# *Electronic Research Infrastructure for Bulgarian Medieval Written Heritage*: history and perspectives

Anna-Maria Totomanova ⓘ🅰ℝ⁸⋆

*Faculty of Slavic Studies, "St. Kliment Ohridski" University of Sofia, Tsar Osvoboditel Blvd. 15, 1504 Sofia, Bulgaria*

**Abstract**

The paper traces the history of the *Histdict* system, which turned into a basis for the new *Electronic Research Infrastructure for Bulgarian Medieval Written Heritage*, which was included into the National Research Roadmap at the end of 2020. Through this act the state declares its support to our resources, that have been so far created and supported by project funding. And of course, it is a big recognition of our efforts and achievements. On the other hand, this act coincided with two other events: the inclusion of RESILIENCE (Research Infrastructure on Religious Studies) in which *Histdict* is taking part, in the European Research Infrastructures Roadmap and the start of the updating and upgrading of the system. Given the situation the Infrastructure is now facing new challenges—not only the successful improvement of the services it offers, but also the inclusion of the Orthodox Cultural Heritage into European research exchange, which will promote and popularize the history and culture of Southeastern Europe.

## 1. Introduction

The *Electronic Research Infrastructure for Bulgarian Medieval Written Heritage* started in 2009 as a doctoral and post-doctoral project BG051PO001-3.3-04-001 ICT Tools for Historical Linguistic Studies, funded by the European Social Fund, OP Human Resources, and was upgraded by three following projects under the same scheme:

- BG051PO001-4.3.04-0004 *E-Medievalia (Electronic Resources for Distant Learning in Medieval Studies)* (2012–2014)
- BG051PO001-3.3.06-0024 *Informatics, Grammar, Lexicography* (2012–2015)
- BG05M2OP001-2-009-0005 *Modern Palæoslavonic and Medieval Studies* (2017–2019)[1]

In the periods when our activities were not covered by European funding, we used other project opportunities such as internal funding in the framework of the University Humanities Complex, in which we participated with the project *Digital Medievalia* (2016–2018). Since the end of 2019 the RI has been supported by the National Research Programme "Cultural Heritage, National Memory and Social Development" funded by the Ministry of Education and Science and coordinated by Sofia University as a leading research institution in the field of Humanities. In the meantime, we were allocated a substantial amount of money for upgrading the existing digital resources under the project *Heritage BG* funded by Science and Education for Smart Growth Operational Programme (*nasledstvo.bg*). Now, after two years spent for unsuccessful bids, we have selected an IT company that is supposed to fulfil this important task.

In 2019 we applied for becoming a part of the National Roadmap for RIs and at the end of 2020 our digital resources were officially accepted as a starting RI.

Since March 2018 we have been involved in two overlapping ESFRI projects aiming at becoming

---

⋆Email address: *atotomanova@abv.bg*.

[1]The projects' goals and results have been constantly reported and popularised, see Ganeva (2018), Totomanova (2012, 2017, 2018), Totomanova-Paneva (2020).

