

DIGITIZAREA, RECUNOAȘTEREA ȘI CONSERVAREA PATRIMONIULUI CULTURAL-ISTORIC

Dr. Elena BOIAN

Dr. Constantin CIUBOTARU

Dr. hab. Svetlana COJOCARU

Dr. Alexandru COLESNICOV

Ludmila MALAHOV

*Institutul de Matematică și Informatică
al AȘM*

Introducere

Problema digitizării și conservării patrimoniului istorico-lingvistic (cultural) reprezintă un domeniu prioritar din agenda digitală pentru Europa. UE evidențiază necesitatea unui efort coordonat în domeniu și întreprinde vaste acțiuni în vederea impulsivării acestui proces, printre care dezvoltarea bibliotecii virtuale *Europeana*, susținută prin rezoluția Parlamentului European din 5 mai 2010 și adoptarea Programului de lucru pentru activități culturale 2011-2014. Menționăm și recomandările Comisiei Europene „Privind digitizarea și accesibilitatea online a materialului cultural și conservarea digitală” din 27 octombrie 2011.

Dezideratele principale ale politicii culturale pentru zonele unde se vorbește limba română țin de studierea, valorificarea și digitizarea patrimoniului cultural-istoric. Procesul de digitizare a patrimoniului 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 număr relativ mic și foarte dispersat de resurse depozitate, o mare diversitate de alfabete folosite la tipărirea lor, în particular câteva „alfabete de tranziție” chirilico-latine. Dificultățile în digitizarea și conservarea acestui tezaur țin de recunoașterea corectă a literelor chirilico-latine, dar și de inexistența unui lexicon adecvat perioadei de tipărire a resursei. O soluție pentru problema lexiconului ar fi alinierea la normele lingvistice contemporane ale textelor vechi [1].

Istoric, limba română a parcurs o cale lungă și bogată de dezvoltare. Există studii care explică apariția foneticii și ortografierii caracteristice etapelor concrete de evoluție a limbii, care sunt necesare atât pentru determinarea alfabetului, cât și a literelor specifice [2,3]. Cunoașterea acestor legități ne permite să construim resurse lingvistice utilizând un instrumentar special elaborat pentru o perioadă istorică concretă.

Prima carte tipărită pe teritoriul românesc a fost *Liturghierul slavon*, îngrijit de către ieromonahul Macarie în anul 1508, iar prima carte tipărită în limba română a fost *Catehismul Românesc* al diaconului Coresi, apărut la Brașov în anul 1535 [3].

Biblioteca Națională a Republicii Moldova deține o colecție de aproximativ 21 000 cărți vechi și rare. Circa 20 de cărți din această colecție sunt tipărite în limba română, în Basarabia (Chișinău și

CULTURAL AND HISTORICAL HERITAGE DIGITIZATION, RECOGNITION AND CONSERVATION

Summary. This article describes digitization of old Romanian texts, problems at their recognition, and motivates the necessity to create specific electronic resources mirroring the history of the standard Romanian language. We provide also statistics of results at recognizing a Romanian text of the 19th century by modern software, and we propose a technology for creation of linguistic lexicon for Moldavian Cyrillic script of 1967–1989, starting from modern (standard) Romanian lexicon. This technology is based on transliteration and parallel texts alignment.

Keywords: digitization, Romanian linguistic resources, text recognition, language technology, Cyrillic script, transliteration, text aligning.

Rezumat. În lucrare se abordează problemele ce apar în procesul de digitizare și recunoaștere a textelor vechi românești, se argumentează necesitatea creării resurselor electronice specifice care caracterizează evoluția limbii române moderne. Se prezintă rezultate statistice obținute la recunoașterea unui text românesc din secolul al XIX-lea, utilizându-se produse program moderne. Se propune o tehnologie în vederea creării lexiconului lingvistic pentru patrimoniul moldovenesc tipărit cu alfabet chirilic în perioada 1967-1989, pornind de la lexiconul românesc modern. Această tehnologie se bazează pe transliterare și pe aliniere paralelă a textelor.

Cuvinte-cheie: digitizare, resurse lingvistice românești, recunoașterea textului, tehnologia limbajului, alfabet chirilic, transliterare, alinierea textelor.

Dubăsari), utilizând alfabetele chirilic și tranzițional [4,5]. Bibliotecile publice din Sankt Petersburg dețin importante mostre de carte românească veche (secolele XVI-XIX). Dintre cele 66 de titluri incluse, spre exemplu, în *Catalogul edițiilor chirilice ale slavilor de sud și ale românilor*, 45 de volume revin slavilor de sud, iar 21 de volume – țărilor românești [6].

