

DIGITAL PRESENTATION OF BULGARIAN LEXICAL HERITAGE TOWARDS AN ELECTRONIC HISTORICAL DICTIONARY

The project: ICT Tools for Historical Linguistic Studies funded by the European Social Fund, OP Human Resources, was designed and carried out with the idea to introduce ICT in such a conservative field as diachronic linguistics. The objective we pursued was twofold:

- to speed up the data collecting from the books created between 10th and 18th cent. and accelerate further data processing;
- to make diachronic linguistics more attractive for young people born in the Computer Age for whom computers are part of their natural habitat.

The Round Table *Interactive Methods in Historical Lexicology and Lexicography* held on 28 V 2010 played a crucial role for the project development. The participants reviewed and summarized the experience in the area of historical lexicography and made the following important decisions:

1. The project should focus on creating software tools for developing a web based Historical Dictionary of Bulgarian, which is the first literary and sacred language of the Slavs with a long written history.
2. Старобългарски речник (Old-Bulgarian Dictionary), created by the Department of History of Bulgarian Language to the Institute for Bulgarian Language, will constitute the foundation for building on a Historical Dictionary of Bulgarian. For this purpose the information it includes will not only be preserved but it will also be enriched and upgraded with materials taken from the Electronic Corpus of Medieval and Early Modern Bulgarian texts.

The project target group participants (PhD and Post-Doc students, young researchers and interns) were assigned individual research tasks in compliance with the decisions made. The Round Table produced a preliminary list of electronic tools for digital processing of the texts. The *Standard of the Dictionary* took shape during the project course based on the decision that we are aiming at designing a *Historical Dictionary of Diachronic Type¹ that should present the history of the Bulgarian words from their first written occurrence until today*. Such a Historical Dictionary has the following features:

¹ The terms *Diachronic* and *Synchronic Historical Dictionaries* were introduced and explained by: Г.А. Богатова, Историческая лексикография как жанр, ВЯ, 1981, p. 83–84.

- **Large chronological span**, starting from the beginning of the Slavonic writing in the 9th cent. up to the modern times;
- **Thematically unlimited text corpus** that includes: literary texts; non-literary texts (geographic and personal names, dialects, vernacular language, inscriptions, graffiti);
- **Open vocabulary** that will be enriched alongside with the corpus building;
- **Diachronic presentation of the lexical material**, which implies the registration of the different meanings of the word and their genetic connection.

The Text Corpus of the Dictionary should include:

- **Bulgarian medieval texts**: opera of the Old-Bulgarian writers; translations from Greek with proven Bulgarian origins (opera of the Holy Fathers, Chronicles, monastic literature, Historical and Apocalyptic texts, juridical texts, miscellanies with stable and mixed content etc.);
- **Non-Literary texts**: notes of the copyists; inscriptions and graffiti; charts;
- **Early Modern Bulgarian texts** (mostly *Damaskins* and *Damaskin miscellanies*);
- **Dialectal texts**.

To create the electronic base of the Historical Dictionary the following electronic tools are needed:

- Digitalized Старобългарски речник;
- Specialized *Diachronic Corpus of Medieval Bulgarian and Early Modern Bulgarian texts*;
- Other specialized corpora, such as *Bulgarian National Corpus* (Български национален корпус)², *dialectal corpora*, *BgSpeech Corpus* (Корпус на българската разговорна реч)³ and so on.

Since the work on the other specialized corpora had already begun, the project team efforts concentrated on creating the Corpus of Medieval and Early Modern Bulgarian texts and on digitalizing the two volumes of Старобългарски речник. The creation of a new Old Bulgarian font was the first step towards the electronic processing of the medieval texts.