a European RI on Religious Studies (ReIReS[2], RESILIENCE[3]). Therefore, we have not been able so far to count on a stable funding on behalf of the state and we had to finance our work under different kinds of project initiatives in order to survive and implement our plans to produce a complex of digital resources for studying and popularising Bulgarian written heritage. Our resources are included in the *Cyrillomethodiana* web portal where you can find detailed information about all our projects related to the RI. The partnership consists of 3 institutions: Sofia University "St. Kliment Ohridski", Cyrillo-Methodian Research Center at Bulgarian Academy of Sciences, BAS Central Library.
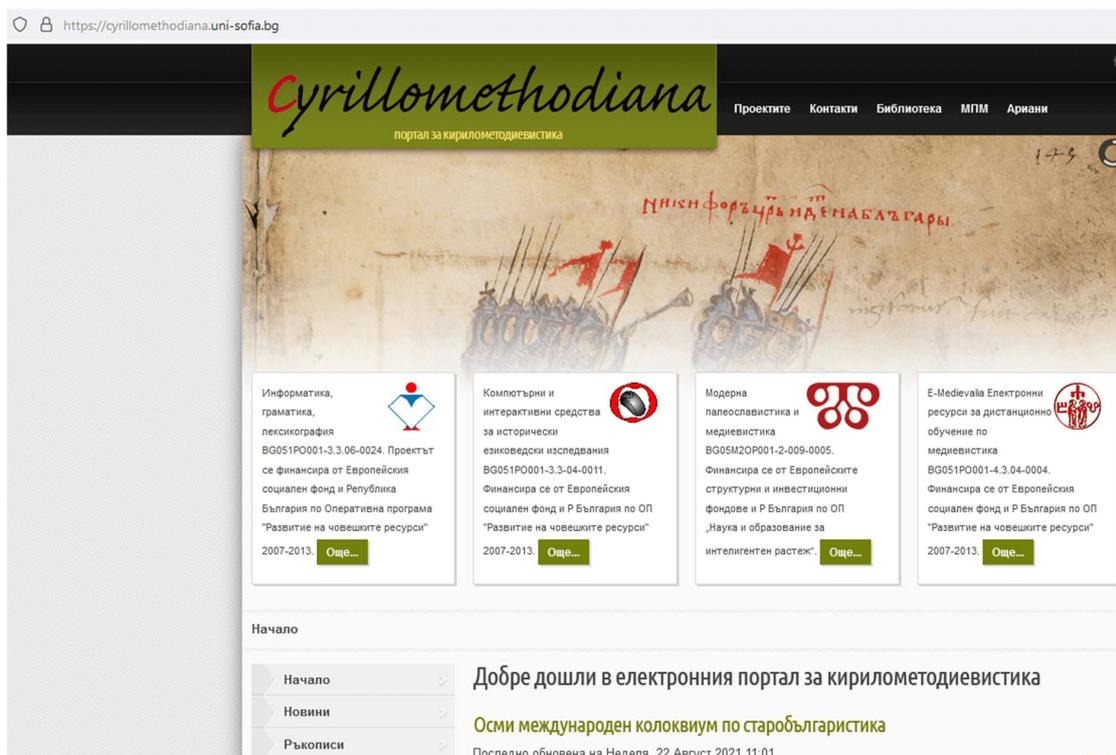


Figure 1: *Cyrillomethodiana* web portal, *cyrillomethodiana.uni-sofia.bg*

**But why the written heritage?**

The written heritage is the most reliable source of history of every nation and country and the basis of the national identity. Our RI provides digital research resources and tools for studying Bulgarian Medieval written heritage in terms of *diachronic linguistics*, *literary history*, *Bulgarian and European history*, *theology*, *cultural history*, and *philosophy*. The Old Bulgarian language (Old Church Slavonic), created by the Holy Brothers Cyril and Methodius, functioned as a sacred language of all Orthodox Slavs in Southeastern Europe during the Middle Ages as well as in Romania and Moldova. That is why researchers in the field of the history and culture of the Slavo-Byzantine Medieval world are interested in using the available digital resources.

**And why electronic?**

 a) Using digital technologies for publishing and processing medieval texts reduces the time of data collection and production of reliable research results. In other words, the introduction of digital technologies in such a conservative field as historical humanities not only optimizes the work of researchers but also creates the conditions for new research initiatives and projects.
 b) Digital tools make the field of historical humanities more attractive for young people born in the digital era.

---

[2]See *reires.eu*.
[3]See *resilience-ri.eu*.

## 2. Our digital resources and tools

As you can see the electronic research resources and tools for processing medieval texts of proven Bulgarian origin have been created in a period spanning over 10 years through a series of scientific and educational projects. The *Histdict* system includes the following digital resources and tools:

### 2.1. Diachronic corpus of Bulgarian language 9ᵗʰ–18ᵗʰ cc.

The Corpus contains more than 150 texts of proven Bulgarian origin from different genres of the 10ᵗʰ–18ᵗʰ cc. and has its own specific software that allows for textual and paleographic annotation. Given the fact that the Bulgarian literature transmitted the Byzantine cultural and literary model to the other Orthodox nations in our part of Europe the Corpus contains both translated and original Medieval Bulgarian texts. The texts are digitally typed and all of them (excluding the works of St. Kliment Ohridski) are reproduced in the orthography, in which they survived (Bulgarian, Serbian or Russian). The Corpus includes also Early Modern Bulgarian texts mostly Damaskini and other compilations as well as some not literary texts such as scribal notes, inscriptions, and juridical documents. Back in 2011 we had only 75 texts whereas now their number has doubled, and we continue uploading new texts. Some of them are provided by colleagues abroad who also use the corpus. One of the largest texts uploaded on to system is the so-called *Chronograph of the Archive* or *Jewish Chronograph* that includes the oldest text of the Octateuch and Kingdoms translated in Bulgaria during the reign of Simeon the Great (893–927). Each text is introduced by the respective archæographic data (title, source, dating, orthography, author etc.). The corpus software needs upgrading to allow for other types of annotation: morphological and content annotation.