Studiile existente explică aspectele legate de dezvoltarea componentelor principale ale limbii: alfabet, lexicon, ortografie cu referire la etapele specifice de evoluție a limbii. Această informație este utilă pentru a crea resurse și instrumente lingvistice racordate la anumite perioade din istoria limbii. Ținând cont de particularitățile fiecărei perioade, vom propune o tehnologie pentru crearea acestor componente. În particular, vom studia problema de digitizare a textelor tipărite cu caractere chirilice în Republica Sovietică Socialistă Moldovenească (RSSM) în perioada 1967-1989.

Lucrarea prezintă un proiect pe termen lung, care abia începe. Pe parcurs ne vom conduce de principiu „din prezent în adâncul secolelor”.

Perioade de evoluție a limbii române

Istoria limbii române cunoaște două epoci în dezvoltarea sa. Prima se referă la formarea dialectului dacoromân, începând cu căderea Sarmisegetuzei (106 A.D.) până în secolul al XV-lea [2]. Se utiliza alfabetul chirilic grație influenței masive a Bisericii Ortodoxe.

Epoca a doua de dezvoltare a limbii române literare (sec. XVI-XX) începe cu apariția primelor texte scrise în limba română și constituie rezultatul

unei îndelungate și complexe evoluții [3]. Procesul de unificare lingvistică este marcat de apariția Bibliei de la București (1688), care a condus ulterior la stabilirea a două mari etape în evoluția lingvistică [7].

Etapa întâi începe cu apariția primelor texte literare românești și se încheie la începutul secolului al XVIII-lea. În cadrul acestei etape pot fi distinse 3 perioade:

- Anii 1532 și 1588, prima fază a limbii literare;
- Anii 1588-1656, faza consolidării principalelor variante ale limbii române literare (muntenească, moldovenească și sud-vest-ardelenească);
- Anii 1656-1715, faza influenței reciproce dintre variantele literare.

A doua etapă se întinde pe un interval între 1715 și 1960. Este epoca de consolidare a limbii unice supradialectale. Procesul de unificare a limbii române literare a cunoscut o evoluție lungă, în cursul a 4 perioade:

1. Anii 1715-1780, momentul primei unificări, aproximativ în 1750;
2. Anii 1780-1836, diversificarea lingvistică;
3. Anii 1836-1881, constituirea principalelor norme ale limbii literare de astăzi;
4. Anii 1881-1960, definitivarea formării normelor limbii române literare contemporane.

Ultima perioadă ne descrie consolidarea stilurilor limbii române literare. În 1904, prin modificările aduse ortografiei, se stabilesc definitiv bazele scrierii fonetice, păstrate, cu unele retușări ulterioare, până în prezent. Vom arăta în Fig. 1-8 exemple de texte tipărite în diverse perioade de evoluție a limbii.

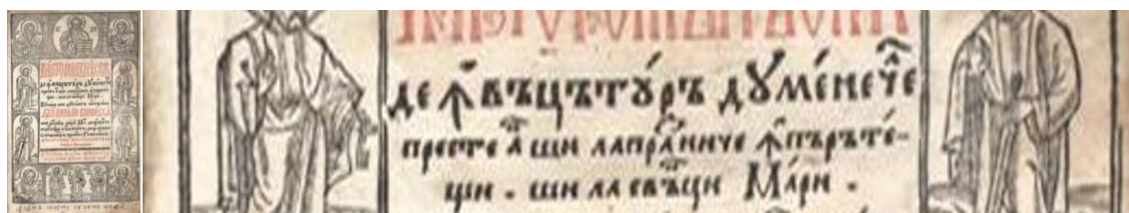


Figura 1. Cazania lui Varlaam, Iași, 1643

<http://tiparituriromanesti.wordpress.com/2011/12/04/cazania-lui-varlaam-iasi-1643/>