In the beginning of 2010 we had already at our disposal a new Old Bulgarian font based on Unicode, that contained more signs than the existing up to that moment Old Bulgarian Unicode fonts. The font has already successfully been used for the digital typing and publishing of some medieval texts. The medieval texts in the last three books of the series “History and Literature” were converted into the new font. The same font is being used for publishing the text of the Bulgarian, Russian and Serbian Synodika for the planned *Brepols* edition *COGD IV*³

² See the description and opportunities of using the BG National corpus on http://www.ibl.bas.bg/BGNC_bg.htm.

³ The corpus was developed as a part of BgSpeech initiative and it is maintained by the Faculty of Slavic Studies to Sofia University at <http://bgspeech.net/>.

as well as for the electronic edition of the so called Архивский хронограф we are preparing under another project. The project team contributed a lot to the improvement of the font functionalities by providing valuable feedback to the software specialists.

The collaboration between the ICT specialists and project participants produced the synergy for the successful use of the font *Cyrillica Bulgarian 10 U* under different type of editing and publishing software and facilitated the Pre-print processing of medieval Slavonic texts. The font was initially elaborated under the project “The Concepts of History across the Orthodox Slavic World” but it was used for the first time and substantially improved under this project. The same font is used by the editorial project for publishing Slavic Synodica as well as by the project *PragmaticFunction Words: A Corpus-Based Description of Variation* run by O. Mladenova at University of Calgary, Canada. The technological development and the mass introduction of so called *web fonts* in the browsers allow the users to read the font without installing it in their own operational systems (fig. 1).

Together with the font a convertor was produced that converts the texts typed with the *Synthesis Soft* fonts into Unicode-based documents. All project participants contributed to the testing and improvement of the convertor and learned how to apply it, converting already typed texts for the diachronic corpus of Bulgarian. By the end of the project the convertor functionalities were expanded to all Synthesis Soft funds plus the Italian Pop-Retkov font, which is of great importance since our Italian colleagues provided us with the digitally typed Alphabetical⁴ and Roman⁵ pateriks (fig. 2). Two additional Unicode fonts were included as well: Cyrillica Ochrid 10 U and Cyrillica Old Style 10 U that is designed for typing Early Modern Bulgarian texts.

The font *Cyrillica Bulgarian 10 U* was used for digitalizing the two volumes of Старобългарски речник, produced by IBL. We express our gratitude to the ICT consultant Mr. Todor Todorov, who developed the font and the convertor and created a second specialized convertor/generator that successfully converted the dictionary containing 11000 entries into a structured XML document without losing a bit of existing information. This second convertor facilitates the process of converting of other already published on paper medieval texts such as Германов *сборник* for example. The software specialists from *Openintegra* elaborated software for editing, expanding and visualizing the dictionary in web environment. It allows an easy and quick access to the media and contributes to popularizing the work of the team all over the world.

⁴ COGD. I–VII. A Special Series of *Corpus Christianorum* by Brepols, 2006 – An International Research Program launched in Bologna and directed by †Giuseppe Alberigo and Alberto Melloni of FSCIRE, Fondazione per le Scienze Religiose Giovanni XXIII, Bologna.

⁵ R. CALDARELLI, *Il Paterik Alfabetico-Anonimo nella traduzione antico-slava*, Roma 1996.

It also enables the data exchange between our institution and other universities since the dictionary is based on the world recognized standard TEI in XML area. The digitalized Old Bulgarian Dictionary is located on the project web page and is accessible for all customers at *histdict.uni-sofia.bg*. We are proud to say that it is the first digitally presented Palaeoslavonic lexicographic manual (fig. 3 and 4).

At the same address *histdict.uni-sofia.bg* one can find also the Diachronic Text Corpus that already contains more than 75 texts of different length and the text collection is constantly filling up. The corpus includes medieval Slavonic texts with proven Bulgarian origins and different orthography (Old Bulgarian – OCS, Middle Bulgarian, Resavian and Russian), Early Modern Bulgarian texts and notes of the medieval copyists. Translations and original works of the Old Bulgarian writers are equally represented in their genre variety – liturgical, exegetical, hagiographic, juridical, chronographic, historical and apocalyptic texts and so on. Some of them have not been published before.

