Vladimir Polomac
Tamara Lutovac Kaznovac

## AUTOMATIC RECOGNITION OF SERBIAN MEDIEVAL MANUSCRIPTS BY APPLYING THE *TRANSKRIBUS* SOFTWARE PLATFORM: CURRENT SITUATION AND FUTURE PERSPECTIVES

The paper investigates the potentials of applying the model for the automatic recognition of (Russian) Church Slavonic manuscripts to Serbian medieval manuscripts written in various types of Cyrillic scripts by employing the Transkribus software platform. The analysis has shown: (a) that the use of the existing generic model for the recognition of Church Slavonic manuscripts can yield rather good results when applied to Serbian medieval manuscripts written in uncial or semiuncial script, (b) that the manuscripts written in the cursive script require the creation of a separate model, and (c) that the creation of a generic model within the Transkribus platform for the Serbian medieval manuscripts would make the process of digitization substantially faster, which in turn would lead to faster realization of tasks within the existing projects related to Serbian historical corpus linguistics and lexicography .

*Key words*: Transkribus, Serbian medieval manuscripts, automatic text recognition, information technology, artificial intelligence, machine learning.

У раду се истражују могућности примене модела за аутоматско рашчитавање (руских) црквенословенских рукописа у оквиру софтверске платформе Transkribus на српске средњовековне рукописе писане различитим типовима ћирилице. Анализа је показала (а) да примена постојећег генеричког модела за аутоматско рашчитавање црквенословенских рукописа може дати веома добре резултате на српским средњовековним рукописима писаним уставом или полууставом, (б) да је за рукописе писане брзописом неопходно креирати посебан модел, и (в) да би се креирањем генеричког модела за српске средњовековне рукописе у оквиру платформе Transkribus, процес дигитализације могао значајно убрзати, што би даље могло водити и убрзању рада на текућим пројектима из српске историјске корпусне лингвистике и лексикографије.

*Кључне речи*: Transkribus, српски средњовековни рукописи, аутоматско рашчитавање текста, информационе технологије, вештачка интелигенција, машинско учење.

**1.** Introduction. The basis for starting the work on this paper is the article by A. Rabus (2019a), on the potentials of the automatic recognition of Church Slavonic manuscripts by using the Transkribus software platform. In the article, the first of its kind within Slavonic studies, the author starts by providing a brief overview of previous attempts to develop the technology for the automatic recognition of medieval manuscripts,[1] including the technology behind the Transkribus

---

[1] The paper primarily focuses on rare occasions of applying the OCR technology to the automatic recognition of old Slavonic printed books and printed editions of Church Slavonic manuscripts (the electronic edition of Bdin collection (*Bdinski sbornik*) and RRuDI corpus), but also emphasizes its complete inapplicability to automatic recognition of old Slavonic manuscripts (cf. Rabus 2019a: 10). The possibility of automatic recognition of old Slavonic manuscripts and printed books by using

platform,[2] reserving the central part of his paper to the creation of the model for the automatic recognition of Church Slavonic manuscripts,[3] accompanied by the quantitative and qualitative analysis of the results obtained by applying the said methods to the manuscripts written in various types of Church Slavonic Cyrillic script. The overall significance of the Rᴀʙᴜs 2019a paper lies in the fact that by using concrete examples the paper convincingly showed that automatic recognition of Church Slavonic manuscripts by applying the Transkribus is indeed a reality, since the first version of the recognized text in the electronic form that has an acceptable error ratio (about 4% of all characters)[4] is automatically achieved, and after the manual correction done by a competent philologist, investing much less time and human and financial resources, the result is the text of the manuscript in electronic form, suitable for further philological and linguistic investigations. In doing so, this paper provides direction on how to significantly expedite the digitization of medieval Slavonic manuscripts by the use of the Transkribus platform, thus providing a significant impetus for the development of diachronically oriented Slavonic studies.[5] The special importance of the Rᴀʙᴜs 2019a paper can be found in the fact that the models for the automatic recognition of Church Slavonic manuscripts are made publicly available through the Transkribus platform, so their potential can be tested on other medieval Slavonic manuscripts as well. The investigation concerning the application of possibilities of such models to medieval Serbian manuscripts written in various types of Cyrillic script represents the overall goal of the current article. Since the models in question are based on

---

the artificial intelligence based on neural networks was for the first time mentioned in Кᴏᴘɴɪᴇɴᴋᴏ – Чᴇᴘᴇᴨᴀɴᴏᴠ – Яsɴɪᴄᴋɪᴊ 2008. This paper is more significant in a theoretical rather than practical sense since the level of character recognition accuracy is around 80%, thus requiring a large amount of time for manual corrections to the text, thus rendering the entire recognition process no more economical than the traditional approach (cf. Rᴀʙᴜs 2019a: 10).

[2] Transkribus is a free access software platform for the automatic recognition and search of manuscripts which was developed within the READ project at the University of Innsbruck. Unlike the traditional approach that focuses on individual letters (OCR technology), Transkribus uses HTR technology based on memorizing and recognizing the entire image of the line from the text. Recently developed and implemented into Transkribus, HTR+ algorithm is based on the artificial intelligence and advanced neural networks and significantly reduces the time required for the training of text recognition models, with a substantially higher accuracy ratio. For more details, see Rᴀʙᴜs 2019a: 10–11.

[3] The functionality of the Transkribus platform is particularly manifested in the potential to train one's own automatic text recognition model, irrespective of the language or script used in the manuscript. The training of the automatic recognition model represents an instance of machine learning based on neural networks in which during the learning process the model compares the manuscript photographs and corresponding letters, words and lines of the text in the diplomatic edition. The successful training of a model requires photographs of the manuscript having the best possible quality and at least 15000 words of previously recognized text. For more details, see Rᴀʙᴜs 2019a: 11–14.

[4] Transkribus possesses the possibility to automatically calculate the ratio of incorrectly recognized letters (CER) by comparing the automatically recognized version of the text and manually corrected version. For more details, see Transkribus Glossary at https://readcoop.eu/glossary/character-error-rate-cer/.

[5] Automatic recognition of Serbian (medieval) manuscripts could significantly expedite the work on the current Serbian historical lexicographic projects (*Dictionary of the 12th–18th Century Serbian Language* and *Dictionary of the Slavonic Serbian Language*), as well as the preparation of the electronic historical corpus of the Serbian language.

Old Church Slavonic and Russian Church Slavonic manuscripts written in uncial or semiuncial Cyrillic scripts, our paper starts from the hypothesis that their application to medieval Serbian manuscripts written in the uncial or semiuncial can yield more or less acceptable results,[6] while the manuscripts written in cursive Cyrillic script should require the creation of a separate recognition model. In structuring the paper we proceeded in line with the stated hypotheses. Accordingly, the Section 2 provides a detailed overview of the existing models for the automatic recognition of medieval Slavonic manuscripts, including the analysis of the results of their application to medieval Serbian manuscripts written in various types of Cyrillic script, while the Section 3 provides concluding remarks and the perspectives for further research.

**2.** Application of existing models for the automatic text recognition to Serbian medieval manuscripts. Two models for the automatic recognition of Church Slavonic Cyrillic manuscripts are available as a part of the Transkribus software platform.[7] The first model, called VMČ_Test_4+, is based on portions of the Russian Church Slavonic manuscript *The Great Reading Menology*, written in semiuncial 16th century Cyrillic script. A total of 173,287 words were used for the training of the model, with CER (Character Error Rates) being 3.82% (more details can be found in Rabus 2019a: 15–19). The other model, dubbed Combined_Full_VKS_2 and based on parts of the Old Church Slavonic *Codex Suprasliensis* (11th century), *The Catecheses of Cyril of Jerusalem* manuscript (11th century) and the Russian Church Slavonic manuscript *Great Reading Menology* (16th century), represents an attempt at creating a generic model suitable for the automatic recognition of different manuscripts written in uncial or semiuncial Cyrillic script. A total of 393,079 words were used for the training of the model, with 3.94% CER (for more details, see Rabus 2019a: 23–27).[8] These models were tested on Serbian medieval manuscripts which are currently the focus of interest of our philological and linguistic investigations: a) on Serbian medieval charters and letters currently being prepared to be used for the development of a specialized electronic corpus (cf. Polomac 2021), as well as b) on the Serbian Church Slavonic manuscript by the name of *Christian Topography of Cosmas Indicopleustes*, which is in the process of preparation for publication in its original graphemic structure together with the associated philological and linguistic studies.[9]

2.1. Automatic recognition of Serbian medieval charters and letters. When choosing the charters and letters to be included in the investigation, we took into account the size of the manuscripts, their philological and cultural significance, state of the manuscripts' preservation and legibility, as well as the availability of

---

[6] The term *semiuncial* denotes the Resavian type of uncial script (cf. Јерковић 1996).

