

Serbian Early Printed Books: Towards Generic Model for Automatic Text Recognition using *Transkribus*

Vladimir Polomac*

* Serbian Language Department, Faculty of Philology and Arts, University of Kragujevac
Jovana Cvijića bb, 34 000 Kragujevac, Serbia
v.polomac@filum.kg.ac.rs

Abstract

The paper describes the process of creating and evaluating a new version of the generic model for automatic text recognition of Serbian Church Slavonic printed books within the *Transkribus* software platform, based on the principles of artificial intelligence and machine learning. The generic model *Dionisio 2.0* was created on the materials of Serbian Church Slavonic books from various printing houses of the 15th and 16th centuries (Cetinje, Venice, Goražde, Mileševa, Gračanica, Belgrade and Mrkša's Church), and, during the evaluation of its performance, it was noticed that CER was about 2–3%. The *Dionisio 2.0* model will be publicly available to all users of the *Transkribus* software platform in the near future.

1. Introduction

The research on creating a model for automatic text recognition of the Serbian Church Slavonic printed books from Venice using a software platform *Transkribus*,¹ presented in Polomac (2022), represents the starting point for this paper. This paper describes the process of transcription and creation of a specific model² for automatic text recognition of *Prayer Book (Euchologion)* printed between 1538 and 1540 in the printing house of Božidar Vuković,³ as well as the process of creating a generic model⁴ for automatic text recognition of other books printed in Venice in the printing house of Božidar Vuković and his son Vičenco.⁵ The most important result of this paper is the creation of the first version of the model *Dionisio 1.0* (named after an Italian pseudonym for Božidar Vuković – *Dionisio della Vecchia*) representing the first publicly available resource for automatic reading of Serbian Church Slavonic manuscripts and printed books within the *Transkribus* software platform (cf. <https://readcoop.eu/model/dionisio-1-0/>).

The *Dionisio 1.0* model structure is shown in Table 1, and its performance is displayed in Table 2.

Book	Word count
<i>Prayer Book (1538–1540)</i>	39,889
<i>Psalter (1519–1520)</i>	10,132
<i>Miscellany for Travellers (1536)</i>	10,618
<i>Festal Menaion (1538)</i>	10,732
<i>Miscellany for Travellers (1547)</i>	10,006
<i>Hieratikon (Liturgikon) (1554)</i>	10,196
Total	91,573

Table 1: *Dionisio 1.0*. Structure and the Amount of Training Data.

Word count	Number of epochs ⁶	CER ⁷ on Train set	CER on Validation set
86,347	100	1.66%	2.09%

Table 2: *Dionisio 1.0* Performance.

¹ *Transkribus* (<https://readcoop.eu/transkribus>) represents an open-access software platform for automatic text recognition and retrieval developed as part of the READ project at the University of Innsbruck. More details about the technological background and operating system cf. Mühlberger et al. (2019).

² 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. For more details see Mühlberger et al. (2019) and Rabus (2019a).

³ Božidar Vuković was a Serbian merchant from Zeta (Podgorica and the area surrounding Lake Skadar). After his arrival at Venice (in 1516 at the latest) he acculturated his Serbian name to the new environment by creating a Latin (*Dionisius a Vetula*) and an Italian pseudonym (*Dionisio della Vecchia*) from his Serbian name and the toponym of Starčeva Gorica (at Lake Skadar), indicating his origin (Lazić, 2018). Books from his printery were

aimed at the Serbian Orthodox Church and its flock under Ottoman rule, yet the motives of his printing business were not only patriotic and religious, but also mercantile and financial (Lazić, 2020b).

⁴ Unlike a specific model that is trained to recognize a single manuscript or printed book, a generic model contains material from different manuscripts or printed books. More details on the possibilities and pitfalls of training generic models can be found in Rabus (2019b).

⁵ After the death of Božidar Vuković, Vičenco Vuković had reprinted several of his father's editions until 1561, and later rented his equipment to other Venetian printers. For more details about his life and work see also Pešikan (1994).

⁶ The term *epoch* in machine learning stands for "one complete presentation of the data set to be learned to a learning machine" (Burlacu and Rabus, 2021).

⁷ The Character Error Rate (CER) is calculated by comparing the automatically generated text and the manually corrected version. See for more details in *Transkribus* Glossary <https://readcoop.eu/glossary/character-error-rate-cer/>.

In the continuation of the research, we aimed at examining the performance of the *Dionisio 1.0* model on Serbian Church Slavonic books created in other printing houses, firstly in Venetian printing houses created after closing Božidar and Vićenco Vuković's printing house, and then in other old Serbian printing houses of the 15th and 16th centuries (Cetinje, Gorazde, Mrkša's Church, Belgrade, Mileševa and Gračanica), thus ultimately offering a generic model for the automatic text recognition of Serbian Church Slavonic printed books as a whole.

2. Applying the *Dionisio 1.0* Model on Books from Other Venetian Printing Houses

In the first experiment, we tested the performance of the *Dionisio 1.0* model on several Serbian Church Slavonic books printed in Venice after closing Božidar and Vićenco Vuković's printing house: *Lenten Triodion* was printed in 1561 by Stefan of Scutari in the Camillo Zanetti's printing house, *Prayer Book (Miscellany for Travellers)* was printed in 1566 by Jakov of Kamena Reka, *Prayer Book (Euchologion)* was created in 1570 in the printing house of Jerolim Zagurović and *Psalter with Appendices* was printed in 1638 in the printing house of Bartol Ginammi (Pešikan, 1994). The starting hypothesis of the paper in the current experiment was that the model trained on the materials of Serbian Church Slavonic books from the printing house of Božidar and Vićenco Vuković would be useful for automatic text recognition of other Venetian editions printed using their printing equipment.

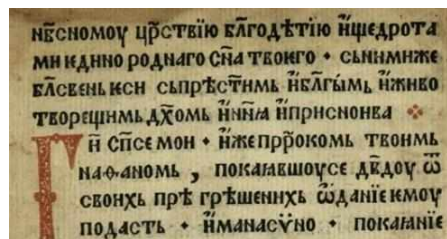
The statistical results of the experiment are shown in the following table.

