

821.163.41.09 Венцловић Стефановић Г.  
091(=163.41):004.352.243  
<https://doi.org/10.18485/kij.2024.71.2.4>

ВЛАДИМИР Р. ПОЛОМАЦ\*  
МАРКО М. МИЛОШЕВИЋ  
Универзитет у Крагујевцу  
Филолошко-уметнички факултет

Оригинални научни рад  
Примљен: 31. 7. 2024.  
Прихваћен: 24. 10. 2024.

## КА (ДИГИТАЛНОМ) ИЗДАЊУ РУКОПИСА ГАВРИЛА СТЕФАНОВИЋА ВЕНЦЛОВИЋА ПОМОЋУ ВЕШТАЧКЕ ИНТЕЛИГЕНЦИЈЕ\*\*

У раду је описан процес креирања и евалуације НТР (енгл. Handwritten Text Recognition) модела за аутоматско рашчитавање српских и српскословенских рукописа Гаврила Стефановића Венцловића помоћу софтверске платформе *Transkribus*, засноване на принципима вештачке интелигенције, машинског учења и напредних неуронских мрежа. Грађа за обуку модела узета је из два рукописа која се чувају у Архиву САНУ: *Слова изабрана*, писана српским језиком и полууставном ћирилицом; *Разглагољник*, писан српскословенским језиком и брзописном ћирилицом; као и из српскословенског *Минеја за октобар* из Архива Српске православне епархије будимске у Сентандреји, писаног полууставном ћирилицом. Најважнији резултат рада јесте креирање прве верзије генеричког НТР модела *Венцл. 1.0* са изузетним перформансама – проценат погрешно препознатих слова износи само 3%. У раду је показано како се овај модел може успешно користити и за аутоматско рашчитавање осталих Венцловићевих српских и српскословенских рукописа, што би коначно могло учинити изводљивим не само критичко издање у оригиналној графици већ и креирање дигиталних претраживих издања, а затим и електронског корпуса Венцловићеве писане заоставштине.

**Кључне речи:** српски језик, српскословенски језик, Гаврил Стефановић Венцловић, машинско учење, вештачка интелигенција, *Transkribus*, аутоматско рашчитавање текста, дигитално издање.

\* v.polomac@filum.kg.ac.rs

\*\* Рад је урађен у оквиру међународног билатералног пројекта *Креирање AI модела за аутоматску обраду српских средњовековних рукописа*, који под покровитељством Министарства за науку, технолошки развој и иновације Републике Србије и Немачке службе за академску размену (DAAD) заједнички реализују Катедра за српски језик ФИЛУМ-а у Крагујевцу и Институт за славистику Универзитета у Фрајбургу (Немачка).

## 1. Увод

Рукописна заоставштина Гаврила Стефановића Венцловића, једног од најзначајнијих српских књижевних стваралаца XVIII века, још увек нам није позната на основу целовитог критичког издања<sup>1</sup>. Овакво стање првенствено је условљено обимом његовог сачуваног рукописног фонда, који обухвата више од двадесет рукописа са око десет хиљада листова<sup>2</sup>. Основни циљ нашег прилога представља покушај да се покаже како се употребом вештачке интелигенције и машинског учења процес рашчитавања Венцловићевих рукописа може више-струко убрзати, што би коначно могло учинити изводљивим не само критичко издање у оригиналној графичкој већ и креирање дигиталних претраживих издања, а затим и електронског корпуса његових српских и српскословенских рукописа.

Полазиште за настанак овога прилога представља претходно истраживање (Поломац и др. 2023) у коме је на примеру Венцловићевог рукописа *Сло-*

<sup>1</sup> Први избор из Венцловићевих рукописа у оригиналној графичкој донео је у другој половини XIX века Г. Витковић (1872, 1887). Почетком XX века Љ. Стојановић (1901) је донео каталожки опис Венцловићевих рукописа из Архива САНУ, а С. Новаковић у својим *Примерима књижевности и језика старог и српско-словенскога* (1904: 180–202) издао у оригиналној графичкој једно Венцловићево слово из рукописа *Великопостник*. Обиман избор из Венцловићевог стваралаштва приредио је савременом ћирилицом у другој половини XX века М. Павић (1966, 1972). Интересовање за приређивање Венцловићевих рукописа поново је оживљено од друге деценије XXI века: Т. Јовановић и Д. Стефановић (2013) приредили су фототипско издање са текстом у оригиналној графичкој Венцловићевог *Сентандрејског буквара* из 1717. године, М. Стефановић Бановић (2016: 270–374) донела је у оригиналној графичкој Венцловићево *Беседе на Благовести* из рукописа *Слова изабрана*, а П. Пенкова (2021: 135–279) *Треће слово против аријанаца Атанасија Александријског* из Венцловићевог рукописа *Разглаголник*.