[7] A. Rabus has also created two publicly accessible models for the automatic recognition of the Glagolitic script: the first model contains approximately 28,000 words from various printed Glagolitic books from Tübingen and Urach (see https://readcoop.eu/model/glagolitic-print/), while the other model comprises approximately 171,000 words from the Breviary of Vid of Omišalj and the Second Beram Breviary (see https://readcoop.eu/model/glagolitic-handwritten-14th-and-15th-century/).

[8] For the potentials and problems in the creation of a generic model for the automatic text recognition within the Transkribus platform, see Rabus 2019b, Hodel et al. 2021.

[9] This edition and study is being prepared for publication by Tamara Lutovac Kaznovac.

images with proper quality. Having in mind these criteria, the expected choice was made by selecting the two most important 14[th] century Serbian charters, written in the uncial type of Cyrillic script: *Banjska Chrysobull* (henceforth: BC) (cf. Трифуновић 2011), and also *Dečani Chrysobull* (third version, henceforth DC III) (cf. Ивић – Грковић 1976: 34–37). The study also includes several other charters and letters from 14[th]–15[th] centuries written in uncial (semiuncial) and cursive Cyrillic script (more on this in 2.1.3). Our methodological approach involved conducting a respective experiment for each manuscript, as well as quantitative and qualitative analyses of the obtained results. Following the process of the automatic recognition of selected manuscript folios (for BC and DC III) or entire manuscripts (in the case of shorter charters and letters) by using the said models, we conducted a manual correction of the text. By comparing the automatically recognized text with the corrected version of the text, we calculated the ratio of unrecognized characters (CER), and this was followed by the qualitative analysis which especially took into account the performance of the models depending on the photographic image quality.

2.1.1. In the first experiment[10] the performance of the VMČ_Test_4+ and Combined_Full_VKS_2 models was tested when applied to the first ten folios of BH (from 5r to 9v). The statistical overview concerning the ratio of incorrectly recognized characters (CER) is given in the following table:
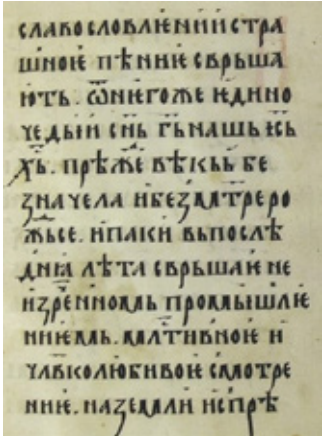
Table 1: CER in BH (sheets 5r–9v)

| Model | VMČ_Test_4+ | Combined_Full_VKS_2 |
|---|---|---|
| Folio | CER | CER |
| 5r | 22,18% | 23.79% |
| 5v | 20,51% | 12,82% % |
| 6r | 26,18% | 16,73%% |
| 6v | 16,95% | 13,90 % |
| 7r | 19,62% | 10,94% |
| 7v | 20,08% | 8,11% |
| 8r | 13,99% | 12,94% |
| 8v | 20,65% | 10,51 % |
| 9r | 20,35% | 18,95% |
| 9v | 21,91% | 8.13 % |
| Mean | 20,24% | 13,68% |

The above overview indicates that the application of the VMČ_Test_4+ model did not produce satisfactory results (mean CER is 20.24%), while the results

___

[10] In all conducted experiments the automatic recognition was performed by employing the linguistic model. For more details, see https://readcoop.eu/transkribus/howto/how-to-train-a-hand-written-text-recognition-model-in-transkribus/.

of applying the generic Combined_Full_VKS_2 model can be estimated as relatively good (mean CER is 13.68%). A comparative representation of the folio 5v photograph and the corresponding automatically recognized text in Table 2 illustratively provides qualitative insight into the performance of both models.

Table 2: VMČ_Test_4+ and Combined_Full_VKS_2 and BH (folio 5v)

| VMČ_Test_4+ | Трифуновић 2011: 20 | Combined_Full_VKS_2 |
|---|---|---|
| славословлени и стра-шное пѣннѣ сврьша-ють. wнегоже единочедыи сн҃ь гб҃ нашь сѫ хъ. прѣже вѣквь без-зна҃ела и без вагреро-жь се. и паки вь послѣ-дн҃алѣта врьшаи не-изренномь пролышли. ни киль. валтивное и г҃лѣколюбивое смотре-нии. на зевали испрѣ- |  | славословлени и стра-шноюе пѣние сврьша-ють· о нѥгоже ѥдино-едыи сн҃ь гб҃ нашь сь хъ. прѣже вѣкьь без-зна҃ела и без мт҃ре ро-жь се. и пакивь послѣ-дн҃а лѣта сврьшаю не и ꙁренномь промышлѥ-ниѥмь· мл҃тивноюе и ч҃лѣколюбивоюе смотре-не. на земли испрѣ- |

Based on the representation given above, we can conclude that both models make recognition errors most frequently when the pajerak mark is involved, which was expected since the mark was not registered in the model training process: instead of the expected страш᾽ное 1/2, без᾽ 6, послѣдн҄а 7/8, неизрѣн᾽номь 8/9, промышлѥниюемь 9/10, мл҄тив᾽ноюе 10, смотрение 11/12, ис᾽ прѣ- 12, the VMČ_Test_4+ model incorrectly recognizes страшное 1/2, без 6, послѣдн҄а 7/8 неизренномь 8/9, пролышли. ни киль. 9/10, валтивное 10, смотрении 11/12, испрѣ- 12, while the generic model recognizes страшноюе 1/2, без 6, послѣдн҄а 7/8, не и ꙁренномь 8/9, промышлѥниюемь 9/10, мл҄тивноюе 10, смотрене 11/12, испрѣ- 12. A large number of mistakes in both models is attributed to the recognition of blanks between words: instead of славословюении 1, без на҄ела 5/6, ис᾽ прѣ- 12, both models recognize славословени и 1, безна҄ела 5/6, испрѣ- 12; instead of ѿ нюегоже 2, мт҄ре рожь 6/7, послѣдн҄а лѣта 7/8, the VMČ_Test_4+ model recognizes wнегоже 2, вагрерожь 6/7, послѣдн҄алѣта 7/8, while instead of паки вь 8, неизрѣн᾽номь 8/9 the generic model recognizes пакивь 8, не и ꙁренномь 8/9. The recognition of the titlo mark and superscript letters also poses a problem for both models, yet the generic models shows somewhat more successful results: instead of х҄ь 5, рож᾽ се 6/7 both models recognize хъ 5, рожь се 6/7; instead of гб҄ь 4, неизрѣн᾽номь 8/9 the VMČ_Test_4+ model renders сѫ 4, неизренномь 8/9, while the generic model renders сь 4, не и ꙁренномь 8/9; concerning the expected мт҄ре рожь 6/7 and мл҄тив᾽ноюе 10 the VMČ_Test_4+ model renders вагрерожь 6/7 and валтивное 10, while the generic model manages to recognize the titlo mark and a superscript letter in these examples.
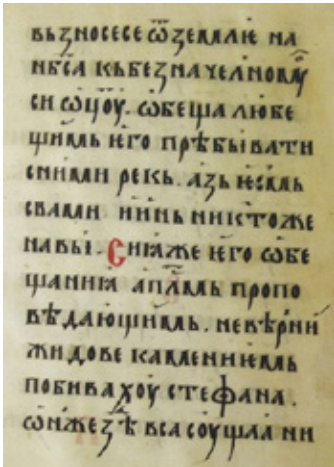
The difference in performance of the two models results from a large number of errors when the VMČ_Test_4+ model attempts to recognize the letters

ѥ and м: instead of славословѥнии 1, страшʼноѥ 1/2, пѣниѥ 2, ѿ нѥгоже 2, ѥдиноѹедыи 3/4, свръшаѥ 9, промышʼлѥниѥмь 9/10, мѫтивʼноѥ 10, ѹлвколюбивоѥ 11, сʼмотрениѥ 11/12, the VMC_Test_4+ model renders славословени и 1, страшное 1/2, пѣннѣ 2, wнѥгоже 2, единоѹедыи 3/4, врьшаи 9, пролышли. ни киль. 9/10, валтивноѥ 10, Глвколюбивоѥ 11, смотрении 11/12; instead of мтⷬе рожь 6/7, промышʼлѥниѥмь 9/10, мѫтивʼноѥ 10, земли 12 there are вагрерожь 6/7, пролышли. ни киль. 9/10, валтивноѥ 10, зевали 12. The remaining letter recognition errors of this model were registered in a small number of examples: the letter т – instead of мтⷬе рожь 6/7 we have вагрерожь 6/7; the letter ь – instead of ісь 4, вѣкьь 5 there is съ 4, вѣквь 5; the letter ѹ – instead of ѹлвколюбивоѥ 11 we have Глвколюбивоѥ 11; the letter ѿ – instead of ѿ нѥгоже 2 there is an incorrect wнѥгоже 2; instead of пѣниѥ 2 we have an incorrect пѣннѣ 2.

