul semnalează utilizatorului care sunt liniile în cauză și, de asemenea, se menționează numărul de caractere ignorate în fiecare dintre cazuri, astfel acesta va putea aprecia cât din informația pe care o analizează este ignorată.



Extragerea informației din corpus se poate face utilizând o listă de *cuvinte cheie*, sau căutând ocurențe în funcție de *prefix*, *sufix*, *întregul cuvânt* sau o *parte oarecare a acestuia*. Toate instanțele elementului căutat vor fi afișate în context, incluzând, la cererea utilizatorului *descriptori statistici*. Descriptorii statistici se întâlnesc nu doar la nivel de text ci și la nivelul întregului proiect. Aceștia descriu proiectul din punctul de vedere al unităților lexicale ce compun textul (*frecvența unităților lexicale*, *număr acestora în text*, *vocabular cumulativ*, *statistici de vocabular/ocurențe*), sau la nivel global cu referire la întregul proiect, dar și la nivelul caracterelor ce compun textul (*frecvență majuscule*, *minuscule și simboluri*).

Foarte interesantă din acest punct de vedere este posibilitatea de a explora corpusul pe baza unor *liste de cuvinte* ce pot fi personalizate de către utilizator. Aceste liste de cuvinte sunt create pornind de la inventarul corpusului cu care se lucrează în acel moment. Opțiunea **keywords** permite atât selecția manuală a cuvintelor cheie cât și utilizarea unor formalime care automatizează acest proces. Aceste liste se pot crea pe baza frecvenței unităților lexicale (se poate stabili între anumite limite prin utilizarea operatorilor <, > sau =), în funcție de structura cuvântului (prefixe, sufixe, sau indiferent de poziția unui element în cuvântul căutat) sau după criteriul de lungime a cuvântului (stabilit între anumite limite prin utilizarea operatorilor <, > sau =). Programul permite atât manipularea listelor de cuvinte precum și importarea sau exportarea acestora în format TXT.

Prezentarea elementelor lexicale recuperate din text se poate face sub forma de **index** sau de **concordanță de tip LINE sau KWIC**, pentru această ultimă opțiune putându-se selecta posibilitatea prezentării elementelor în funcție de contextul de stânga sau de dreapta al cuvântului.

SCP 4.9 Simple Concordance Program 4.09 (build 3), May 2006 - [gawain.scp]

File View Tools Help

KeyWords 834 words: --,a,abloy, Ascending Alphabetic Order  Invert

Concordance Word List Statistics

Exact Keys un  2  As it occurs    
 KeyWords  prefix  suffix  word  anywhere As it occurs   
Left context  Index  Concordance  
Right context

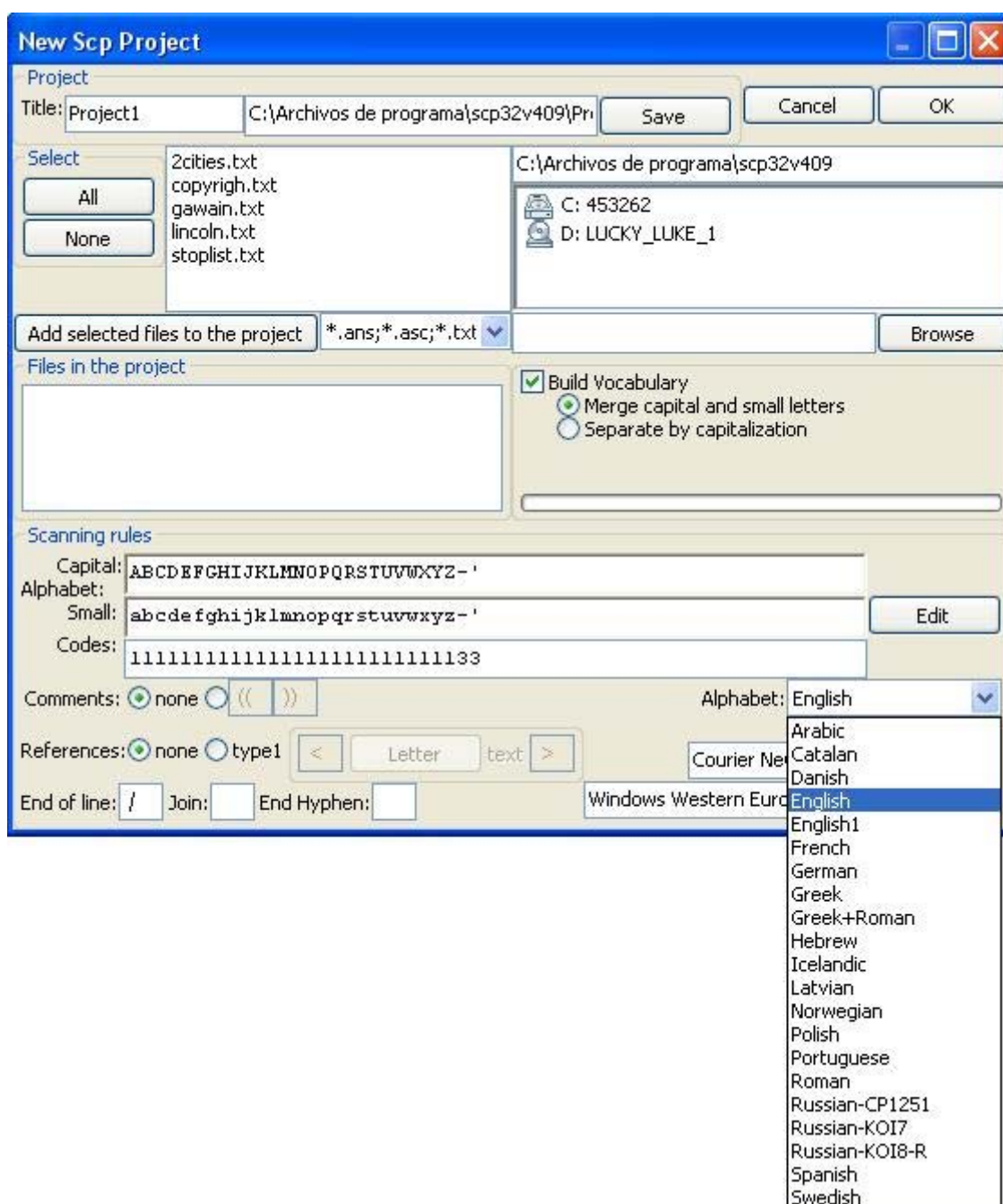
Concordance: 1 key, 8 citations FLPV

```

un
3 1140 201  couthe cowpled hor houndes, // Unclosed the kenel dore and calde
3 1161 202      of arwes; // At uche wende under wande wapped a flone, // That
3 1181 203  daylyght lemed on the woves, // Under covertour ful clere, cortyned
3 1201 203  and to-hir-warde torned, // And unlouked his yye-lyddes and let
3 1210 204      gay lady, // `Ye ar a sleper unslyye, that mon may slyde hider
3 1245 205  as ye reherce here / I am wywe unworthy, I wot wel mysyven. /
3 1327 208      were, // And didden hem derely undo as the dede askes. // Serched
3 1352 209  thay hyyes, // Bi the bakbon to unbynde. // Bothe the hede and the