<sup>2</sup> Фонд Венцловићевих српскословенских рукописа нешто је обимнији од оног на српском народном језику и чине га углавном рукописи богослужбене намене: 1) *Служабник* (РГБ, Собрание Н. П. Румянцеве Ф. 256 № 401), 1711–1716. године; 2) *Сентандрејски буквар* (САНУ 1 (141)), 1717. године (Јовановић/Стефановић 2013); 3) *Минеј за октобар* (СА 70 (8)), прва четвртина XVIII века; 4) *Минеј за новембар* (СА 71 (9)), прва четвртина XVIII века; 5) *Минеј за децембар* (СА 72 (10)), прва четвртина XVIII века; 6) *Минеј за јануар* (СА 73 (11)), прва четвртина XVIII века; 7) *Минеј за фебруар* (СА 74 (12)), прва четвртина XVIII века; 8) *Минеј за март и април* (СА 75 (13)), прва четвртина XVIII века; 9) *Минеј за мај* (СА 76 (14)), прва четвртина XVIII века; 10) *Минеј за август* (СА 78 (16)), 1716–1717. године; 11) *Часови и богородичник* (САНУ 78 (132)), 1725. године; 12) *Црквени зборник* (САНУ 31 (140)), 1725/30. године; 13) *Разглаголник* (САНУ 95 (135)), 1734. године; 14) *Пресађеница* (САНУ 96 (133)), 1735. године; 15) *Молитве, акатисти и др.* (САНУ 77 (134)), 1739. године; 16) *Каноник* (САНУ 71 (138)), 1739. године. Венцловићеву писану заоставштину на српском језику чине текстови непосредно упућени православним верницима – беседе, слова, поуке, у које понекад инкорпорира друге жанрове, попут житија: 1) *Поученија и слова разлика* (САНУ 94 (271)), 1732. године; 2) *Мач духовни (прва књига)* (САНУ 92 (267)), 1733/34. године; 3) *Мач духовни (друга књига)* (САНУ 93 (268)), 1733/34. године; 4) *Великопостник* (САНУ 97 (136)), 1740/41. године; 5) *Слова изабрана* (САНУ 101 (137)), 1743. године; 6) *Пентикости* (САНУ 98 (272)), 1743. године; 7) *Житија, слова и поуке* (САНУ 84 (270)), 1744/45. године; 8) *Поученије изабрано (прва књига)* (САНУ 99 (139)), 1745. године; 9) *Поученије изабрано (друга књига)* (САНУ 100 (269)), 1746. године. Поред сведочанстава о Венцловићевим рукописима који су нам остали непознати (нпр. трећа књига Барановичевих проповеди *Мач духовни* (Павић 1972: 98), овде треба поменути и претпоставку о Венцловићу као преписачу *Богородичника* РР128 (датиран између 1710. и 1720. године) из Библиотеке Матице српске у Новом Саду, која се износи у Грбић и др. 1999: 115–120.

ва изабрана (САНУ 101 (137))<sup>3</sup> показано како се помоћу софтверске платформе *Transkribus* (засноване на принципима вештачке интелигенције, машинског учења и напредних неуронских мрежа)<sup>4</sup> може вишеструко убрзати рад на рашчитавању грађе и креирању електронског корпуса за потребе израде историјског речника српског језика. Најважнији практични резултат рада Поломац и др. 2023 представља креирање првог НТР (енгл. Handwritten Text Recognition) модела<sup>5</sup> за аутоматско рашчитавање Венцловићевих рукописа на српском језику. Помоћу овога НТР модела аутоматски се могу добити прве верзије рашчитаних текстова Венцловићевих рукописа на српском језику са веома ниским процентом погрешно препознатих карактера (око 4,5–6%)<sup>6</sup>, који варира у зависности од рукописа (детаљније у Поломац и др. 2023: 304, 307–308). Квалитативне перформансе овога модела биле су лошије од квантитативних будући да *Transkribus* није увек био у стању да препозна наредна слова која су у складу са начелима преношења грађе за историјски речник српског језика била означена као текстуални таг у суперскрипту (Поломац и др. 2023: 311). Додатни проблем представљала је чињеница да је процес обуке модела био заснован на фонту *BeogradPro* који није у складу са Unicode стандардом. У ослонцу на ове резултате, у наставку истраживања желели смо да креирамо велики генерички НТР модел са унапређеним квантитативним и квалитативним перформансама способан за аутоматско рашчитавање Венцловићевих рукописа не само на српском већ и на српкословенском језику. Овај модел, заснован на фонту *Bukyvede* који је усаглашен са Unicode стандардом<sup>7</sup>, представљао би кључни алат за вишеструко убрзавање рада на рашчитавању Венцловићевих рукописа и припреми њиховог (дигиталног) издања и електронског корпуса.