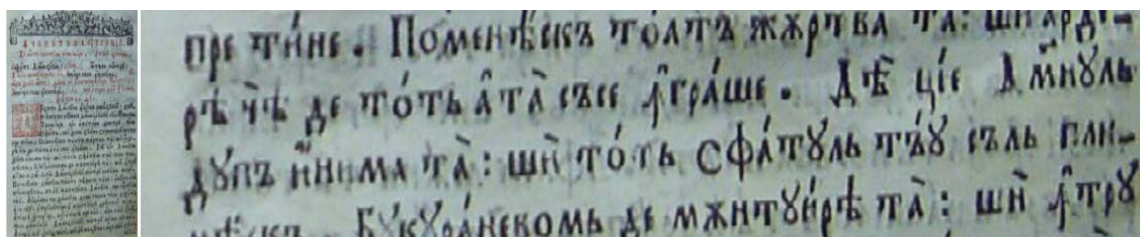


Figura 2. Ceaslov, 1748. <http://muzeu.reintregirea.ro/index.php?cid=carte-53>

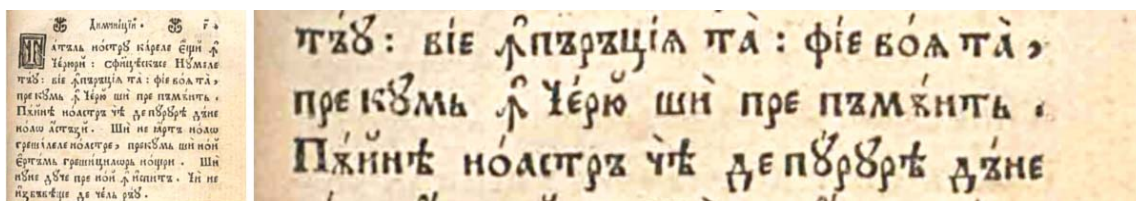


Figura 3. *Татъл ностру*. Acathist cu multe alease rugăciuni pentru evlaviaa fieștecăruia creștin. Acum a treia oară tipărit. Blaj: Tipografia Seminariului, 1786. http://documente.bcuculuj.ro/web/bibdigit/patrimoniul/BCUCLUJ_FCS_BRV497.pdf



Figura 4. *Letopisiștile Țării Moldovii* publicate pentru întâiași dată de Mihail Kogălniceanu. Tom I. Iașii. La toate libreriile. 1852. <http://tiparituriromanesti.wordpress.com/2012/03/24/miron-costin-carte-pentru-descalecatul-dintaiu-a-tarii-moldovii-si-neamul-moldovinesc/>

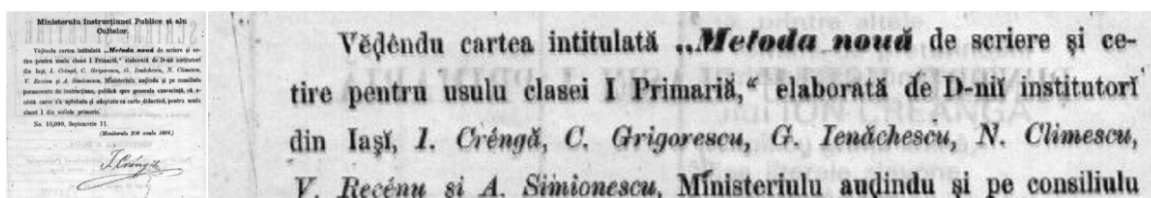


Figura 5. *Metodă nouă de scriere și citire: Pentru usulu clasei I, Primariă* / I. Crengă, C. Grigorescu, G. Ienăchescu, Ed. a II-a. – Iassy: Tipografia H. Goldner, 1868. – 71 p. <http://www.scribd.com/doc/70357520/Carte-rara-in-limba-romană-din-colecțiile-bibliotecii-Contribuții-bibliografice-Fascicula-2>

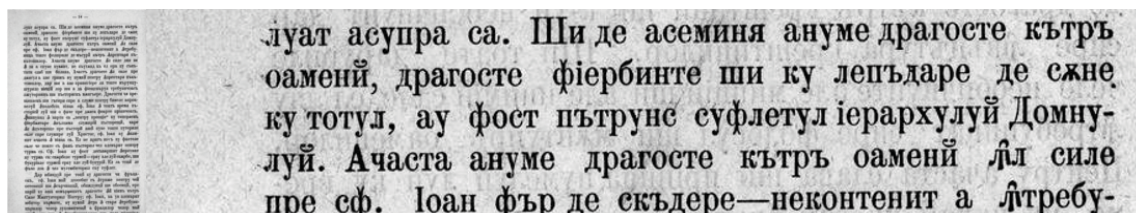


Figura 6. O pagină din revista „Луминѣторул”, 1908, Nr. 1. http://upload.wikimedia.org/wikipedia/ru/b/b1/Rumynskaja_Kirillica_Grazhdanskij_Shrift.jpeg

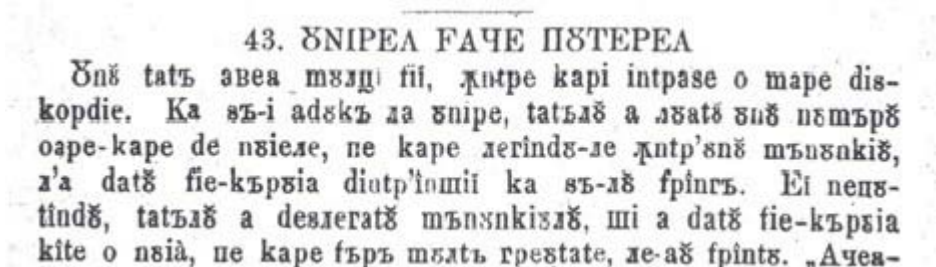


Figura 7. Una din variantele alfabetului tranzițional utilizat în alfabetul lui I. Creangă.