The qualitative analysis shows that the performance of the generic model is significantly better that it appears to be judging solely by the CER ratio. As has already been observed, the largest number of errors is associated with the recognition of the pajerak mark and blanks between words. The examples in which the generic model incorrectly recognizes letters are not numerous. This model also has problems recognizing the ligature ѥ in several examples: instead of славословѥнии 1, свръшаѥ 9, сʼмотрениѥ 11/12 there are incorrect славословени и 1, свръшаю 9, смотре не 11/12. The remaining letter recognition errors were registered as single instances: the letter а – instead of паки вь 8 we have the incorrect плкивь 8; the letter ѿ – instead of ѿ нѥгоже 2 there is the erroneous о нѥгоже 2; the letter и – instead of сʼмотрениѥ 11/12 we have смотре не 11/12; the model omitted letters in two examples – ѹ in ѥдиноедыи 3/4 (from the expected ѥдиноѹедыи 3/4) and н in смотре не 11/12 (from expected сʼмотрениѥ 11/12).

The potential of the generic model can further be illustrated when applied to the folio 7v in which CER is only 8.11% (Table 3).

Table 3: Combined_Full_VKS_2 and BH (folio 7v)

| Трифуновић 2011: 24 | Combined_Full_VKS_2 | Ground Truth |
|---|---|---|
|  | вьзно се се ѿ землю на нбслкь безнаѥлному си оцбоу. обеца любе-цимь ѥго прѣбывати с ними рекь. азь ѥсми свлми. и инь никстоже навы. Сиа же ѥго обе-цаниа. апⷧлмь пропо-вѣдаюциⷨмь. невѣрни жидове камениѥмь побⷪивахоу стефана. он же зⷮѣ вса соуцаа ни- | вьзносе се ѿ землѥ на нбса кь безнаѥлʼному си wцⷪоу. wбеца любе-цимь ѥго прѣбывати с ними рекь. азь ѥсми с вами. и инь никʼтоже на вы. Сиа же ѥго wбе-цаниа. апⷧлмь пропо-вѣдаюциⷨмь. невѣрʼни жидове камениѥмь побивахоу стефана. wнⷪ же зⷮѣ вса соуцаа ни- |

The largest number of errors in this folio is again attributed to the recognition of the pajerak mark and blanks between words: instead of безнач҄ел'ному 2, ник҃тоже 6, невѣр́ни 9, ѡн' 12 there is the incorrect безначелному 2, никтоже 6, невѣр́ни 9, он 12; instead of вьзносе 1, нб҃са кь 2, с вами 6, на вы 7 the generic model renders вьзно се 1, нбслкь 22, свлми 6, навы 7. Expectedly, there are errors connected with the titlo mark and the use of superscript letters: instead of нб҃са кь 2, безнач҄ел'ному 2, повивахоу 11 we have the incorrect нбслкь 2, безначелному 2, повивахоу 11. A small number of errors in letter recognition can most frequently be associated with the letters ѡ and ѧ: instead of ѡц҃оу 3, ѡвещѧниѧ 7/8, ѡн' 12 the generic model incorrectly renders оцоу 3, обещаниа 7/8, он 12; instead of нб҃са кь 2, с вами 6 there are incorrect нбслкь 2, свлми 6. In one example the model makes an error when recognizing the ligature ѥ and the letter є: instead of землѥ 1 and камениѥмь 10 we have the erroneous землю 1 and камениємь 10.

2.1.2. In the next experiment we tested both models on a portion of DC III (in the folios 8r–10v, as well as the folio 76v, which is the initial part of *The Nun Evgeniya's Charter to the Dečani Monastery* (see Младеновић 2007: 391–406). The statistical overview of the incorrectly recognized character ratio (CER) is given in the following table.

Table 4: CER in DC III (8r–10v, 76v)

| Model | VMČ_Test_4+ | Combined_Full_VKS_2 |
|---|---|---|
| Folio | CER | CER |
| 8r | 22.46% | 13.45% |
| 8v | 29.64% | 15.84% |
| 9r | 24.82% | 13.23% |
| 9v | 23.26% | 15.84% |
| 10r | 24.97% | 12.79% |
| 10v | 27.83% | 13.79% |
| 76v | 19.34% | 9.90% |

The generic model yielded results almost twice better on average than VMČ_Test_4+ model. The successfulness of the generic model is well evidenced by the comparative representation of the folio 76v and the automatically recognized text (Table 5).

Table 5: Combined_Full_VKS_2 and DC III (folio 76v)

| Младеновић 2007: 399 | Combined_Full_VKS_2 |
| --- | --- |
|  | ...ρѣ пости моа. и гѝ възлюбихь блголѣпиѥ домоу твоѥго· й мѣсто вьселенїа славы твоі ѥ́ ρε̑ꙋ бжтвьный дѣдь. свѣтло бо нниꙗликьствоуѥть веселиѥ красоующи се бжть. вна̀ цр̑квы днь. и хо̑ꙋ ρꙋ̑ающїй се въпиеть, прѣиспьцрена дх̑а блг̑тию, крѣпо̑ моа ѥси и пѣнїе. Сего ρа̑ вьспоюте все дн̑їи живота моѥго. ибо прильпе дш̑а моа по тебѣ, мене же приѥть десница твоа. нн̑а бо блгоцвьтоущїи витсе. и ꙗко цр̑ко ꙏ дѣꙗн̑наа порфіроꙑ добрꙑ дѣль о множениемь и ꙗко лоза плодовитаа Въ стρа̑ на домоу бж̑иꙗ. сн̑ове ꙗко лѣтораслихла слыннїи. ва̑ и ино сеце блгодѣꙗниꙗ. и прѣвѣньства чинови прѣвьсходеце. и добродѣтелию дроугь дроугь дроугꙿа ретоуюце. цр̑ию Благоꙋѣстиѥмь· блгопокорениимь силниценно наꙗелници оукрашениимь. блгоговѣꙏниимь сцен̑ници. чиномь всачьскꙑи оукрашенна. по мѣρѣ кꙑꙗко по дарованию дх̑а добродѣтели пло̑ деце. Съвоꙋзолюбве съдрьжеце се̑, въ |

Based on the given representation, it can be concluded that the generic model most frequently confuses letters и and ѥ in different positions: твои ѥ̑ 3, ва̀ и ино сеце 13, благопокорениимь 16, оукрашениимь 17 and блгоговѣꙏниимь 17/18 instead of твоѥе 3, ваниѥ носеце 13, блгопокорениѥмь 16, оукрашениѥмь 17 and блгоговѣниѥмь 17/18. In two examples the letters ь and ъ are confused as a part of a preposition and a prefix: възлюбихь 1, Съвоꙋзолюбве 20 instead of вьзлюбихь 1, сьвоꙋзо̑ любве 20. Other errors in the recognition of the letters amount to single instances: твоꙗ 9 instead of твоꙗ 9, ꙏ дѣꙗннаа 10 instead of ꙍдѣꙗннаа 10, о множениемь 11 instead of Ꙋмножениѥмь 11, хласлыннїи 13 instead of маслыннїи 13, цр̑ию 15 instead of цр̑иѥ 15, ценно 17 instead of Сцен̑но 17. Along with these errors, in a large number of examples we registered the failure to recognize pajerak mark, superscript letters and the titlo (most commonly superscript с under the titlo: днь 5, хоꙋ 5, крѣпо̑ 6, цркою 10 instead of дн̑ь 5, хо̑ꙋ 5, крѣпо̑ 6, цркою 10), and the blanks between words. Certain errors occur because the model recognizes accent marks as superscript letters: for instance, твои ѥ̑ 3, красоуюцꙵй 4, ρꙋ̑ающїй се 5 and прѣвьсходеце 14 instead of твоѥе 3, красоуюци 4, ρꙋ̑аюцїи се 5 and прѣвьсходеце 14.

A relatively good result is also recorded by the generic model in the folios 8r–10v DC III (CER is 14.15% on average). This special mention is due to the fact that the generic model yields approximately the same results in the folio 76v written in the Serbian Church Slavonic and the folios 8r–10v which mostly consist of the proper names of people from the monastery grounds.