```

C:\Archivos de programa\scp32v409\gawain.scp gawain 1 file, 313 lines, 2045 words, 834 word vocabulary



Limitele programului sunt vizibile în momentul în care limba textelor dintr-un anumit corpus presupune existența în texte a unor caractere speciale, cum este cazul limbii române. **SCP** permite utilizarea unui set restrâns de caractere (ANSI / ASCII) deși numărul limbilor cu care se poate lucra este destul de mare. Din păcate această listă nu cuprinde și limba română iar fonturile nu sunt adaptate pentru limbă noastră. Acest fapt are ca rezultat dificultăți în explorarea unui corpus în limba română, totuși, acest lucru nu este imposibil.

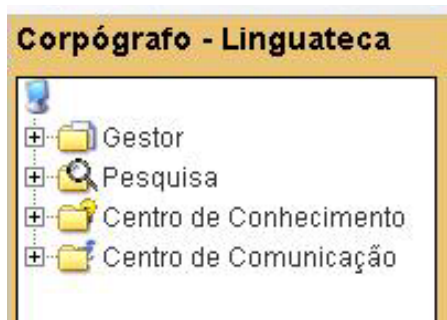
### Corpografo

Situându-se la un nivel superior în ceea ce privește complexitatea și funcțiile pe care le oferă, **Corpografo** este prezentat ca instrument de analiză și exploatare de corpus textual pentru elaborare de instrumente lingvistice în format digital. Scopul cu care a fost construit este crearea de surse lingvistice complexe cum ar fi: *dicționare generale, dicționare-tezaur, glosare specializate,*

corectoare ortografice, aplicații de recuperare a informațiilor sau de traducere automată.

Proiect portughez creat de o echipă formată din: Luís Sarmiento, Ana Sofia Pinto, Luís Miguel Cabral, Débora Oliveira, Belinda Maia, Diana Santos<sup>7</sup>. Este vorba despre un proiect lingvistic cu acces gratuit on-line, însă necesitând obținerea unei chei pentru fiecare utilizator și având de întâmpinat toate dificultățile accesului on-line la o bază de date aflată la mare distanță, **Corpografo** poate fi accesat la URL: [www.linguateca.pt/corpografo](http://www.linguateca.pt/corpografo). Cu o interfață ușor de utilizat, care însă prezintă dificultatea unui meniu de comenzi în portugheză, este un proiect în curs de dezvoltare ce dorește să integreze într-un mediu electronic unice diverse instrumente de analiză textuală ce sunt considerate în mod obișnuit ca entități individuale. Este destinat cercetătorilor, studenților, traducătorilor, etc.

Foarte bine organizat din punctul de vedere al sarcinilor pe care le poate îndeplini, **Corpografo** permite crearea și gestionarea de corpusuri personalizate, analiza, căutarea și extracția terminologică pe baza acestor corpusuri. Interfața acestuia prezintă patru meniuri de lucru reprezentate prin categoriile: **Gestor**, **Pesquisa**, **Centro de Conhecimento** și **Centro de Comunicação**. Fiecare dintre aceste categorii corespunde unei etape de lucru în cercetarea terminologică a unui corpus. Le vom prezenta pe fiecare pe scurt în cele ce urmează.



**Gestor** reprezintă un spațiu virtual de stocare a fișierelor care contribuie la formarea unui corpus ce stă la baza cercetării terminologice. Crearea corpusului corespunde unei prime etape a cercetării, de modul în care se realizează această etapă de investigare depinzând calitatea și eficiența interogărilor ulterioare a corpusului creat. Astfel, o mare atenție trebuie acordată, în utilizarea **Corpografo**, unei etape premergătoare de „pregătire” a textelor înainte de explorarea corpusului.

De asemenea, o mare importanță este acordată organizării informației în cadrul corpusului. Pentru o mai bună gestionare și un control efectiv al informației conținute într-un corpus, există posibilitatea de a clasifica fișierele din interiorul acestuia în funcție de *tematica* pe care o tratează textul, iar în cadrul acestei categorii există descriptori ca: *mediul*, *domeniul* și *subdomeniul de specializare*. Se poate menționa *sursa textului* (prin menționarea organizației, instituției sau a editurii) precum și *autorul*, în cazul în care aceste informații sunt repetitive, există, posibilitatea de a relaționa între ele textele care au aceeași sursă sau aparțin aceluiași autor. Alte informații care însoțesc textul sunt: *numele fișierului*, *titlul documentului*, *limba textului*, *reguli de acces public*, *data de introducere* și *de modificare*, precum și o *descriere* a textului.