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

Figura 8: Text tipărit cu alfabet chirilic, perioada 1967-1989. Utilizat în prezent în Transnistria (M. Eminescu, „Epigonii”)

Perioade de evoluție a alfabetului român

În secolul al XVII-lea, tiparul românesc utiliza un alfabet chirilic cu 47 de litere, majoritatea dintre ele fiind împrumutate din alfabetul bisericii slavone. S-au adăugat câteva litere grecești, în mare parte pentru redarea adecvată a numelor proprii, dar și litere originale românești. De exemplu, litera **Ѡ** utilizată pentru a reda prefixul (prepoziția) **în**, **îm**, sau litera **î** modernă la începutul cuvântului. Acest alfabet a fost utilizat la tipărirea *Cazaniei* lui Varlaam la Iași în anul 1643 (Fig.1). Primul abecedar românesc a fost tipărit în 1699 la Bălgrad (Alba Iulia), iar prima gramatică românească a fost tipărită în 1757 de Dimitrie Eustatievici.

Începând cu anul 1830 și până la adoptarea oficială a alfabetului latin român în 1862, nu exista un alfabet stabil, astfel în această perioadă au fost utilizate cel puțin șapte modificări ale așa-numitului „alfabet de tranziție”, chirilico-latin, care conținea atât litere latine, cât și litere chirilice (Fig. 4, 7). De exemplu, **e - e** (1830) - **ε** (1846); **к - k**; **ш - шт**; **с - дз - dz - ċ** (1846).

Utilizarea grafiei latine în România nu a influențat activitatea tipografică din Basarabia. După alipirea Basarabiei la Imperiul Rus în 1812, limba oficială la Chișinău devine limba rusă. În anul 1833 limba română a fost exclusă din circuitul oficial, dar a continuat să fie utilizată în activitățile eparhiale. Astfel, pe parcursul anilor 1867-1871 apărea versiunea română a monitorului eparhiei Chișinău tipărit cu caractere chirilice. Tipografia bisericească din Chișinău a fost sistată pe perioada 1883-1890, procesul fiind reluat la începutul secolului XX.

Spre deosebire de alfabetul chirilic utilizat pentru scrierea limbii române din secolele XIV-XV până în anul 1862, alfabetul chirilic folosit în Republica Autonomă Sovietică Socialistă Moldovenească (RASSM) începând cu anii 1930 și, ulterior,

în Republica Sovietică Socialistă Moldovenească (RSSM) și Transnistria în prezent, este de fapt o adaptare a alfabetului chirilic rusesc. De menționat că în perioada 1932-1938 în RASSM a fost utilizat alfabetul latin. În Republica Moldova alfabetul chirilic a fost utilizat până în 1989.

Vom prezenta mai jos (Tab.1) perioadele de evoluție a alfabetului român începând cu *Cazania* lui Varlaam. Pe lângă alfabet, există și alți factori care caracterizează evoluția limbii, precum ortografierea și lexiconul.

Recunoașterea textelor tipărite

Procesul de digitizare și de recunoaștere pentru manuscrise este destul de complicat, deoarece necesită efectuarea unor operații suplimentare, de exemplu, ajustarea contrastului, „curățirea imaginii”, segmentarea textului. De asemenea, trebuie elaborate algoritmi speciali de recunoaștere și lexicoane specializate. Procesul de digitizare și recunoaștere e constituit din următoarele etape (Figura 9):

- Digitizarea textului pentru obținerea copiei electronice grafice;
- Recunoașterea cu metode standardizate, adică utilizarea nemijlocită a OCR (Optical Character Recognition) [8], sau prin instruirea lui. În caz contrar, se vor folosi proceduri ale Inteligenței Artificiale, așa-numitul proces de conversie. Transliterarea textului se va efectua ținând cont de literele specifice utilizate în textul inițial.
- Verificarea textului recunoscut se produce utilizând resursele lingvistice reutilizabile specializate pentru perioada de timp respectivă.

Digitizarea textelor constă în scanarea lor și obținerea variantei electronice în formă de imagine. Pentru recunoașterea textelor din imagine se aplică OCR. Sistemele standard OCR utilizează diferite metode de recunoaștere a textelor. Am cercetat posi-

Tabelul 1

Evoluția alfabetului român începând cu anul 1642

România	Basarabia
1642 – 1710 (alfabet chirilic)	
1710 – 1830 (alfabet chirilic modificat)	1710 – 1814 (alfabet chirilic modificat)
1830 – 1862 (alfabet tranzițional, mixt chirilico-latin)	1814 – 1880 (alfabet chirilic bazat pe alfabetul rus și cel slavonic bisericesc, ocazional alfabet tranzițional și latin)
1862 – 1904 (alfabet latin)	1880 – 1905 (n-a existat tipar românesc) 1905 – 1918 (alfabet chirilic bazat pe alfabetul civil rus)
1904 – 1960 (alfabet latin modificat)	1919 – 1940, 1941 – 1944 (alfabet latin modificat) 1940 – 1941 (alfabet chirilic bazat pe alfabetul rus) [Vezi mai sus în text situația din Transnistria]
1960 – 1993 (alfabet latin modificat)	1944 – 1989 (alfabet chirilic bazat pe alfabetul rus; din 1967 apare litera ж)
1993 – prezent (alfabet modern român bazat pe alfabetul latin)	1989 – prezent (alfabet modern român bazat pe alfabetul latin) [Vezi mai sus în text situația din Transnistria]

