

INSTRUMENTAR PENTRU DIGITIZAREA ȘI TRANSLITERAREA TEXTELOR TIPĂRITE ÎN LIMBA ROMÂNĂ CU CARACTERE CHIRILICE

Svetlana COJOCARU, Constantin CIUBOTARU,
Alexandru COLESNICOV, Ludmila MALAHOV,
Tudor BUMBU

Abstract: The paper discusses some of the problems regarding the digitization of Romanian Cyrillic printings of the 17th – 20th centuries. An application was created to accompany ABBYY Fine Reader OCR engine and facilitate the process of digitization. The proposed tools and solutions involved were already successfully used at the re-edition of a 20th century book in mathematics in the modern Latin Romanian script, and at research in Romanian philology.

Keywords: cultural heritage, digitization, OCR, Romanian language, Romanian Cyrillic script, 17th – 20th centuries.

Introducere

Problema digitizării și conservării patrimoniului istorico-lingvistic reprezintă un domeniu prioritar din agenda digitală pentru Europa. Dezideratele principale ale politicii culturale pentru spațiile unde se vorbește limba română țin de studierea, valorificarea și digitizarea acestui patrimoniu. Procesul de digitizare necesită soluționarea unui șir de probleme legate de recunoașterea, editarea, traducerea, interpretarea, circularea și recepționarea textelor tipărite atât în limba română, cât și în alte limbi moderne.

Soluționarea acestor probleme pentru patrimoniul istorico-lingvistic românesc se confruntă cu dificultăți și aspecte specifice: un număr mare de perioade în evoluția limbii, un volum mic de resurse depozitate foarte dispersat, o mare diversitate de alfabetele folosite la tipărirea lor, în particular câteva

„alfabete de tranziție” chirilico-latine, lipsa instrumentarului pentru recunoașterea corectă a literelor chirilice din diferite perioade istorice, precum și inexistența unui lexicon adecvat perioadei de tipărire a resursei.

Instrumentarul descris mai jos integrează o serie de componente software, atât existente, cât și dezvoltate de autori, care formează o platformă pentru procesarea preliminară a imaginii, recunoașterea textului și transliterarea lui în grafie latină modernă. Este prezentată evoluția alfabetelor utilizate pentru tipăriturile de limbă română în România și pe teritoriul actual al Republicii Moldova, descrise componentele instrumentarului și aplicarea lui pentru recunoașterea și transliterarea textelor din diferite perioade istorice.

Evoluția alfabetelor românești

Utilizarea scrisului chirilic pentru tipăriturile românești a cunoscut o evoluție atât în timp, cât și în spațiul geografic. Pe teritoriul României aceasta înregistrează variații de alfabete (sau elemente ale alfabetelor) chirilice vechi (alfabete chirilice românești – ACR) de la primele texte tipărite până în 1830, când începe perioada de tranziție, caracterizată prin utilizarea alfabetelor mixte (alfabete de tranziție – AT). Sunt înregistrate 17 astfel de variante, perioada respectivă finalizându-se în 1862 prin trecerea completă și definitivă la grafia latină (alfabet latin pentru limba română – ALR).

Pe teritoriul actual al Republicii Moldova situația a fost diferită. În 1924, odată cu formarea Republicii Autonome Sovietice Socialiste Moldovenești (în componența Republicii Sovietice Socialiste Ucrainene) a fost impusă grafia chirilică, care prelua, de fapt, alfabetul limbii ruse cu excluderea a trei litere (ѐ, ѓ și њ). În anii 1932-1938 s-a revenit la grafia latină, care a fost din nou înlocuită cu cea chirilică (alfabetul chirilic moldovenesc), ultima fiind stabilită și pe întregul teritoriu al Republicii Sovietice Socialiste Moldovenești (RSSM), formate în anul 1940. Alfabetul din 30 de litere rusești a fost completat în 1967 prin introducerea literei ж pentru redarea sunetului dž. Grafia chirilică fost înlocuită cu cea latină prin votarea de către Sovietul Suprem al RSSM pe 31 august 1989 a celor mai importante compartimente din Legea cu privire la funcționarea limbilor vorbite pe teritoriul RSS Moldovenești, această zi fiind declarată ulterior sărbătoare națională – Ziua limbii. În Tab. 1 este prezentată o schemă comparativă a etapelor principale ale evoluției alfabetelor pe teritoriul României și Basarabiei.