Book	CER
<i>Lenten Triodion</i> (1561)	9.41%
<i>Miscellany for Travellers</i> (1566)	11.63%
<i>Prayer Book (Euchologion)</i> (1570)	13.67%
<i>Psalter with Appendices</i> (1638)	16.04%

Table 3: Application of the *Dionisio 1.0* model on publications from other Venetian printing houses.

The unexpectedly high CER does not necessarily indicate poor performance of the *Dionisio 1.0* model. The largest number of errors in text recognition is the result of the fact that in these books accent marks are used differently than in the books from the printing house of Božidar and Vićenco Vuković, which were used to train the *Dionisio 1.0* model. This fact is especially evident in *Prayer Book (Euchologion)* from the printing house of Jerolim Zagurović (1570) and *Psalter with Appendices* from the printing house of Bartol Ginammi (1638) in which only *spiritus lenis* with an *oxia* over the initial vowel grapheme was used.

To illustrate this claim, we shall use a comparative presentation of a photograph of a part of sheet 2b *Prayer Book (Euchologion)* (1570) and an automatically read text using the *Dionisio 1.0* model.

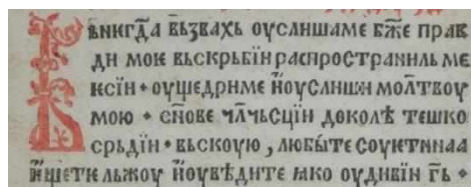


1-1 нѣсномоу црѣствію бѣгодѣтію ѿщедрота
1-2 ми ѣдїнорѣднѣго снѣ твоѣго · съ нїми же
1-3 бѣсвєнь ієси съ прѣстїмь ѿ бѣгїмь ѿ живо-
1-4 творещїмь дхѣомь ѿ нїна ѿ прїсноїва ·
1-5 Гї спїсе мой · ѿже пррѣкомь твоїмь
1-6 надѣномь, покaлѣвшоу се дѣдоу ѿ
1-7 своїхъ прѣгрѣшенїхъ ѿданїе їемоу
1-8 подѣсть · ѿ мѣнасвнѣ · покaлѣнїе

Figure 1: The Automatically Read Text of a Segment of Sheet 2b *Prayer Book (Euchologion)* from 1570.

The greatest number of errors in text recognition refers to cases in which the model outputs accent marks in accordance with the material on which it was trained, although in the text of *Prayer Book (Euchologion)* these marks were not used: so instead of *щедротами* 1/2, *твоѣго* 2, *нїми* 2, *бѣсвєнь* 3, *животворещїмь* 3/4, *дхѣомь* 4, *прїсноїва* 4, *мой* 5, *твоїмь* 5, *наданомь* 6, *своїхъ* 7, *прѣгрѣшенїхъ* 7, *їемоу* 7, *подѣсть* 8, *манасвнѣ* 8, *покaлѣнїе* 8 the model outputs *щедротами* 1/2, *твоѣго* 2, *нїми* 2, *бѣсвєнь* 3, *животворещїмь* 3/4, *дхѣомь* 4, *прїсноїва* 4, *мой* 5, *твоїмь* 5, *наданомь* 6, *своїхъ* 7, *прѣгрѣшенїхъ* 7, *їемоу* 7, *подѣсть* 8, *манасвнѣ* 8, *покaлѣнїе* 8. Along with the accent marks, the model incorrectly reads a *pajerak* mark in two examples only: instead of *ѣдїнорѣднѣго* 2, *покaлѣвшоу* 6 there is the incorrect *ѣдїнорѣднѣго* 2, *покaлѣвшоу* 6. In one example, instead of *oxia* there is an incorrect double *circumflex*: instead of *бѣгїмь* 3 there is the incorrect *бѣгїмь* 3.

The same problem is exhibited by the comparative presentation of the photograph of a part of sheet 5b *Psalter with Appendices* (1638) and the automatically read text.



1-4 ѿнегда възвѣхъ оуслїшиша ме бѣже прѣвѣ-
1-5 дї моє въ скрѣбїи распространїлї мѣ-
1-6 їєсї · оущѣдри ме ѿ оуслїшиша молѣтвоу
1-7 мою · снѣве члѣтєсїи до колѣтѣ тѣшко-
1-8 срдїи · вьскоуїю, любїте соудїтїнаа
1-9 ѿцїтїе лѣжоу ѿ оуѣдїте їако оудївїи гѣ ·

Figure 2: The Automatically Read Text of a Part of Sheet 5b *Psalter with Appendices* from 1638.

Here, too, the largest number of errors refers to cases in which the *Dionisio 1.0*. model outputs accent marks according to the patterns of their use in the Venetian books that served for its training, although in the text of Ginammi's *Psalter with Appendices* these marks were not used. Thus instead of *възвахъ* 4, *оуслиша* 4, *правди* 4/5, *моѣ* 5, *скръбїи* 5, *распространиль* 5, *ме* 5, *ієсїи* 6, *оушѣдри* 6, *оуслиши* 6, *мою* 7, *сїнове* 7, *до колѣ* 7, *тешкосръдїи* 7, *вскоуюю* 8, *любыте* 8, *соуіетннаа* 8, *льжоу* 9, *оуѣдїте* 9, *іако* 9, *оудївїи* 9 the model incorrectly outputs *възвахъ* 4, *оуслиша* 4, *прѣвди* 4/5, *моѣ* 5, *скръбїи* 5, *распространїль* 5, *мѣ* 5, *ієсїи* 6, *оушѣдри* 6, *оуслиши* 6, *моѡ* 7, *сїнове* 7, *до колѣ* 7, *тѣшкосръдїи* 7, *вскоуюѡ* 8, *любыте* 8, *соуіетннаа* 8, *льжоу* 9, *оуѣдїте* 9, *іако* 9, *оудївїи* 9. Here, as well, the other types of errors are confirmed by isolated examples: *pajerak mark*: instead of *правди* 1/2 there is the incorrect *прѣвди* 1/2; space between words: instead of *ме* 5 the incorrect *мѣ* 5; initials: instead of *Вънегда* 4 the incorrect *ѡнегда* 4; incorrect accent recognition: instead of *ї* 6 there is the incorrect *й* 6.