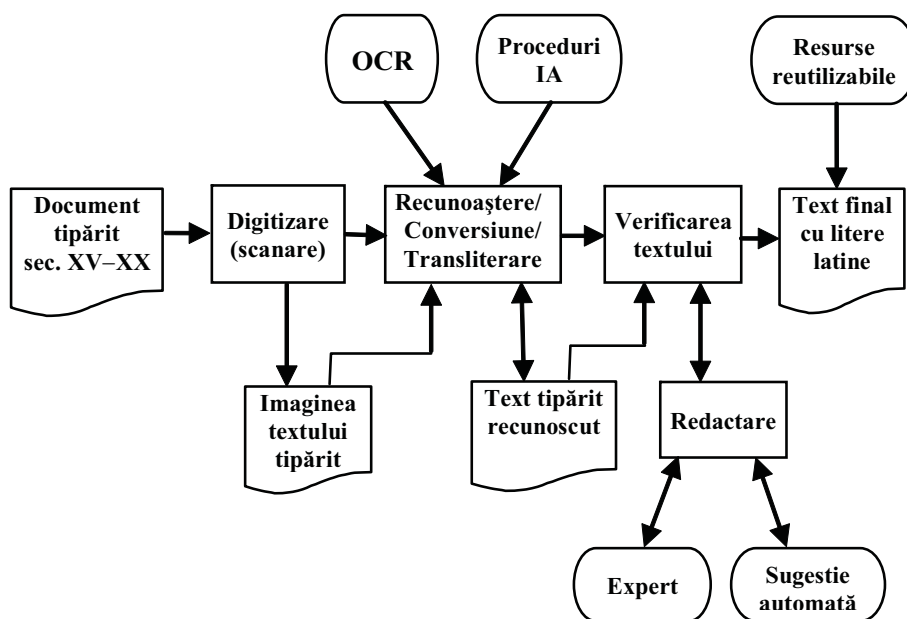


Figura 9. Etapele tehnologice de recunoaștere a textelor tipărite

bilitățile a două sisteme: IRIS și ABBY FineReader. Rezultatele experiențelor de recunoaștere a textului tipărit în sec. al XIX-lea sunt expuse în continuare. S-a determinat că sistemul IRIS, în procesul de instruire, nu poate selecta orice fragment din imaginea textului și de aceea acest sistem nu satisface scopurile noastre de recunoaștere a textului tipărit cu scrisul vechi român.

În continuare se vor folosi tehnici de recunoaștere a formelor pentru identificarea individuală a caracterelor unei pagini de text, inclusiv semnele de punctuație, spațiile și sfârșitul de linie. Textul recunoscut se va prezenta ca un fișier editabil.

Transliterarea este un proces strict individual ce depinde de perioada examinată. În funcție de textul inițial, se vor utiliza programe care conțin informație despre caracterele specifice întâlnite în text. Transliterarea presupune stabilirea unei relații bidirecționale univoce între două sisteme de scriere astfel, încât un cunoscător să poată reconstitui textul original din varianta transliterată. Procesul de transliterare se va folosi doar la necesitate.

Verificarea textului se efectuează cu aplicații special elaborate [9], care utilizează resursele reutilizabile specifice pentru perioada istorică a textului tipărit. Totodată, cuvintele noi obținute se vor introduce în lexiconul corespunzător.

Procesarea textelor tipărite cu alfabet chirilic în RASSM și RSSM

Perioada inițială de utilizare a alfabetului chirilic în Republica Autonomă Sovietică Socialistă Moldovenească (RASSM) se referă la anii 1924-

1940 și se asociază cu utilizarea unui lexicon foarte specific, caracterizat de:

- utilizarea cuvintelor rusești (de exemplu, совет, указ, словарь) în locul echivalentelor românești (consiliu, decret, dicționar);
- excluderea neologismelor românești, ele fiind considerate „burghezisme”;
- utilizarea lexiconului local (Transnistrean);
- introducerea unor neologisme auto-inventate pentru unele noțiuni abstracte neatestate în limbajul basarabean. De exemplu, амувремник (amuvremnic) în loc de contemporan;
- utilizarea particularităților accentului local (Transnistrean), de exemplu, ди (di) în loc de de, мержи (merji) în loc de merge, сунити (suniti) în loc de sunete etc.