Tab. 1. Etapele principale ale evoluției alfabetelor utilizate în tipar pe teritoriul României și Basarabiei

România	Basarabia
1642 – 1797 (alfabet chirilic, până la 47 litere)	
1797 – 1830 (alfabet chirilic, 43 litere)	1710 – 1814 (alfabet chirilic, 43 litere)
1830 – 1862 (alfabete de tranziție, mixt chirilic-latin)	1814 – 1880 (alfabet chirilic bazat pe alfabetul rus și cel slavon bisericesc; ocazional alfabete de tranziție)
1862 – 1904 (alfabet român bazat pe alfabetul latin, versiunea întâia)	1880 – 1905 (nu a existat tipar românesc) 1905 – 1918 (alfabet chirilic bazat pe alfabetul civil rus)
1904 – prezent (alfabet modern român bazat pe alfabetul latin)	1919 – 1940, 1941 – 1944 (alfabet modern român bazat pe alfabetul latin) 1940 – 1941 (alfabet chirilic bazat pe alfabetul rus)
	1944 – 1989 (alfabet chirilic bazat pe alfabetul rus; în 1967 apare litera Ж)
	1989 – prezent (alfabet modern român bazat pe alfabetul latin)

Astfel, problema digitizării și transliterării textelor românești scrise cu caractere chirilice poate fi divizată în trei compartimente majore: alfabet chirilic românesc vechi, alfabete mixte și alfabet chirilic contemporan moldovenesc.

O particularitate importantă a alfabetelor menționate o constituie existența aplicației univoce într-o direcție, anume: ACR→AT→ACM→ALR. Aceasta permite aducerea tuturor cuvintelor la alfabetul modern, folosindu-l drept o reprezentare generală a literelor din diferite epoci, precum și utilizarea mijloacelor existente de procesare a limbajului natural și (într-o anumită măsură) a resurselor lingvistice moderne.

Etapele principale ale procesării textelor

Principalele componente ale procesului de digitizare și transliterare sunt următoarele:

- ☞ Recunoașterea optică a caracterelor din tipărituri românești din sec. XVII-XX;
- ☞ Transliterarea rezultatelor obținute în grafia modernă latină;
- ☞ Crearea modelelor (pattern-urilor) pentru reprezentarea caracterelor și îmbinărilor de caractere;
- ☞ Crearea alfabetelor și dicționarelor specifice pentru anumite perioade, spații geografice, tipografii;
- ☞ Transliterarea inversă (din grafie latină în chirilică).

În linii generale abordarea noastră este ilustrată în Fig. 1.

