

Trotter, David/Bozzi, Andrea/Fairon, Cédric (éd.) (2016). *Actes du XXVII<sup>e</sup> Congrès international de linguistique et de philologie romanes (Nancy, 15-20 juillet 2013). Section 16 : Projets en cours ; ressources et outils nouveaux*. Nancy, ATILF : <http://www.atilf.fr/cilpr2013/actes/section-16.html>.

## CLRE. Corpus lexicographique roumain essentiel. 100 dictionnaires de la langue roumaine alignés au niveau de l'entrée

### 1. Introduction

Le projet *CLRE, Corpus lexicographique roumain essentiel. Les dictionnaires de la langue roumaine alignés au niveau de l'entrée* (2010–2013), continue une série d'initiatives qui visent l'informatisation de la lexicographie roumaine et qui se sont déroulées au sein de l'Institut de Philologie Roumaine « A. Philippide », de l'Académie Roumaine de Iasi. Ce dernier projet, CLRE, a été financé par le Conseil National de la Recherche Scientifique (CNCS) de Roumanie. Le collectif de recherche a été formé par des lexicographes de l'Institut de Philologie Roumaine « A. Philippide », de l'Académie Roumaine de Iasi (c'est-à-dire Elena Tamba, Ana Catană-Spenchiu, Marius Clim) et par des informaticiens de l'Université « Alexandru Ioan Cuza » de Iasi (Marius Răschip et Mădălin Ionel Patrașcu). Pendant cette période on a envisagé de créer un outil pour aider en particulier les lexicographes qui rédigent le *Dictionnaire-trésor de la langue roumaine*. CLRE représente un important outil pour les spécialistes et facilitera beaucoup le travail des lexicographes, parce que ce corpus offre surtout une perspective historique de la lexicographie roumaine et, implicitement, de la langue roumaine. Cette perspective signifie que le lexicographe aura la possibilité de voir comment un mot du vocabulaire roumain a évolué, s'il y a une variation sémantique, orthographique, ortho-épique ou morphologique. Ce corpus inclut divers types de dictionnaires : des dictionnaires généraux, des dictionnaires auxiliaires (de néologismes, d'étymologie, d'orthographe, orthoépique et de morphologie) et des dictionnaires spéciaux (c'est-à-dire des encyclopédies ou d'autres dictionnaires spécialisés, qui sont de grande importance pour la perspective diachronique sur la langue).

Mais cet outil sera également disponible au grand public qui aura accès à un grand nombre de dictionnaires, qui ont été inclus dans la base de données et alignés au niveau de l'entrée. De la Bibliographie de DLR (*Dictionnaire-trésor de la langue roumaine*) on a choisi 100 dictionnaires – 150.000 pages de dictionnaire. Cela constitue une première pour la langue roumaine, parce que ce corpus est à présent le plus grand et on peut toujours ajouter des nouveaux dictionnaires.

Ainsi, ce projet a eu pour buts :

- a) la réalisation d'une base de données qui comprenne les dictionnaires essentiels de la Bibliographie du DLR, alignés au niveau de l'entrée. La Bibliographie du *Dictionnaire de la langue roumaine* contient plus de 4000 titres des textes représentatifs pour la culture roumaine et elle est utilisée pour extraire des citations pour exemplifier les sens des mots-titre ;
- b) l'élaboration des logiciels qui permettent la consultation interactive de ce corpus, qui peut constituer un cadre moderne de travail et de recherche lexicographique ;
- c) la réalisation d'une liste de mots quasi-exhaustive, pour la langue roumaine, à partir du corpus aligné ;
- d) l'augmentation de la visibilité de l'activité de recherche des professionnels linguistes et informaticiens, dans le domaine de la langue roumaine, pour promouvoir les moyens informatiques de traitement linguistique créés dans le projet.

## 2. La digitalisation des dictionnaires

En ce qui concerne les dictionnaires qui font partie du corpus CLRE, ils ont subi un processus de digitalisation, qui signifie qu'ils ont été numérisés à l'aide d'un scanner en V et puis traités du point de vue informatique, à l'aide d'un logiciel de reconnaissance optique des caractères (Abby Fine Reader 9) pour reconnaître les entrées et pour ensuite indexer les résultats dans une base de données XML.

Cette étape a généré de nombreux problèmes à cause du type d'alphabet utilisé, mais également à cause du format du dictionnaire.

Pour la langue roumaine, les dictionnaires ont été écrits en utilisant trois types d'alphabets : l'alphabet cyrillique, l'alphabet de transition ou l'alphabet latin. Les premiers deux types ont été un véritable défi pour les informaticiens, vu que le logiciel de reconnaissance optique a été presque inefficace pour ces textes. On verra, dans les deux figures suivantes, des captures de ces dictionnaires.