<sup>3</sup> Детаљније о садржају рукописа в. у Стојановић 1901: 155–171.

<sup>4</sup> Више о софтверској платформи *Transkribus* в. у Милбергер и др. 2019. За досадашње резултате у примени ове платформе на словенске и српске средњовековне рукописе и старе штампане књиге в. Рабус 2019, Бураку/Рабус 2021, Поломац/Лутовац Казноац 2021, Поломац 2022а, 2022б, 2023, Рабус/Томпсон 2023.

<sup>5</sup> Софтверска платформа *Transkribus* свим корисницима омогућава обуку сопственог НТР модела за аутоматско рашчитавање текста, независно од времена настанка, језика или писма. Обука НТР модела 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. Минимална количина података неопходна за успешно тренирање модела за старе штампане књиге износи око 5 000 речи, а за рукописе око 15 000 речи.

<sup>6</sup> Процент погрешно препознатих карактера (енгл. Character Error Rate, скраћено CER) представља основни квантитативни показатељ успешности НТР модела и добија се поређењем аутоматски рашчитаног и ручно коригованог текста. НТР модел може се сматрати изузетно успешним уколико је CER испод 5%.

<sup>7</sup> В. <https://kodeks.uni-bamberg.de/aks/schrift/bukyvede.htm>.

## 2. Креирање и евалуација генеричког НТН модела за Венцловићеве рукописе

Процес креирања генеричког НТН модела идеално је осмишљен за потребе аутоматског рашчитавања Венцловићевих српских и српскословенских рукописа писаних полууставном и брзописном ћирилицом. Највећи методолошки проблем за остварење овога циља представљало је одсуство јавно доступних дигиталних снимака Венцловићевих рукописа, као и ограничена финансијска средства која су нам била на располагању за потребе њихове дигитализације. Грађа за обуку модела узета је из два рукописа која се чувају у Архиву САНУ: *Слова изабрана* (САНУ 101 (137)), писана српским језиком и полууставном ћирилицом<sup>8</sup>; *Разглагољник* (САНУ 95 (135)), писан српскословенским језиком и брзописном ћирилицом<sup>9</sup>; као и из српскословенског *Минеја за октобар* (СА 70 (8)) из Архива Српске православне епархије будимске у Сентандреји<sup>10</sup>, писаног полууставном ћирилицом.

Креирању генеричког НТН модела претходило је креирање посебних НТН модела за сва три наведена Венцловићева рукописа. У овом процесу примењен је степенати приступ, који је посебно погодан за рашчитавање великих рукописа, а који се састоји из неколико фаза: најпре, ручно рашчитавање мањег дела рукописа (првих петнаестак хиљада речи) и обука прве верзије модела, затим, аутоматско рашчитавање следећег дела рукописа (наредних петнаестак хиљада речи), ручна корекција аутоматски рашчитаног текста и обука нове верзије модела са побољшаним перформансама. Други корак је поновљен све док квантитативне и квалитативне перформансе модела нису постале изузетне. Поступак се може илустровати на примеру рукописа *Слова изабрана* (САНУ 137):

Табела 1: Креирање НТН модела за рукопис *Слова изабрана* (САНУ 137)

НТН модел	Број страна за обуку	Број речи за обуку	CER
<i>Венцл. САНУ 137.1</i>	67	14 102	5,9%
<i>Венцл. САНУ 137.2</i>	135	29 065	5%
<i>Венцл. САНУ 137.3</i>	203	44 145	4%
<i>Венцл. САНУ 137.4</i>	270	59 091	3,51%
<i>Венцл. САНУ 137.5</i>	437	94 926	3%

<sup>8</sup> Располагали смо само дигиталним снимцима *Беседа на Божић* из овога рукописа (укупно 475 страна).

<sup>9</sup> Из овога рукописа били су нам на располагању само дигитални снимци *Беседа на Шестоднев* (укупно 144 страна).

<sup>10</sup> Располагали смо снимцима целог рукописа (укупно 440 страна) захваљујући љубазности Његовог Преосвештенства Епископа будимског, г. Лукијана.