Figure 2: Diachronic corpus, *histdict.uni-sofia.bg*

### 2.2. Digitized Old Bulgarian dictionary (10 500 dictionary entries)

Our main goal under the first project was to compile an electronic *Historical dictionary of Bulgarian language* and we decided that the easiest way to achieve it was to take the *Old Church Slavonic Dictionary*[4] compiled and published by Institute for Bulgarian Language to the BAS as a basis for this dictionary. That is why we started digitizing the *Old Church Slavonic Dictionary* but in two years we managed to digitize only 10 000 of its entries, and the remaining 500 entries were digitized later under the second project since they do not comply with the entire format of the other entries or contained different types of errors and had to be processed one by one. The full electronic version of this dictionary is now available online and contains the lexis of the Classical OCS manuscripts of 10ᵗʰ–11ᵗʰ cc.

---

[4]Старобългарски речник. Т.I, 1999. София, Валентин Траянов, 1027 с. Т. II, 2009. Валентин Траянов, 1325 с.

Figure 3: Old Bulgarian Dictionary online, *histdict.uni-sofia.bg*

## 2.3. Reverse Greek–Old Bulgarian dictionary

The structured XML document that included all entries of the Old Bulgarian Dictionary was used to produce a reversed Greek-Old Bulgarian Dictionary, which is also available online and is widely used by colleagues who deal with the metaphrastic practices of Old Bulgarian men of letters. The reverse dictionary was a spin-off result of E-Medievalia project that produced an interactive teaching-learning platform in Medieval Studies. The platform contains 24 interactive courses for both under- and postgraduate students in the area of Philology, Philosophy, Theology, History and Arts, among them there are courses of Old Bulgarian (ocs), Old Bulgarian literature and History of Bulgarian, which are available also in English. The platform is combined with a virtual classroom, that we used for workshops and conferences. When the pandemic started this digital resource allowed us to start immediately the on-line teaching.