S	
Anașajpa, f. pl. i. baș Țerâgjen, Țerâgjen.	Animer, m. % . f. pl. pî. 70. part. von anikt.
Anașajropia, m. oșpe, f. pl. i. e. baș Țerâgjen, Țerâgjen; nișli; Țerâgjen.	Animș, oșpe, oșp. v. 1) anșiften, baș Țerâgjen, Țerâgjen; nișli; Țerâgjen.
Anașajr, m. % . f. pl. pî. 70. part. von anașpece.	Aniu, m. -pl. i. baș Țerâgjen, Țerâgjen.
Anașajropa, f. pl. i. f. anașajpa, Țerâgjen, Țerâgjen; nișli; Țerâgjen.	Anrepa, m. pl. n. ein langes Țerâgjen.
Anașajka, f. pl. e. Țerâgjen, Țerâgjen; nișli; Țerâgjen; Țerâgjen.	Anrepa, f. pl. e. baș Țerâgjen, Țerâgjen.
Anașajka, adv. baș Țerâgjen, de anașajka, nișli; Țerâgjen.	Anrepa, adv. 1) nișli; Țerâgjen. 2) baș Țerâgjen, Țerâgjen. 3) nișli; Țerâgjen.
Anașajka, oșpe, oșp. v. 1) baș Țerâgjen, Țerâgjen. 2) nișli; Țerâgjen. 3) nișli; Țerâgjen. 4) nișli; Țerâgjen.	Anrepa, f. pl. i. 1) Țerâgjen, Țerâgjen. 2) Țerâgjen, Țerâgjen.
Anașajka, f. a. e. Țerâgjen.	Anrepa, m. % . f. pl. pî. 70. 1) part. von anașpece. adv. geiș Țerâgjen, Țerâgjen.
Anașajka, m. oșpe, f. Țerâgjen.	Anrepa, adv. vor geiș Țerâgjen.
Anașajka, m. pl. i. baș Țerâgjen.	Anrepa, adv. nișli; Țerâgjen, Țerâgjen; nișli; Țerâgjen.
Anașajka, f. pl. e. geiș Țerâgjen.	Anrepa, oșpe, oșp. v. 1) nișli; Țerâgjen, Țerâgjen. 2) baș Țerâgjen, Țerâgjen. 3) nișli; Țerâgjen, Țerâgjen.
Anașajka, f. pl. e. geiș Țerâgjen.	Anrepa, f. pl. i. 1) Țerâgjen, Țerâgjen. 2) Țerâgjen, Țerâgjen. 3) nișli; Țerâgjen, Țerâgjen.
Anașajka, f. pl. e. geiș Țerâgjen.	Anrepa, m. % . f. pl. pî. 70. part. von anașpece.
Anașajka, f. pl. e. geiș Țerâgjen.	Anrepa, m. pl. e. baș Țerâgjen, Țerâgjen.
Anașajka, f. pl. e. geiș Țerâgjen.	Anrepa, m. pl. e. baș Țerâgjen, Țerâgjen.

Figure 1: Andrei Iser, *Vocabular românesc-nemțesc*. Brașov, I. Gött, 1850.

45

Блажд, з. адж. saust, mltb, jahm.

Блаждець, ф. die Saustmuth, Wiltde; Zahmheit.

Боа, ф. die Niesenschlange.

Боавъ, е. ф. das Korn; die Beere.

Боавъ, е. ф. die Krankheit.— de anz, die Wassersucht. — лмоаскъ, die syphilitische Krankheit. — конилоп, die Gallsucht, Epilepsie; — аавъ, der weiße Fluß; — скалъ, die Durrsucht der Kinder.

Боавъ, е. ф. die Beere. podxa katorna kopъchel. — die Bombe.

Боар, i. б. der Ochsentreiber, Ochsenhüter, Ochsenhirt. —, der Wärenthüter (ein Gestirn), o zodie de 93 stele.

Боб, spl. n. die Bohne, Saubohne.

Бовиць, е. ф. die Punkte auf einem Zeuge, Stoffe. —, das Körnchen, Beerlein.

Бовок, чi. б. die Knoche. —, das Gänßchen, die junge Ente.

Боворникъ, че. ф. der Ehrenpreis (Pflanze). o еспianъ.

Бовотазъ, ф. das Dreikönigsfest, das Fest der heiligen drei Könige.

Борат, з. адж. reich, wohlhabend, begütert.

Бораче, i. ф. der Lebuchen, Pfefferkuchen; der Käseuchen.

Боръgie, иl. ф. der Reichthum, die Wohlhabenheit.

Бодаспеѝ, i. б. der Laucher. o uacepe.

Бозиѝ, в. der Altich, en cad.

Болва, e.l.e. ф. die Farbe.

Болвангепie, иl. ф. die Färberei.

Болвангѝ, i. н. der Färber.

Болер, i. б. der Bojar.

Болерекъ, оакъ. адж. bojarisch, herrschaftlich.

Болери, в. abeln, in den Adelsstand, Bojarenstand erheben.

Болерие, ф. die Bojarenwürde, der Adel. —, die venerische Krankheit.

Figure 2 : S. Petri, *Vocabular portativ românesc-nemțesc și nemțesc-românesc*. Sibiu, G. de Clozius, 1861.