Table 6: Combined_Full_VKS_2 and DC III (folio 8r)

| Ивић – Грковић 1976: 143 | Combined_Full_VKS_2 |
| --- | --- |
|  | трикь. тврьдоѥ абра-<br>леньмоу хранко́ и драико. ра-<br>викь. богоѥ голіа. мирко рада новикь. ме-<br>доѥ абра тмоу милошьлмць имь добро-<br>славь. срьдакь абра́ тмоу прибць и милета.<br>веселко и радосла вьлдѣдь имь дражоѥ ни<br>зоуклинь. богꙗнь каменарь. николаписко-<br>тикь. смиль абратмоу милеша и добро-<br>славьлмць имь братославь. доброславькоу-<br>манови. смиль абра тмоу и вань а щц̈ь им<br>....рославь.<br>богоѥ ко влꙋ̈ꙗ сн̈ьмоу мирославь. храни сла-<br>вь везили ꙗ а щц̈ь моугюргь. доброу и абра-<br>тмоу богань· Дее ꙋрьвено брѣжа не за селкь<br>дѣꙋаньскни. милко а братмоу ни но славь-<br>истлико. приби славь милкови и богоѥ и<br>мирославь. милошьни послалилдѣдь имь<br>хоунко. добрьꙋинь абратвлоу радославь кра-<br>дмоужь а щц̈ь имь добрень. прибоѥ а сн̈<br>моу доброславь. боудоѥ абратва оура икоа дѣ-<br>дь имь щбрадь. храноѥ адѣ дмоу щ зрина.<br>ра и коа братмоу богоѥ и боже та ими лѣн |

Based on the comparative representation of the folio 8r and corresponding automatically recognized text (Table 6) it can be concluded that the generic model most frequently does not recognize the letter ꙗ (confuses it with л): леньмоу 2, лмць 3, 9, ко влꙋ̈ꙗ 12, лдѣдь 5, 17, истлико 16 instead of ꙗ сн̈ моу 2, а щц̈ь 3, 9, коваꙗь 12, а дѣдь 5, 17, и стаико 16, as well as the letter ѡ̈ (confuses it with the letters ѡ and о): щц̈ь 9, щц̈ь 10, 13, 19 instead of ѡ̈ц̈ь 9, 10, 13, 19, letter м: ....рославь 11, абратвлоу 18, абратва оура икоа 20 instead of мирославь 11, а братъ моу 18, а брат моу раико а 20, and the letter є: леньмоу 2, Дее 14 instead of ꙗ сн̈ моу 2, Ꙗ се 14. Other errors in the recognition of letters are restricted to individual instances: Дее 14 instead of Ꙗ се 14 (capital letter), милошьни послалилдѣдь 17 instead of милошь нинослаꙗи а дѣдь 17, ѡ̈ зрина 21 instead of ѡзрина 21. One example can be singled out: крадмоужь 20/21 instead of радмоужь 20/21, in which the initial letter is crossed out. Like the folio 76v, the folios 8r–10v contain the largest number of errors that are associated with the recognition of pajerak mark (it occurs quite frequently), superscript letters and titlo mark. What is characteristic of the folios 8r–10v is that the model does not recognize blanks between words, probably because in the course of training there was no opportunity for it to gain insight into the specific onomastic vocabulary of the charter.

2.1.3. Encouraged by the relatively good results of the generic model applied to DC III, we conducted additional experiments with several charters and letters from the 14th–15th centuries written in different types of Cyrillic script.

A statistical overview of the results when applying the generic model to the letters from the 14th–15th centuries written in cursive script is given in the following table.
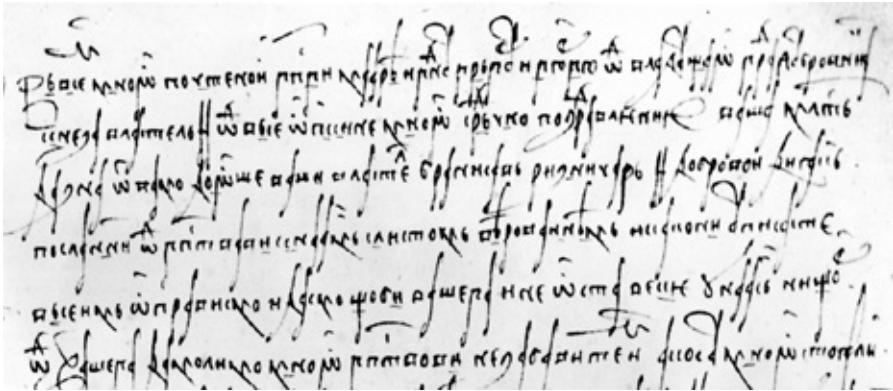
Table 7: Combined_Full_VKS_2 and cursive scripts of 14th–15th century[11]

| Letter | CER |
| --- | --- |
| Emperor Stefan Uroš V's Letter to Dubrovnik (around 1358) | 51.15% |
| Letter of Jerusalem Metropolitan Mihailo to Dubrovnik (1386) | 43.46% |
| King Tvrtko I Kotromanić's Letter to Dubrovnik (may 1389) | 38.44% |
| Letter of Dubrovnik to Lady Mara Branković and her Sons (1402) | 55.20 % |
| Letter from Dubrovnik about the Settlement of Nikša Sorkočević's Debt (1419) | 50.48% |
| Letter of Turkish Sultan Murad II to Dubrovnik (1431) | 56.33% |

The data above confirm the initial hypothesis that the Combined_Full_ VKS_2 generic model is not suitable for the automatic recognition of the charters and letters written in cursive Cyrillic script, which is expected since the model was trained exclusively by using the manuscript material written in uncial or semiuncial. An extraordinarily high percentage of errors indicates that it is necessary to train a separate model for the automatic recognition of manuscripts written in cursive script.

The lack of usefulness of the generic model in the process of recognizing cursive Cyrillic script is illustratively evidenced by the comparative representation of the photographs of the letter which was sent to Lady Mara Branković and her sons in 1402 from the Dubrovnik office (cf. Стојановић 1929: 146–147) and the corresponding automatically recognized text in the following table.

Table 8: Cursive Cyrillic script and Combined_Full_VKS_2

| Letter from Dubrovnik to Lady Mara Branković and Her Sons (1402) |
| --- |

| Combined_Full_VKS_2 |
|---|

Велєнно ѡно ꙋтено и инини порѣ и нн҃ѣ иренѡ҃ и ꙋѳно ѿ плаꙋено ѡ прорⷣа. проꙋни
се нєдоꙋ поꙋтельнꙋ҃ все ѿ псеи немно ѿ сроꙋно посрѣвлени, помохⷶ҃ть
ндⷶ҃ хѡ҃. но аще и ции властѣ҃ срⷩици рⷩа ни ꙋтрь. и о пробои. ꙋинъ.
посланни нанесенїю листовь вѣрованнѣхь и скони оннѡꙑꙗнє-
вьсеихь ѡ прⷢци и схо искоꙵови оꙋшенїи не ѡсновенїе но се҃. нициѡⷹ. аı҃.
ѿвашего. холихохноѡрнилова не ревнꙵ теи. ꙗкѡ҃ сонѡ ѿ ст҃оель

The statistical overview of the results of applying the generic model to the 14th–15th century charters written in uncial script has been given in the following table.

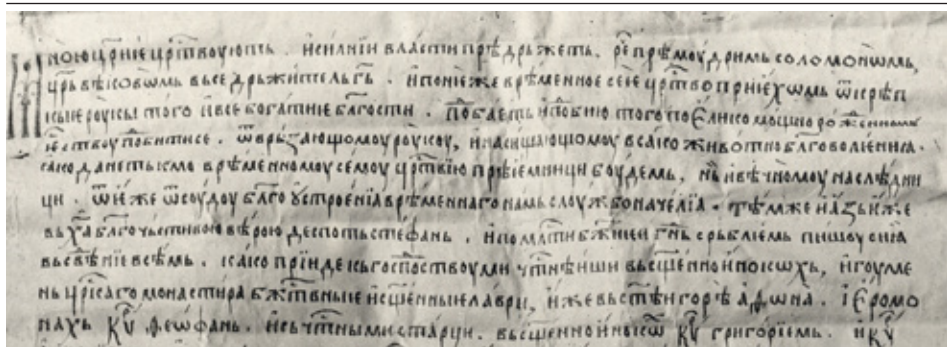Table 9: Combined_Full_VKS_2 and charters from 14th–15th century written in uncial script

| Charter | CER |
|---|---|
| Dečani Chrysobull (first version, 1330–1331) (lines 1–19) | 17.55% |
| Dečani Chrysobull (first version, 1330–1331) (lines 20–49) | 24.46% |
| Emperor Dušan's Charter to the St. Sava Cell in Karyes (1348) (lines 1–35) | 12.80% |
| Emperor Dušan's Charter to the St. Sava Cell in Karyes (1348) (lines 36–64) | 13.91% |
| Despot Stefan Lazarević's Charter to Despotess Yevpraksia (1404–1405) | 9.99% |
| Lady Mara Branković and her Sons' Charter to Dubrovnik (1405) (lines 1–27) | 18.80% |
| Lady Mara Branković and her Sons' Charter to Dubrovnik (1405) (lines 28–43) | 27.36% |
| Despot Stefan Lazarević's Charter to the Mileševa Monastery (1413) | 11.20% |
| Despot Stefan Lazarević's Charter to the Great Lavra Monastery (1414–1415) | 6.78% |

The ratio of unrecognized characters to other charters of 14th–15th centuries written in uncial Cyrillic script when the generic model is applied is mostly positioned within the ranges recorded in the BC and DC III folios. The exceptions are the first version of the Dečani Chrysobull and Lady Mara Branković and her Sons' Charter to Dubrovnik (1405), where the CER is higher than in the other charters. Along with the expected errors related to the recognition of the pajerak mark, superscript letters and titlo mark in both charters, the high CER in the first version of the Dečani Chrysobull can also be explained by the lower quality of the photograph, which in turn led to problems in recognizing entire portions of the text, while in the case of Lady Mara Branković and her Sons' Charter to Dubrovnik problems can be attributed to popular vocabulary which the model had no opportunity to familiarize itself with during the training, and, just like in the case of DC III, this led to a greater number of errors concerning the blanks between words.