Figure 4: Greek–Old Bulgarian Dictionary online, *e-medievalia.uni-sofia.bg*

## 2.4. *Historical dictionary of the Bulgarian language (a developing resource of diachronic type)*

The design of the historical dictionary software turned out to be the biggest challenge for the project team and especially for our ICT specialists, since as mentioned above, we decided to edit and upgrade the digitized version of the Old Bulgarian Dictionary, complementing it with new entries and new meanings in the old entries. Neither linguists nor ICT specialists were aware of how difficult this task might be. Only by the end of 2014, a few months before the end of the project "Informatics, Grammar, Lexicography" we finally received a software product that allowed for both editing old entries and adding new ones and started compiling the Historical Dictionary of Bulgarian. The Christian Terminology that had not been studied properly since Fr. Miklosich published his famous study in 1876 wa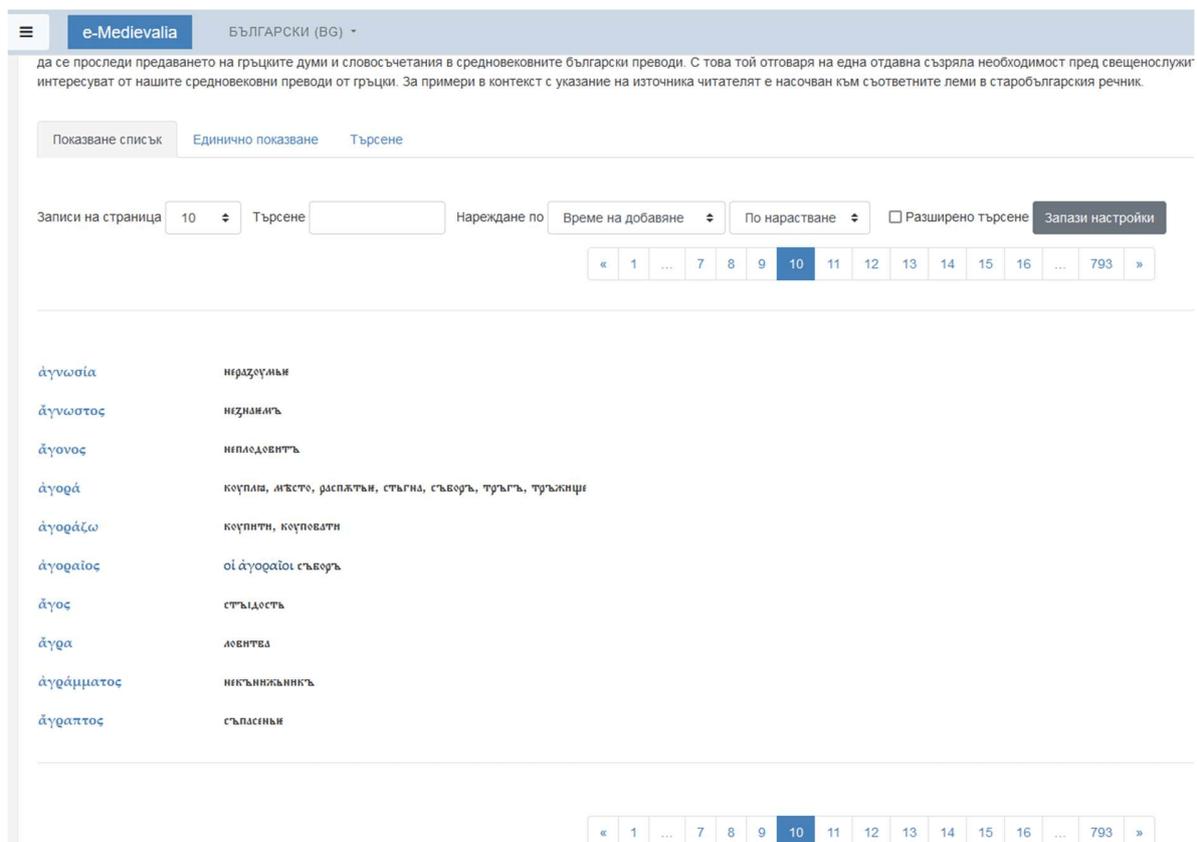s the first lexical group to be included into the Historical dictionary (Miklosich, 1876). The constantly developing version of this dictionary is also available online and the new and edited entries appear in colors—green or blue. Recently the historical dictionary has been complemented with the entries or specific meanings taken from the synchronous historical dictionaries that contain the lexis of a single writer or a group of writers or some specialized lexis.



Figure 5: Historical Dictionary, *histdict.uni-sofia.bg*

## 2.5. *Synchronous historical dictionaries*

Trying to expand and test the use of digital tools under the project "Modern Palæoslavonic and Medieval Studies" we produced two dictionaries of this type: *Patriarch Euthymius' Language dictionary* and *Terminological dictionary of John the Exarch*. Both dictionaries also have paper versions[5] and the second part of *Patriarch Euthymius' Language dictionary* was financed by the National Research Programme.

For compiling the dictionary entries an alphabetical list of all word forms in Patriarch Euthymius' writings contained in the corpus was produced electronically. The list shows the title of the opus and the exact place the forms occur and allows a quick work lemmatization.

A similar list for the terminological dictionary was also made available, but given the complexity of the task we decided to limit the linguistic material and explore the terminology of John the Exarch's *Hexaemeron* and *Theology*. As a result, around a 1000 of terms were included into this specialized dictionary.

---

[5]Речник на езика на Патриарх Евтимий. Част I. **А–Н**, Част II. **О–Ѧ**. УИ „Св. Климент Охридски", София, 2019, 2020; Терминологичен речник на Йоан Екзарх. София, 2019.

Figure 6: *Patriarch Euthymius' Language dictionary*, *histdict.uni-sofia.bg*



Figure 7: *Terminological dictionary of John the Exarch*, *histdict.uni-sofia.bg*

The specialized dictionary software for writing and editing dictionary entries seemed to be too complicated for the project participants—especially PhD and Postdoc students that had no experience under the previous projects. That is why we compiled the new dictionary entries in Word environment and included them into the respective specialized dictionaries through a converter. For this purpose, a special entry format was elaborated and observed. The converter for St. Patriarch Euthymius' dictionary turned out to be a great finding because it enabled us to convert more than 400 dictionary entries in few seconds. The online version of *Patriarch Euthymius' Language dictionary* is slightly different from the paper version, the latter reporting up to three examples under each identified meaning. The electronic dictionary, on the other hand, displays all occurrences of the lexemes in the text and allows quick access to the respective text in the corpus where the user can find the word and see and copy the respective contexts. It is done by a simple click on the signature.

