

UTILIZAREA RESURSELOR REUTILIZABILE ALE LIMBII ROMÂNE ÎN PROCESUL DE RECUNOAȘTERE A TEXTELOR

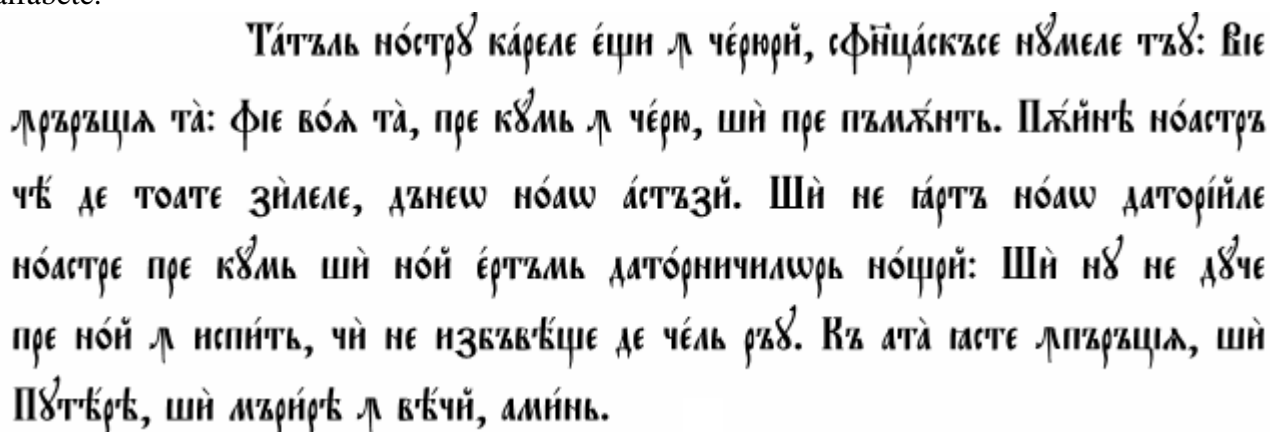
E. Boian, C. Ciubotaru, S. Cojocaru, A. Colesnicov, L. Malahov
Institutul de Matematică și Informatică al AȘM, Moldova
lena@math.md

Abstract. *The paper describes the technology for text digitalization and converting of scanned Romanian texts written with the Cyrillic letters. Full-text digitalization is performed using lexicon Reusable Romanian Linguistic Resources (RRLR) and spell checker RomSP. This technology can be applied for other kind of resources like some specific Romanian records written in the Cyrillic script.*

Cuvinte cheie: *digitalizarea textului, dechirilizarea textului, recunoașterea textului, lexicon lingvistic.*

I. Introducere

Pentru a avea acces la tezaurul cultural, opere literare vechi, materialele vechi din arhive, creații ale autorilor români și moldoveni, resurse folclorice specifice, este necesar să se creeze biblioteci virtuale care să includă aceste surse. O mare parte din aceste surse pot fi manuscrise sau texte tipărite cu litere chirilice. Există mai multe versiuni de alfabet chirilic folosit la scrierea textelor românești [1]. Vom menționa doar alfabetul folosit pe tot teritoriul românesc în secolele XIV-XV (până în anul 1862) și cel folosit în Republica Sovietică Socialistă Autonomă Moldovenească începând cu anul 1924 (cu o întrerupere între 1932-1938) și, ulterior, în Republica Sovietică Socialistă Moldovenească. În figurile 1 și 2 vedem exemple de texte scrise cu aceste alfabete.



ТѢТЪЛЪ НОСТРЪ КАРЕЛЕ ЕЦИН Л ЧЕРЮРИ, СФНЦАСКЪСЕ НЪМЕЛЕ ТЪЪ: ВІЕ
ЛРЪРЪЦІА ТѢ: ФІЕ ВОЛ ТѢ, ПРЕ КЪМЪ Л ЧЕРЮ, ШІН ПРЕ ПЪМЪНТЬ. ПЪИИНЪ НОСТРЪ
ЧЪ ДЕ ТОАТЕ ЗІЛЕЛЕ, ДЪНЕУ НОЛЪ АСТЪЗІ. ШІН НЕ ІАРТЪ НОЛЪ ДАТОРІІЛЕ
НОАСТРЕ ПРЕ КЪМЪ ШІН НОЙ ЕРТЪМЪ ДАТОРНИЧІЛОРЪ НОЦІРІ: ШІН НЪ НЕ ДЪЧЕ
ПРЕ НОЙ Л ІСПІТЬ, ЧІН НЕ НЗЪВЪКІЩЕ ДЕ ЧЕЛЪ РЪЪ. КЪ АТѢ ІАСТЕ ЛПЪРЪЦІА, ШІН
ПЪТЪБЪРЪ, ШІН МЪРІРЪ Л ВЪЧІЙ, АМІНЪ.

Fig.1. Rugăciunea „Tatăl nostru”.