Le processus de reconnaissance optique des caractères n'a pas offert des résultats utilisables, parce que chaque dictionnaire en cyrillique a son propre type de caractères. Cette diversité nécessitera une adaptation du logiciel de reconnaissance optique pour chaque œuvre lexicographique. Et aussi, comme on peut y voir dans les deux captures ci-dessus, vu qu'il s'agit de dictionnaires bilingues, on a des caractères spécifiques pour chaque langue, ce qui rend la tâche du logiciel ABBYY, de reconnaissance des caractères, presque impossible. Dans le cas de dictionnaires en alphabet cyrillique et de transition, on a opté pour la translittération manuelle des entrées. De cette manière on peut aligner les entrées de ces dictionnaires avec les autres dictionnaires. Mais l'utilisateur spécialisé aura également à sa disposition la capture numérisée de l'entrée du mot choisi. Ainsi il pourra utiliser toutes les informations contenues dans ces dictionnaires.

Le format des dictionnaires a constitué tout un autre défi pour les informaticiens, qui ont dû créer un logiciel de délimitation des entrées assez flexible.



Figure 3 : Nifon Bălășescu, *Dictionariu româno-frances*. Volum I. Fascicul I. Bucuresci, Tipografia Mitropolitului Nifon, 1859.

Dans ce cas on remarque que l'auteur du dictionnaire a soumis neuf entrées à un autre mot-titre. Mais à cause de cela, ces neuf entrées manquent d'information morphologique, à une seule exception, et on a aussi des problèmes à les délimiter.

AFLUIREA REZERVISTILOR	30
<p><b>AFLUIREA REZERVISTILOR.</b> Deplasarea rezerviștilor la unitățile la care sînt repartizați potrivit planului de mobilizare a cetățenilor cu obligații militare. Poate fi : directă, cînd deplasarea se face din comune sau de la organizațiile socialiste la unități ; prin puncte de adunare, cînd cetățenii respectivi se grupează mai întîi în anumite locuri dinainte stabilite, de unde se trimit la unități.</p>	
<p><b>AGENT DE TRANSMISIUNI.</b> Persoană prin care se transmit diferite mesaje, verbale sau scrise, între diferiți comandanți (șefi) militari, realizînd, în acest fel, legătura între aceștia. Un A.d.T. poate fi militar sau luptător din formațiunile de apărare și acționează pe jos sau transportat (cu motocicletă, auto, avion, elicopter).</p>	
<p><b>AGENTIE DE CAMPANIE A BĂNCII NAȚIONALE A REPUBLICII SOCIALISTE ROMÂNIA.</b> Instituție bancară organizată de Banca Națională, cu atribuții de executare de casă pentru unități și M.U., potrivit dispozițiilor tehnice primite de la filiala de campanie și dispozițiilor organului financiar al ordonatorului de credite pe lingă care este repartizată.</p>	
<p><b>AGRAFĂ.</b> Material genistic din sîrmă curbată și ascuțită la ambele capete, folosit la prinderea sîrmei (ghimpate) pe pari sau țărui, pentru realizarea gardurilor și rețelelor de sîrmă.</p>	
<p><b>AGRESIUNE ARMATĂ.</b> După textul „Definiției agresiunii armate“ adoptat de Comitetul special al O.N.U. la 15 aprilie 1974 și însușit de Adunarea Generală în același an agresiunea este „folosirea forței armate de către un stat împotriva suveranității, integrității teritoriale sau independenței politice a unui stat sau în orice mod incompatibil cu Carta Națiunilor Unite“. Statul care a recurs primul la forța armată, în contradicție cu Carta a O.N.U., este calificat agresor. Se consideră A.A. acel act comis cu intenție, de o anumită gravitate, care se deosebește de așa-numitele „incidente minore“ (de pildă, incidente sporadice de frontieră). Cazurile tipice de A.A. sînt : invadarea sau atacarea teritoriului unui stat, bombardarea, blocarea porturilor sau coastelor sale de către</p>	

Figure 4: *Lexicon militar*. București, Editura Militară, 1980.

Dans cette figure on remarque une autre spécificité lexicographique, c'est-à-dire qu'il n'y a pas un seul mot-titre, mais un syntagme, une construction fixe ou une expression. Dans ce cas il est impossible d'englober tous les mots dans une entrée et le lexicographe choisira le mot-titre pour pouvoir aligner ce dictionnaire de cette façon.



Figure 5: *Dicționar de informatică*. București, Editura Științifică și Enciclopedică, 1981.

Dans ce dictionnaire il y a une différenciation entre deux types d'entrées. La première est écrite avec des minuscules, ce qui dénote qu'on a une expression ou un syntagme qui utilise un mots-titre déjà défini au-dessus. Dans la deuxième situation, le mot-titre est utilisé aussi dans le texte de la définition de la même manière formelle. Vu que le mot-titre est écrit plusieurs fois, cela pose des problèmes lors de la délimitation des entrées, le logiciel pouvant les considérer comme des entrées distinctes.