<sup>7</sup> Pentru mai multe detalii asupra rolului fiecărui membru al echipei în cadrul proiectului, vezi [www.linguateca.pt/corpografo](http://www.linguateca.pt/corpografo)

Pentru realizarea operațiunii de „pregătire” a textului, este necesară accesarea opțiunii de *editare a textelor* din meniul **Texto**. Se consideră necesară eliminarea tuturor elementelor irelevante din text, cum ar fi: *caractere speciale, referințe bibliografice, erori de ortografie, titluri, subtitluri, note bibliografice* existente în text, precum și *formulele* ce pot apărea. O opțiune importantă pentru o primă luare de contact cu textul este **Ver dicționar**, ceea ce permite obținerea unei liste generale a *atomilor* unui text. O altă opțiune utilă este cea de comparare a două texte.

**Corpógrafo - Linguateca**

**Gestor**

- Ficheiros
- Corpora
- Autores
- Organizações
- Repositório multimédia
- Pesquisa
- Centro de Conhecimento
- Centro de Comunicação

---

**Operações**

- \* Criar pasta
- \* Renomear texto
- \* Apagar
- \* Mover para ...

---

**Adicionar ficheiros**

- \* Do meu computador
- \* A partir de um URL

---

**Informações**

**Lista de Ficheiros**

	FICHEIROS:\ [texto]		
1	? Sony Ericsson.doc	552 átomos	
2	? ordonanta 34.doc	8081 átomos	
3	? HG 348.doc	2881 átomos	
4	? Legea 510.doc	4734 átomos	
5	? legea 591.doc	4073 átomos	
6	? ordonanta 34.doc	8081 átomos	
7	? ordonanta 79.doc	20288 átomos	
8	? 2_TNABF 0-30 MHz.pdf	9919 átomos	
9	? 3_TNABF 30-1000 MHz.pdf	14002 átomos	

Odată încheiată etapa de *introducere, descriere și pregătire* a materialului lingvistic se trece la crearea corpusului, după care se poate începe explorarea acestuia.

Funcțiile destinate explorării corpusului sunt grupate în meniul **Pesquisa** (căutare). **Pesquisa** corespunde etapei de cercetare lingvistică propriu-zisă, iar instrumentele pe care le pune la dispoziție permit *căutarea, studierea și extragerea* de informații dintr-un corpus. O opțiune importantă se consideră a fi căutarea de secvențe de cuvinte consecutive de diverse lungimi ce se pot stabili de către utilizator (**N-gramas**) și care apar frecvent într-un text. Aceasta permite observarea unor structuri și combinații sistematice de cuvinte, foarte utilă în identificarea termenilor de specialitate.

Foarte importantă la nivelul explorării textului este opțiunea de realizare a diverse tipuri de concordanțe: **concordanța la nivel de frază**<sup>8</sup>, **concordanța Janela**<sup>9</sup> și **concordanță KWIC**<sup>10</sup>. În afară de aceste posibilități de explorare de corpus, **Corpografo** prezintă avantajul de a putea efectua căutări și în funcție de expresii regulate, care, de asemenea, pot constitui baza de pornire în crearea de concordanțe.

<sup>8</sup> Tip de concordanță în care contextul se limitează la fraza în care apare cuvântul căutat.

<sup>9</sup> Tip de concordanță în care se poate defini numărul de cuvinte care să constituie contextul de stânga și dreapta al termenului căutat.

<sup>10</sup> Tip de concordanță în care se permite definirea contextului în funcție de numărul de cuvinte sau de caractere.



Corpógrafo - Linguateca

Gestor  
Pesquisa  
Concordância Frase  
Estudo de N-Gramas  
Concordância Janela  
Concordância KWIC  
Centro de Conhecimento  
Centro de Comunicação

Expressão de pesquisa: electronic  
Corpus pesquisado: telecomunicatii (34965 átomos : 7 ficheiros)

Nº de concordâncias obtidas: 3 distribuídas por 3 ficheiros  
Frequência da totalidade das instâncias da concordância: 0.00 %

<< <<< >>> >>

#	f	Frase onde ocorre a concordância	Info
1	-1	7 ( 1 ) Finançarea cheltuielilor curente de capital ale IGCTI se asigură din venituri proprii , care provin din următoarele surse: a ) sumele încasate de la titularii dreptului de utilizare a spectrului de frecvențe radio pentru utilizarea acestuia; b ) tariful de înregistrare a Sistemului electronic de achiziții publice; c ) tariful de reînnoire a înregistrării în Sistemul electronic de achiziții publice; d ) tariful de utilizare a Sistemului electronic de achiziții publice; e ) tariful de participare în Sistemul electronic de achiziții publice; f ) tariful de eliberare a certificatului digital de înregistrare în Sistemul informatic pentru atribuirea electronică a autorizațiilor de transport rutier internațional de marfă; g ) tariful de reînnoire a certificatului digital de înregistrare în Sistemul informatic pentru atribuirea electronică a autorizațiilor de transport rutier internațional de marfă; h ) donații , legate de sponsorizări în condițiile legii; i ) credite interne de externe contractate în condițiile legii; j ) sumele provenite din amenzi administrative aplicate conform art. 56 din Ordonanța de urgență a Guvernului nr .	
2	-1	operarea Sistemului electronic de achiziții publice , a Sistemului electronic național de Sistemului informatic pentru atribuirea electronică a autorizațiilor de transport internațional rutier de marfă; 18 .	
3	-1	Ele pot viza: a ) contribuții financiare pentru susținerea serviciului universal; b ) plata tarifului de monitorizare anual; c ) interoperabilitatea serviciilor de interconectare rețelelor; d ) disponibilitatea resurselor de numerotație pentru utilizatorii finali , inclusiv condiții impuse în temeiul legislației speciale privind serviciul universal; e ) cerințe privind protecția mediului , planurile de urbanism de amenajare a teritoriului , precum și cerințe de condiții legate de acordarea dreptului de acces pe proprietăți , colocarea și utilizarea partajată a resurselor , inclusiv , dacă este cazul , garanțiile de ordin financiar sau tehnic necesare pentru a asigura executarea corespunzătoare a lucrărilor de infrastructură; f ) obligații privind retransmisia serviciilor de programe prin rețelele de comunicații electronice , în conformitate cu prevederile legislației privind audiovizualul; g ) prelucrarea datelor cu caracter personal de protecția vieții private; h ) protecția consumatorilor; i ) restricții privind transmiterea conținutului ilegal de video și audiovizual , în conformitate cu prevederile legale aplicabile în domeniul comerțului electronic de audiovizual; j ) informații care trebuie furnizate în temeiul art. 4 alin .	