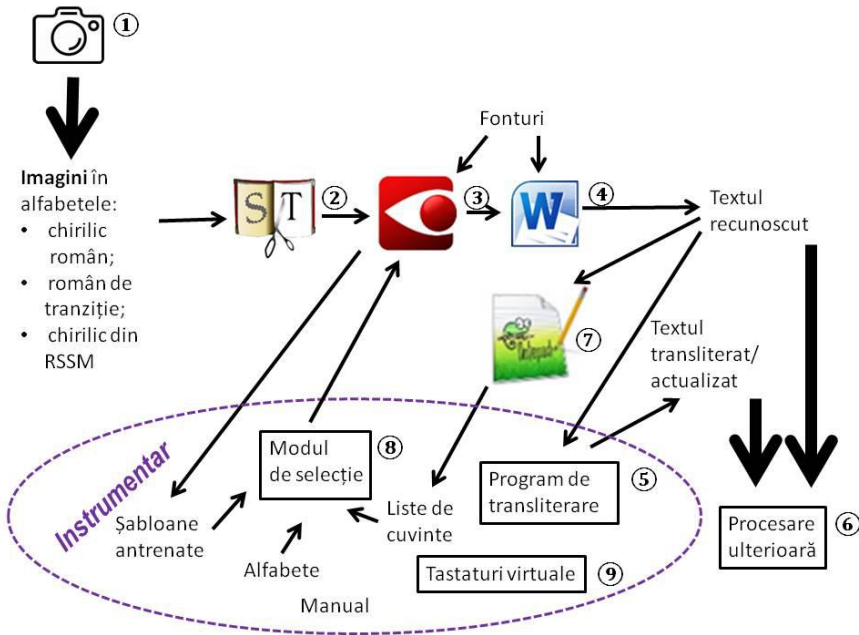


Fig. 1. Schema digitizării și transliterării textelor românești tipărite cu caractere chirilice

Operațiunile principale (notate cu cifrele 1-6) includ următoarele procedee:

1. Obținerea imaginii prin scanare, cu utilizarea softului din dotarea scannerului. Calitatea dorită este de 600 dpi sau mai mult.
2. Pregătirea imaginilor pentru OCR. În cazul nostru am utilizat ScanTailor – un soft utilitar gratuit, care efectuează corecții automate masive ale defectelor imaginii, spre exemplu, corectează unghiul de înclinare a paginii sau curăță unele pete mici. Desigur, există mai multe instrumente de acest fel. ABBYY Finereader (AFR), pe care noi îl aplicăm la pasul următor, conține și el un propriu editor de imagini, chiar mai performant, capabil să alinieze linii ondulate sau să corecteze distorsiuni trapezoidale. Însă aceste corecții ar trebui aplicate manual pentru fiecare imagine (pagină) în parte, ceea ce este destul de laborios la procesarea volumelor mari de text.
3. Recunoașterea optică a caracterelor (OCR) este efectuată cu ABBYY Finereader (AFR). Acest program operează cu texte din circa 140 de limbi, ultimele versiuni incluzând și caractere din alfabetul chirilic vechi. Funcționează pe rețele neurale, grație modelării intrinseci a limbajului dă dovadă de performanță și acuratețe.
4. Textul recunoscut este salvat ca un document Microsoft Word, care ne oferă o colecție bogată de fonturi pentru a reda caracterele cât mai aproape de cele originale, precum și vaste posibilități de formatare și editare.
5. În funcție de scopurile urmărite de utilizator, textul obținut în una din cele trei clase de alfabet chirilic, descrise mai sus, poate fi transliterat în grafie latină modernă. Pentru alfabetul chirilic românesc vechi, precum și pentru alfabetele de tranziție există o corespondență destul de constantă în raport cu alfabetul latin românesc. În cazul alfabetului chirilic moldovenesc ne confruntăm cu mai multe iregularități cauzate de anumite litere, dar și de scrierea cuvintelor de origine străină și a substantivelor proprii. Modul de soluționare a acestor probleme va fi expus în secțiunile ce urmează.
6. Procesarea manuală sau automatizată a rezultatului obținut pentru corectarea finală a textului.

Modulele adiționale (7-9) oferă niște servicii pentru facilitarea operării și pentru îmbunătățirea rezultatului. AFR utilizează liste de cuvinte pentru soluționarea ambiguităților și eliminarea cratimelor. Aceste liste trebuie să

conțină câte un cuvânt per linie utilizând codurile UTF-8. În cadrul abordării noastre folosim, de regulă, editorul Notepad++, care se distribuie gratuit și operează atât cu UTF-8, cât și cu alte codificări. Utilizând pluginul TextFX, aplicația Notepad++ poate crea liste de cuvinte sortate unice, adică fără repetări, adică anume în formatul solicitat de Finereader. De asemenea, Notepad++ permite selecția fonturilor necesare pentru editarea cuvintelor românești în caractere chirilice, inclusiv din alfabetele vechi.