### 3. La validation de la segmentation des dictionnaires

Ce processus a été réalisé par les lexicographes, mais également avec l'aide de quelques volontaires qui possédaient des connaissances philologiques assez solides.

<p>45</p> <p><input checked="" type="checkbox"/> AMERICANIZA</p> <p><input checked="" type="checkbox"/> AMBULACRĂR, <i>-Ă, ambulacrar, -e adj.</i> (fr. <i>ambu-lacraire</i>) Cu privire la ambulacru, de ambulacru.</p> <p><input checked="" type="checkbox"/> AMBULĂCRU, <i>ambulacre</i> s.n. (fr. <i>ambulacre</i>) Tub subțire situat pe fața inferioară a corpului echino- dromilor și terminat printr-o ventuză, care servește la locomotie, respirație și pipăit.</p> <p><input checked="" type="checkbox"/> AMBULANT, <i>-Ă, ambulanti, -le adj.</i> (fr. <i>ambulant</i>, lat. <i>ambulans, -ntis</i>) Care se deplasează dintr-un loc în altul, care nu are un loc de reședință fix.</p> <p><input checked="" type="checkbox"/> AMBULANȚA, <i>ambulante</i> s.f. (fr. <i>ambulance</i>) Mijloc de transport al bolnavilor, al ranților etc.; sa!-vare (3).</p> <p><input checked="" type="checkbox"/> AMBULATORIU, <i>-IE, ambulatorii</i> (fr. <i>ambula-toire</i>, lat. <i>ambulatorius</i>) 1. Adj. (Despre tratamente medicale) Care nu necesită spitalizare. 2. S.n. Dis- pensar.</p> <p><input checked="" type="checkbox"/> AMBUSCADA, <i>ambuscade</i> s.f. (fr. <i>ambuscade</i>) Mar nevră de atac prin surprindere asupra unui inamic în mișcare.</p> <p><input checked="" type="checkbox"/> AMBUSCĂT, <i>-Ă, ambuscăți, -te adj.</i>, s.m. și f. (fr. <i>ambuscade</i>) (Militar) scutit de obligațiile periculoase din timpul războiului. <input checked="" type="checkbox"/> AMBUSURI, <i>ambusuri</i> s.f. (fr. <i>embouchure</i>) Parte a unui instrument muzical prin care se suflă aerul cur-gura.</p> <p><input checked="" type="checkbox"/> AMBUTEIĂ, <i>ambuteiaje</i> vb. I (fr. <i>embouteiller</i>) A blo-ca o cale (rutiera, maritimă etc.) cu vehicule, na-ve etc.</p> <p><input checked="" type="checkbox"/> AMBUTEIAJ, <i>ambuteiaje</i> s.n. (fr. <i>embouteillage</i>) 1. Imbutelire. 2. Blocare a circulației rutiere sau navale din cauza aglomerației.</p> <p><input checked="" type="checkbox"/> AMBUTISĂ, <i>ambutitoz</i> vb. I (cf. fr. <i>embouter</i>) A su-pune un metal unor socuri mecanice pentru a-l da o anumită formă.</p>	<p>45</p> <p><b>AMERICANIZĂ</b></p> <p><b>AMBULACRĂR, -Ă, ambulacrar, -e adj.</b> (fr. <i>ambu-lacraire</i>) Cu privire la ambulacru, de ambulacru.</p> <p><b>AMBULĂCRU, ambulacre</b> s.n. (fr. <i>ambulacre</i>) Tub subțire situat pe fața inferioară a corpului echino-dromilor și terminat printr-o ventuză, care servește la locomotie, respirație și pipăit.</p> <p><b>AMBULĂNT, -Ă, ambulanti, -le adj.</b> (fr. <i>ambulant</i>, lat. <i>ambulans, -ntis</i>) Care se deplasează dintr-un loc în altul, care nu are un loc de reședință fix.</p> <p><b>AMBULANTORIU, -IE, ambulatorii</b> (fr. <i>ambula-toire</i>, lat. <i>ambulatorius</i>) 1. Adj. (Despre tratamente medicale) Care nu necesită spitalizare. 2. S.n. Dis-pensar.</p> <p><b>AMBUSCĂDĂ, ambuscade</b> s.f. (fr. <i>ambuscade</i>) Mar-nevră de atac prin surprindere asupra unui inamic în mișcare.</p> <p><b>AMBUSCĂȚ, -Ă, ambuscăți, -te adj.</b>, s.m. și f. (fr. <i>ambuscade</i>) (Militar) scutit de obligațiile periculoase din timpul războiului.</p> <p><b>AMBUSURĂ, ambusuri</b> s.f. (fr. <i>embouchure</i>) Parte a unui instrument muzical prin care se suflă aerul cu-gura.