Concordâncias  
\* Voltar ao menu  
\* Analisar Frequências  
\* Contexto Janela 1-1  
\* Contexto Janela 2-2  
Imprimir!

Informações

**Centro de Conhecimento** reprezintă spațiul în care se pot sistematiza și organiza informațiile obținute în urma explorării unui corpus. Este vorba despre informații de tip *lexical, morfologic, sintactic și semantic*, ce pot permite crearea de materiale lingvistice cum ar fi: *liste de cuvinte, glosare, reguli de căutare, tipare, relații semantice*, etc. Cu ajutorul **Corpografo** se poate crea o reprezentare formală a conceptelor și informațiilor lingvistice asupra unui domeniu specializat. Acest aspect are aplicații importante în cercetarea terminologică, formalizarea cunoștințelor dintr-un domeniu specializat fiind un element fundamental. La nivel tehnic această secțiune a **Corpografo** permite gestionarea de baze de date terminologice. Se insistă asupra faptului că **Corpografo** nu doar permite gestionarea unor simple liste de cuvinte ci utilizatorul are posibilitatea de a stabili *relații* între termeni, ceea ce îi permite crearea de *rețele conceptuale multidimensionale*, **Corpografo** dovedindu-se a fi un instrument foarte flexibil din acest punct de vedere.

Corpógrafo - Linguateca

Gestor  
Pesquisa  
Centro de Conhecimento  
BD Terminológicas  
Gestor de Relações  
Centro de Comunicação

acces restricționat

1) Dados Gerais

idioma	?
tipo	entrada terminológica
estado administrativo	admitido
registro	técnico
frequência de utilização	usado com frequência
proveniência	empréstimo interlinguístico

2) Autores

Dados relativos aos autores indisponíveis

3) Fontes

1:

4) Morfologia

gênero	N
número	S
animacidade	I
morfologia	NC AJ

5) Definições

Tip de acces personalizat pentru fiecare utilizator, prin care fiecare utilizator i se alocă un drepturi de acces în rețea.

Ficha do termo  
\* Dados Gerais  
\* Autores  
\* Fontes  
\* Morfologia  
\* Definições  
\* Relações Semânticas  
\* Equivalentes Tradução  
\* Media associado  
\* Estatísticas  
\* Eliminar Termo  
\* Listar Termos  
Imprimir!

Informações

**Centro de Conhecimento** permite de asemenea gestionarea și editarea bazelor de date terminologice, căutarea și stabilirea de relații între termeni. În ceea ce privește căutarea termenilor într-un corpus, utilizatorul are posibilitatea de a stabili lungimea unei secvențe textuale, trecerea de la *forma flexionată* a acestuia la *forma normalizată*, de asemenea se pot omite termenii care deja sunt introduși în baza de date. După introducerea tuturor termenilor selectați de către utilizator în baza de date, aceștia vor fi descriși, în baza de date existând descriptori pentru *limbă*, *descriere morfologică*, *autor*, *referință bibliografică*, de asemenea este posibilă *căutarea unei definiții a unității lexicale* în cauză în corpusul de texte ce se explorează, căutarea de eventuale *relații semantice* între termenii bazei de date, *căutarea de echivalente de traducere*, *asocierea de elemente multimedia* pentru termenul respectiv, *consultarea de statistici* cu referire la un anumit termen existent în corpus.

Ultimul meniu, **Centro de Comunicaçao**, se referă la partea de documentație asupra **Corpografo** și la posibilitatea de a contacta administratorul **Corpografo** pentru schimb de mesaje.

**Corpografo - Linguateca**

**Publicações**

**The Corpografo - a Web-based environment for corpora research**  
**Autores:** Luís Sarmento, Belinda Maia & Diana Santos.  
**Descrição:** Artigo incluído nos Proceedings of LREC 2004 . Lisboa, Portugal, 25 May 2004.  
**Tipo:** Artigo  
[\[ Descarregar Artigo \]](#)

**The pedagogical and linguistic research implication of GC to on-line parallel and comparable corpora**  
**Autora:** Belinda Maia  
**Descrição:** Apresentado na conferência CP3A 2003: Corpora Paralelos, Aplicações e Algoritmos Associados, 3 de Junho de 2003.  
**Tipo:** Artigo e Apresentação  
**Nota:** O Corpografo era anteriormente denominado "Gestor de Corpora" ou simplesmente "GC"  
[\[ Descarregar Artigo \]](#) | [\[ Descarregar Apresentação \]](#)

