

# Generation of the Romanian Cyrillic lexicon for the period 1967 – 1989

Constantin Ciubotaru, Valentina Demidova, Tudor Bumbu

## Abstract

This paper is dedicated to the issue of generation of the Romanian Cyrillic lexicon used between 1967 and 1989. The rules for transliteration of words in the modern Romanian lexicon in their equivalents written in Cyrillic and vice versa are specified and argued. The respective algorithms have been developed and implemented, which is an expert tool in the lexicon generation process. The activity of the expert is reduced to the verification of the transliterated variants and the modification of the transliteration rules.

**Keywords:** lexicon, transliteration, Moldovan Cyrillic, cyrilization, romanization of Cyrillic, morpho-syntactic tagsets.

## 1 Introduction

The problem of digitizing and preserving the linguistic historical heritage is a priority domain of the digital agenda for Europe. The digitization process requires solving a series of problems related to the recognition, editing, translation, and interpretation of printed texts. The solution of these problems has to do with specific difficulties: a large number of periods in the evolution of the language, a small volume of resources widely distributed, a large diversity of alphabets used in their printing.

The paper addresses the issues related to the digitization and transliteration of the historical linguistic heritage printed with Cyrillic alphabet in the period 1967–1989 on the territory of the Moldovan Soviet Socialist Republic following the linguistic norms of the modern Romanian language. The alphabet used in this period (the Moldovan Cyrillic alphabet, AlfCYR) is, in fact, the alphabet of the Russian language from which the letters “ѐ”, “иц” and “ѡ” were excluded and extended in 1967 by the introduction of the letter “ж”. Specific difficulties: lack of resources in electronic format, fragmentary grammatical descriptions that allow ambiguous interpretations.

According to the definition of dexonline [1] *transliteration* is “the transcription of a text from one alphabet to another, rendering the letters by their equivalents, without regard to the phonetic value of the signs”.

The process of transliterating the Romanian words in their written equivalents with characters of the AlfCYR alphabet we will call *cyrillization*. For instance, *puiului* ⇒ *пуюлуй* (chicken), *fiului* ⇒ *фуюлуй* (son), *cenușiu* ⇒ *ченушуй* (greyish), *viermi* ⇒ *верми* (worms), *vierii* ⇒ *виепуї* (boars).

The inverse procedure for cyrillization will be called *romanization of Cyrillic*, e.g. *пуюлуй* ⇒ *puiului* (chicken), *бѡем* ⇒ *biet* (poor), *боер* ⇒ *boier* (boyar), *nenm* ⇒ *piept* (chest).

Taking into account the lack of Romanian Cyrillic resources for this period, the following lexicon generation algorithm is proposed. Sets of rules are defined after which the algorithms for cyrillization and romanization of Cyrillic are constructed. The cyrillization algorithm is applied to the Romanian lexicon (we used the lexicon developed at the “Al. I. Cuza”, Iasi [2]), noted by LexROM. There will be obtained a variant of the Cyrillic lexicon which can be further processed by the algorithm for the romanization of Cyrillic, thus obtaining a new variant for the Romanian lexicon. The ideal situation would be for these two lexicons to coincide. If mismatches occur, the expert(s) intervenes, who can modify the rules of cyrillization \ romanization of Cyrillic, repeat the whole process or intervene with corrections on the constructed

Cyrillic lexicon. The diagram of this process is inserted in Figure 1. The accuracy of the obtained Cyrillic lexicon largely depends on the qualification of the expert(s).