The given examples of the most common errors show that, despite the high percentage of incorrectly recognized characters, after the automatic post-correction of the transcripts which would include accent marks removal using the *Search/Replace chosen chars in transcript* option, the *Dionisio 1.0*. model can also be very efficient in recognizing Serbian Church Slavonic books created in the printing houses of Jerolim Zagurović and Bartol Ginammi during the 16th and 17th centuries.

The greatest number of errors in the automatic recognition of the text *Lenten Triodon* (1561) by Stefan of Scutari and *Prayer Book (Miscellany for Travellers)* (1566) by Jakov of Kamena Reka also refers to the recognition of accent marks. However, what distinguishes these books from the books from the printing houses of Jerolim Zagurović and Bartol Ginammi is that accent marks are actually used, yet in different positions compared to the books from the printing house of Božidar and Vićenco Vuković on which the *Dionisio 1.0*. model was trained. To illustrate this claim, we will first use a comparative presentation of a part of sheet 3a *Lenten Triodon* (1561) by Stefan of Scutari and the automatically read text using the *Dionisio 1.0*. model.

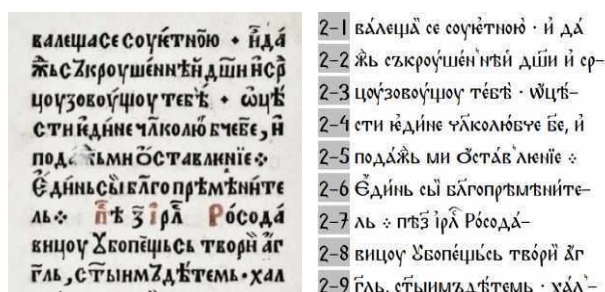


Figure 3: The Automatically Read Text of a Part of Sheet 3a *Lenten Triodon* from 1561.

Errors in accent mark recognition: instead of *валеца* 1, *соуіетноѡ* 2, *сзкроушєннѣи* 2, *срцоу* 2/3, *тєвѣ* 3, *ѡцѣ* 3, *ѡстѣвлєнїє* 5, *пєцїь сѣтворї* 8, *хал* 9/10 the model

incorrectly outputs *валеца* 1, *соуіетноѡ* 1, *сзкроушєннѣи* 2, *срцоу* 2/3, *тєвѣ* 3, *ѡцѣ* 3, *ѡстѣвлєнїє* 5, *пєцїсь твѡрї* 8, *хал* 9/10. Errors in recognizing spaces between words are also of high frequency: instead of *да* 1, *цоузоуѡцѡу* 3, *пѣ* 7, *ѡбѡпєцїсь твѡрї* 8, *а҃г* 8, *с҃тїнмздѣтемь* 9 the model incorrectly outputs *да* 1, *цоузоуѡцѡу* 3, *пѣз* 7, *ѡбѡпєцїсь твѡрї* 8, *а҃г* 8, *с҃тїнмздѣтемь* 9. In a fewer number of examples, errors in recognizing *pajerak mark*, superscript letters and *titlo mark* can be found: instead of *сзкроушєннѣи* 2, *ѡстѣвлєнїє* 5, *хал* 9/10, *ср* 2, *пѣ* 7 the model incorrectly outputs *сзкроушєннѣи* 2, *ѡстѣвлєнїє* 5, *хал* 9/10, *ср* 2, *пѣ* 7.

A comparative presentation of a part of sheet 7a *Prayer Book (Miscellany for Travellers)* from 1566 and the automatically read text using the *Dionisio 1.0*. model displays similar errors.

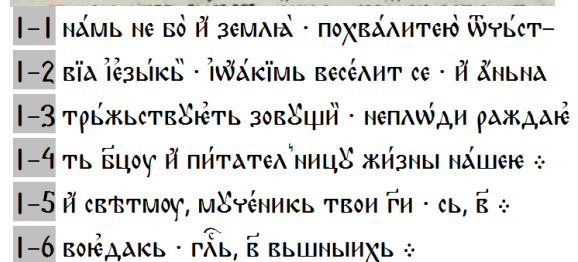
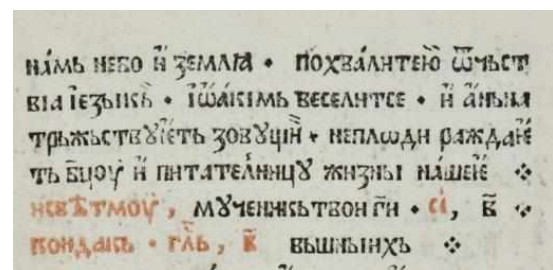


Figure 4: The Automatically Read Text of a Part of Sheet 7a *Prayer Book (Miscellany for Travellers)* from 1566.

Errors in recognizing accent: instead of *нево* 1, *земля* 1, *похвалїте* ю 1, *ѡцѣствїа* 1/2, *їєзыкѣ* 2, *весєлїт* 2, *трѣжьствѣїєть* 3, *неплѡди* 3, *раждаїє* 3, *пїтѣтелнїцѣ* 4, *жїзны* 4, *нѣшєє* 4, *и* 5, *мѣчєнїкѣ* 5, *кондакѣ* 6 the model incorrectly outputs *не бо* 1, *земля* 1, *похвалїтеѡ* 1, *ѡцѣствїа* 1/2, *їєзыкѣ* 2, *весєлїт* 2, *трѣжьствѣїєть* 3, *неплѡди* 3, *раждаїє* 3, *пїтѣтелнїцѣ* 4, *жїзны* 4, *нѣшєє* 4, *ї* 5, *мѣчєнїкѣ* 5, *воїєдакѣ* 6. A certain number of errors is connected to recognizing spaces between words: instead of *нево* 1, *похвалїте ю* 1, *раждаїє* 3, *свѣт моу* 5 the model incorrectly outputs *не бо* 1, *похвалїтеѡ* 1, *раждаїє* 3, *свѣтмоу* 5. Several errors in recognizing letters may perhaps be related to poor quality of the photograph: instead of *сї* 5, *кондакѣ* 6 the model incorrectly outputs *сь* 5, *воїєдакѣ* 6.