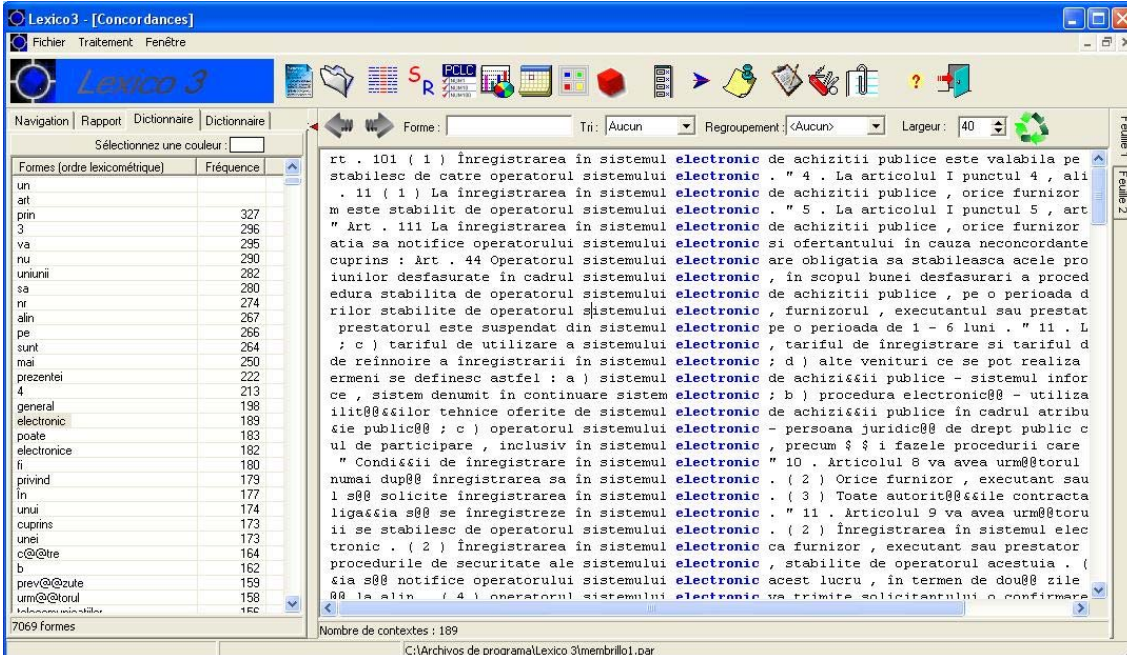
**Gestor de Corpora - Um ambiente Web para Linguística Computacional**  
**Autor:** Luís Sarmento  
**Descrição:** Apresentado na conferência CP3A 2003: Corpora Paralelos, Aplicações e Algoritmos Associados, 3 de Junho de 2003.  
**Tipo:**Artigo e Apresentação  
**Nota:** O Corpografo era anteriormente denominado "Gestor de Corpora" ou simplesmente "GC"  
[\[ Descarregar Artigo \]](#) | [\[ Descarregar Apresentação \]](#)

Din punctul de vedere al limbii române, **Corpografo** prezintă limite, deoarece limba română nu face parte dintre limbile propuse pentru studiu și nu are implementate caracterele speciale din limba română. Totuși utilizarea textelor în limba română este posibilă, dezavantajul fiind afișarea caracterelor speciale sub formă de coduri .HTML. În cazul în care ar exista un interes pentru acest instrument, contactarea membrilor proiectului și propunerea implementării caracterelor speciale pentru limba română și a limbii române ca limbă de lucru în acest program, ar fi cea mai simplă soluție. Avantajul cert pe care îl oferă față de celelalte instrumente electronice de explorare de corpus ce fac obiectul acestui articol, este faptul că este foarte flexibil în ceea ce privește formatul fișierelor care pot să compună corpusul. **Corpografo** acceptând simultan fișiere .DOC, .RTF, .PDF, .TXT, .PS, .HTML ceea ce reprezintă de departe cea mai largă gamă de formate de documente text acceptate până în prezent de un asemenea instrument electronic de acest tip.

## Lexico3

Ultimul dintre instrumentele electronice care marchează evoluția actuală în lingvistica aplicată, și pe care ne-am propus să îl prezentăm este **Lexico3**, un complex de programe de statistică textuală, după cum îl prezintă autorii săi. Elaborat de către o echipă a Universității Paris 3– Sorbonne Nouvelle, din care fac parte Cédric Lamalle, William Martinez, Serge Fleury și André Salem. Este un instrument complex, cu distribuție gratuită în scopul cercetării și testării, care poate fi descărcat de la URL <http://lexico3.no-ip.org/>.

Cu o interfață transparentă și ușor de manevrat, modul de lucru cu **Lexico3** se reduce la introducerea unui corpus de texte în format .TXT într-o bază, după care se poate trece direct, fără alte etape intermediare, la explorarea textului prin utilizarea a diverse instrumente de investigare, analiză și statistică pe care **Lexico3** le pune la dispoziția utilizatorului. Avantajele certe pe care le permite programul în această primă etapă de încărcare a corpusului este faptul că permite o etichetare a corpusului, precum și faptul că utilizatorului i se cere confirmarea utilizării unui inventar de delimitatori textuali cum ar fi: `.,:;!/?_-'"/[]{}$%` și care nu vor fi analizați ca și componente ale textului ci vor avea un rol important în fragmentarea acestuia. Listă care poate fi modificată de către acesta, în funcție de tipul de text pe care dorește să îl supună analizei statistice.