## 2 Cyrillicization

If the problem of digitizing and recognizing the printed text is solved relatively simply, then the problem of cyrillization is more difficult. The rules of transliteration according to cyrillization can be divided into two categories: general rules and context-sensitive rules. The general rules establish the situations when the transliteration directly substitutes letter with a letter(s) ignoring the contextual dependencies. Thus, the letter *b* will be always transliterated into **б**, we will note this through  $b \Rightarrow \mathbf{б}$ . The same action will be performed for the following pairs:  $a \Rightarrow \mathbf{а}$ ,  $\check{a} \Rightarrow \mathbf{э}$ ,  $\hat{a} \Rightarrow \mathbf{ы}$ ,  $c \Rightarrow \mathbf{к}$ ,  $d \Rightarrow \mathbf{д}$ ,  $e \Rightarrow \mathbf{е}$ ,  $f \Rightarrow \mathbf{ф}$ ,  $g \Rightarrow \mathbf{г}$ ,  $h \Rightarrow \mathbf{х}$ ,  $i \Rightarrow \mathbf{и}$ ,  $\hat{i} \Rightarrow \mathbf{и}$ ,  $j \Rightarrow \mathbf{ж}$ ,  $k \Rightarrow \mathbf{к}$ ,  $l \Rightarrow \mathbf{л}$ ,  $m \Rightarrow \mathbf{м}$ ,  $n \Rightarrow \mathbf{н}$ ,  $o \Rightarrow \mathbf{о}$ ,  $p \Rightarrow \mathbf{п}$ ,  $r \Rightarrow \mathbf{р}$ ,  $s \Rightarrow \mathbf{с}$ ,  $\check{s} \Rightarrow \mathbf{ш}$ ,  $t \Rightarrow \mathbf{т}$ ,  $\check{t} \Rightarrow \mathbf{ч}$ ,  $u \Rightarrow \mathbf{у}$ ,  $v \Rightarrow \mathbf{в}$ ,  $x \Rightarrow \mathbf{кс}$ ,  $z \Rightarrow \mathbf{з}$ .

For instance, *chirilizarea*  $\Rightarrow$  **кхирилизация** (cyrillization). Of course, that *chirilizarea*  $\Rightarrow$  **кирилзация** would be correct. In this case, the contextual dependencies intervene. Because such dependencies are many and quite complicated, we will only expose a few. More information on this topic can be found in [4]. From the example above we notice that the transliterations *chi*  $\Rightarrow$  **ки** and *ea*  $\Rightarrow$  **я** are correct. A few more contextual rules: *gi*  $\Rightarrow$  **жи**, *ghi*  $\Rightarrow$  **ги**, *ci*  $\Rightarrow$  **чи**, *ci*  $\Rightarrow$  **ч**, *iu*  $\Rightarrow$  **ю**, *iu*  $\Rightarrow$  **у** (the examples in paragraph 1).

Many difficulties arise when transliterating the letter *i* at the end of the word. All the possible transliterations:  $i \Rightarrow \mathbf{и}$  (*citi*  $\Rightarrow$  **чити** (read)),  $i \Rightarrow \mathbf{й}$  (*pui*  $\Rightarrow$  **пуй** (chicken)),  $i \Rightarrow \mathbf{ь}$  (*arici*  $\Rightarrow$  **аричь** (hedgehog), plural),  $i \Rightarrow$  (*arici*  $\Rightarrow$  **арич**, singular). To make the correct decisions, contextual rules can sometimes be supplemented with morpho-syntactic information. We use MSD (morpho-syntactic-description) tags present in the LexROM lexicon [2]. For example, for infinitive verbs the letter *i* at the end of the word will be translated into **и**, the masculine nouns in the

plural nominative-accusative case the articulated form will end in *uî*, and the transliteration  $i \Rightarrow \mathfrak{b}$  at the end of the word is characteristic for nouns in the plural dative-genitive case, and also for verbs, the second person present, past, more than perfect.

The order of the two categories of transliteration rules is very important. Contextual dependencies always take precedence over general rules.

### CYRILLIZATION ALGORITHM

#### 0. Start

1. The lexicon of the modern Romanian language is given [2] (we will note it  $\text{LexROM}_1$ ) and the rules of cyrillization (general and context sensitive).

\\* We will build the Romanian Cyrillic lexicon for the period 1967–1989 (we will note it by  $\text{LexCYR}$ ) \*\

2. Initially  $\text{LexCYR} = \emptyset$

3. **For all** words  $wrom$  from  $\text{LexROM}_1$ :

3.1. Apply on  $wrom$  the context-sensitive rules for cyrillization. Note the result by  $wcyr_1$ .

3.2. Apply on  $wcyr_1$  general rules for cyrillization. Note the result by  $wcyr$ .

3.3. Include  $wcyr$  in  $\text{LexCYR}$ .

4. Stop

## 3 Romanization of Cyrillic

Romanization of Cyrillic is facing the same problems as cyrillization. The general and contextual rules for this procedure are also defined. The general rules are relatively simple, for example,  $a \Rightarrow a$ ,  $p \Rightarrow r$ ,  $\mathfrak{io} \Rightarrow iu$ ,  $\mathfrak{b} \Rightarrow i$ . If only the general rules apply to transliteration, we obtain, for instance,  $\mathfrak{nyrolyû} \Rightarrow puiului$ ,  $\mathfrak{b\text{ь}em} \Rightarrow biet$ ,  $\mathfrak{boep} \Rightarrow boer$ ,  $\mathfrak{nenm} \Rightarrow pept$ . The last two transliterations are incorrect and their correct variants are:  $\mathfrak{boep} \Rightarrow boier$ ,  $\mathfrak{nenm} \Rightarrow piept$ . In this case contextual