The potential of the generic model in relation to the recognition of charters written in the Serbian Church Slavonic and uncial Cyrillic script is best illustrated through the example represented by Despot Stefan Lazarević's Charter to Great Lavra Monastery (1414–1415). A comparative representation of the image of the first ten lines of the charter and the automatically recognized text is given in the following table.

Table 10: Despot Stefan's Charter to Great Lavra Monastery (1414–1415) (lines 1–10)

| Младеновић 2007: 283 |
| --- |



| Combined_Full_VKS_2 |
| --- |

оунȣюцр̄нїец̄ртвоують. и силнїи власти пр́ѣдръжеть, р́ечъ пр́ѣмоудримъ соломонъмь цр̄ь вѣковъмь вьсєдръжитель г̄бъ. и понеже вр́ѣменное сеи цр̄тво прѣхȣмь ѿ крѣпкые роукы того и всєбогатіи бл̄гости. п́обаеть и п́обию того по єлико моцно рожен̄номȣ єствоуп́обити се. Ѿврьзаюцюмоу роукоу, и ѹасицаюцюмоу всако животно бл̄говолениа. іако да не тъкмо вр́ѣменномоу семоу цр̄твїю прѣемници боудемь, нъ и вѣчномоу наслѣдни ци. ѿ єже ѿ соудоу бл̄гоѵстроєнїа вр́ѣменнаго намь слоужбо наѵелїа. тѣм́ же и азь иже вь х́а бл̄гоѵьстиною вѣрою деспоть стефань. и по м́ати бж̄нїе и г̄нⷯьсрьблемь пишоу сиа вь свⷮѣнїе всѣмь. іако прїиде і́с господствоу миѵт̄нѣиши вь сц̄ен̄но и нокѵхь, и гоуме- нь цр̄каго мона стира бж̄твныи и сц̄енныи лаври, иже вьстⷡѣи горⷮѣ адⷩина. І́єромо- нахь ѵ дефⷩань. и сь ѵтⷩными старци. вь сц̄ен̄но и нокⷲ. вь григорїемь. И

The largest number of errors is attributed to the failure to recognize the ligature ꙗ: цр̄нїе 1, понеже 2, сеи 2, крѣпкые 2/3, всє богатии 3, єствоу 4, бл̄говолениа 4, прѣемници 5, еже 6, бж̄нїе 7, срьблемь 7, бж̄твныи 8, сц̄енныи 8 instead of цр̄ие 1, понеже 2, сеи 2, крѣпкые 2/3, всєбогатие3, єствоу 4, бл̄говолꙗниа 4, прѣꙗмници 5, ꙗже 6, бж̄иꙗ7, срьблꙗмь 7, бж̄тꙗныи8, сц̄енныꙗ 8. In a smaller number of instances the generic model confuses the letters ь and ъ: прѣдръжеть 1, тъкмо 5, нъ 5 instead of прѣдрьжеть 1, тькмо 5, нь 5, the letters н and и: цр̄нїе 1, бж̄нїе 7 instead of цр̄ие 1, бж̄ие 7, the letters ѵ and н: ѹасицаюцюмоу 4 instead of насицаюцюмоу 4, as well as н and в: бл̄гоѵьстиною 7 instead of бл̄гоѵьстивою 7. The model was quite successful in recognizing superscript letters, the titlo mark and blanks between words: cf. цр̄твоують 1, цр̄тво 2, цр̄твїю 5, бл̄гоѵстроєнїа 6, etc. Certain errors represent an instance of hypercorrection[12]: п́обаеть и п́обию 3, п́обити се 4 instead of п́обаеть и п́обию 3, п́обити се 4, рожен̄номȣ 3 and тѣм́ же 6 instead of рожєнномȣ 3 and тѣмже 6.
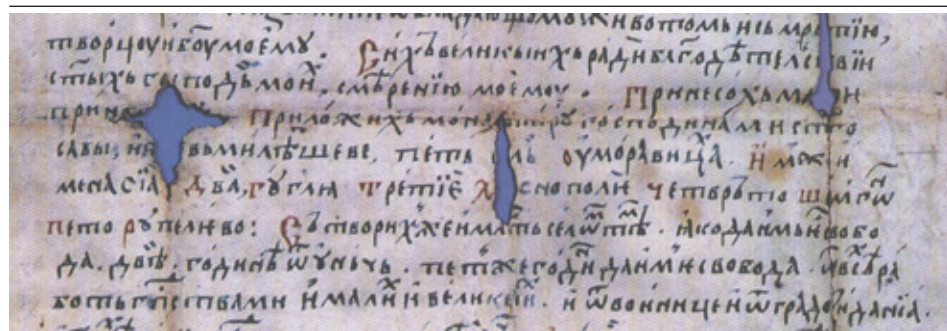
An excellent result was also recorded with the generic model used in the other two charters of Despot Stefan Lazarević: to Despotess Yevpraksia (1404–1405) (CER 9.99%) and to Mileševa Monastery (1413) (CER 11.20%). When com-

---

[12] For more information on hypercorrection as an atypical error in the process of applying the generic model, see Rabus 2019b: 12.

pared to to Despot Stefan's Charter to the Great Lavra Monastery (1414–1415), the somewhat higher CER can be explained by a larger number of errors in less legible places where the charter was folded or errors which are due to damaged parts of the text. As an illustration for this claim we can use the example from Despot Stefan's Charter to the Mileševa Monastery from the following table 11.

Table 11: Despot Stefan's Charter to the Mileševa Monastery (line 20–27)

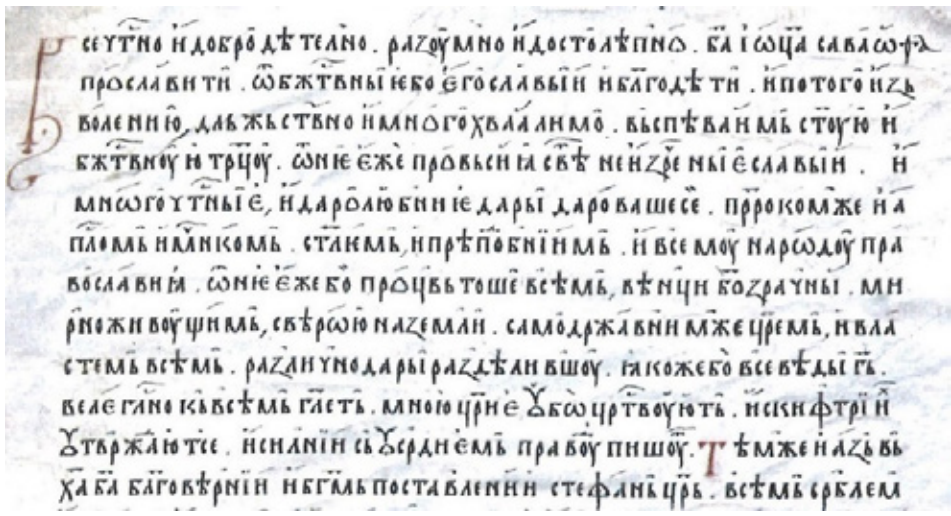| Младеновић 2007: 432 |
| --- |



| Combined_Full_VKS_2 |
| --- |

творцоу и бо҃у моемꙋ. Сихъ великыихь ради бл҃годѣтелствїи
ст҃ыхь бо подь мой, смѣренїю моемоу. Принесохь мате
при на приложихь· мона и и рꙋгосподи нами ст҃го
савы, нивьми лѣшеве. петь вль оуморавица. Имжїи
мена сїардвагꙋ гл҃а. третїе сно поле четь ро̑ тош мгѡ
петорꙋпелиево. въ творих́ же и мать селѡ тѣ. ꙗко да имь свобо-
да. двѣ̈. годинъ ѿꙋньꙋь. пет́ же годи да ими свобода. и всѣ̈ ра̑"
боть гн҃ства ми и мали и великый. и ѡ воиницѣ и ѡ градоу и данїа.