The illustrated examples of the most frequent errors in *Lenten Triodon* (1561) and *Prayer Book (Miscellany for Travellers)* (1566) show that the *Dionisio 1.0*. model can be used for obtaining transcripts that can, after appropriate manual correction, be used for creating specific models for automatic text recognition of the aforementioned two books.

3. Applying the *Dionisio 1.0*. Model on Books from Other Serbian Printing Houses of the 15th and 16th Centuries

In the second experiment, the performance of the *Dionisio 1.0*. model was tested on selected books from other printing houses of the 15th and 16th centuries (Cetinje, Goražde, Gračanica, Mileševa, Belgrade and Mrkša'a Church). During the research, we started from the hypothesis that the model trained on the material of books from the Venetian printing house Vuković will be useful for books from other printing houses, since there are not many orthographic variations in Serbian Early Printed Books as there are in medieval manuscripts.

The results of the experiment are shown in the following table.

Book (Printed House, Year)	CER
<i>Octoechos, mode 1–4</i> (Cetinje, 1495)	8.24%
<i>Psalter with Appendices</i> (Goražde, 1519)	6.44%
<i>Octoechos, mode 5–8</i> (Gračanica, 1539)	11.11%
<i>Prayer Book (Euchologion)</i> (Mileševa, 1546)	5.43%
<i>Tetraevangelion</i> (Belgrade, 1552)	11.28%
<i>Tetraevangelion</i> (Mrkša's Church, 1562)	12.06%

Table 4: Application of the *Dionisio 1.0*. model on publications from other printing houses in the 15th and 16th centuries.

Based on the previous table, it can be concluded that the *Dionisio 1.0*. model achieved the best results in the automatic recognition of the text of *Prayer Book (Euchologion)* (1546) from the printing house of the Mileševa monastery and *Psalter with Appendices* (1521) from the Goražde printing house. These results can be explained by the fact that *Prayer Book (Euchologion)* (1546) had been printed in Mileševa with the same typographic characters as *Psalter with Appendices* (1521) from Božidar Vuković's printing house, as well as by the fact that *Psalter with Appendices* (1519) was printed in Goražde using the typographic equipment imported from Venice (Lazić, 2020a).⁸

To illustrate the efficiency of the *Dionisio 1.0*. model we may firstly use the comparative presentation of the photograph of a part of sheet 5b *Prayer Book (Euchologion)* (1546) from the printing house of the Mileševa monastery and the automatically read text in Figure 5.

In this book, as well, the greatest number of errors refers to accent marks recognition: instead of и́мени 4, и́стинный 4, и́дино́рѡднаго 4/5, сѣ́аго 5, и́ ме 7, сподо́бльшаго 7, ѿ́ноу 10 the *Dionisio 1.0*. incorrectly outputs и́мени 4, и́стинный 4, и́дино́рѡднаго 4/5, сѣ́аго 5, и́ ме 7, сподо́бльшаго 7, ѿ́ноу 10. Other errors are fewer in number and relate to recognizing initials, spaces between words and *pajerak*

mark: instead of ѿ́ноу и́мени 4, по́дъ 9 и дрѣ́внѡю 10 the model incorrectly reads и́мени 4, по́ дъ 9 и дрѣ́внѡю 10.

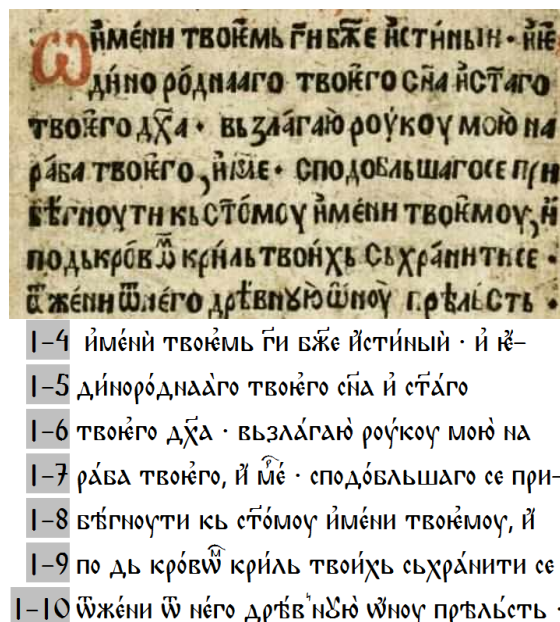


Figure 5: The Automatically Read Text of a Part of Sheet 5b *Prayer Book (Euchologion)* from 1546.

Similar errors are indicated by the comparative illustration of the photograph of a part of sheet 35a *Psalter with Appendices* (1519) from the Goražde printing house and the automatically read text using the *Dionisio 1.0*. model.

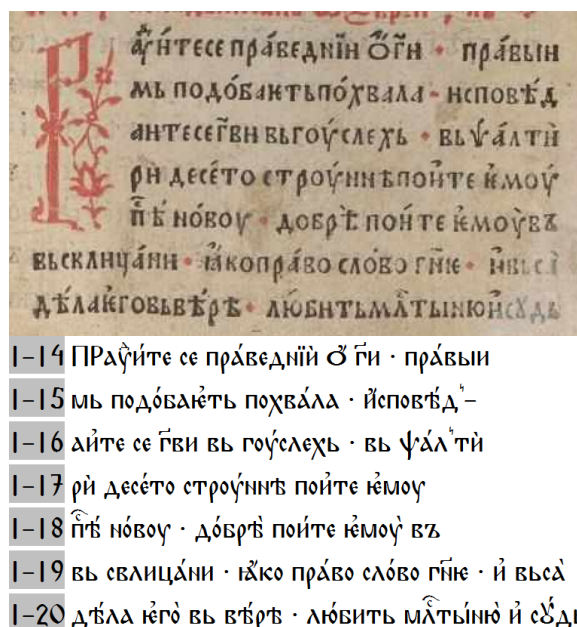


Figure 6: The Automatically Read Text of a Part of Sheet 35a *Psalter with Appendices* from 1519.