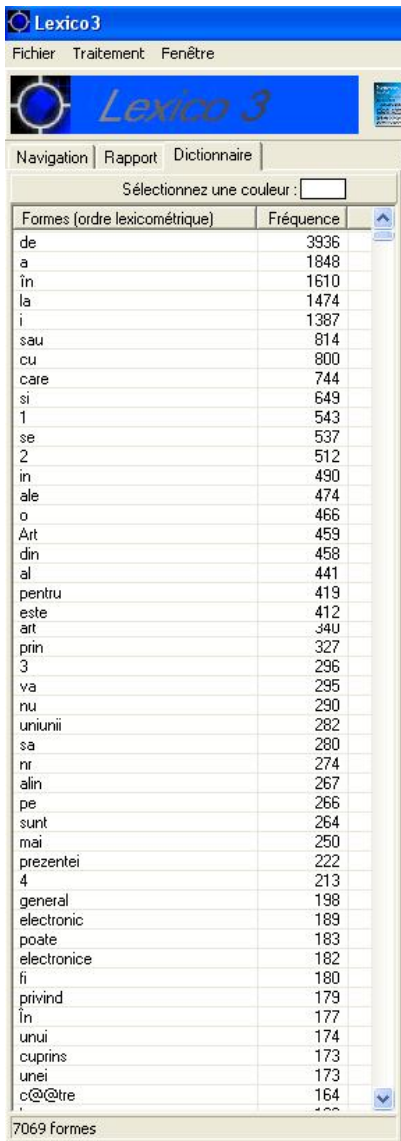
The screenshot shows the Lexico3 software interface. On the left, there is a table with two columns: 'Formes (ordre lexicométrique)' and 'Fréquence'. The table lists various forms and their frequencies. On the right, there is a text snippet with several instances of the word 'electronic' highlighted in blue. The interface also includes a menu bar, a toolbar, and a status bar at the bottom.

Formes (ordre lexicométrique)	Fréquence
un	327
art	296
prin	295
3	290
va	282
nu	280
uniunii	274
sa	267
nr	266
alin	264
pe	250
sunt	222
mai	213
prezentei	198
4	189
general	183
electronic	182
poate	180
electronice	179
fi	177
privind	174
in	173
unui	173
cuprins	164
unei	162
c@atre	159
b	158
prev@zute	158
urm@torul	158
telecomunicati	158

În ceea ce privește etichetarea corpusului, ea poate fi de orice tip, de la cea mai simplă la cea mai complexă, în funcție de ceea ce dorește utilizatorul, etichetele sunt similare cu cele XML, pot fi definite cu cea mai mare libertate, urmând ca apoi să i se menționeze programului care sunt etichetele importante și ce rol au în text. De exemplu, corpusul de texte demonstrativ al **Lexico3** se referă la presa din timpul Revoluției Franceze și prezintă următoarea etichetare: `<mois=01><quinzaine=11> <semaine=111> <Sda=1793> <numero=260> <edito=0> <Epg=1><Sat=0>`, destul de transparentă în prima parte, mai puțin transparentă în ceea ce privește ultimele patru elemente de etichetare. Utilitatea acestei etichetări se remarcă în momentul în care pentru o ilustrare grafică a distribuției unei ocurențe într-un text se cere delimitarea textului. Ca

delimitatori se pot folosi atât semnele de punctuație cât și etichetele care au rolul de descriptori într-un text.

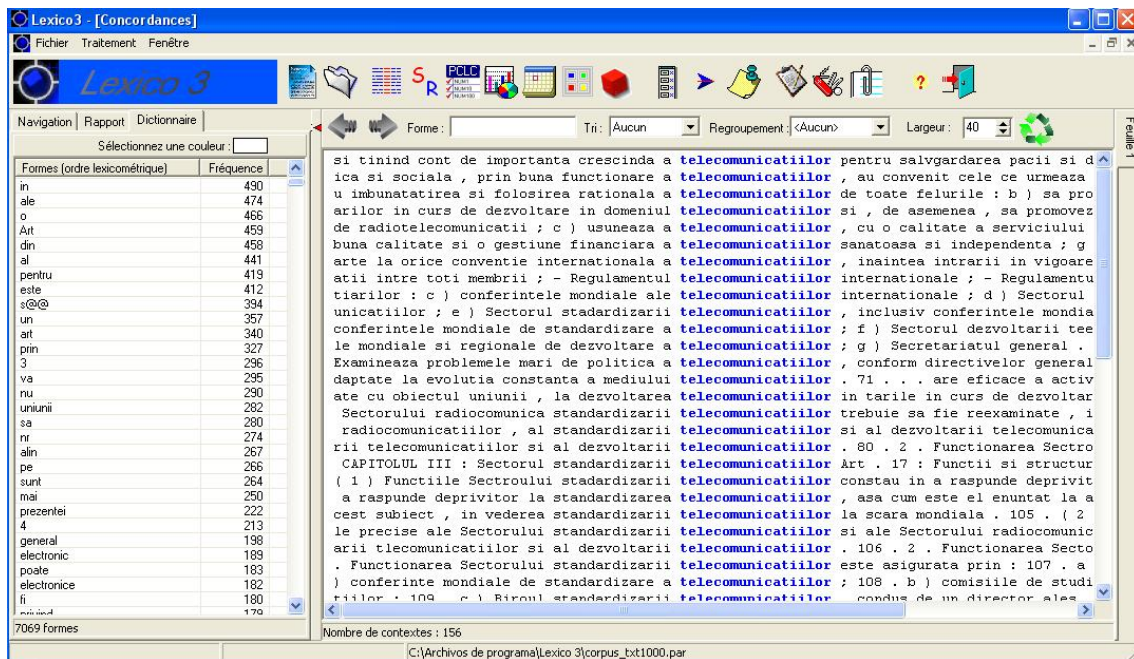
Prima operațiune care se poate efectua după introducerea corpusului în bază este consultarea dicționarului, lucru ce permite identificarea unităților lexicale care constituie textul și frecvența acestora în text.