</p> <p><b>AMBUTEIA, ambuteiaje</b> vb. I (fr. <i>embouteiller</i>) A blo-ca o cale (rutieră, maritimă etc.) cu vehicule, na-ve etc.</p> <p><b>AMBUTEIAJ, ambuteiaje</b> s.n. (fr. <i>embouteillage</i>) 1. Imbutelire. 2. Blocare a circulației rutiere sau navale din cauza aglomerației.</p> <p><b>AMBUTISĂ, ambutitoz</b> vb. I (cf. fr. <i>embouter</i>) A su-pune un metal unor socuri mecanice pentru a-l da o anumită formă.</p> <p><b>AMENAJĂ, amenajă</b> vb. I (fr. <i>aménager</i>) Așeză, găsi un organig (un interior, o grădini etc.) în vederea unui scop, a unei utilizări.</p> <p><b>AMENAJABIL, -Ă, amenajabil, -e adj.</b> (fr. <i>aménagement</i>) Care se poate amenaja.</p> <p><b>AMENAJAMENT, amenajamente</b> s.n. (fr. <i>aménagement</i>) Organizare în scopul ameliorării condițiilor de viaține (a exploatareii forestiere, a turismului etc.).</p> <p><b>AMENDĂ, amendă</b> vb. I (<i>amendă</i> = -) A aplica o amendă; a penaliza.</p> <p><b>AMENDĂ, amendă</b> vb. I (fr. <i>amender</i>) 1. A amelo-rea, a aduce îmbunătățiri, a curăți (un text, lege). 2. A îmbunătăți totul prin încorporarea unor sub-stanțe nutritive.</p> <p><b>AMENDABIL, -Ă, amendabil, -e adj.</b> (fr. <i>amenda-ble</i>) 1. Care se poate amenda. 2. Care este modificat prin amendare.</p> <p><b>AMENDAMENT, amendamente</b> s.n. (fr. <i>amendement</i>) 1. Propunere de îmbunătățire a unui proiect de lege. 2. Operațiune de ameliorare a solului. 3. Sub-stanță necesară ameliorării condițiilor de existență a plantelor.</p> <p><b>AMENDĂ, amenajă</b> s.f. (fr. <i>amende</i>) 1. Pedepșă în bani. 2. <i>amendă onorabilă</i> = recunoașterea în public a unei greșeli.</p> <p><b>AMENDABIL, -Ă, amendabil, -e adj.</b> (fr. <i>amenda-ble</i>) 1. Care se poate amenda. 2. Care este modificat prin amendare.</p> <p><b>AMENDAMENT, amendamente</b> s.n. (fr. <i>amendement</i>) 1. Propunere de îmbunătățire a unui proiect de lege. 2. Operațiune de ameliorare a solului. 3. Sub-stanță necesară ameliorării condițiilor de existență a plantelor.</p> <p><b>AMENDĂ, amenajă</b> s.f. (fr. <i>amende</i>) 1. Pedepșă în bani. 2. <i>amendă onorabilă</i> = recunoașterea în public a unei greșeli.</p> <p><b>AMENDABIL, -Ă, amendabil, -e adj.</b> (fr. <i>amenda-ble</i>) 1. Care se poate amenda. 2. Care este modificat prin amendare.</p> <p><b>AMENDAMENT, amendamente</b> s.n. (fr. <i>amendement</i>) 1. Propunere de îmbunătățire a unui proiect de lege. 2. Operațiune de ameliorare a solului. 3. Sub-stanță necesară ameliorării condițiilor de existență a plantelor.</p> <p><b>AMENAJĂ, amenajă</b> vb. I (<i>amendă</i> = -) A aplica o amendă; a penaliza.</p> <p><b>AMENDĂ, amendă</b> vb. I (fr. <i>amender</i>) 1. A amelo-rea ca fiind încălț, dispune pe o încă, care atâră ca un cîmpac, mla, mîpăc.</p> <p><b>AMENAJĂCĂ, amenajăcă</b> s.f. (fr. <i>amenajăcă</i>) (La pl.) Familii de plante terose care se înfruntesc în formă de arbuști (și la sg.) plantați care fac parte din rezervația forestieră.</p> <p><b>AMERICANISM, americanisme</b> s.n. (fr. <i>américanisme</i>) 1. Ansamblu de trăsături proprii americanilor; loc: trăsături a comportamentului americanilor. 2. Cu-șuș, expresie din engleza americană.</p> <p><b>AMERICANIZĂ, americanizat</b>, <i>-ate</i> s.m. și f. (fr. <i>américanisation</i>) 1. Specialist în americanizare. 2. Adapt al americanizării (1).</p> <p><b>AMERICANIZETĂ, s.f.</b> (germ. <i>Amerikanizetät</i>) Simbol simbolic al culturii și al civilizației ameri-cane.</p> <p><b>AMERICANIZĂ, americanizat</b> vb. I (fr. <i>américani-sation</i>) A da stă la un caracter american; a se com-porta ca american.</p>
--	--