Except for the example бл҃годѣтелствїи 20, in which the model succeeded in reconstructing the letter т from a damaged part of the charter, in other places where the charter was damaged or folded errors occurred as expected: бо подь 21 instead of господь 21, мате 21 instead of малоѥ 21, мона и и рꙋгосподи 22 instead of монастирꙋ господина 22, нивьми лѣшеве 23 instead of иже вь милѣшеве 23, вль 23 instead of сель 23, Имжїи 23 instead of Им꙼жи 23, тош мгѡ 24 instead of шемгѡ̈ 24, и всѣ̈ 26 instead of ѿ всѣ̈ 26, градоу и данїа 27 instead of градозиданїа 27.

The qualitative analysis concerning the performances of the generic model applied to Emperor Dušan's Charter to the St. Sava Cell in Karyes (book number: Hil 31) (Skopje, 1348) also reveals excellent recognition results, despite CER in the first part of the charter (lines 1–35) being 12.80%, and 13.91% in the second part. To illustrate this claim, we can compare the photographic image of the first twelve lines of the charter and the automatically recognized text in the table that follows.

Table 12: Emperor Dušan's Charter to the St. Sava Cell in Karyes (lines 1–12)

| Живојиновић 2008: 59–70 |
| --- |



| Combined_Full_VKS_2 |
| --- |

се ꙋ҆тно и добродѣ҆телно. разоумно и достолѣпно. бꙋ҇ ї ѿца сава҇ ѿ фа-
прославити. ѿ бж҇твныи бо его славы и и блг҇одѣти. и по тог҇о изь-
волению, дльжьство и много҇ хваа ан҇м҇о. вь спѣваимь стꙋ҆ю и
бж҇твноую трꙋцоу. ѿ не еже провьси ꙗ҆ свѣ҇ неизре҇ны є҆ славы й҆. и
мнꙍгоꙋтны є҆ й дароюбнниє҆ дарь дароваше се· пророком же й а҆-
пломь и мнкомь. стаиемь и прѣп҇о҇бнїимь· и всемоу нарꙍдоу пра-
вослави ꙗ. ѿне є҆ же бо процвьтоше всѣмь, вѣнци бо зраꙋны. мн
рно живоꙋщимь свѣрꙍю на земли. самодржа в҇ ним҇ же цꙋремь· й влꙋа-
стемь всѣмь· разлиꙋно дарь раздѣлив҇шоу· ꙗ҆коже бо все вѣды гꙋь
велеглꙗно къ всѣмь глꙋеть. мною цꙋри є҆ овꙋбꙋ цꙋртвоꙋють. и скифтрꙋй
отврꙋжаютсе· й сианїй сꙋ о҆срдиемь правоу пишоꙋ· Тꙋѣм҇же й а҆зꙋ вь
хꙋꙋ бꙋꙋ блг҇овѣрнїи и бꙋгꙋмь поставлении стефань цꙋрь. всѣмь срꙋблем

The above representation indicates that the largest number of errors are result from the failure to recognize the pajerak mark: добродѣтелно 1, разоумно 1, достолѣпно 1, дльжьство 3, пророком же 5, прѣпо҇бнїимь 6, всѣмь 7, 10, вѣнци 7, бо зраꙋны 7, мирно 7/8, земли 8, самодржа в҇ ним҇ же 8, разлиꙋно 9, раздѣлив҇шоу 9, все 9, Тꙋѣм҇же 11, блг҇овѣрнїи 12, поставлении 12, etc. – instead of добродѣтел҇но 1, разоум҇но 1, достолѣп҇но 1, дльжьств҇но 3, пророком҇ же 5, прѣпо҇б҇нїимь 6, в҇сѣмь 7, 10, вѣн҇ци 7, бо҇зраꙋ҇ны 7, мир҇но 7/8, зем҇ли 8, самодржав҇ним҇ же 8, разлиꙋ҇но 9, раз҇дѣлив҇шоу 9, в҇се 9, Тꙋѣм҇же 11, блг҇овѣр҇нїи 12, постав҇лении12, etc. The CER level is also affected by a large number of examples in which it is necessary to remove superscript accent marks since the model often renders them hypercorrectly: й 5x2, 8, 11x2, дароюбнниє҆ 5, а҆пломь 5/6, цꙋремь 8, ꙗ҆коже 9, сианїй 11, сꙋ о҆срдиемь 11, а҆зꙋ 11, вь҆ 11, etc. The hypercorrectness is also evident in the process of

recognizing the superscript letters: сава̏ ѿ фа 1, того̏ 2, мно̏го̏ хваа ан҃мо̏ 3, тр̏цоу 4, провьси ꙗ̈ 4, неизре҆ꙝны є̈ 4, й 4, прѣпⷪбⷪнїим 6, ѽнꙁ є̈ 7, процвьтошѐ 7, цр̏и є̈ 10, ср̏блем 12 instead of савамда 1, того 2, многохваалимо 3, тр̏цоу 4, провьсиꙗ4, неизрѣные 4, и 4, прѣпⷪбⷪ'нїимь 6, ѽнкеꙗже 7, процвьтоше 7, цр̏ие 10, ср̏блем 12. Errors in the recognition of letters are mostly associated with the ligature ѥ: instead of бжтвные 2, ѽнкеꙗже 4, 7, there is бжⷮтвный 2, ѽне є̈ же 4, ѽне є̈ же 7, aslo with letter ꙋ: instead of ꙋбⷡ 10, ꙋтвр'жают' се 11, ꙋсрдиꙗемь 12 there is оⷡвбⷡ 10, отвръжаюⷮтсе 11, ꙝсрдиꙗемь 12, letter л: instead of хваалимо 3, стⷧ҆лемь 6, сиⷧ҆нїи 11 there is хваа ан҃мо̏3, станемь 6, сианїй 11, letter и: instead of хваалимо 3, мир'но 7/8 there is хваа ан҃мо̏ 3, мнрно 7/8 and letter ы: instead of даⷬы 4, 8 there is даⷬь 4, 8. Other errors: instead of савамда 1 there is сава̏ ѿ фа 1, instead of велегл҃но 10, кь 10 there is велеглⷩ҆но 10, къ 10.

2.2. Aᴜᴛᴏᴍᴀᴛɪᴄ ʀᴇᴄᴏɢɴɪᴛɪᴏɴ ᴏꜰ ᴛʜᴇ *Cʜʀɪꜱᴛɪᴀɴ Tᴏᴩᴏɢʀᴀᴩʜʏ* ᴏꜰ Cᴏꜱᴍᴀꜱ Iɴᴅɪᴄᴏᴩʟᴇᴜꜱᴛᴇꜱ (1649). The Serbian Church Slavonic translation of Cosmas Indicopleustes' *Christian Topography* (1649) is preserved as a part of a more extensive manuscript which is located in the vault of the Holy Trinity Monastery near Pljevlja (Montenegro), under the book register number Pljevlja 79. The first part of the manuscript (folios 1–101) contains the *Hexameron* by John the Exarch, while the *Christian Topography* (henceforth CT) can be found in the second part of the manuscript, to folio 240.[13] The manuscript was copied in uncial Cyrillic script in 1649 by Gavrilo Trojičanin, the most renowned calligrapher of the time, using a Russian model, in the Holy Trinity Monastery near Pljevlja, and it was ornamented by miniatures by Andrija Raičević, an icon painter and miniaturist (cf. Jᴀɢɪħ 1922: 1; Pᴀᴋɪħ 2016: 407). For the purposes of this research, we used the photographic images stored on microfilm from the Department for Archeography of the National Library of Serbia.[14] Although the quality of the photographs was not ideal, the application of the special VMČ_Test_4+ model and the Combined_Full_VKS_2 generic model rendered substantially better results than those obtained with Serbian medieval charters and letters. The experiment involved 9 folios, in which the average CER value was 6.43% for the VMČ_Test_4+ model and 5.34% for the Combined_Full_VKS_2 generic model. If CER calculations are applied only to unrecognized letters, then the results become even more impressive: CER is reduced to 2.67% with the VMČ_Test_4+ model, and to 1.42% with the Combined_Full_VKS_2 generic model. It is important to mention that the CER values do not vary substantially from one folio to another, meaning that both models yield consistent results regardless of the analysed folio. The variations in CER values depend mostly on the frequency of superscript letters in the individual folios. In other words, the higher the use of the superscript letters, the higher the CER value. What particularly needs to be mentioned are various types of errors in the process of recognizing superscript letters: a) complete omission of an superscript letter (in a large number of instances): обра instead of обра̏, сло instead of слⷪ, въпросившем instead of въпросившеⷨ, хⷬтїанскы instead of хⷬтїанскы̏, новы instead of новы̏, etc.;
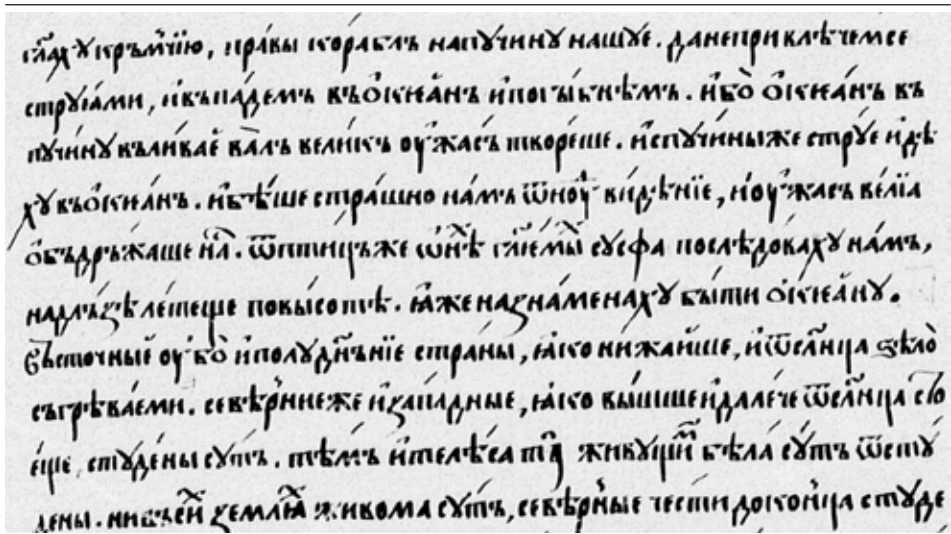
---