Прва верзија модела *Венцл. САНУ 137.1* са одличним квантитативним показатељима (процент погрешно препознатих карактера износио је 5,9%) обучена је на ручно рашчитаној грађи првих 67 страна рукописа (14 102 речи). Помоћу овога модела аутоматски је рашчитано, а затим и ручно кориговано наредних 68 страна рукописа, како би се добила додатна грађа за обуку друге верзије модела *Венцл. САНУ 137.2* са побољшаним перформансама: процент погрешно препознатих карактера пада на 5%. Овај поступак поновљен је неколико пута, тако да је пета верзија модела *Венцл. САНУ 137.5* обучена на 437 страница текста (скоро 95 хиљада речи) са свега 3% погрешно препознатих карактера. Табела 1 потврђује претходна наша истраживања о томе да више грађе за обуку модела директно утиче на његову успешност (на мањи процент погрешно препознатих слова) (Поломац 2022а: 16). С друге стране, на основу ове табеле, али и претходних наших искустава у обуци модела<sup>11</sup>, може се закључити како даље повећање количине грађе за обуку не би нужно водило и даљем значајнијем унапређењу перформанси модела<sup>12</sup>.

Степенасти приступ приликом припреме грађе за креирање модела за аутоматско рашчитавање рукописа *Слова изабрана* (САНУ 137) примењен је и на Венцловићевим српскословенским рукописима. У табели 2 приказане су перформансе модела за аутоматско рашчитавање рукописа *Разглаголник* (САНУ 135), а у табели 3 за *Минеј за октобар* (СА 8).

Табела 2: Креирање НТР модела за рукопис *Разглаголник* (САНУ 135)

НТР модел	Број страна за обуку	Број речи за обуку	CER
<i>Венцл. САНУ 135.1</i>	54	15 020	5 %
<i>Венцл. САНУ 135.2</i>	137	38 669	4%

Табела 3: Креирање НТР модела за рукопис *Минеј за октобар* (СА 8)

НТР модел	Број страна за обуку	Број речи за обуку	CER
<i>Венцл. СА 8.1</i>	110	23 883	4,8%
<i>Венцл. СА 8.2</i>	418	90 909	4%

Оба модела креирана су тако што је најпре ручно рашчитана мања количина рукописа (око 15 хиљада речи за *Разглаголник*, око 24 хиљаде речи за *Минеј за октобар*), на основу које су креиране прве верзије модела (*Венцл. САНУ 135.1* и

<sup>11</sup> В. модел *Dionisio 2.0* који је обучен на око 180 хиљада речи са процентом погрешно препознатих слова од 2,44% (Поломац 2022б: 159).

<sup>12</sup> Максималне перформансе софтвера *PyLaia* у оквиру платформе *Transkribus* у обуци модела за словенске средњовековне рукописе износе око 2–3% погрешно препознатих слова.

Венцл. СА 8.1). Помоћу ових модела аутоматски је рашчитан остатак текста, који је након ручне корекције искоришћен за обуку друге верзије са унапређеним перформансама. Поређење модела Венцл. САНУ 135.2 и Венцл. СА 8.2 показује како су исте перформансе (CER од 4%) постигнуте са различитом количином грађе за обуку. У случају *Минеја за октобар* за обуку модела је искоришћен цео рукопис, док је у случају *Разглаголника* могуће претпоставити да би још грађе за обуку додатно унапредило перформансе модела<sup>13</sup>.

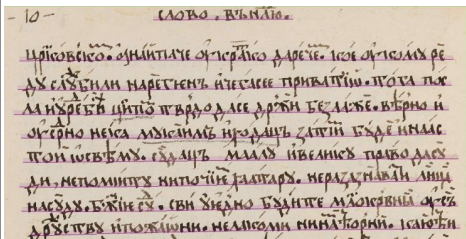
Након (полу)аутоматског рашчитавања наведена три Венцловићева рукописа могли смо да на основу њих обучимо прву верзију генеричког НТР модела под називом Венцл. 1.0. Параметри и перформансе овога модела приказане су у следећој табели:

Табела 4: Параметри и перформансе генеричког НТР модела Венцл. 1.0

Модел	Број страна за обуку	Број речи за обуку	Број епоха <sup>14</sup>	CER на скупу за обуку	CER на скупу за проверу
Венцл. 1.0	973	215 314	130	2%	3%

Права представа о успешности овога модела може се добити тек квалитативном анализом, односно поређењем фотографија рукописа и аутоматски рашчитаног текста. У табели 5 приказан је део листа 10а рукописа *Слова изабрана* и аутоматски рашчитан текст помоћу модела Венцл. 1.0.