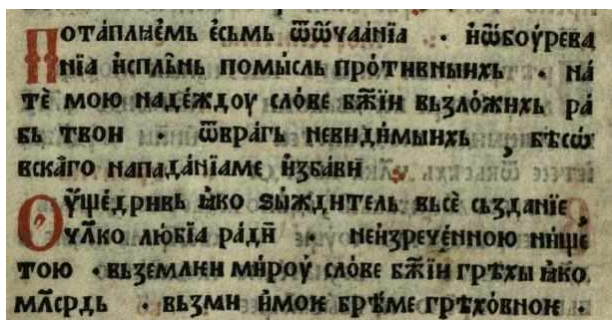
⁸ Scholars likewise claim that *Psalter with Appendices* (1519) and *Prayer Book (Euchologion)* (1544) from Goražde printing house could have been printed in Venice, as well, which

corresponds to the widespread practice of the time to place a counterfeit place of printing on the colophonies of editions (Lazić, 2020a).

The previous illustration demonstrates how the *Dionisio 1.0* model makes the most frequent errors while recognizing accent marks: instead of прѣведнѣи 14, подѡбаѣтъ 15, похвала 15, исповѣданте се 15/16, ѡуѣдѣи 16/17, ѣмоу 17, доврѣ 18, ѣго 19, млтвину 19, сѣдь 19 the model incorrectly outputs прѣведнѣи 14, подѡбаѣтъ 15, похвала 15, ѣсповѣдѣайте се 15/16, ѡуѣдѣи 16/17, ѣмоу 17, доврѣ 18, ѣго 19, млтвину 19, сѣдь 19. The other errors pertain to recognizing spaces between words, *pajerak* mark and initials: instead of прѣвыи- 14, ѡуѣдѣи- 16, десѣтостроуѣнѣ 17 the model incorrectly reads: прѣвыи 14, ѡуѣдѣи 16, десѣто строуѣнѣ 17; instead of исповѣданте се 15/16, ѡуѣдѣи- 16 there is the incorrect ѣсповѣдѣайте се 15/16, ѡуѣдѣи 16; instead of Рауѣите се 14 there is the incorrect ПРАуѣите се 14. There is merely one example of an incorrectly recognized letter: instead of въсклицани 19 the model incorrectly reads въ свлицани 19.

The *Dionisio 1.0* model also shows a similar performance during the automatic recognition of the text of the oldest printed Serbian Church Slavonic book – *Octoechos, mode 1–4* (1495) from the Cetinje printing house. The percentage of unrecognized characters is somewhat higher than in the previous two books due to poor photo quality and issues with recognizing certain letters and punctuation marks.

To illustrate the efficiency of the model, we will use a comparative presentation of a part of sheet 33b and the automatically read text in the following figure.



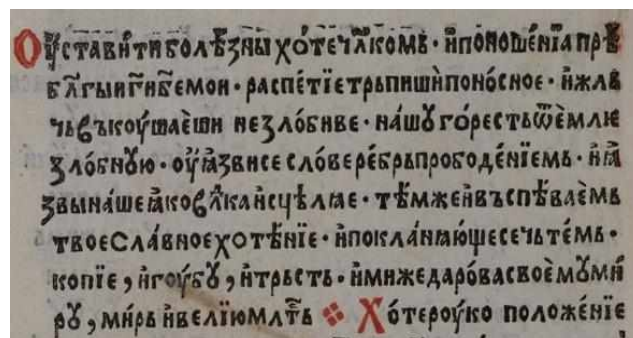
- 1-8 Потѡпаѣемъ ѣсѣмъ ѡ ѡуѣдѣи . ѣ ѡуѣдѣи
- 1-9 нѣа ѣсплѣнь помысль прѡтивныхъ . на
- 1-10 тѣ моѡ на дѣждѡу слѡве бѣжѣи възлѡжихъ рѡ-
- 1-11 въ твоѡ ѡ врагѣ невидѣмыхъ . ѡ бѣсѡ
- 1-12 вѣскаго нападѣнѣа ме ѣзбѡви
- 1-13 Оуѣдѣривъ ѣко ѡждѣитель въсѣ сѣздѣнѣе .
- 1-14 ѡуѣколюбѣа рѣдѣ . неѣзрѣченноу ницѣ
- 1-15 тоѡ . въ зѣмлѣи мѣроу слѡве бѣжѣи грѣхы ѣко
- 1-16 млсѡдѣ . възми ѣ моѣ брѣме грѣхѡвнѡе .

Figure 7: The Automatically Read Text of a Part of Sheet 33b *Octoechos, mode 1–4* from 1495.

In this book, too, the largest number of errors in the automatic text recognition occurs with accent marks: instead of ѣсѣмъ 8, ѣсплѣнь 9, на 9, моѡ 10, твоѡ 11, бѣсѡвскаго 11/12, ѣзбѡви 12, ѣко 13, сѣздѣнѣе 13,

неѣзрѣченноу 14, ницѣтоѡ 14/15, зѣмлѣи 15, ѣко 15, възми 16, ѣ 16 the *Dionisio 1.0* model incorrectly reads: ѣсѣмъ 8, ѣсплѣнь 9, на 9, моѡ 10, твоѡ 11, бѣсѡ вѣскаго 11/12, ѣзбѡви 12, ѣко 13, сѣздѣнѣе 13, неѣзрѣченноу 14, ницѣтоѡ 14/15, зѣмлѣи 15, ѣко 15, възми 16, ѣ 16. The issues with recognizing spaces between words and *pajerak* mark can be illustrated by the following examples: instead of ѡуѣдѣи- 8, на дѣждѡу 10, бѣсѡ- 11 there is the incorrect ѡуѣдѣи 8, на дѣждѡу 10, бѣсѡ 11; instead of бѣсѡвскаго 11/12 there is the incorrect бѣсѡ вѣскаго 11/12. In this book, as we have already mentioned, the *Dionisio 1.0* model likewise incorrectly recognizes certain letters and punctuation marks: instead of ѡ 8, ѡждѣитель 13, млсѡдѣ 16 there is the incorrect ѡ 8, ѡждѣитель 13, млсѡдѣ 16; instead of невидѣмыхъ, 11, ѣзбѡви :12 there is the incorrect невидѣмыхъ . ѡ 11 ѣзбѡви :12.