Figure 6: La validation de la segmentation en entrées d'une page d'un dictionnaire

Dans cette étape, les informaticiens ont créé un logiciel qui essaie de reconnaître chaque entrée d'après ses caractéristiques formelles. Ce logiciel signale les entrées qui sont considérées comme véritables, mais il souligne également les autres mots qui peuvent être envisagés comme des entrées. Le lexicographe doit valider chaque entrée à son tour. Par cette démarche on valide les entrées correctes, mais, en même temps, chaque entrée est corrigée pour pouvoir les aligner facilement. Pour chaque page validée, le lexicographe a accès également à la page numérisée, pour pouvoir corriger les erreurs. S'il y a des mots qui ne peuvent pas être validés comme des entrées, alors le lexicographe peut choisir de valider partiellement la page respective. Après ça, l'informaticien traitera de nouveau cette page pour reconnaître les entrées restantes.

#### 4. La segmentation des entrées de dictionnaire

Un autre élément distinct et définitoire pour le projet CLRE c'est le fait qu'il met à disposition des utilisateurs la variante numérisée du dictionnaire, pour pouvoir vérifier chaque entrée. Ainsi, chaque entrée est segmentée de façon automatique après l'opération de validation réalisée par les lexicographes.

Mais l'exactitude de la segmentation est liée, bien sûr, au format des dictionnaires. Et pour cela les informaticiens ont dû adapter leur logiciel pour toute une variété des formats. En général, le validateur statistique utilise les suivants éléments pour faire la segmentation : l'alignement à gauche et à droite du texte et la police (le style de la police – italique, souligné, gras –) et les types des caractères (majuscules, minuscules).

răsărit. Ortus. Oriens.	răgăitură. Eructatio.
3730 răsăd. Planta.	Râm. Roma.
rășină. Pix. Resina.	3775 rămă. Lumbricus terrae.
rășipă. Dissipatio. Ruina.	rămător. Porcus.
rășipitură. Idem.	rămuitor. Idem.
rășipitor. Dissipator. Dirutor.	rămuiesc. Verto.
3735 rășipită. –	răpă. Ripa.
rășpândă. –	3780 răpesc. –
rășpund. Respondeo.	răpiti. -ă. –
răspuns. Responsum.	răpitor. –
rășpesc. Dissipo. Diruo.	răs. Ritus.
3740 resteu. –	răs. –
răstescu-mă. Comminor.	3785 rât. Nasus porcinus.
răstesc. –	râu. Rivus. Fluvius.
răstignesc. Crucifigo.	răușor. Rivulus. Fluviolus.
răstignitură. Crucifixio. (122)/	răză. Tela viilis et lacera.
3745 răstignitor. Crucifixor.	răzos. Lacer.
răstorn. Inverso.	3790 rob. Captivus. Mancipium.
răsturnat. Inversus. Inversio.	robesc. Captivo.
răsturnătură. Inversio.	robie. Captivitas.
răstunchiat. –	robică. Captiva.
3750 răsură. Rosa campestris.	robotar. Laborator.
rătedz. –	3795 robotor. Captivator.
rătedz. –	rod. Rodo.
rătedzat. –	rod. Fructus. (124)/
rătedzătură. –	roadă. Proles.
3755 rătăiesc. Aberro.	rodesc. Fructifico.
rătăicură. Aberratio. Error.	3800 rodesc. Creo.
rătund. -ă. Rotundus. -a.	roditor. Fructifer.
revenală. Humor. Humiditas.	roditor. Creator.
revărs. Variegatio.	rodzătură. Rosio.
3760 revărsat. Variegatus.	roună. Ros.
revărsătură. Variegatio.	3805 rouărendză. Rorat.
răbesc. –	rog. Rogo. Oro.
răzbitură. –	rogoz. Carectum. Săg.
război. Praelium. Pugna.	rogojină. Storea.
3765 războială. –	roi. Examen ap<i>-um.
războiesc. Praelior. Pugno.	3810 roiesc. Examen emitto.
răzbuună. Serenat.	roc. Finis mundi. Consummatio
răzmiriță. Bellum.	seculi.
răzmirițescu-mă. – (123)/	rochie. –
3770 răd. Rideo.	Romn. Roma.
rădzător. Risibilis.	romon. –
răgăiesc. Eructo.	3815 romonită. –
	ros. -ă. Corrosus. -a.

Figure 7: *Dictionarium Valachico-Latinum. Primul dicționar al limbii române.*  
 Studiu introductiv, ediție, indici și glosar de Gh. Chivu. București,  
 Editura Academiei Române, 2008 [cca 1650]

Dans ce dictionnaire, considéré comme le premier dictionnaire de la langue roumaine, l'auteur a réalisé aussi une numérotation des mots introduits et a marqué cela en faisant une note tous les cinq mots et en inscrivant le nombre correspondant. Le résultat est que le logiciel de segmentation reconnaît les chiffres, et non pas les mots, comme étant les entrées.