We used a similar converter for the *Terminological dictionary* as well and now we are trying to elaborate another converter for digitizing the reverse Greek–Old Church Slavonic dictionary produced by Christov (2019). The dictionary contains 14 625 Greek entries, to which correspond 33 307 Greek–Slavonic parallels.

The converted dictionary entries allow for online redaction using the specialized software for the historical dictionary. However, given the complexity of this software, we intent to produce a converter for uploading the entries in the Historical dictionary as well, thus speeding up the work on this important resource.

## 2.6. *Search engine and virtual keyboard*

The difficulties to produce a proper search engine came from the fact that Old Bulgarian is a language with complicated inflectional morphology especially in the nominal paradigm (six cases plus Vocative, three genders, three grammatical numbers, simple and compound forms of the adjectives and declinable participles), which further was reduced to 2 to 5 forms, depending on grammatical gender. So, it turned out that in order to design and produce a reliable search engine we had to have a tagger (electronic tool that

allows for morphological annotation of the inflectable words). And to produce one we needed to create a grammatical dictionary of ocs taking into account all possible representations of a single form. Yet we needed a quick search tool in order to facilitate the work on the historical dictionary and the use of the *Corpus*. The temporary solution came from our main software specialist who digitized the ocs dictionary and developed the software for the Historical dictionary. The search engine is installed in the system and one can search on all its entries (dictionaries, the *Corpus*, and the *Chronograph*) or choose one or more of them. The machine shows the list of the texts, in which the form, we look for, is found and to find it together with the respective context one has to use the browser internal search functionalities. Therefore, one of the main tasks of the ICT company we selected will be to design and install a real search engine that will enable us to search without using browser functionalities and according to different parameters: beginning or the end of the word, letter strings contained within the word, grammatical forms. The virtual keyboard, that might need some fine tuning as well, is also in place.



Figure 8: Search engine, *histdict.uni-sofia.bg*

## 2.7. Grammatical dictionary and semi-automatic morphological analyzer

The work on the grammatical dictionary started with compiling the tagset of ocs, which contains 2200 tags and describes the complicated ocs grammar (Totomanova *et al.*, 2015). The tagset was followed by a grammatical dictionary, which includes paradigms of the inflectable words taking into account different phonetic and orthographic variations (Totomanova *et al.*, online). The paradigms (rules) were ascribed to all inflectable words in the historical dictionary and as a final step the grammatical dictionary was installed into it. Now by clicking the sign + located next to the lemma the user can see the whole paradigm of the respective word. Based on that our ICT specialists created a prototype of the morphological annotator, which is also an open access tool on our site.

Yet this is another temporary solution because our goal is however to create an automatic tagger. For this purpose, we are upgrading the grammatical dictionary and editing the grammatical rules. Entering the uninflectable words (adverbs, prepositions, conjunctions, particles, interjections) as well as pronominal forms into grammatical dictionary turned out to be a big challenge. The inclusion of the participial forms seems to represent even a bigger challenge, given the fact that ocs verbs possess 5 participles, out of which four are declinable and have also determined counterparts. We developed rules for automatic generation of the participial paradigms depending on the verb type and so far, have produced the participial forms of the verbs of the first two ocs conjugations (i.e. 894 verbs with over 3 million forms). The grammatical dictionary upgrade made us realize the need to edit the lemmas and the grammatical definitions in the

Historical Dictionary, which seems to turn into a constant process and is also one of the tasks the ICT company is supposed to fulfil.



Figure 9: Grammatical dictionary and tagger, *histdict.uni-sofia.bg*

### 2.8. Old Bulgarian Unicode fonts and converter

For compiling the Corpus and digitizing the Old Bulgarian Dictionary we needed not only the respective software but also specialized Unicode (UTF-8) OCS fonts readable online. To avoid the retyping of already digitally typed texts the new fonts were installed in a converter that allowed for their conversion to the new fonts. Now we use the third version of the converter, which includes 3 specialized fonts with different design *CyrillicaBulgarian10U*, *CyrillicaOchrid10U* and *CyrillicaOldStyleU*. The last one is meant for typing the Early Modern Bulgarian texts. In the beginning we were able to convert only the fonts of the families *Cyrillica Bulgarian*, *Cyrillica Ohrid* and *Cyrillica Shafarik*, developed by Synthesis Soft company and widely spread among the Palæoslavonic research community, to which gradually other OCS and Ancient Greek fonts were added: the Italian *PopRetkov*, Unicode font *BukyVede*, which is also based on Synthesis Soft design, *TimesGreekClassic* and *TimesGreekOld*, also developed by the same company, convert now into *Palatino*, as well as all varieties of the modern font *TimesCyrillic* in *Times New Roman*. Our OCS fonts and the converter spread throughout Europe since they were compatible with the editorial software programs and a series of books and periodicals in both Bulgaria and abroad are using them.