In the rest of the books listed in Table 4, (*Octoechos, mode 5–8* (1539) from Gračanica, *Tetraevangelion* (1552) from Belgrade and *Tetraevangelion* (1562) from Mrkša's Church), the CER is slightly higher, around 11–12%. The categories in which the *Dionisio 1.0* model outputs errors are mostly the same in all three books, so we will only take a comparative presentation of a part of sheet 27b *Octoechos, mode 5–8* (1539) from Gračanica and the automatically read text as an illustration.



- 1-1 Оуѣстѡвити бѡлѣзнь хѡте ѡкомѣ . ѣ попошѣнѣа прѣ-
- 1-2 бѣгынѣи бѣ моѡ . распѣтѣе трѣпиши понѡсноѣ . ѣ жѣв
- 1-3 ѡвъѡкуѣшѣи незлѡбѣе . нашѡ горѣсть ѡемлѣ
- 1-4 злѡбнѣю . ѡуѣзѡви се слѡве рѣбрь прѡводѣнѣемъ . ѣ ѣ-
- 1-5 зѡвы нашѣ ѣко вѣка ѣцѣлѣе . тѣмъ же ѣ въспѣѡемъ
- 1-6 твоѣ слѡвноѣхѡтѣнѣе . ѣ поклѡняѡще се ѡ тѣмъ .
- 1-7 копѣе , ѣ гоуѣдѣ , ѣ трѣсть . ѣ мѣже дарѡва своѣ мѡмѣ-
- 1-8 рѡ , мѣръ ѣ вѣлѣю млѣ . Хѡте роуѣко положѣнѣе

Figure 8: The Automatically Read Text of a Part of Sheet 27b *Octoechos, mode 5–8* from 1539.

The greatest number of errors is related to the recognition of accent marks: instead of бѡлѣзнь 1, ѣ 1, 2, 5, 6, 8, моѡ 2, трѣпиши 2, понѡсноѣ 2, ѡвъѡкуѣшѣи 3, ѡемлѣ 3, прѡводѣнѣемъ 4, ѣзѡвы 4/5, ѣко 5, ѣцѣлѣе 5, въспѣѡемъ 5, твоѣ 6, слѡвноѣхѡтѣнѣе 6, поклѡняѡще се 6, ѣмѣже 7, своѣмѡ мѣ- 7, вѣлѣю 8 the *Dionisio 1.0* model incorrectly outputs бѡлѣзнь 1, ѣ 1, 2, 5, 6, 8, моѡ 2, трѣпиши 2, понѡсноѣ 2, ѡвъѡкуѣшѣи 3, ѡемлѣ 3, прѡводѣнѣемъ 4, ѣзѡвы 4/5, ѣко

5, ѿсцѣлаѣ 5, въспѣваѣмъ 5, твоѣ 6, славноѡхотѣнїѣ 6, покланїающе се 6, ѿ мїже 7, своѣ мѣмї- 7, вѣлію 8. Recognizing spaces between words represents the problematic issue in a multitude of cases: instead of жль-2, ѡ въвкоушаѣши 3, славноѡ хотѣнїѣ 6, ѡтѣмъ 6, ѿмїже 7, своѣмѣ мї- 7, роукоположенїѣ 8 the model incorrectly outputs жлѣ 2, ѡвъвкоушаѣши 3, славноѡхотѣнїѣ 6, ѡтѣмъ 6, ѿ мїже 7, своѣ мѣмї- 7, роукоположенїѣ 8. The other errors pertain to the recognition of superscript letters and *pajerak mark*, as well as regular letters in a few examples: instead of жль-2, мѣтъ 8 the model outputs жлѣ 2, мѣтъ 8; instead of тѣмже 5 there is the incorrect тѣмъ же 5; instead of поношенїа 1, жль- 2, ѡѣмѣиѣ 3 the model reads поношенїа 1, жлѣ 2, ѡѣмѣиѣ 3.

The quantitative and qualitative analysis conducted in this chapter demonstrates that the *Dionisio 1.0* recognizes the text of the Serbian Church Slavonic books created in other printing houses of the 15th and 16th centuries with varying degrees of success. The quantitative analysis shows that the lowest CER was recorded in books from Mileševa and Goražde printing houses, which is expected considering the fact that these books were printed using the typographic printing equipment from Venice. An acceptable CER was noted during the recognition of *Octoechos, mode 1–4* (1494) from the Cetinje printing house, while this percentage exhibited in books from other printing houses (Belgrade, Gračanica, Mrkša’s Church) underscores the need for training a new version of the generic model with improved performance. The qualitative analysis showed that the *Dionisio 1.0* model usually makes errors when recognizing accent marks, but also when recognizing spaces between words. The errors in recognizing superscript letters, *pajerak mark*, initials and regular letters are far less common.

4. Creation and evaluation of the generic model *Dionisio 2.0*.

When creating a new version of the model, we started from the transcripts of Serbian Church Slavonic books listed in Table 4 obtained using the *Dionisio 1.0* model. By means of the manual correction of the transcripts, the Ground Truth⁹ data was obtained for training the generic model *Dionisio 2.0*. In accordance with our findings on the interdependence of model success and the amount of training data (Polomac, 2022), as well as similar findings for Church Slavonic books from the Berlin State Library (Neumann, 2021), the goal was set to provide a critical mass of at least 10000 words for each printed book in order to train the generic model *Dionisio 2.0*. While training the generic model *Dionisio 2.0* we used the Ground Truth data prepared for the *Dionisio 1.0* model (see Table 1 here), as well as the new Ground Truth data from Serbian Church Slavonic books printed in other printing houses of the 15th and 16th centuries listed in the following table.