Ш’ачел реже-ал поезией, вечник тынэр ши фериче,
Че дин фрунзе ыць дойнеште, че ку флуерул ыць зиче...

Fig. 2. Mihai Eminescu, „Epigonii”.

Pentru a le include în biblioteci virtuale ele trebuie să fie digitalizate (scanate, recunoscute și reprezentate ca texte cu litere latine). În cele ce urmează acest proces îl vom numi dechirilizare. În

lucrare se propune o tehnologie pentru dechirilizarea textelor, inclusiv a manuscriselor. Efectuarea acestor lucrări va permite unificarea, omogenizarea și integrarea mediului cultural național în circuitul internațional, va consolida statutul limbii române ca limbă de comunicare pe continentul european.

II. Tehnologia de dechirilizare a textelor

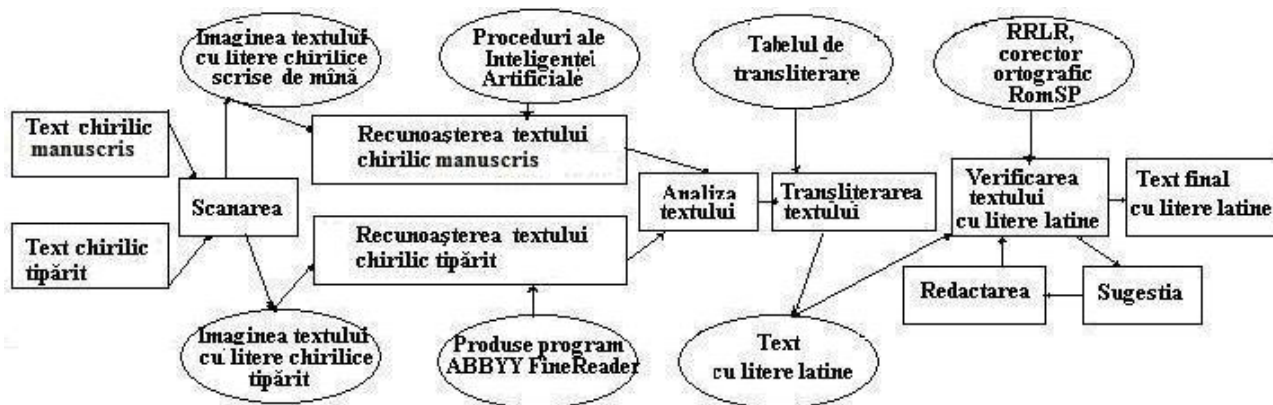


Fig.3. Procesul de dechirilizare a textului.

Procesul de dechirilizare a unui text (manuscris sau tipărit) este prezentat în Fig.3 și este alcătuit din următoarele etape:

- scanarea, clasificarea și indexarea textului;
- segmentarea și recunoașterea textului chirilic;
- transliterarea textului;
- verificarea textului cu litere latine utilizând corectorul ortografic RomSP [2] și resursele reutilizabile ale limbii române (RRLR) [2, 3].

II.1. Scanarea, clasificarea și indexarea textelor

Conversia documentelor în format electronic constă din mai multe etape. La prima etapă documentele sunt scanate. Următoarea etapă este clasificarea, care permite regruparea documentelor cu ajutorul unor șabloane predefinite (de exemplu, scrisori, articolele de ziar, contracte sau facturi, etc.). După scanare și clasificare se va face indexarea documentului electronic, care constă în generarea unor câmpuri cu informații referitor la text: data, numele autorului, titlul, etc., Clasificarea și indexarea pot fi efectuate manual sau semiautomat. Regimul semiautomat presupune utilizarea în regim interactiv a unor reguli speciale, care contribuie la sporirea acurateții și accelerează procesul de completare.

Exportul datelor extrase din diverse aplicații, de exemplu, business, baze de date, arhive electronice, se poate face în diverse formate electronice. Pentru arhivele electronice, de regulă, se utilizează formatul pdf.

Următorul pas este recunoașterea imaginii și transformarea ei în format textual.

II.2. Recunoașterea textului chirilic

Imaginea textului chirilic scris poate fi transformată în format textual utilizând rețele neuronale pentru recunoașterea literelor chirilice. Aceste rețele se bazează pe metoda de auto-organizare dezvoltată de T. Kohonen [4, 5]. Abilitatea de recunoaștere a rețelei este rezultatul unui proces de învățare pornind de la exemple de recunoaștere corectă. Ca date de intrare se utilizează descrierea sintetică a unui obiect (de exemplu, descrierea grafică a unei litere, sau ansamblul caracteristicilor acesteia). În majoritatea situațiilor descrierea este un vector de valori numerice obținute printr-o prelucrare prealabilă (preprocesare) a informațiilor brute. În calitate de date de ieșire se obține un indicator al clasei căreia îi aparține obiectul (de exemplu, numărul de ordine al literei în cadrul alfabetului).