The open access Histdict system is a complex of unique resources and tools for publication and research of medieval Slavonic texts and in a European and global context it corresponds to *Thesaurus Linguæ Græcæ* and *Perseus*, representing the classical written heritage. Compared to them our RI has a big advantage—it is based on the link between the diachronic corpus and the respective dictionaries, that are being elaborated using the lexical material form the corpus.

## 3. Conclusions: future plans and development

The main objective of RI is to maintain, develop and upgrade the electronic research resources and tools through the creation of new functionalities and research capabilities. The following activities are envisaged:

- Complementing the corpus with new texts in order to collect all Bulgarian written heritage in digital format

- Creating a parallel corpus with translations into Modern Bulgarian language so that the works of Old Bulgarian writers will be accessible to the general public as well
- Continuous updating of the historical dictionary
- Creating synchronous electronic dictionaries on medieval texts
- Maintaining and upgrading the existing software
- Research
- Training of doctoral and postdoctoral students and young scholars
- Publishing: research and publications of texts and dictionaries

On June 30th, 2021, the ESFRI Forum included RESILIENCE in the Research Infrastructure Roadmap 2021. This means that RESILIENCE will take its place in the strategic Research Infrastructures for the European Research Area and that it can work on the further development of the Research Infrastructure for Religious Studies. Therefore, the situation now is different. As a part of RESILIENCE we will expand the circle of potential users of our digital resources including also the learning platform E-Medievalia with its courses in Old Church Slavonic. But the most important result of our participation in RESILIENCE opens for us new horizons and provides the great opportunity to include the cultural heritage of Southeastern Europe related to the Christianity into European research exchange. At the same time, it represents another big challenge for us. And to face this challenge decently we are looking for new partners...

## Bibliography

Christov, I. (2019). *Гръцко-църковнославянски речник*. Съставен от Иван Христов въз основа на *Речника на църковнославянския език* от архимандрит д-р Атанасий Бончев. Редактор А. Тотоманова, Зографски манастир "Света гора".

Ganeva, G. (2018). *Electronic Diachronic Corpus and Dictionaries of Old Bulgarian*, in "Studia Ceranea", **8**, p. 111–119, Crossref.

Miklosich, Fr. (1876). *Die christliche Terminologie der Slavischen Sprachen: Eine sprachgeschichtliche Untersuchung von Franz Miklosich*. Denkschriften der kaiserlichen Akademie der Wissenschaften. Philosophisch-Historische Klasse. Band 24. Wien.

Totomanova, A. (2012). *Digital Presentation of Bulgarian Lexical Heritage. Towards an Electronic Dictionary*, in "Studia Ceranea", **2**, p. 221–234, Crossref.

Totomanova, A. (2017). *Диахронный корпус болгарского языка. Состояние и перспективы*, in "Filologija", **68**, p. 223–242, Crossref.

Totomanova, A. (2018). *OCS Biblical Language in the Era of the New Technologies*, in *Дигиталатни и аналитични подходи към писменото наследство. Материали от 7-мата международна конференция El'Manuscript "Писменото наследство и информационните технологии"*, Sofia, p. 250–265.

Totomanova, A., Slavova, T. & Ganeva, G. (2015). *Морфосинтактичен тагсет на старобългарския книжовен език*, in *Информатика, граматика, лексикография BG051-3.3-06-0024/2012. Сборник доклади и материали от заключителната конференция*, Sofia, 29-30.06.2015, p. 17–117.

Totomanova, A., Slavova, T. & Ganeva, G. (online). *Граматически речник на старобългарския език*, [online].

Totomanova-Paneva M. (2020). *Дигитални ресурси за подготовка на докторанти, постдокторанти и млади учени*, in *Девети международен есенен научно-образователен форум. Съвременният учител и предизвикателствата на информационното общество*, Sofia, p. 91–94.