Табела 5: Квалитативна анализа модела Венцл. 1.0 на рукопису *Слова изабрана*

Слова изабрана (САНУ 137), део листа 10а	Аутоматски рашчитан текст, Венцл. 1.0
	<ul style="list-style-type: none"> <li>I-1 слово. въ илю.</li> <li>I-2 цифров'скѣ. а наплате оу крако да рече. ко е оу комѣ ре-</li> <li>I-3 дѣ слѣбни нарегиъ и чега се е приватив. тога пос-</li> <li>I-4 ла и дрѣвъ цифѣ твр'де да се др'жи безлаже. вѣрно и</li> <li>I-5 оу ср'но нека мѣксаниъ и ходоцѣ за тѣи бѣде и наас-</li> <li>I-6 тон и свѣмѣ. сѣдацѣ маалѣ и великѣ право да сѣ</li> <li>I-7 ди. непоминѣ ни поучѣ хаатарѣ. неразазнавати лица</li> <li>I-8 на сѣдѣ. вѣже сѣ. сви бѣд'но бѣдите мѣокр'вны оуез-</li> <li>I-9 дрѣствѣ и пожатни. не лакоми ни нѣторни. каюти</li> </ul>

Наведени узорак показује да модел Венцл. 1.0 најчешће греша приликом препознавања размака међу речима: уместо оу<sup>к</sup>рѣко 2, слѣ<sup>б</sup>ни ли 3, вѣз лаже 4, мѣкс

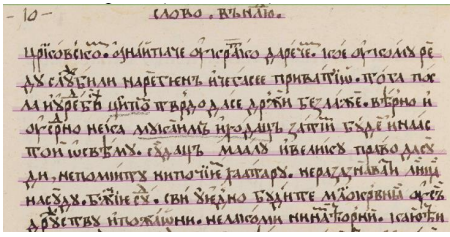
<sup>13</sup> Нажалост, нисмо располагали финансијским средствима за прибављање више снимака овога рукописа, тако да ову тезу нисмо успели да проверимо.

<sup>14</sup> Под термином *epocha* у машинском учењу подразумева „one complete presentation of the data set to be learned to a learning machine” (Бурлаку/Рабуц 2021: 1).

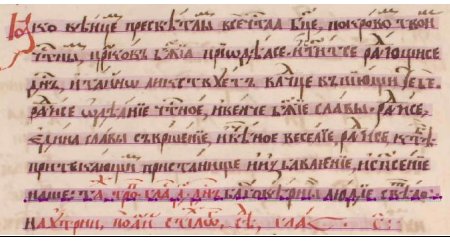
нимъ 5, сѡ- 6, не по митѡ 7, не разазнавати 7, ни по чїѡ 7, оу сѡ- 8 модел погрешно доноси оу крако 2, сѡбѡли 3, безлаже 4, мѡклимъ 5, сѡ 6, непомитѡ 7, неразазнавати 7, ни по чїѡ 7, оу сѡ- 8. Грешке у другим категоријама долазе само изоловано: размак међу речима и надредно слово представљају проблем у примеру оу ср'но 5 (треба оуср'но 5); погрешно донет пајерак у примеру наас'- 5 (треба наас- 5); погрешно донето надредно слово у примеру ѡѡ 8 (треба сѡ 8).

На сличне закључке упуђује и квалитативна анализа аутоматског рашчитавња Венцловићевих српскословенских рукописа. У табелама 6 и 7 дат је приказ дела листа 5а рукописа *Разглагольник*, односно дела листа 10б *Минеја за октобар*, и аутоматски рашчитаног текста помоћу модела *Венцл. 1.0*.

Табела 6: Квалитативна анализа модела *Венцл. 1.0* на рукопису *Разглагольник*

<i>Разглагольник</i> , део листа 5а	Аутоматски рашчитан текст, <i>Венцл. 1.0</i>
	<ul style="list-style-type: none"> <li>1-1 слово, въ нѡю.</li> <li>1-2 црков'скѡ, а наипаче оу крако да реѡе, ко е оу комѡ ре-</li> <li>1-3 дѡ сѡбѡли нареченъ и чегѡ се е приватїи, тога пос-</li> <li>1-4 ла и зрѡеѡ цїиѡ твердо да се држи безлаже, вѡрно и</li> <li>1-5 оу ср'но нека мѡклимъ и ходоцъ за тїѡ бѡде и наас'-</li> <li>1-6 тон в свѡмѡ, сѡдацъ маалѡ и великѡ право да сѡ</li> <li>1-7 ди, непомитѡ ни по чїѡ хадтарѡ, неразазнавати лица</li> <li>1-8 на сѡдѡ, бѡже сѡ, сви ѡед'но бѡдите мѡкр'вни оу сѡ-</li> <li>1-9 дрѡеѡѡ и пожажни, не лакомн ни нѡфор'ни, каюти</li> </ul>