The screenshot shows the Lexico3 software interface. At the top, there is a menu bar with 'Fichier', 'Traitement', and 'Fenêtre'. Below the menu bar is a title bar with the Lexico3 logo and the text 'Lexico 3'. The main window has a tabbed interface with 'Navigation', 'Rapport', and 'Dictionnaire' tabs. Below the tabs is a text input field labeled 'Sélectionnez une couleur :'. The main area contains a table with two columns: 'Formes (ordre lexicométrique)' and 'Fréquence'. The table lists various word forms and their corresponding frequencies. At the bottom of the table, it indicates '7069 formes'.

Formes (ordre lexicométrique)	Fréquence
de	3936
a	1848
în	1610
la	1474
i	1387
sau	814
cu	800
care	744
și	649
1	543
se	537
2	512
in	490
ale	474
o	466
Art	459
din	458
al	441
pentru	419
este	412
art	340
prin	327
3	296
va	295
nu	290
uniunii	282
sa	280
nr	274
alin	267
pe	266
sunt	264
mai	250
prezentei	222
4	213
general	198
electronic	189
poate	183
electronice	182
fi	180
privind	179
în	177
unui	174
cuprins	173
unei	173
c:@@tre	164
.	164

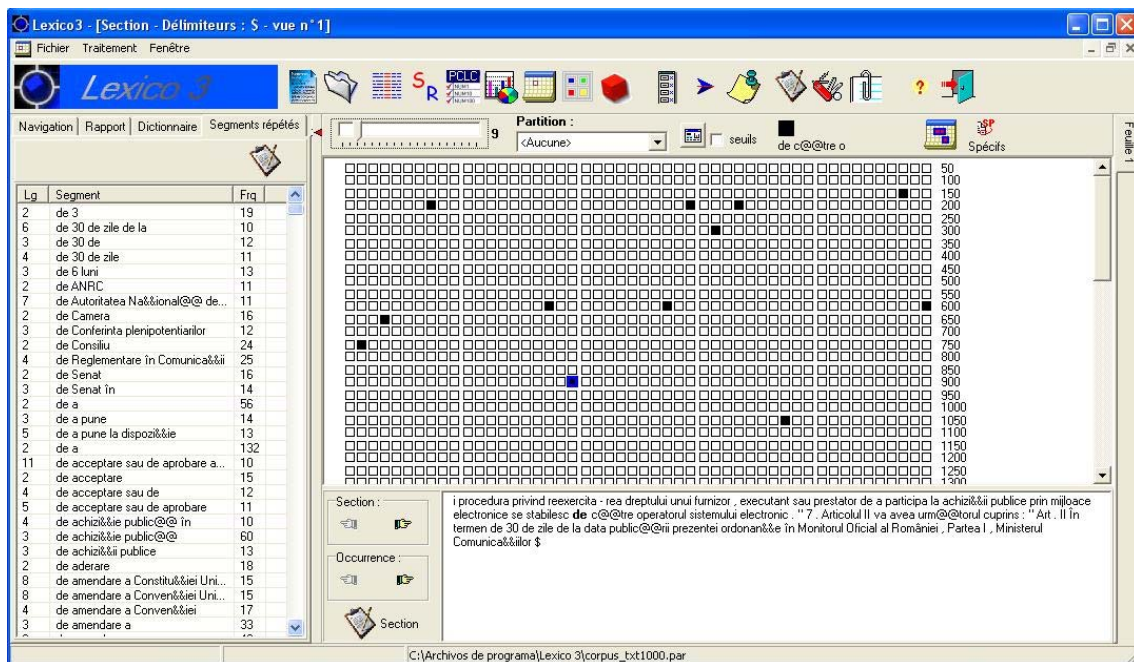
Pasul următor este crearea unei concordanțe. Spre deosebire de instrumentele prezentate anterior, interfața grafică a Lexico3 permite deplasarea elementelor lexicale cu ajutorul mouse-ului, deoarece are implementată funcția **drag and drop**. Concordanța pe care o poate realiza este una simplă, singurul element deosebit fiind faptul că se poate stabili de către utilizator dimensiunea contextului iar ordonarea ocurențelor se poate face în ordinea apariției în text sau în ordine alfabetică în funcție de contextul de stânga sau de dreapta.



Există de asemenea o opțiune ce permite utilizatorului să afișeze toate formele unei unități lexicale prezente în text, acest lucru putând permite identificarea de câmpuri lexicale, familii de cuvinte sau forme în flexiune, împreună cu frecvența apariției lor în text.



Dacă este vorba despre un corpus etichetat, se poate obține o reprezentare grafică a segmentelor care compun textul, așa cum au fost ele marcate de către utilizator. Astfel, se poate vedea, de exemplu, distribuția anumitor elemente lexicale, lucru destul de important în cazul analizei unui discurs sau a comparării a două texte. De asemenea pentru a putea studia distribuția și uzul unei unități/segment lexical repetitiv, în partea inferioară a ecranului, se poate afișa contextul de utilizare.



O altă opțiune foarte utilă se referă la crearea unei liste a segmentelor care se repetă într-un text. Sunt posibile, de asemenea, analiza elementelor specifice unui fragment al corpusului sau o analiză factorială a corespondențelor în text.



Rezultatele analizei se pot păstra într-un raport care poate fi citit cu **Netscape** sau **Internet Explorer**.