Instrumentarul conține un modul de selectare a grupului de date pregătite a priori pentru AFR (alfabet, listă de cuvinte, set de șabloane pentru recunoaștere). Selecția se efectuează în funcție de perioada de timp, regiune geografică și tipografie. Acest modul apelează programul AFR incluzând datele selectate. Alfabetele și listele de cuvinte sunt pregătite manual, patternurile pentru antrenare sunt stocate de către AFR. Pentru secolul XVII au fost incluse circa 3500 de patternuri și o listă de peste 2600 de cuvinte, pentru secolul XVIII – circa 1800 cuvinte și peste 4000 de patternuri. Fig. 2 ilustrează șabloanele introduse în AFR pentru antrenarea recunoașterii literei „a” în tipărituri din secolul XVII. Pentru culegerea de pe ecran a caracterelor românești vechi în scopul introducerii sau redactării textelor respective este creată o tastatură virtuală.

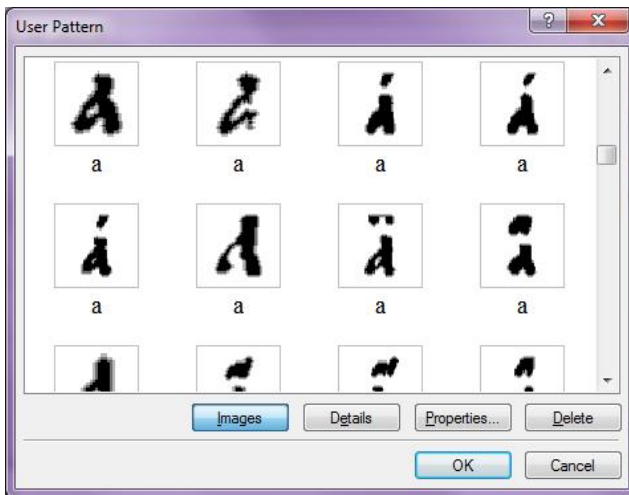


Fig. 2. Șabloane AFR pentru recunoașterea literei „a”

Evident că evoluția unei limbi nu se reduce doar la modificarea scrisului, ci se exprimă, în primul rând, în dezvoltarea lexiconului și a ortografiei. Aceste subiecte nu sunt tratate în studiul nostru, exceptând posibilitatea de adaptare la ortografierea contemporană a unor cuvinte transliterate din alfabetul chirilic moldovenesc, efectuată în cazul unei solicitări venite de la utilizator.

Procesarea textelor din secolul XVII

Fiecare perioadă din cele trei enumerate mai sus își are specificul său de procesare. Pentru textele din secolul XVII este caracteristic un nivel mai avansat de degradare, zgomote, caractere alipite sau rupte. În recunoașterea lor ne confruntăm cu o serie de probleme, dintre care vom evidenția următoarele:

1. Suprascrierea literelor (slovelor), având drept scop economia de spațiu și de eforturi la scriere / citire;
2. Omiterea unor litere;
3. Multitudinea diverselor semne (title, diacritice etc.) plasate deasupra liniei;
4. Prescurtări (în denumiri de funcții, luni ale anului etc.);
5. Scrierea numeralelor atât cu cifre, cât și cu litere;
6. Diversitatea fonturilor.

Suprascrierea și prezența semnelor deasupra rândului conduc spre separarea eronată a unei linii în două. Pentru soluționarea acestei probleme, la sugestia autorilor Finereader-ului, s-a recurs la majorarea formală a densității imaginii.

În urma procesării cu instrumentarul nostru literele suprascrise sunt incluse în cuvânt, iar celelalte semne supralinie sunt omise.

Tipografiile din secolul XVII foloseau o varietate de fonturi pentru tipărirea cărților și a documentelor, dar dintre acestea putem distinge, în general, două seturi substanțial diferite, atât după stil, cât și după utilizarea caracterelor. În Fig. 3 sunt prezentate fragmente a două pagini din două cărți tipărite în perioade destul de apropiate (1648 și 1679), care au fonturi distincte.