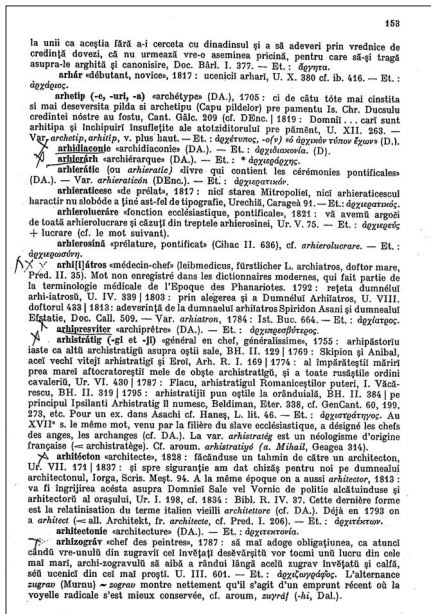


Figure 8: Găldi Ladislav, *Les mots d'origine néogrecque en roumain à l'époque des Phanariotes*. Budapest, 1939.

Il y a aussi des dictionnaires, en particulier des dictionnaires anciens, qui ne sont pas faciles à trouver et qui en plus présentent de « traces » d'utilisation. Dans ce cas, c'est presque impossible, pour le logiciel, de reconnaître une entrée, si devant l'entrée il y a toute sorte de notations faites par de divers utilisateurs. Dans ce cas, c'est le validateur humain qui décide où commence l'entrée et qui corrige les éventuelles erreurs.

Les dictionnaires spécialisés posent de divers problèmes ayant trait à leur spécificité. Il y a toute sorte de marques d'orthographe, de morphologie et bien sûr, l'accent.

Après la validation, on essaye d'améliorer les résultats du logiciel de reconnaissance optique des caractères en passant par un nouveau processus de reconnaissance des caractères des dictionnaires analysés.

Les deux dernières étapes du projet ont envisagé l'alignement des entrées des dictionnaires et les possibilités d'interroger cette base de données.

En ce qui concerne l'alignement des dictionnaires, après avoir fini le processus de segmentation des entrées à l'aide de l'interface spécialement créée, on a réalisé un alignement primaire de tous les mots des dictionnaires et ensuite fait un alignement avec l'eDTLR<sup>1</sup>. Les entrées alignées sont incluses dans la base de données qui

<sup>1</sup> Le projet national eDTLR. *Le Dictionnaire (Trésor) de la Langue Roumaine en format électronique* a été financé par le Conseil National du Management des Projets (CNMP), pour la période 2007-2010 et il représente l'un de plus importants projets de type collaboratif des

permettra faire des recherches ultérieures, en utilisant des divers critères. Ainsi, on peut chercher dans un seul dictionnaire ou dans tous les dictionnaires introduites en CLRE, mais on peut effectuer également une recherche selon d'autres critères : on peut chercher les mots commençant avec la même lettre ou qui incluent un certain groupe de lettres. Le résultat de la recherche permet la visualisation de l'entrée du dictionnaire, tout d'abord en format image. Sur la même page, on peut voir aussi la variante reconnue par le logiciel ABBYY, mais elle ne peut pas être utilisée dans le cas des dictionnaires en alphabet cyrillique ou de transition. L'utilisateur peut voir également les informations concernant le dictionnaire qui a servi de source pour le mot en question.

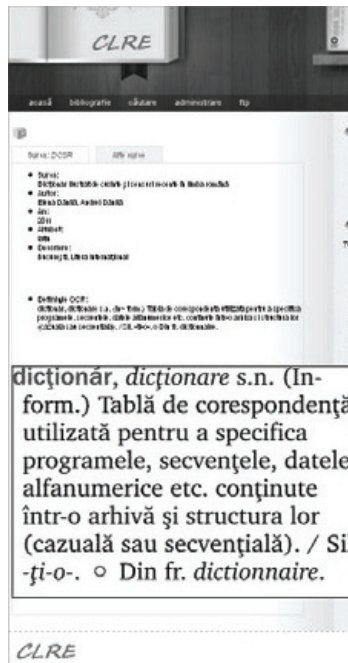


Figure 9: Une capture d'écran de l'article « dictionnaire » d'un lexique écrit en alphabet latin.

dernières années, qui entremêle l'expérience des lexicographes et des informaticiens pour créer des outils électroniques. Le projet a été coordonné par dr. Dan Cristea de la Faculté d'Informatique de l'Université « Alexandru Ioan Cuza » de Iasi.

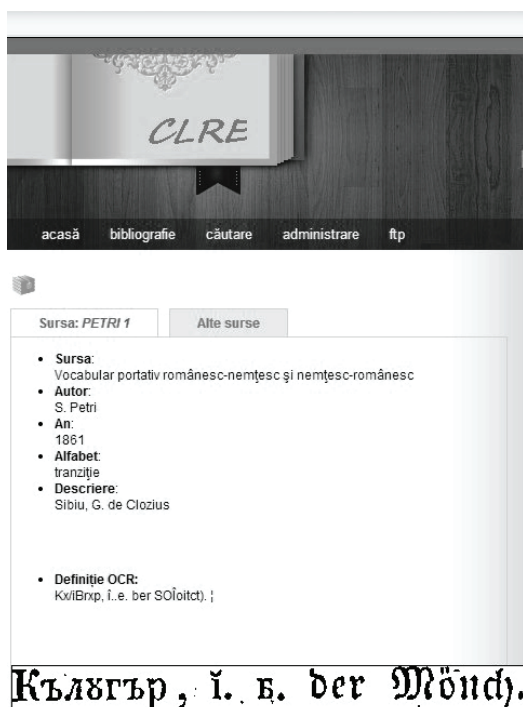


Figure 10: Une capture d'écran de l'article « moine » d'un dictionnaire écrit en alphabet cyrillique.