⁹ The term Ground Truth Data in machine learning refers to completely accurate data used to train the model. In our case, these would be exact transcripts of digital photographs of the

Book (Printed House, Year)	Word count
<i>Octoechos, mode 1–4</i> (Cetinje, 1495)	15,667
<i>Psalter with Appendices</i> (Goražde, 1519)	16,445
<i>Octoechos, mode 5–8</i> (Gračanica, 1539)	15,179
<i>Prayer Book (Euchologion)</i> (Mileševa, 1546)	15,003
<i>Tetraevangelion</i> (Belgrade, 1552)	15,333
<i>Tetraevangelion</i> (Mrkša’s Church, 1562)	15,733

Table 5: The *Dionisio 2.0* model – Ground Truth data from other printing houses of the 15th and 16th centuries.

The performance of the generic model *Dionisio 2.0* is shown in the following table.

Word count	Number of epochs	CER on Train set	CER on Validation set
176,481	200	2.03%	2.44%

Table 6: Performance of the generic model *Dionisio 2.0*.

In order to compare the performance of the two models, we tested them on ten sheets from *Psalter with Appendices* (1495) from the Cetinje printing house and *Hieraticon* (1521) from the Goražde printing house, the latter two representing Serbian Church Slavonic books that did not form the material for training the model. The results of the experiments are shown in the following table.

Book (Printed House, Year)	<i>Dionisio 1.0</i> CER	<i>Dionisio 2.0</i> CER
<i>Psalter with Appendices</i> (Cetinje, 1495)	5.71%	1.50%
<i>Hieraticon</i> (Goražde, 1521)	9.38%	4.61%

Table 7: Comparing the Performance of the Two Models on Books from Cetinje and Goražde Printing Houses.

As can clearly be seen from the previous table, the *Dionisio 2.0* model displays significantly better results compared to the *Dionisio 1.0* model. To illustrate the exceptional efficiency of the *Dionisio 2.0* model we provide a comparative presentation of a part of sheet 3b *Psalter with Appendices* (1495) from Cetinje printing house and the automatically read text in the figure 9.

As we can see in the figure, the *Dionisio 2.0* model errors only in a few examples in which the *spiritus lenis* and *perispomena* are insufficiently clearly differentiated: instead of ѿнїи 8, поїотъ 9, истїноїѡ 10, наклѡзують 11 the model incorrectly outputs ѿнїи 8, поїотъ 9, истїноїѡ 10, наклѡзують 11. There is a single example of the model mixing *spiritus lenis* and *oxia*: instead of ѡнїѡ 13 there is the incorrect ѡнїѡ 13. The space between words was also

manuscript. For more details on this term, see Transkribus Glossary at <https://readcoop.eu/glossary/ground-truth/>.

incorrect in one example solely: instead of мнѣти 9 there is the incorrect мнѣ ти 9. In the other examples on the shown part of sheet 3b the *Dionisio 2.0* model regularly recognizes letters, spaces between words, *titlo* and accent marks. The exceptional efficiency of the *Dionisio 2.0* model in recognizing *Psalter with Appendices* (1495) from the Cetinje printing house, especially compared to *Hieraticon* (1521) from the Gorazde printing house, has resulted from the fact that there are no superscript letters in *Psalter with Appendices* (1495), while accent marks are given in expected positions.

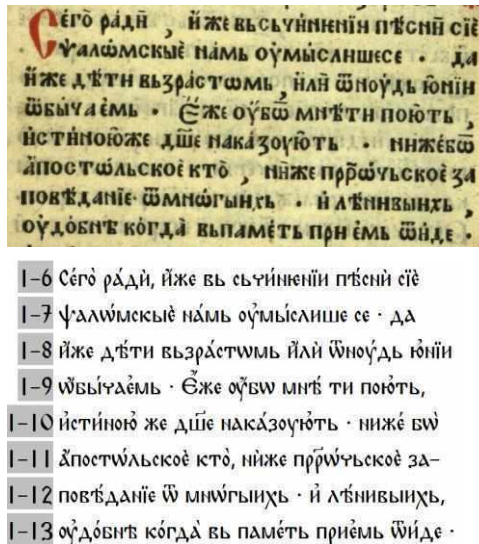


Figure 9: The Automatically Read Text of a Part of Sheet 3b *Psalter with Appendices* (1495).

On the other hand, superscript letters, as well as accent marks, found frequently in unexpected positions, are present in *Hieraticon* (1521) from the Gorazde printing house, which definitely affects a somewhat less efficient CER in this book. To illustrate the aforementioned, we shall use the comparative presentation of a part of sheet 9b and the automatically read text in the following figure.

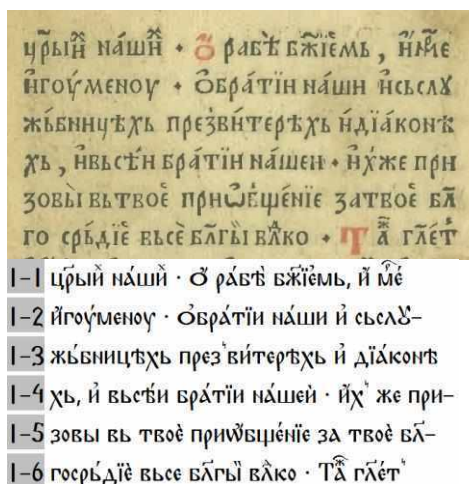


Figure 10: The Automatically Read Text of a Part of Sheet 9b *Hieraticon* (1521).