Fig. 3. Două pagini tipărite cu utilizarea diferitelor fonturi; caracterele atipice sunt evidențiate sub text

Este evident, că aceste două texte sunt diferite după stil, având fonturi accentuat distincte în cazul literelor “т” și “з”. Litera “т” din primul text este tipărită ca “т” standard, iar în cel de al doilea apare în forma scrierii de mână, adică “m”. Același lucru îl observăm și în cazul literei “з”, al cărei mod de scriere diferă substanțial. Dacă am aplica procedura de recunoaștere a unui text de al doilea tip utilizând modele din primul, rezultatul ar avea o rată de eroare destul de mare. Deci este necesar să avem un instrument cu ajutorul căruia utilizatorul va putea alege cel mai potrivit set de modele pentru cartea sau documentul din secolul XVII tipărit în grafie chirilică românească.

Utilizatorul trebuie să fie familiarizat cu modelele existente, el va analiza vizual o pagină din cartea care urmează să fie recunoscută și va alege modelul cel mai potrivit, orientându-se după anumite caractere distincte (precum ar fi scrierea diferită a literelor “т” și “з” în exemplele de mai sus). În caz că modelul cel mai potrivit nu poate fi ales vizual, se recomandă recunoașterea unei pagini prin câteva modele diferite, iar cel care va prezenta un rezultat mai bun va fi utilizat pentru recunoașterea întregii lucrări.

O interfață specială este creată pentru selecția regiunii geografice, unde s-a tipărit textul. Putem alege una din următoarele variante: Iași, București, Târgoviște, Bălgrad (Alba Iulia), Uniev (Cernăuți), Sas Sebeș, Snagov sau Buzău. În cadrul unei regiuni avem posibilitatea selectării tipografiei, spre exemplu, pentru București sistemul este antrenat în recunoașterea fonturilor din Tipografia Domnească și cea a Scaunului Mitropoliei Bucureștilor (Fig. 4).

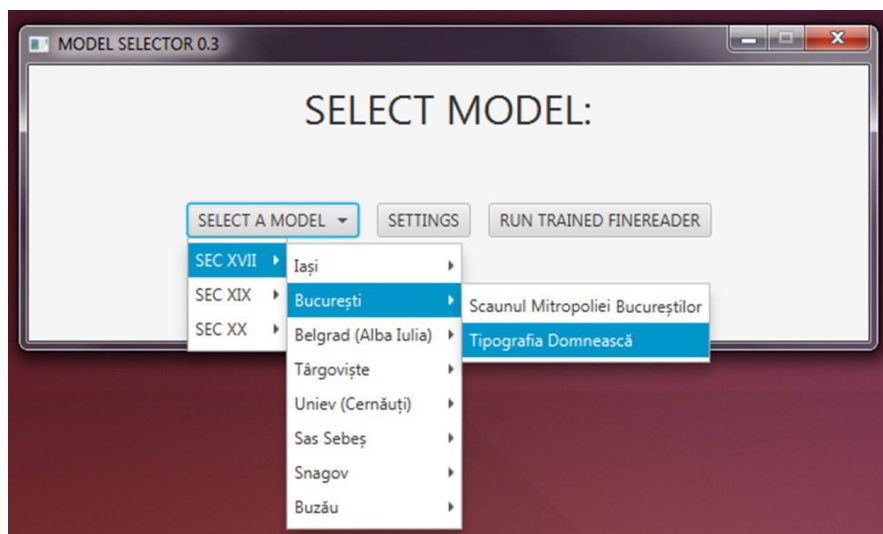


Fig. 4. Interfața de selectare a modelului OCR

Transliterarea în majoritatea cazurilor reprezintă o aplicație de tip „literă → literă”, excepție făcând șapte cazuri (Г, К, Ч, И, Ъ, А, А) când este necesară o analiză simplă a contextului de dreapta cu lungimea de 1-2 simboluri. Spre exemplu, litera А trece în а la începutul cuvântului și după Ъ, И, în е - după litera Ч, și în еа după orice consoană, la sfârșitul cuvântului.