Ne vom referi în continuare la perioada 1967-1989 de utilizare a alfabetului chirilic în Republica Sovietică Socialistă Moldovenească (RSSM). Pentru procesarea OCR a textelor apărute în această perioadă, este necesar să fie instruit sistemul OCR pentru a recunoaște litera adițională ж și pentru a crea lexiconul respectiv. Existența unui lexicon caracteristic acestei perioade ar permite automatizarea procesului de verificare și validare a cuvintelor recunoscute prin procedeele expuse mai sus. Acest lexicon poate fi creat: a) manual, b) prin transliterarea cuvintelor românești scrise cu caractere latine în varianta corectă scrisă cu caractere chirilice și c) prin alinierea variantelor de text tipărite în paralel cu caractere latine și caractere chirilice.

Prin transliterare vom înțelege transcrierea unui cuvânt din limba română în forma echivalentă scrisă

cu caractere chirilice conform normelor lingvistice acceptate pe perioada 1967-1989 în RSSM. Metoda transliterării s-ar potrivi ideal în cazul dacă se reușește formalizarea tuturor regulilor de transcriere. Un studiu prealabil arată că acest proces este anevoios și nu poate fi automatizat în totalitate din cauza iregularităților legate de discordanța dintre fonologia, morfologia și sintaxa limbii române și normele lingvistice acceptate în RSSM. Parțial acest proces poate fi automatizat implicând elemente de formalizare a regulilor de transcriere, de intervenție manuală și de aliniere.

Dificultăți evidente apar la transliterarea cuvintelor de proveniență străină. Dacă în limba română aceste cuvinte se scriu, de regulă, ca și în limba originală, atunci transcrierea lor cu caractere chirilice se face conform pronunțării. De exemplu, *design* – *дизайн*, *cowboy* – *ковбой*, *watt* – *ват*, *charleston* – *чарлстон*. Aceste cuvinte pot fi transliterate doar în regim manual.

Pentru lexiconul original românesc procesul respectiv poate fi parțial automatizat. În acest scop s-au stabilit reguli de transcriere a literelor și îmbinărilor de litere. Inserăm mai jos câteva astfel de reguli.

1) Reguli de transcriere „literă → literă”. De exemplu, *a* → *а*, *ă* → *э*, *b* → *б*, *d* → *д*, *f* → *ф*, *l* → *л*, *m* → *м*, *n* → *н*, *r* → *р*, *ș* → *ш*, *t* → *т*, *ț* → *ц*, *v* → *в*, *z* → *з* (*bardă* – *бардэ*, *zarvă* – *зарвэ*, *măr* → *мэр*).

2) Reguli de transcriere pentru literele *î* și *â*. Îmbinările *âi*, sau *îi* se vor transcrie în *и* pentru cuvintele *mâine*, *râine*, *câine* și derivatele lor (*mâine* → *мыне*, *râine* → *ныне*, *câine* → *кыне*, *mîine* → *мыне*, *rîine* → *ныне*, *cîine* → *кыне*). În alte situații se va aplica regula „literă → literă”: *â* → *ы*, *î* → *ы* (*român* → *ромын*, *întâi* → *ынтый*).

3) Reguli pentru *ea* și *ia*. Se transcriu în *я*, cu o singură excepție: pronumele *ea* se transcrie *еа*; în același timp, verbul *ia* se transcrie ca *я*;

4) Transcrierea lui *i* prin trei litere diferite: *и*, *й*, *ь*. Concomitent, menționăm existența cazurilor când litera *i* este omisă (*ierpure* → *енуре*), sau trecută în *ы* (*introduce* → *ынтродучере*). Reguli de transcriere pentru litera *c*.

a) *c* → *к*, dacă după *c* urmează una din vocalele *a*, *â*, *î*, *o*, *u*, sau o consoană diferită de *h* (*încrețit* → *ынкрецим*, *clocot* → *клокот*, *casă* → *касэ*, *cisoș* → *кукош*; *câmp* → *кымп*).

b) Combinațiile *che*, *chi* se vor transcrie în *ке* și, respectiv, *ки* (*cheltuială* → *келтуялэ*, *chihlimbar* → *кихлимбар*, *chibzui* → *кибзуи*).

c) Dacă după îmbinarea *ce* nu urmează *a*, atunci se aplică regula *ce* → *че* (*cercel* → *черчел*, *сер* → *чен*).

d) *cea* → *ча* (*ceară* → *чарэ*, *ceas* → *час*, *ceață* → *чацэ*, *ceașcă* → *чаишкэ*). Excepție pentru articolul demonstrativ *cea* (*acea*) → *чя* (*ачя*).