The previous illustration points to the fact that the *Dionisio 2.0* model makes errors almost exclusively during accent marks recognition. Thus, instead of рабѣ 1, бжїемъ 1, мѣ 1, игоуменоу 2, и 3, нашей 4, и хъ 4, призовы 4/5, твоѣ 5x2, приѡбщєніе 5, блгосрбдїе 5/6, все блгы 6 the model incorrectly reads рабѣ 1, бжїемъ 1, мѣ 1, игоуменоу 2, и 3, нашей 4, и хъ 4, призовы 4/5, твоѣ 5x2, приѡбщєніе 5, блгосрбдїе 5/6, все блгы 6. Along with the aforementioned errors, there are a few examples of incorrect recognition of spaces between words: instead of ѡ братїи 2, съ слѡ-2, дїаконѣ-3 все блгы 6 the model reads ѡ братїи 2, съ слѡ-2, дїаконѣ 3 все блгы 6.

5. Concluding Remarks

The research showed how the *Transkribus* software platform, based on the principles of machine learning and artificial intelligence, could be used to create efficient models for automatic text recognition of Serbian Church Slavonic printed books from the end of the 15th to the middle of the 17th century. Having in mind the limitations of the *Dionisio 1.0* model in the automatic recognition of the text of the Serbian Church Slavonic books printed outside Venice, the paper describes the process of creating a generic model *Dionisio 2.0*, capable of recognizing Serbian Church Slavonic printed books as a whole. The generic model *Dionisio 2.0* was trained on the material of the Serbian Church Slavonic books printed in various Serbian printing houses of the 15th and 16th centuries: Cetinje, Venice, Gorazde, Gračanica, Mileševa, Belgrade and Mrkša's Church. The quantitative analysis of the performance of this model showed that it could be used to automatically obtain transcripts with a minimum percentage of incorrectly recognized characters (about 2-3%). Most frequently, CER depends on the quality of the photo of the book, the frequency of use of accent marks and superscripts, as well as the correct use of accent marks in the appropriate positions. Using the *Dionisio 2.0* model transcripts of Serbian Church Slavonic printed books can be obtained automatically, which, after being edited by a competent philologist, can be used for further philological and linguistic research, primarily for creating searchable digital editions of books, as well as electronic corpora, thus creating opportunities for diachronic research of Serbian early modern literacy on a large quantity of data. In the near future, the generic model *Dionisio 2.0* will become publicly available to all users of the *Transkribus* software platform, which will enable further improvement of its performance, which could ultimately lead to the creation of a generic model for automatic text recognition of Church Slavonic printed books as a whole.

6. Acknowledgment

The research conducted in the paper was financed by the Ministry of Education, Science and Technological Development of the Republic of Serbia, contract no. 451-03-68/2022-14/ 200198, as well as by the German Academic Exchange Service (DAAD) within the project *Automatic Text Recognition of Serbian Medieval*

Manuscripts and Early Printed Books: Problems and Perspectives.

7. References

- Constața Burlacu and Achim Rabus. 2021. Digitising (Romanian) Cyrillic using Transkribus: new perspectives. *Diacronia*, 14:1–9.
- Miroslav Lazić. 2018. Od Božidara Vukovića do Dionizija dela Vekije: identitet i pseudonim u kulturi ranog modernog doba. In: Anatolij A. Turilov et al., eds., *Scala Paradisi*, pages 165–185, SANU, Beograd.
- Miroslav Lazić. 2020a. Inkunabule i paleotipi: srpskoslovenske štampane knjige od kraja 15. do sredine 17. veka. In: Vladislav Puzović and Vladan Tatalović, eds., *Osam vekova autokefalije Srpske pravoslavne crkve*, Vol. 2, pages 325–344. Sveti arhijerejski sinod Srpske pravoslavne crkve–Pravoslavni bogoslovski fakultet, Beograd.
- Miroslav Lazić. 2020b. Between an Imaginary and Historical Figure: Božidar Vuković’s Professional Identity. *Ricerche Slavistiche*, 43:141–156.
- Vladimir Neumann, 2021. Deep Mining of the Collection of Old Prints *Kirchenslavica Digital*. *Scripta & e-Scripta* 21: 207–216.
- Vladimir Polomac. 2022. Serbian Early Printed Books from Venice. Creating Models for Automatic Text Recognition using *Transkribus*. *Scripta&e-Scripta*, 22 [in print].
- Günther Mühlberger, L. Seaward, M. Terras, S. Oliveira Ares, V. Bosch, M. Bryan, S. Colluto, H. Déjean, M. Diem, S. Fiel, B. Gatos, A. Greinoecker, T. Grüning, G. Hackl, V. Haukkovaara, G. Heyer, L. Hirvonen, T. Hodel, M. Jokinen, P. Kahle, M. Kallio, F. Kaplan, F. Kleber, R. Labahn, M. Lang, S. Laube, G. Leifert, G. Louloudis, R. McNicholl, J. Meunier, J. Michael, E. Mühlbauer, N. Philipp, I. Pratikakis, J. Puigcerver Pérez, H. Putz, G. Retsinas, V. Romero, R. Sablatnig, J. Sánchez, P. Schofield, G. Sfikas, C. Sieber, N. Stamatopoulos, T. Strauss, T. Terbul, A. Toselli, B. Ulreich, M. Villegas, E. Vidal, J. Walcher, M. Wiedermann, H. Wurster, and K. Zagoris. 2019. Transforming scholarship in the archives through handwrittn text recognition. *Journal of Documentation*, 5 (75):954–976.
- Mitar Pešikan. 1994. Leksikon srpskoslovenskog štamparstva. In: Mitar Pešikan et al., eds., *Pet vekova srpskog štamparstva 1494–1994: razdoblje srpskoslovenske štampe XV–XVII*, pages 71–218, Narodna biblioteka Srbije–Matica srpska, Beograd.
- Achim Rabus. 2019a. Recognizing Handwritten Text in Slavic Manuscripts: a Neural-Network Approach using *Transkribus*. *Scripta & e-Scripta*, 19:9–32.
- Achim Rabus. 2019b. Training Generic Models for Handwritten Text Recognition using *Transkribus*: Opportunities and Pitfalls. In: *Proceeding of the Dark Archives Conference*, Oxford, 2019b, in print.