## 5. Conclusions

Tant le spécialiste que l'utilisateur non initié peuvent consulter virtuellement, à l'aide du CLRE, une multitude de dictionnaires qui, autrement, ne seraient disponibles que dans les bibliothèques.

Ce projet est lié avec d'autres projets en ligne réalisés en Roumanie (eDTLR et la troisième édition du dictionnaire RDW<sup>3</sup>, version informatisée) mais également avec le processus de réalisation du *Dictionnaire de la Langue Roumaine*, soutenu par l'Académie Roumaine. Il y a aussi une liaison entre CLRE et d'autres projets internationaux : par exemple avec DÉRom<sup>2</sup> (le CLRE étant considéré comme source de documentation) ou avec *ENeL: European Network of e-Lexicography* qui est une action de la ligne de financement COST et qui cherche à créer un réseau paneuropéen des ressources lexicographiques pour l'usage académique et pour le public. Ces collaborations servent à confirmer l'importance de cette ressource électronique pour la langue roumaine.

<sup>2</sup> *Dictionnaire Étymologique Roman* (DÉRom) ([www.atilf.fr/DERom/](http://www.atilf.fr/DERom/)).

L'information accessible par l'intermédiaire du CLRE peut être comparée avec celle offerte du modèle des corpus lexicographiques créés pour d'autres langues: *Le rayon des dictionnaires*, <http://www.atilf.fr/> – collection de dictionnaires informatisés français, du XVI<sup>e</sup> jusqu'au XX<sup>e</sup> siècle ou *Das Wörterbuchnetz*; *Nuevo tesoro lexicográfico de la lengua española*, <http://buscon.rae.es/ntlle/SrvltGUILoginNtll> – base de données qui inclut les versions facsimilées de tous les dictionnaires édités et publiés par la Real Academia Española; <http://germazope.uni-trier.de/Projects/WBB/> – réseau de dictionnaires de langue allemande, créée à l'Université Trier d'Allemagne, etc.

De cette façon, CLRE offre aux spécialistes intéressés par l'étude de la langue roumaine, aux romanistes, mais aussi aux spécialistes informaticiens qui désirent analyser le traitement du langage naturel à partir de la langue roumaine, la désambiguïsation sémantique, etc. – un instrument de travail extrêmement utile.

L'Académie Roumaine, Filiale de Iasi	Marius-Radu CLIM
L'Académie Roumaine, Filiale de Iasi	Elena TAMBA,
L'Institut de Philologie Roumaine « A. Philippide »	Ana CATANĂ-SPENCHIU
L'Institut de Philologie Roumaine « A. Philippide »	Mădălin PĂTRAȘCU

## Références bibliographiques

- Atkins, B.T.S., Rundell, M. 2008. *The Oxford Guide to Practical Lexicography*. Oxford, Oxford University Press.
- Clim, Marius/Dănilă, Elena/Haja, Gabriela, 2008. «Premise ale informatizării cercetării lexicografice academice românești», in: *Limba română. Dinamica limbii, dinamica interpretării*, București, Editura Universității din București, 585-591.
- Dănilă, Elena, 2013. «Corpus lexicographique roumain essentiel. Les dictionnaires de la langue roumaine alignés au niveau de l'entrée», in: Herrero, Casanova/ Calvo Rigual, Cesareo (ed.), 2013. *Actas del XXVI Congreso Internacional de Lingüística y de Filología Románicas 6-11 septiembre 2010, Valencia* (6-11 septembre 2010), Berlin/New York, La maison d'édition Walter de Gruyter, volum VIII, 125-134.
- Dănilă, Elena/ Clim, Marius-Radu/ Catană-Spenchiu, Ana, 2011. «Towards a Romanian Lexicographic Corpus», *Philologica Jassyensia*, An VII, Nr. 2 (14), 191-198.
- Tamba, Dănilă Elena/ Clim, Marius-Radu/ Pătrașcu, Mădălin/ Catană-Spenchiu, Ana, 2012. «The Evolution of the Romanian Digitalized Lexicography. The Essential Romanian Lexicographic Corpus», in: Fjeld, Ruth Vatvedt, Torjusen, Julie Matilde, (ed.), 2012. *Proceedings of the 15 th EURALEX International Congress, 7-11 august 2012, Oslo*, Press Reprosentrales, UiO, 1014-1017 ; on-line : [http://www.euralex.org/proceedings-toc/euralex\\_2012/](http://www.euralex.org/proceedings-toc/euralex_2012/) .