[13] For more details on the manuscript consult Mᴏшин 1958: 254; Cᴛᴀнᴋᴏвиħ 2003: 26. For more information about the paleographic and linguistic particularities of the manuscript, see Jᴀɢɪħ 1922. On the miniatures see Pᴀᴋɪħ 2016: 407–417.

[14] The image has been obtained with the blessing of His Grace Atanasije (Rakita), whom we hereby express our sincerest gratitude.

b) recognizing an superscript letter as the one which is lowered into the text line: e.g. ρе҃ instead of ρе҃, пѫ҃ instead of пѫ҃, etc., c) incorrect recognition of a superscript letter: e.g. слѡ҃ instead of слѡ҃, быва́ющй instead of быва́ющй, свой instead of свой, живѫ҃щй instead of живѫ҃щй, нарица́е instead of нарица́е, etc., and d) hypercorrection, that is the occurrence of a superscript letter in the positions where it would be expected, but where it was not found in the text itself: e.g. ιавлιа́е instead of ιавлιа́е, силны́ instead of силны, оутврь́жа́е instead of оутврь́жа́е, сътвори́ instead of сътвори, etc. Along with the errors related to the recognition of superscript letters, both models make errors when it comes to the recognition of blanks between words, while the errors related to the recognition of the titlo, punctuation, and regular letters occur infrequently.

The extraordinary performance of the generic model when performing the recognition in CT is evidenced by the comparative representation of the first ten lines from the folio 116r (in which the CER is only 3.78%) and the automatically recognized text in the table 13.

Table 13: CT (folio 116r, line 1–10) and Combined_Full_VKS_2



Combined_Full_VKS_2

гл҃ахѹ крьмтїю. правы корабл҃ на пѹ чинѹ наш҃е. да не привлѣчем се стрѹιами, и вь падемь вь оккань и погыбнѣмь. ибо оккань вь пѹ чинѹ вь ливае̑ валь великь оуж҃аствореше. испѹчины же трѹе̑ и дѣ хѹ вь окиань. и бѣше страшно намь ѿноу̑ видѣнїе, и оужась велïа обьдрьжаше на҃ . ѿ птиць же ѡнѣ гл҃емы сѹсфа послѣдовахѹ намь. на длъзѣлетеще̑ по высотѣ. ιаже назнаменахѹ быти окïеанѹ. вьсточные оубо и полѹднꙑнїе страны, ιако нижайше, и ѿ сл҃нца ѕѣло сьгрѣваеми. се вѣрни еже и западные, ιако вышше̑ и далече ѿ сл҃нца сто- еще. стѹдены сѹ҃ть. тѣмь и телѣ сата живѹщй бѣла сѹ҃ть ѿстѹ дены. ни вь сй землιа. живома сѹ҃ть, се вѣрные чести до конца стѹде

Judging from the above representation, it can be concluded that the generic model most frequently manifested problems with the recognition of the ligature ѥ: окканъ 2 instead of окѥанъ 2, окианъ 4 instead of окѥанъ 4, гл҃ємы 4 instead of гл҃ѥмы 4, окїеанꙋ 6 instead of окѥанꙋ 6. When errors concerning letters are considered, only one individual instance was recorded: the failure to recognize the letter с in трꙋѣ 3 instead of стрꙋѣ 3. Among other errors the most numerous ones are those connected with the recognition of blanks between words: пꙋ чинꙋ 1, въ падемъ 2, въ ливаѥ 3, оужасътвореше 3, и дѣхꙋ 3/4, се вѣрни еже 8, телѣ сата 9, ѿстꙋдены 9/10, въ си 10 instead of пꙋчинꙋ 1, въпадемъ 2, въливаѥ 3, оужасъ твореше 3, идѣхꙋ 3/4, севѣрние же 8, телѣса та 9, ѿ стꙋдены 9/10, въси 10. The extraordinarily low CER level in this folio is especially affected by the infrequency of superscript letters. The following errors were recorded: a) a superscript letter was not recognized: гл҃ємы 4 instead of гл҃ѥмы̏ 4, землꙗ 10 instead of землꙗ̏ 10, b) incorrect superscript letter: живꙋцїй 9 instead of живꙋ̏цїй 9, c) hypercorrection: въливаѥ̏ 2 instead of въливаѥ 3, трꙋѣ 3 instead of стрꙋѣ 3, длъзъвлетецїе̏ 6 instead of длъзъвлетецїе 6. Along with these examples, we need to mention the instances in which the model performed the recognition of the superscript letters correctly: ѿноꙋ̏ 4, на̏ 5, wнѣ̏ 5 and сꙁ̏ 10. The model made errors in all three examples that contained the pajerak mark in the photograph: крꙋмчїю 1, севѣрние 8, конца 10 instead of крꙋм'чїю 1, севѣр'ние 8, кон'ца 10, and in one example there was an occurrence of hypercorrected вышꙋше 8 instead of вышше 8. The titlo used for the purpose of abbreviation was mostly well recognized: гла̃хꙋ 1, гл҃ємы 5, сл҃нца 7, 8, in contrast to one example where it was omitted: полꙋднънїе 7 instead of полꙋдн҃нїе 7. Punctuation (comma and full stop) was also recognized excellently, with only a single instance of a full stop instead of a comma at the end of the fifth line.

**3.** Concluding remarks. The conducted research has confirmed the initial hypothesis that the application of the existing models for the automatic recognition of Church Slavonic Cyrillic manuscripts can also be quite successful in an overall sense when applied to Serbian medieval manuscripts written in the uncial or semiuncial script, while the application to Serbian medieval manuscripts written in cursive script renders transcripts which are not usable, thus showing that there is a need to create a special model for the recognition of the Serbian cursive Cyrillic script. Among the investigated manuscripts written in uncial or semiuncial scripts, the best results were obtained by applying the existing models to the Serbian Church Slavonic manuscript Cosmas Indicopleustes' *Christian Topography* (1649), which can probably be attributed to the fact that both models for the automatic recognition contain materials from chronologically close manuscripts written in the same script. The success of the application of the existing models to Serbian medieval charters written in the uncial or semiuncial script often varies depending on the quality of the images and the preservation state of the manuscript. While the application of the VMČ_Test+ model mostly produced unsatisfactory results, the application of the generic model resulted in quite usable transcripts. A special mention is due to the manuscript DC III, in which the generic model manifested approximately the same recognition results both in the parts of the manuscript written in the Old Serbian and the parts written in the Serbian Church Slavonic language. Although the ratio of unrecognized characters (CER) was above 10% in

many charters, the qualitative analysis has shown that the benefit of the transcripts obtained by applying the generic model can be deemed to be quite satisfactory, especially if we take into account the fact that the model was unable to gain insight into some of the special characters (such as the pajerak mark) during the training process and that the recognition errors are most frequently related to the blanks between words, superscript letters and titlos, and much more rarely to individual letters. The benefit of the generic model is especially evident when used for the automatic recognition of the more voluminous manuscripts written in uncial or semiuncial script. The transcripts of the parts of the manuscripts obtained by applying the generic model can be manually corrected and then used for subsequent training of the generic model for the purposes of enhancing its performance or for the purposes of training a special model for the recognition of the remainder of a voluminous manuscript. (e.g., BC, DC III or CT)[15]. By employing this procedure (cf. RABUS 2019b: 13), in a relatively short period of time we can obtain considerable amounts of data (photographs and corresponding transcripts of Serbian medieval manuscripts), which can be used to create special models for individual voluminous manuscripts, and ultimately a generic model for Serbian medieval manuscripts, thus significantly expediting the work on the current projects involving the Serbian historical corpus linguistics and lexicography.