Corectitudinea recunoașterii depinde de calitatea imaginii, de cantitatea șabloanelor introduse și de mărimea dicționarului. Acuratețea la recunoaștere fără antrenare este destul de mică – 35%, în urma învățării supervizate se poate ajunge la circa 70% din cuvinte recunoscute corect. Oricum, o corectare manuală rămâne necesară.

Procesarea textelor din secolele XIX - XX

Diverse surse atestă circa 17 versiuni de alfabet de tranziție. Mai mult, varietatea scrierii poate fi întâlnită chiar și în cadrul uneia și aceleiași lucrări, unde unele pagini sunt tipărite doar cu caractere chirilice, altele mixte sau pur latine.

În recunoașterea textelor tipărite cu alfabet mixte am utilizat două abordări. În cadrul primei textul scanat este reprodus în urma digitizării în glifele sale originale. Acest lucru este posibil prin configurarea și antrenarea AFR, precum și prin dotarea lui cu un dicționar de epocă. Corectitudinea recunoașterii este de circa 93%. Cea de a doua abordare a fost propusă pentru a soluționa

problema varietății alfabetelor. AFR permite atât obținerea rezultatului în glife originale, cât și substituirea oricărei glife printr-o secvență de litere din alfabetul selectat. Finereader propune această metodă pentru ligaturi, dar ea poate fi utilizată în mod mai general pentru o substituție arbitrară. În cazul alfabetelor de tranziție a fost construită o versiune generalizată de alfabet de intrare în care se stabilește o corespondență univocă cu un singur alfabet de ieșire, indiferent de modul de prezentare a literei de intrare. Spre exemplu, atât **т** (chirilic) cât și **t** (latin) vor fi recunoscute drept **t**. S-a constatat, că în abordarea bazată pe ligaturi rata de acuratețe crește până la 97% din cuvinte sau 0.6% caractere eronate față de 1.5% în cazul primei abordări.

Ca și în cazul alfabetului chirilic vechi pentru transliterarea textelor din alfabete de tranziție sunt suficiente două tipuri de reguli: „literă → literă” sau substituție dependentă de context.

Digitizarea textelor tipărite cu caractere chirilice moderne (alfabet chirilic moldovenesc) se caracterizează printr-un grad mai înalt de acuratețe, atât grație calității mai bune a imaginii, cât și faptului că AFR este antrenat din start pentru recunoașterea alfabetului rusesc, fiind necesară doar adăugarea literei **ж**. Corectitudinea recunoașterii după antrenare depășește 98%.

Însă o problemă mai dificilă pentru textele din această perioadă o constituie transliterarea. Ca și în cazul alfabetului chirilic vechi și al celor de tranziție, pentru mai multe caractere există corespondența univocă „literă → literă”, pentru câteva litere problema se rezolvă prin analiza contextului, însă avem și cazuri, când o soluție exactă nu poate fi identificată. Din ele face parte transliterarea caracterului **я**, care poate fi reprezentat prin **ia**, **ea** sau **a**, spre exemplu: **тряз**→**treaz**, **амязэ**→**amiază**, **абрeвия**→**abrevia**. Circa 20 de reguli pentru această literă, atât euristice cât și statistice, au permis să fie rezolvate mai multe cazuri, dar nu s-a ajuns la o soluție completă.

O altă problemă parvine de la cuvintele de origine străină, care în scrierea cu caractere latine respectă forma din limba originală, pe când în alfabetul chirilic sunt aplicate principiile fonetice, spre exemplu, cuvântul **дизайн**, în transliterare directă ar fi **dizain**, pe când ortografierea corectă este **design**. Aceeași situație apare și în cazul substantivelor proprii preluate din alte limbi. Numele **Шекспир** (Shakespeare) în transliterare directă apare ca **Şexpir**, etc. Toate aceste cazuri formează o listă de excepții, care sunt procesate prioritar, urmând apoi aplicarea regulilor.