Imaginea textului chirilic tipărit poate fi recunoscută cu ajutorul produsului program ABBYY FineReader [6], care convertește orice fișier în format editabil, obținut prin scanarea documentelor. Acestea pot fi transformate în formate editabile ale pachetului Microsoft Office (Microsoft Word, Excel), *rtf*, *html*, *txt* etc. Recunoașterea optică a caracterelor (OCR) transformă imaginile în text. OCR face posibilă editarea și reutilizarea textului inserat în imaginile scanate. OCR utilizează o formă de recunoaștere a imaginii, cunoscută sub numele de recunoaștere a modelului pentru identificare individuală a caracterelor într-un text dintr-o pagină, inclusiv semnele de punctuație, spațiile și sfârșitul de linie.

II.3. Dechirilizarea textului

Procesul de dechirilizare presupune recunoașterea caracterelor chirilice, segmentarea cuvintelor, transcrierea cu caractere latine și recunoașterea ulterioară utilizând resursele lingvistice existente și algoritmi pentru sugestii (corectarea semiautomată).

Pentru transliterare se va utiliza tabelul cu reguli de transliterare a textelor cu litere chirilice în cele cu litere latine. Pentru verificarea textului obținut se vor utiliza RRLR și corectorul ortografic RomSP. RRLR conține o bază de date cu informație lingvistică pentru limba română la nivel de cuvânt și un set de programe de gestionare a bazei de date. Ca volum RRLR conține circa un milion de cuvinte.

Verificarea textului se efectuează la nivel de cuvânt cu RomSP folosind RRLR. Textul se verifică la nivel de cuvânt în modul următor. Dacă cuvântul selectat este găsit în RRLR, rezultă că el este corect și se selectează următorul cuvânt din text. În caz contrar (cuvântul este greșit sau nu este inclus în RRLR) se inițiază un dialog, care propune următoarele opțiuni:

- a corecta cuvântul manual,
- a solicita prin intermediul sugestiei o listă de cuvinte din RRLR similare cuvântului greșit și a înlocui cuvântul greșit cu unul corect din lista de sugestii;
- a cere eliminarea cuvântului sau marcarea acestuia pentru a-l prelucra ulterior consultând un dicționar de specialitate sau un expert lingvist;
- pentru verificarea corectitudinii unui cuvânt care lipsește în dicționar, cum ar fi nume, prenume, abreviere, denumire geografică etc. acest cuvânt se va declara ca unul corect pentru restul sesiunii de lucru cu textul;
- totodată există posibilitatea includerii cuvântului nou în RRLR (folosind un set de programe de gestionare cu funcția de flexionare).

Vom menționa, că aplicarea acestei tehnologii pentru procesarea textelor din surse vechi necesită și completarea lexiconului cu termenii specifici perioadei și domeniului respectiv.

III. Concluzii

Tehnologia descrisă permite digitalizarea și recunoașterea resurselor specifice cu utilizarea lexiconului RRLR și corectorului ortografic RomSP. Această tehnologie poate fi propusă pentru digitalizarea resurselor de altă natură, cum ar fi resursele folclorice specifice, scrise cu litere chirilice. Resursele digitalizate vor reprezenta manuscrise și înregistrări specifice stocate într-o bază de date și protejate prin metode originale special elaborate. Pentru facilitarea accesului utilizatorilor la aceste resurse va fi elaborată o interfață specială și un set de instrumente care vor permite recunoașterea textelor. Crearea unei platforme de tipul e-learning ar permite utilizarea resurselor digitalizate ca materiale didactice de instruire.

Crearea resurselor specifice, plasarea acestor resurse pe Internet pentru acces public vor extinde aria și posibilitățile de cercetare, inclusiv și în domeniul științelor umanitare, prin modificarea mediului comunicativ informațional. Aceasta se manifestă prin posibilitatea aplicării tehnologiilor propuse pentru alte limbi, de exemplu, limba sîrbă.

Această tehnologie va putea fi utilizată și la completarea resurselor lingvistice RRLR cu cuvinte noi, extrase din materialele digitalizate, atestate de experți lingviști.

IV. Referințe

1. Bărbulescu Ilie. Fonetica alfabetului cirilic în textele române din vécul XVI și XVII. Anul 1904, București.
2. Burlaca O., Ciubotaru C., Cojocaru S., Colesnicov A., Magariu G., Malahov L., Petic M., Verlan T. Applications based on reusable linguistic resources. In *Multilinguality and interoperability in language processing with emphasis on Romanian*, Editors: D. Tufiș, C. Forăscu, București, 2010.
3. <http://www.math.md/elrr/>
4. Egmont-Petersen, M., de Ridder, D., Handels, H. (2002). "Image processing with neural networks - a review". *Pattern Recognition* 35 (10), pp. 2279–2301.
5. Kohonen T. 1988. *Self-organization and associative memory*. 2d ed. New-York, Springer-Verlag.
6. <http://www.abbyy.com/>