e) Dacă după îmbinarea *ci* nu urmează una din vocalele *a*, *o*, *u*, atunci se aplică regula *ci* → *чи* (*ciment* → *чимент*, *ciclu* → *чиклу*, *cimbrișor* → *чимбришор*). Dacă cuvântul se termină în *ci*, atunci poate fi aplicată una din regulile: ca excepție *ci* → *ч* (*arici* → *арич*, *beci* → *беч*, *prichici* → *прикич*); *ci* → *чь*, pentru plural (*arici* → *аричь*, *saci* → *сачь*, *maci* → *мачь*); *ci* → *чи*, alte situații (*aci* → *ачи*, *răci* → *рэчи*, *înveșnici* → *ынвешничи*).

f) *cio* → *чо* (*ciorbă* → *чорбэ*, *cioacărlie* → *чокырлие*, *cioban* → *чобан*, *cocioabă* → *кочоабэ*).

g) *ciu* → *чу* (*ciuperchi* → *чуперчь*, *ciubotă* → *чуботэ*, *bucium* → *бучум*).

Utilizând astfel de reguli (lista cărora poate fi prelungită), procesul de transliterare se transformă într-o acțiune de trecere prin „ciur și prin dârmon”. Pornind de la lexiconul contemporan al limbii române [11,12] se stabilește un set de filtre, fiecare filtru având un coeficient de prioritate, care depinde de probabilitatea obținerii unui rezultat corect la aplicarea regulilor acestui filtru. Mai întâi, se vor aplica acele filtre care exclud, sau minimizează, intervenția manuală. Cuvintele filtrate se exclud din lexicon și asupra lexiconului rămas se aplica alte filtre. Din păcate, toate aceste etape de filtrare necesită un anumit grad de intervenție manuală.

Procesarea textelor tipărite cu alfabet latin și litere adiționale

Pentru ilustrarea tehnologiei descrise vom cerea procesul de recunoaștere și verificare a unui text digitalizat din cartea [10], tipărită în anul 1894, (Fig.10).

Textul din Fig.10 a fost recunoscut cu sistemul OCR IRIS. Ca urmare au rămas nerecunoscute cuvintele ortografiate cu litere specifice secolului al XIX-lea. De exemplu, se obține **tnsălbătăcitu** în loc de **însălbătăcitu**.

Acest rezultat nu poate fi îmbunătățit, deoarece IRIS nu posedă capacitatea de a selecta fragmente arbitrare din imagine. Utilizarea unui lexicon modern permite să se recunoască **avutū** ca **avută**, varianta corectă pentru acest context fiind **avut**. Cuvintele specifice lexiconului secolului al XIX-lea nu pot fi recunoscute corect, deoarece pentru aceasta sunt necesare dicționare corespunzătoare perioadei date care, în cazul nostru, ar conține cuvintele **remasū**, **vieța**, **împêratū** etc.

Dacă în textul recunoscut se vor restabili literele specifice și textul obținut se va verifica cu ajutorul

Românii, deși au avut o mie de ani se suferă în-
vasiunile barbare, care au distrus toate operele mărețe
ale arhitecturii romane, în câtu acestu faptu a rămasu
până ađi în ȃicerea populară „n'a rămasu pétră pe
pétră“, totuși nici moravurile nici sufletulu loru nu s'a
însălbătăcită. Ei au păstratū o adăncă intimitate și do-
ioșie în viéța familiară. Căsătoria este încungiurată de-o
mulțime de ceremonii când grave, când vesele. Mirésa
este „o fată de împăratu“, mirele „ficiorū de împăratū“,
ceea ce indică respectū și fericire. Căsătoria este „pe
viéță și mórte“, pentru aceea și jelirea la mórtea u-
nuia dintre soți este adăncă și lungă. In cealaltă lume
însă ér' se întênescu pentru a trăi împreună. Cultulu
moșilorū (sufletele răposăților) este în fôrte mare o-
nóre până ađi. Anumite sêrbătóri peste anū suntă con-
sacrate acestu cultū.

Figura 10. Text digitizat, 1894 (Densușianu, 1984, p. 130)

corectorului ortografic RomSp [9], care posedă un lexicon al limbii române moderne de circa un milion de cuvinte, vom constata că 57 la sută din cuvintele textului sunt recunoscute drept corecte. Acestea sunt cuvintele, ortografierea cărora a rămas intactă față de perioada secolului al XIX-lea, de exemplu, **sufere, acesta, fericire**. Cuvinte „suspicioase” sunt cele afectate de modificările ortografice, de exemplu, **cealaltă (cealaltă), doioșie (duioșie), miie (mie), avut (avut), ađi (azi)**.

Pentru recunoașterea corectă a textului trebuie de instruit sistemul OCR ca să recunoască literele și să completeze lexiconul cu cuvinte noi, specifice secolului al XIX-lea. De exemplu: **avutū, o miie, invasiunile, pétră, nici, sufletulu, lorū, însălbătăcitū, doioșie, viéța, ficiorū, împăratū, miresa**.