Табела 7: Квалитативна анализа модела *Венцл. 1.0* на рукопису *Минеја за октобар*

<i>Минеја за октобар</i> , део листа 10б	Аутоматски рашчитан текст, <i>Венцл. 1.0</i>
	<ul style="list-style-type: none"> <li>1-1 Іако вѡице пресѡвѡѡи всеѡгда бїѡе, покровѡ твоѡ</li> <li>1-2 ѡтчи, црковъ вѡжа прїводѡ се, и тїт се рѡвнѡци се</li> <li>1-3 дїѡ, и танїѡ ликетвѡеѡѡ вѡлѡе вѡпѡци теѡе.</li> <li>1-4 рѡви се ѡдѡѡанїѡ ѡтноѡ, и венѡе бѡже славы, рѡви се,</li> <li>1-5 єдина славы єѡвршенїѡ, и вѡвноѡ веселѡе, рѡви се, к тѡѡ</li> <li>1-6 прїтекаюци прїстанїѡци и изѡѡѡленїѡ, и єпїенїѡ</li> <li>1-7 наше: тѡ, трѡ, глѡд, ѡ: дїѡ бѡгов'вѡрнї ѡлѡдїє єѡѡло :</li> <li>1-8 на ѡтрни, по дї єтїѡѡ, сѡ, глас: Г :</li> </ul>

На основу узорка из табеле 6 може се закључити како се грешке у аутоматском рашчитавњу *Разглагольника* најчешће односе на размак међу речима: уместо невїма вѡ 1, дами 2, по водѡ 2, вѡсеѡѡно 3 модел *Венцл. 1.0* погрешно доноси не вїмаѡѡ 1, да ми 2, поводѡ 2, вѡсе кѡно 3. Остале малобројне грешке односе се на писање надредног слова и титле: уместо несѡна 6, оѡже 8 модел погрешно доноси несѡна 6, оѡже 8. Грешком смо сматрали и писање пајерка у двама примерима:

кон'ци 7, њстаџ'никъ 8, изнад којих се на снимку види надредни знак који обликом подсећа на тачку или акценатски знак.

У узорку текста из *Минеја за октобар*, који је наведен у табели 7, најчешће грешке односе се на рашчитавање надредног слова: уместо т'џит 2, т'џинџ 3, пристанице 6 модел погрешно доноси: т'џит 2, т'џинџ 3, пристанице 6. Пајерак је изостављен у примеру с'звр'шене 5 (треба с'звр'шене 5), а титла је погрешно донета у примеру д'нџ 3 (треба д'нџ 3). У једном примеру погрешно је рашчитано слово љ: уместо веселје 5 модел погрешно доноси веселџе 5.

### 3. ПРИМЕНА МОДЕЛА *ВЕНЦЛ. 1.0* НА ДРУГЕ ВЕНЦЛОВИЋЕВЕ РУКОПИСЕ

Претходно истраживање (Поломац и др. 2023: 307–310) показало је да се помоћу модела обученог на материјалу Венцловићевог рукописа *Слова изабрана* аутоматски могу добити рашчитани текстови и других Венцловићеких рукописа на српском језику са одличним квантитативним и квалитативним перформансама. У наставку истраживања интересовале су нас квантитативне перформансе генеричког НТР модела *Венцл. 1.0* на другим Венцловићевим рукописима писаним српским и српскословенским језиком. За ове потребе креиран је експеримент у коме су помоћу модела *Венцл. 1.0* аутоматски рашчитани мањи узорци (првих десет страна) српских рукописа *Великопостник* (САНУ 136) и *Поученије изабрано (прва књига)* (САНУ 139), као и српскословенских рукописа *Пресађеница* (САНУ 133) и *Минеј за новембар* (СА 9). Након ручне корекције аутоматски рашчитаних текстова, израчунат је проценат погрешно препознатих карактера (CER) за све рукописе, приказан у следећој табели:

Табела 8: Примена модела *Венцл. 1.0* на друге рукописе

Рукопис	CER
<i>Великопостник</i> (САНУ 136)	3,13%
<i>Поученије изабрано</i> (САНУ 139)	2,66%
<i>Пресађеница</i> (САНУ 133)	4,29%
<i>Минеј за новембар</i> (СА 9)	3,80%

Наведена табела показује како се НТР модел *Венцл. 1.0* може ефикасно користити и за аутоматско рашчитавање других Венцловићеких српских и српскословенских рукописа. Изузетне перформансе модела на рукопису *Поученије изабрано* (САНУ 139) чак су и незнатно боље него на рукописима на којима је модел обучен, док су на осталим рукописима сличне или тек незнатно слабије од процента погрешно препознатих карактера на скупу података за проверу мо-



дела. Овакви квантитативни показатељи очекивани су с обзиром на чињеницу да су наведени рукописи по језику и писму веома слични рукописима на којима је модел обучен<sup>15</sup>. Квалитативна анализа перформанси модела *Венцл. 1.0* на рукописима из претходне табеле није указала на битнија одступања у односу на стање наведено уз табеле 6 и 7. И на овим рукописима модел најчешће греша приликом препознавања размака међу речима, ређе приликом препознавања надредних слова, титле и пајерка.

#### 4. Закључне напомене

Истраживање проведено у раду показало је како се помоћу софтверске платформе *Transkribus*, засноване на принципима вештачке интелигенције, машинског учења и напредних неуронских мрежа, може креирати генерички НТР (енгл. Handwritten Text Recognition) модел са изузетним показатељима (процент погрешно препознатих карактера износи свега 3%) у аутоматском рашчитавању обимне писане заоставштине Гаврила Стефановића Венцловића. Овај генерички НТР модел, назван *Венцл. 1.0*, а обучен на грађи од 973 стране и око 215 хиљада речи коју чине део рукописа *Слова изабрана* (српски језик, полуустав), део рукописа *Разглагољник* (српскословенски језик, брзопис) и цео рукопис *Минеј за октобар* (српскословенски језик, полуустав), може се ефикасно употребити и за аутоматско рашчитавање других Венцловићевих српских и српскословенских рукописа. Квалитативна анализа перформанси овога НТР модела показала је да се најчешће грешке односе на препознавање размака међу речима, ређе на препознавање надредних слова, титле и пајерка, а само изузетно на препознавање појединачних слова у реду. Помоћу НТР модела *Венцл. 1.0* процес рашчитавања Венцловићевих рукописа може се десетоструко убрзати, чиме се коначно стварају претпоставке за припрему целовитог критичког издања у оригиналној графичкој његове писане заоставштине не само у традиционалном штампаном већ и у дигиталном формату. Предуслов за реализацију овога прворазредног националног и културног задатка јесте завршетак дигитализације Венцловићевих српских и српскословенских рукописа у Архиву САНУ.

<sup>15</sup>Интересантно би било истражити перформансе овога модела на Венцловићевим српским рукописима писаним брзописном ћирилицом (нпр. *Мач духовни*) будући да они нису били обухваћени обуком модела. Претпоставка о нешто лошијим перформансама у односу на рукописе из табеле 8 није могла бити проверена услед одсуства дигиталних снимака рукописа.

## ЛИТЕРАТУРА

**Бурлаку/Рабус 2021:** С. Burlacu, A. Rabus, Digitising (Romanian) Cyrillic using Transkribus: new perspectives, *Diacronia*, 14, 1–9.

**Витковић 1872:** Г. Витковић, О књижевном раду јеромонаха Гаврила Стефановића, *Гласник Српског ученог друштва*, 24, 151–177.

**Витковић 1887:** Г. Витковић, *Прошлост, установа и споменици угарских краљевих шајкаша: од 1000 до 1872*, У Београду: У Штампарији Краљевине Србије.

**Грбић и др. 1999:** Д. Грбић, К. Шкорић, М. Гроздановић-Пајић, *Тирилске рукописне књиге Библиотеке Матице српске. Књ. 7: Акатисти, Стихологије, Богородичници*, Нови Сад: Библиотека Матице српске.

**Јовановић/Стефановић 2013:** Т. Јовановић, Д. Стефановић, *Венцловићев Сентандрејски буквар: 1717*, Будимпешта–Београд: Радионица „Венцловић”–Арте.

**Милбергер и др. 2019:** G. 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, J. Pratikakis, J. Puigcerver Pérez, H. Putz, G. Retsinas, V. Romero, R. Sabltnig, 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, K. Zagoris, Transforming scholarship in the archives through handwritten text recognition, *Journal of Documentation*, 5/75, 954–976.