Programul de transliterare funcționează în două moduri, direct și actualizat. În cel de al doilea caz cuvintele sunt aduse la normele moderne de scriere. Spre exemplu, la scrierea cu grafie chirilică într-o serie de cuvinte litera *i* lipsește sau este utilizată scrierea cu *î* în loc de *i*: **требье (trebue), ынтродучере (introducere)** etc. În funcție de cererea utilizatorului este păstrată ortografia originală sau efectuată actualizarea, în exemplele de mai sus cuvintele respective fiind transliterate în **trebuie și introducere**.

Resurse chirilice pentru OCR

Calitatea recunoașterii textelor ar putea fi îmbunătățită având la dispoziție vocabulare cu caractere chirilice din perioada respectivă. În cazul secolului XX aceste vocabulare pot fi obținute prin transliterarea inversă, adică din grafie latină în cea chirilică. Există mai multe resurse lingvistice, care pot fi supuse acestei transformări, de exemplu, DEX online¹ sau ELRR². În urma analizei lor am considerat mai potrivit lexiconul elaborat la Universitatea „A. I. Cuza”, Iași³, care conține circa 1 milion de intrări, este bine structurat și însoțit de tag-uri morfologice. Transliterarea inversă se confruntă cu propriile dificultăți, în particular, la procesarea literei *i*, care în grafie chirilică poate avea trei reprezentări: **и, й, ъ** sau poate fi omisă. Exemple: arici → арич (la singular, litera *i* este omisă), arici → аричь (la plural, substituția *i* → **ь**), [a] cheltui → келтуи (infinitiv, substituția *i* → **и**), [eu] cheltui → келтуй (prezent, singular, persoana I, substituția *i* → **й**). O problemă similară de lipsă a aplicației univoce apare și la transliterarea inversă a diftongilor și triftongilor, spre exemplu, există trei variante posibile pentru diftongul **ia**: **я, ия** și **иа**. Ca și în cazul transliterării directe, rămâne problema cuvintelor de origine străină. Nu în toate cazurile a fost posibilă stabilirea de reguli formale, care ar permite automatizarea completă a procesului, pentru soluționarea ambiguităților fiind necesare și intervenții manuale.

Procesarea începe cu selectarea cuvintelor incluse în dicționarul de excepții, asupra celor rămase se aplică o serie de filtre, care transliterează separat prefixele, sufixele, diftongii și triftongii, apoi se efectuează filtrul final, care se reduce la aplicația „literă → literă”.

¹ <https://dexonline.ro/>

² <http://www.math.md/elrr/>

³ <http://nlptools.info.uaic.ro/WebPosRo/resources/posDictRoDiacr.txt>

Concluzii

Instrumentarul propus a demonstrat posibilitatea aplicării lui la digitizarea și transliterarea textelor chirilice românești din diferite perioade istorice. După cum e și firesc, cele mai dificile pentru recunoaștere sunt tipăriturile vechi, unde gradul de corectitudine este mai redus și intervenția manuală este mai solicitată. Mai puțin laborioase sunt transliterările textelor din secolul XX. Cu ajutorul tehnologiei propuse a fost digitizată și transliterată cartea de matematică⁴. Procesarea unui volum de 224 de pagini a servit drept test pentru validarea instrumentarului. S-a constatat o calitate foarte bună a recunoașterii, dar au fost semnalate mai multe erori de transliterare, majoritatea din ele fiind cauzate de ambiguități la interpretarea literei **я**. Nu în toate cazurile s-a reușit actualizarea ortografiei, fiind necesare intervenții manuale.

Svetlana COJOCARU
E-mail: svetlana.cojocaru@math.md
Constantin CIUBOTARU
Alexandru COLESNICOV
Ludmila MALAHOV
Tudor BUMBU

Institutul de Matematică și Informatică, Academia de Științe
din Republica Moldova, Chișinău

⁴ V. Andrunachievici, I. Chitoroagă. *Numere și ideale*, Chișinău, „Lumina”, 1979.