Limitele **Lexico3** se referă la aceeași problemă a acceptării limbii române ca limbă de lucru. Dar putând procesa fișiere în format .TXT, **Lexico3** prezintă avantajul de a putea lucra cu fișiere text ce conțin caractere UTF-8. Afișarea lor pe ecran lasă mult de dorit însă informațiile sunt accesibile utilizatorului. O altă limită a programului, care însă poate fi datorată faptului că pe Internet este doar o versiune demo, este limitarea corpusului de texte la 200 de pagini, ceea ce reprezintă un eșantion suficient pentru a studia posibilitățile de lucru pe care le oferă programul **Lexico3** dar care în contextul real al cercetării lingvistice și al analizei de corpus este insuficient.

## Concluzii

Prezentarea acestor instrumente electronice de explorare și exploatare a unui corpus de texte poate forma cercetătorului o idee generală asupra tendințelor actuale de evoluție a instrumentelor electronice cu aplicație în domeniul lingvistic. Se remarcă astfel o trecere de la instrumentele simple care rezolvau o problemă la nivel punctual<sup>11</sup> la crearea de instrumente complexe, cu aplicație multidisciplinară, care unifică mai multe programe într-o suită ce prezintă avantajul fiabilității și flexibilității. De asemenea, la nivel lingvistic se constată tendința de a elabora instrumente independente de limbă<sup>12</sup> astfel același instrument putând fi folosit pentru analiza oricărui text. Așa cum se vede din prezentarea de față, această tendință are anumite limite iar limba română prezintă o problemă complexă în contextul lingvisticii aplicate. Există astfel diverse probleme la nivelul utilizării acestor instrumente de analiză și recuperare a informației pe teren lingvistic românesc. Totuși, este de remarcat faptul că nu este imposibil de a utiliza aceste instrumente cu rezultate bune, în momentul în care se conștientizează care sunt aceste probleme.

De asemenea, se constată o tendință de a oferi acces direct utilizatorului unui instrument electronic de analiză textuală, dar un acces controlat. Astfel, **Corpografo** nu este doar un instrument pus la dispoziția comunității științifice pentru analiză și cercetare lingvistică. Prin permiterea unui acces on-line la **Corpografo** se pot obține informații asupra interesului comunității științifice asupra acestui instrument de lucru, asupra opțiunilor celor mai des utilizate și a scopului în care acest instrument este utilizat, a limbilor de lucru, a domeniilor de interes pentru analiză lingvistică, astfel acest instrument se convertește într-un instrument de măsură a necesităților existente în domeniul cercetării terminologice în special. De asemenea posibilitatea pe care acesta o oferă de a intra în contact direct cu cercetătorii care au proiectat instrumentul, asigură de asemenea feed-back-ul necesar pentru a dezvolta și îmbunătăți acest instrument electronic.

Lipsa unui asemenea instrument de analiză textuală și de explorare de corpus, dedicat limbii române sau a soluțiilor de incorporare a acesteia în cadrul altor instrumente de lucru deja existente, ne poate da o idee despre tendințele actuale în cercetarea lingvistică la nivel mondial și, de asemenea, limitele și carențele cercetării în acest domeniu la nivel local pentru limba română. Lipsa unor astfel de instrumente stă, în mod cert, la baza lipsei din peisajul lingvistic românesc a unor surse lingvistice de o calitate comparabilă cu cea a celor ce descriu alte limbi, materiale cum ar fi: *dicționare generale, dicționare specializate, dicționare-tezaur, glosare de termeni, corectoare automate, etc.* Considerăm că o îmbunătățire a acestui aspect al cercetării lingvistice și o extindere a utilizării instrumentelor de analiză și recuperare a informației bazată pe corpus poate avea ca rezultat o mai bună gestionare a informației lingvistice și, în consecință o mai bună calitate în ceea ce privește elaborarea de resurse și instrumente lingvistice ce descriu limba română.

---

<sup>11</sup> De exemplu nu se poate afirma că Corpografo sau Lexico3 sunt concordancier-e. Sunt și concordancier-e dar mai au și alte funcționalități care le fac să iasă din această clasificare foarte îngustă.

<sup>12</sup> Există instrumente lingvistice cu mai mare tradiție care sunt dedicate doar unei limbi (ex: *The British National Corpus* -<http://www.natcorp.ox.ac.uk/>) sau unui text (ex: concordanța textelor biblice sau cea a Constituției europene).

## Bibliografie

1. **Simple Concordance Program**, <http://www.textworld.com/>
2. **Corpografo**, [www.linguateca.pt/corpografo](http://www.linguateca.pt/corpografo)
3. **Lexico3**, <http://lexico3.no-ip.org/>
4. **The British National Corpus**, <http://www.natcorp.ox.ac.uk>
5. **Concordance biblique**, [http://www.lueur.org/bible/bible\\_rechercher.php](http://www.lueur.org/bible/bible_rechercher.php), pagina web a Bisericii baptiste protestante din Angers.
6. André Salem, **Approches quantitatives des corpus textuels**, conferință IULA, Barcelona, 2006.
7. Luigi Sansonetti, **Exploration textuelle d'interactions verbales entre un adulte et un enfant avec Lexico3**,
8. Andrea Kuncova, Aude Mansondieu, **Outils de statistique textuelle. Manuel d'utilisation abrégé (Dix premiers pas avec Lexico3)**, SYLED-CLA2T, Université de la Sorbonne Nouvelle – Paris 3
9. Belinda Maia, Luís Sarmento, **Gestor de Corpora – Um ambiente Web integrado para. Linguística baseada em Corpora**, [www.linguateca.pt/corpografo](http://www.linguateca.pt/corpografo)
10. Belinda Maia, Luís Sarmento, Diana Santos, **The Corpógrafo – a Web-based environment for corpora research**, [www.linguateca.pt/corpografo](http://www.linguateca.pt/corpografo)
11. Luís Sarmento (2004), **Relatório Técnico sobre o Corpógrafo**, <http://poloclup.linguateca.pt/docs/cg/>.