**Новаковић 1904:** С. Новаковић, *Примери књижевности и језика старога и српско-словенскога*, Београд: Краљевска-српска државна штампарија.

**Павић 1966:** М. Павић, *Црни биво у срцу*, Београд: Просвета.

**Павић 1972:** М. Павић, *Гаврил Стефановић Венцловић*, Београд: Српска књижевна задруга.

**Пенкова 2021:** П. Пенкова, *Буквалният превод на Трето слово против арианите от Атанасий Александрийски по ръкописа на Гаврило Стефанович Венцлович в зборника Разглаголник 1732–1734. г.: Изследване и издание на текста*, София: Издателство „Валентин Траянов”.

**Поломац/Лутовац Казновац 2021:** V. Polomac, T. Lutovac Kaznovac, Automatic Recognition of Serbian Medieval Manuscripts by Applying the Transkribus Software Platform: Current State and Future Perspectives, *Зборник Матице српске за филологију и лингвистику*, LXIV/2, 7–26.

**Поломац 2022а:** V. Polomac, Serbian Early Printed Books from Venice: Creating Models for Automatic Text Recognition using Transkribus, *Scripta&e-Scripta*, 22, 11–29.

**Поломац 2022б:** V. Polomac, Serbian Early Printed Books: Towards Generic Model for Automatic Text Recognition using Transkribus, in: D. Fišer, T. Erjavec (eds.), *Proceedings of the Conference on Language Technologies and Digital Humanities*, Ljubljana: Institute for Contemporary History, 154–161.

**Поломац 2023:** V. Polomac, Macarius: A HTR model for Romanian Slavonic Early Printed Books, *Slavistica Vilnensis*, 68/2, 10–23.

**Поломац и др. 2023:** V. Polomac, M. Kurešević, I. Bjelaković, A. Colić Jovanović, S. Petrović, Digitizing Cyrillic Manuscripts for the Historical Dictionary of the Serbian language using Handwritten Text Recognition Technology, *Slověne*, 12/1, 295–316.

**Рабус 2019:** A. Rabus, Recognizing Handwritten Text in Slavic Manuscripts: a Neural-Network Approach using Transkribus, *Scripta&e-Scripta*, 19, 9–32.

**Рабус/Томпсон 2023:** A. Rabus, W. Thompson, Performance of Generic HTR Models on Historical Cyrillic and Glagolitic: Comparison of Engines, *Scripta&e-Scripta*, 23, 11–34.

**Стефановић Бановић 2015:** М. Стефановић Бановић, *Беседе, слова и поуке на Благовести у преводу и преради Гаврила Стефановића Венцловића*. Докторска дисертација у рукопису, Београд: Филолошки факултет.

**Стојановић 1901:** Љ. Стојановић, *Каталог рукописа и старих штампаних књига Српске краљевске академије*, Београд: Српска краљевска академија.

---

Vladimir R. Polomac  
Marko M. Milošević

## TOWARDS A (DIGITAL) EDITION OF GAVRILO STEFANOVIĆ VENCLOVIĆ'S MANUSCRIPTS USING ARTIFICIAL INTELLIGENCE

### Summary

The paper describes the process of creating and evaluating an HTR (Handwritten Text Recognition) model for the automatic text recognition of Serbian and Serbian Church Slavonic manuscripts by Gavrilo Stefanović Venclović using the Transkribus software platform, based on the principles of artificial intelligence, machine learning, and advanced neural networks. The training material for the model was taken from two manuscripts held in the Archives of the Serbian Academy of Sciences and Arts: *Slova izabrana*, written in the Serbian language and semi-uncial Cyrillic script; *Razlagolnik*, written in the Serbian Church Slavonic language and cursive Cyrillic script; and the Serbian Church Slavonic *Menaion for October* from the Archives of the Serbian Orthodox Diocese of Buda in Szentendre, written in semi-uncial Cyrillic script. The most important result of the work is the creation of the first version of the

generic HTR model *Vencl. 1.0* with exceptional performance – the percentage of incorrectly recognized characters is only 3%. The paper demonstrates how this model can be successfully used for the automatic transcription of other Serbian and Serbian Church Slavonic manuscripts by Venclović, which could finally make feasible not only a critical edition in the original script but also the creation of searchable digital editions, and subsequently, an electronic corpus of Venclović's written legacy.

**Keywords:** Serbian language, Serbian Church Slavonic, Gavrilo Stefanović Venclović, machine learning, artificial intelligence, Transkribus, automatic text recognition, digital edition.