Ținând cont de faptul că sistemul OCR ABBYY Fine Reader este înzestrat cu facilități de instruire, am mai efectuat un experiment. Sistemul a fost instruit în mod special ca să poată recunoaște literele specifice secolului al XIX-lea. Iată câteva astfel de litere:

- **ü** (literă finală, mută sau citită),
- **é** (é se pronunță ca diftongul ea),
- **ó** (ó se pronunță ca diftongul oa),
- **đ** (se citea z sau dz),
- **ê** (se folosea ca â).

Unele rezultate ale experimentelor sunt relatate în Tabelul 2.

Pentru a obține rezultate mai performante la verificarea textelor tipărite este necesar ca pentru perioada istorică corespunzătoare:

- să fie instruit scannerul pentru a recunoaște caracterele specifice;
- să fie elaborat un lexicon cu cuvinte și fraze uzuale specifice perioadei;
- să fie extinse facilitățile corectorului ortog-

rafic (spellchecker) pentru a utiliza alfabetul și lexiconul elaborat.

Procesarea textelor tipărite cu alfabet tranziționale

Există cel puțin șapte versiuni ale alfabetului tranzițional (mixt chirilico-latin). Majoritatea literelor acestor alfabet pot fi recunoscute de ABBYY Fine Reader prin evidențierea codurilor respective din setul Unicode. O singură literă specifică pentru aceste alfabet lipsește în Unicode - **Ѡ**. În acest caz urmează să fie inclusă o variantă de literă echivalentă (de exemplu, o săgeată **↑**, sau slavonica „yus” **ѡ**) și instruit sistemul pentru recunoașterea acestei variante grafice.

Tabelul 2

Rezultatele experimentelor OCR cu texte din secolul al XIX-lea

Mod de recunoaștere	Cuvinte corecte	Cuvinte suspecte
IRIS	57%	43%
ABBYY FR, fără instruire	63%	37%
ABBYY FR, cu instruire și dicționar pentru o pagină	98%	2%
ABBYY FR, cu instruire, mai multe pagini, aceeași carte	95%	5%
ABBYY FR, cu instruire, pagini din altă carte	95.4%	4.6%

Concluzii

Resursele digitizate sunt înregistrări specifice stocate într-o bază de date postată pe Internet. Tehnologia propusă se axează pe soluționarea cu succes, pentru fiecare perioadă din evoluția limbii, a două probleme majore: 1. Elaborarea (dezvoltarea) algoritmilor pentru recunoașterea literelor specifice perioadei; 2. Elaborarea instrumentarului și interfe-

țelor necesare pentru crearea resurselor lingvistice (lexiconului) corespunzătoare perioadei în scopul eficientizării procesului de recunoaștere a cuvintelor și de aliniere la normele lingvistice contemporane.

La trecerea de la o perioadă la alta, în limitele posibilităților, se vor utiliza instrumentarul și resursele deja elaborate, materializând astfel principiul „din prezent în adâncul secolelor”.

Resursele electronice create pot fi plasate în Internet pentru acces public, contribuind la dezvoltarea mediului de comunicare informațională pentru limba română. În plus, aceste resurse ar constitui un suport esențial pentru cercetători, iar convertite în text literar contemporan ar putea fi utilizate ca materiale didactice în procesul de instruire.

Bibliografie

1. M. Moruz, A. Iftene, A. Moruz, D. Cristea, *Semi-automatic alignment of old Romanian words using lexicons*, In: Proceedings of the 8-th International Conference „Linguistic resources and tools for processing of the Romanian language”, Iași, Editura Universității „A.I. Cuza”, 2012, p. 119-125.

2. G. Ivănescu, *Istoria limbii române*, Iași, 1980. [G. Ivănescu, History of the Romanian language, Iași, 1980.

3. Ștefan Munteanu și Vasile Țăra, *Istoria limbii române literare*, Editura Didactică și Pedagogică, București, 1978.

4. *Cartea Moldovei (sec XVII – înc. sec XX). Ediții cu caractere chirilice (sec XVII – înc. sec XX)*, Catalog general, Chișinău, 1992.

5. Zamfira Mihail, *155 cărți într-o carte*, Editura Prometeu, Chișinău, 2010, 532 p.

6. Valori Bibliofile-2008, Rev. Gazeta bibliotecarului, Iunie-Iulie 2008, nr. 6-7, p.1 <http://87.248.191.115/bnrm/publicatii/files/3/93.pdf>

7. Gheție I., *Istoria limbii române literare*, București, 1978.

8. Optical Character Recognition (OCR) Technology.

9. 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, p.461-476.

10. Densușianu, A., *Istoria limbii și literaturii române*, Iași, 1894, <http://ru.scribd.com/doc/123035210/Istoria-limbii-si-literaturii-romane>.



Idel Ianchelevici. *Perennis perdurat poeta*, 1972, bronză