rules are also needed. E.g,  $z\omega \Rightarrow gh\omega$ , if  $\omega \in \{e, u, я, ю, в\}$  and  $z\omega \Rightarrow g\omega$ , if  $\omega \notin \{e, u, я, ю, в\}$  (*гeоргинэ*  $\Rightarrow$  *gheorghinэ* (dahlia), *гогоашэ*  $\Rightarrow$  *gogoaшэ* (donut)). Rules for letter *я*: *я*  $\Rightarrow$  *ia* (usually at the beginning of the word), *ия*  $\Rightarrow$  *ia* (usually at the end of the word), *я*  $\Rightarrow$  *ea*. It is very difficult to make the right decision for the presence of the letter *я* inside the word. In such situations the algorithm will use the rule *я*  $\Rightarrow$  [*ia*][*ea*], leaving the correct decision to the expert.

Another difficult problem to mention is the transliteration of the letter *ы*, which can be substituted either by *î* or by *â*. Our algorithm follows the recommendations of the Romanian Academy regarding this spelling.

## ROMANIZATION OF CYRILLIC ALGORITHM

### 0. *Start*

1. The Romanian Cyrillic Lexicon for the period 1967–1989 (LexCYR) and romanization of Cyrillic rules (general and context sensitive) are given

\\* We will build the modern Romanian lexicon (we will note it by LexROM<sub>2</sub>) by applying the transliteration method \*\

2. Initially LexROM<sub>2</sub> =  $\emptyset$

3. **For all** words *wcyr* from LexCYR:

3.1. Apply on *wcyr* the context-sensitive rules for romanization of Cyrillic. Note the result by *wrom<sub>1</sub>*.

3.2. Apply on *wrom<sub>1</sub>* general rules for romanization of Cyrillic. Note the result by *wrom<sub>2</sub>*.

3.3. Apply on *wrom<sub>2</sub>* rules of transliteration for letter *ы*. Note the result by *wrom*.

3.4. Include *wrom* in LexROM<sub>2</sub>.

4. *Stop*

## 4 Lexicon generation technology

As mentioned above, the electronic resources for the period 1967–1989 are completely absent, complete exposure of the grammar used is missing, many of the interpretations of the transliterated words are ambiguous. Therefore, the expert plays a major role in the lexicon generation process. The proposed technology aims to automate this process. With the available cyrillization and romanization of Cyrillic algorithms, but also access to the formalized rules, the lexicon generation process can take a few iterations. At each iteration, the expert(s) intervene(s) to modify the set of rules and, possibly, directly, the Cyrillic lexicon. This scheme is described in detail in Figure 1.

## 5 Conclusion

The paper proposes a technology for the generation of the Romanian Cyrillic lexicon for the period 1967–1989 applying the transliteration method. Starting from the lexicon of the modern Romanian language elaborated at the University “Al. I. Cuza”, Iasi [2] (1.096.674 words) the cyrillization and romanization of Cyrillic algorithms are applied consecutively. The intermediate results are available to the experts, who can modify or extend the set of rules applied to transliteration, but also directly correct the obtained Cyrillic lexicon. The final lexicon will be a result of performing several such iterations. The main problems that need to be solved by the expert(s) are the ambiguities that arise as a result of cyrillization\romanization of Cyrillic. For example, *iu* ⇒ [*io*] [*uy*], *eie* ⇒ *e* [*e*] [*üe*] [*ue*] – at cyrillization, *я* ⇒ [*ia*] [*ea*] – at romanization of Cyrillic.