REFERENCES

HODEL, Tobias, David SCHOCH, Christa SCHNEIDER, Jake PURSELL. General Models for Handwritten Text Recognition: Feasibility and State-of-the-Art. German Kurrent as an Example. *Journal of Open Humanities Data* 7: 13 (2021), 1–10. http://doi.org/10.5334/johd.46
POLOMAC, Vladimir. Towards Fundamental Principles for Creating the Electronic Corpus of Serbian Medieval Charters and Letters. *Scripta&e-Scripta* 21 (2021): 41–53.
RABUS, Achim. Recognizing Handwritten Text in Slavic Manuscripts: a Neural-Network Approach Using Transkribus. *Scripta&e-Scripta* 19 (2019a): 9–32.
RABUS, Achim, Training Generic Models for Handwritten Text Recognition Using Transkribus: Opportunities and Pitfalls. *Proceeding of the Dark Archives Conference*, Oxford, 2019b, in print.
TRANSKRIBUS GLOSSARY: https://readcoop.eu/glossary/<15.07.2021>

*

Живојиновић, Драгић. Скопска хрисовуља цара Душана за келију Св. Саве Јерусалимског у Кареји (Хил. 31). Стари српски архив 7 (2008): 59–70.
Ивић, Павле, Милица Грковић. Дечанске хрисовуље. Нови Сад – Београд: Институт за лингвистику – БИГЗ, 1976.
Јагић, Ватрослав. Козма Индикоплов по српскому рукопису г. 1649-е. Палеографско-језичка студија. Споменик СКА XLIV (1922): 1–39.
Јерковић, Вера. Полуустав у српским повељама од краја XIV века и током XV века. Зборник Матице српске за филологију и лингвистику XLII (1996): 89–113.
Корниенко, Сергей Иванович, Федор Михайлович Черепанов, Леонид Нахимович Ясницкий. Распознавание текстов рукописных и старопечатных книг на основе нейросетевых технологий. Валерий Дмитриевич Соловьев, Виктор Аркадьевич Баранов (ред.). Современные информационные технологии и письменное наследие: от древних текстов к электронным библиотекам. *El' Manuscript 08. Материалы Международной научной конференции*. Казань: Издательство Казанского государственного университета, 2008, 155–156.

---

[15] Special models trained for the recognition of each of the more voluminous manuscripts can achieve extremely high accuracy, with CER under 1% (cf. Rabus 2019b: 2–3).

Младеновић, Александар. Повеље и писма деспота Стефана. Београд: Чигоја штампа, 2007.

Мошин, Владимир. Ћирилски рукописи манастира Св. Тројице код Пљеваља. Историски записи XIV/1–2 (1958): 235–260.

Ракић, Зоран. Шестоднев Јована Егзарха и Хришћанска топографија Козме Индикоплова. Душан Оташевић, Зоран, Ракић, Ирена Шпадијер (ур.). Свет српске рукописне књиге (XII–XVII век). Београд: Галерија Српске академије наука и уметности, 2016, 407–417.

Станковић, Радоман. Рукописне књиге манастира Свете Тројице код Пљеваља: водени знаци и датирање. Београд: Народна библиотека Србије, 2003.

Стојановић, Љубомир. Старе српске повеље и писма, књ. 1. Београд: Српска краљевска академија, 1929.

Трифуновић, Ђорђе. Повеља краља Милутина манастиру Бањска. Књига прва: Фототипија изворног рукописа. Београд: ЈП Службени гласник, 2011.

Владимир Поломац
Тамара Лутовац Казновац

АУТОМАТСКО РАШЧИТАВАЊЕ СРПСКИХ СРЕДЊОВЕКОВНИХ РУКОПИСА ПОМОЋУ СОФТВЕРСКЕ ПЛАТФОРМЕ *TRANSKRIBUS*: СТАЊЕ И ПЕРСПЕКТИВЕ

Р е з и м е

Софтверска платформа *Transkribus* (https://readcoop.eu/transkribus/), недавно развијена на Универзитету у Инсбруку (Аустрија), представља алат за ручно и аутоматско рашчитавање и претраживање старих рукописа и штампаних књига, независно од времена настанка, језика или писма. Кључна предност Транскрибуса у односу на друге сродне апликације огледа се у могућности да корисник сам креира сопствени модел за аутоматско рашчитавање текста. Тренирање модела за аутоматско рашчитавање текста представља пример машинског учења заснованог на напредним неуронским мрежама у коме модел упоређује фотографије рукописа и одговарајућа слова, речи и линије текста у дипломатичком издању. За успешно тренирање модела неопходно је обезбедити што квалитетније фотографије рукописа и најмање 15000 речи рашчитаног текста. За аутоматско рашчитавање старословенских и црквенословенских ћириличких рукописа у оквиру Транскрибуса доступна су два модела која је развио немачки слависта А. Рабус: први модел, назван VMČ_Test_4+, заснован на деловима рускословенског рукописа Великие Минеи-Четьи, писаног полууставном ћирилицом XVI века; други модел, назван Combined_Full_VKS_2, заснован на деловима старословенског Супрасаљског кодекса (XI век), рукописа Катихизиса Кирила Јерусалимског (XI век) и рускословенског рукописа Великие Минеи-Четьи (XVI век), представља покушај креирања генеричког модела за аутоматско рашчитавање различитих црквенословенских рукописа писаних уставном или полу- уставном ћирилицом.

Основни циљ нашег рада представља истраживање могућности примене Рабусових модела за аутоматско рашчитавање српских средњовековних рукописа писаних различитим типовима ћирилице. Наведени модели тестирани су на српским средњовековним рукописима који су тренутно у фокусу наших филолошких и лингвистичких истраживања: на српским средњовековним повељама и писмима који се приређују за потребе изградње специјализова- ног електронског корпуса, као и на српскословенском рукопису Хришћанске топографије Козме Индикоплова (1649) који се припрема за објављивање у оригиналној графији уз пратеће филолошке и лингвистичке студије. Проведено истраживање потврдило је почетну хипотезу да примена постојећих модела за аутоматско рашчитавање црквенословенских ћириличких рукописа у начелу може бити веома успешна и на српским средњовековним рукописима писаним уставом или полууставом, док се применом на српске средњовековне рукописе писане брзописом добијају неупотребљиви транскрипти, што упућује на потребу креирања специјал- ног модела за рашчитавање српске брзописне ћирилице. Међу испитиваним рукописима писаним уставом или полууставом најбоље резултате добили смо применом постојећих модела на српско- словенски рукопис Хришћанске топографије Козме Индикоплова, што је вероватно у вези са чињеницом да оба модела за аутоматско рашчитавање садрже материјал из хронолошки блиских

рукописа писаних истим писмом. Успех примене постојећих модела на српске средњовековне повеље писане уставом или полууставом најчешће варира у зависности од квалитета фотографије и очуваности рукописа. И док је примена модела VMČ_Test+ углавном дала незадовољавајуће резултате, применом генеричког модела на свим испитиваним повељама добили смо веома употребљиве транскрипте. Иако је проценат непрепознатих карактера (CER) у многим повељама био изнад 10%, квалитативном анализом показано је да се вредност транскрипата добијених применом генеричког модела може сматрати веома задовољавајућим, посебно ако се има у виду чињеница да модел током тренинга није имао прилике да види неке специфичне карактере (нпр. пајерак) и да се грешке у рашчитавању најчешће односе на размак међу речима, надредна слова и титле, а знатно ређе и на појединачна слова. Вредност генеричког модела посебно долази до изражаја приликом аутоматског рашчитавања обимнијих рукописа писаних уставом или полууставом. Транскрипти дела рукописа добијени помоћу генеричког модела могу се ручно кориговати, а затим искористити за поновно тренирање генеричког модела у циљу побољшања његових перформанси или за тренирање специјалног модела за рашчитавање остатка обимног рукописа. Овим поступком у релативно кратком временском року можемо доћи до великих количина података (фотографија и одговарајућих транскрипата српских средњовековних рукописа) помоћу којих можемо креирати специјалне моделе за појединачне обимније рукописе, а у крајњем исходу и генерички модел за српске средњовековне рукописе, што би могло значајно убрзати рад на текућим пројектима из српске историјске корпусне лингвистике и лексикографије.

University of Kragujevac
Faculty of Philology and Arts
Department of Philology
Jovana Cvijića bb, 34000 Kragujevac, Serbia
*v.polomac@filum.kg.ac.rs*
*tamara.kaznovac@filum.kg.ac.rs*