For all the words in LexROM starting with the letter “c”, in total 171846 words, 6381 ambiguities were found at the first iteration, which represents 3.7 %. Two iterations were needed to overcome these ambiguities. Of course, the accuracy depends considerably on the qualification of the expert. Also here, at the first iteration comparing the lexicon LexROM<sub>1</sub> with LexROM<sub>2</sub> we obtained 8960 words that do not

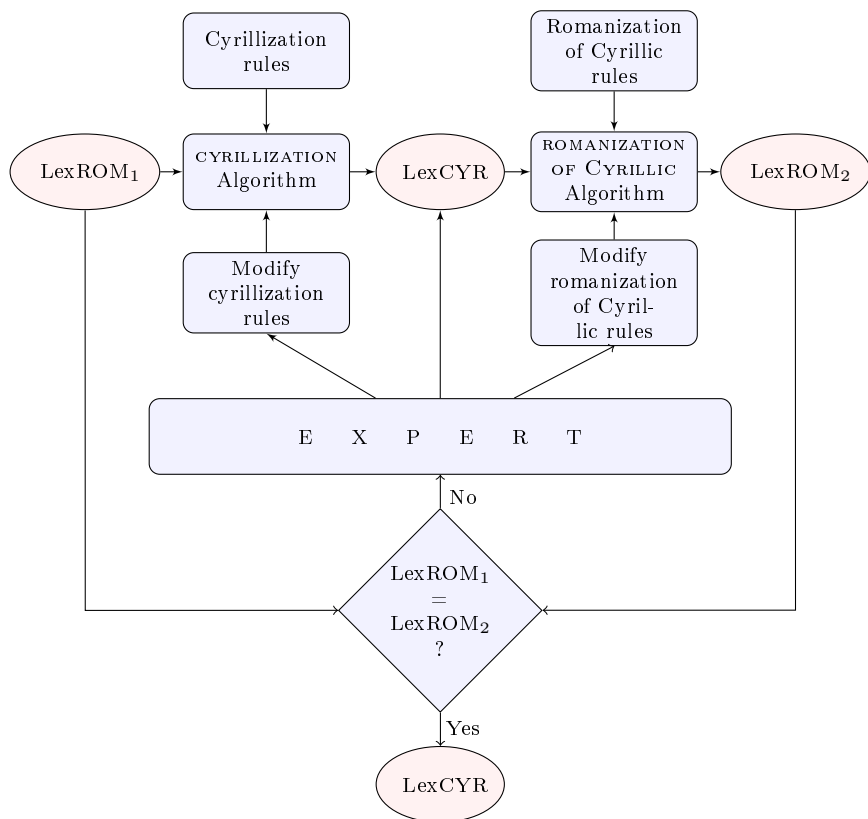


Figure. 1. The scheme for generating the Romanian Cyrillic lexicon

coincide, which represents about 5.21%. Most of these mismatches are incorrect general rules for the romanization of the Cyrillic. The proposed technology allows returning to the previous intermediate variants, thus revising the lexicon. To better understand the role of the expert and contextual dependencies, we applied on LexROM lexicon only the general rules (paragraph 2) of cyrillization. As a result, we got only

42.2% correct words.

## Bibliography

- [1] <https://dexonline.ro/definitie/translitera%C8%9Bie>.
- [2] <http://nlptools.info.uaic.ro/WebPosRo/resources/posDictRoDiacr.txt>
- [3] S. Cojocaru, E. Boian, C. Ciubotaru, A. Colesnicov, V. Demidova, L. Malahov. *Regeneration of printed cultural heritage: challenges and technologies*. Chişinău: The Third Conference of Mathematical Society of the republic of Moldova, 19-23 August, 2014, pp. 481–489.
- [4] V. Demidova. *Particular Aspects of the Cyrillicization Problem*. Chişinău: The Third Conference of Mathematical Society of the republic of Moldova, 19-23 August, 2014, pp. 493–498.
- [5] V. Demidova. *The Peculiarities of the Romanization of Cyrillic*. *Studia universitatis Moldaviae*, 2015, no. 2(82), Seria “Ştiinţe exacte şi economice”, pp. 16–20 (in Romanian).

Constantin Ciubotaru<sup>1</sup>, Valentina Demidova<sup>1</sup>, Tudor Bumbu<sup>1,2</sup>

<sup>1</sup>Vladimir Andrunachievici Institute of Mathematics and Computer Science

<sup>2</sup>The State University "Dimitrie Cantemir"

E-mails: [chebotar@gmail.com](mailto:chebotar@gmail.com), [demidova@math.md](mailto:demidova@math.md), [bumbutudor10@gmail.com](mailto:bumbutudor10@gmail.com)