Most project participants actively committed themselves to the workshop held on 20 XI 2011 that was dedicated to the digital presentation of the medieval texts in the corpus. To our great satisfaction in two weeks all interested parties – the project team, target group representatives, tutors and ICT specialists – all together managed to enter into the corpus a bigger number of texts than it was initially planned. The ICT specialists from *Openintegra company* supported our team helping to alleviate errors that occurred during the testing while entering texts and added new functionalities to the corpus software suggested by the team. We consider that to be an enormous success, given the fact that this is the first diachronic corpus based on Slavonic material connected to the elaboration of a historical dictionary and provided with a program for linguistic annotation.

The software we developed is *user friendly* and very easy to use. The electronic tools for text commentaries (both paleographic and codicological) as well as for visualizing variant readings, create new opportunities for adequate presentation of the medieval Slavonic texts that will be applied to the digital edition of the Chronograph of Archive, planned under the project “The Concepts of History across the Orthodox Slavic World” and to other electronic publications (fig. 6–11 show the Corpus functionalities)⁶.

The software is fully transferable and could be used for digital processing of texts, for creating corpora and dictionaries of different languages. That is why the software developers and the team have the intention to publish it as an Open source material so that our colleague from abroad might access it. In return we hope to receive from them some ideas about its further improvement and application.

⁶ К. Диди, Патерик Римский. Диалоги Григория Великого в древнеславянском переводе, Москва 2001.

The corpus itself turned out to be a wonderful tool for digital presentation of the Bulgarian lexical heritage on diachronic scale. The openness and accessibility of the data it contains provide opportunities for its expanding through adding new meanings and lexemes. The text uploading is very simple and the copyright of the authors is preserved through the introduction of different access levels.

The corpus is also a study tool and could be easily installed into the teaching-learning process in area of Palaeoslavonic and medieval studies as well as in diachronic linguistics.

The corpus is supplied with a *Search engine* that allows searching the texts by metadata (author, genre, orthography etc.) as well as directly in the text content.

A programme for editing the articles of the digitalized Старобългарски речник was developed to make the dictionary the basis for creating the Historical Dictionary of Bulgarian. We have already started adding new lexemes that are not registered in the Old Bulgarian manuscripts and developed a number of new dictionary units using the experience and methodology of the authors of Старобългарски речник (fig. 5).

Yet the real work on the dictionary is only about to start. For this purpose we have to focus our efforts on the following directions: Developing new dictionary entries.

Expanding the chronological coverage of the existing dictionary entries.

Editing the units/articles of the Historical Dictionary.

In order to solve these problems we have to establish a connection between the Corpus and the Historical Dictionary, that shall allow us to discover both the missing lexemes and the new unregistered so far meanings. Producing glossaries and lists of lexemes for lexicographically non explored texts from the corpus will be one of the project spin-off results. I do not think however that we should overlook the materials that can be found in already published lexicographic manuals. Adding new dictionary entries and new meanings in the existing ones will require a careful editing of Старобългарски речник entries, since the Historical Dictionary will rather focus on tracking the word meaning development throughout the centuries than on the exhaustive presentation of the lexical material. But we are still at the beginning and expect to gain valuable experience in this regard.

The set of electronic tools for creating corpora and dictionaries on medieval Bulgarian text material seems to be the most impressive and important project result. I am deeply convinced that the free access to both the corpus and its digital version will attract to our work many followers from both the country and abroad who will contribute to this extremely important lexicographic project.

The Diachronic Corpus of Bulgarian we created is the first of this kind since it is connected to a dictionary and supplied with respective electronic tools for text processing. The electronic source might have many applications since it could be used for:

1. Producing e-based lexicographic manuals of different types:

- Diachronic Historical Dictionaries;
- Historical Dictionaries of synchronic type (Dictionaries of Literature or of different authors, different periods etc.);
- Glossaries;
- Thematic dictionaries;
- Etymological dictionaries.

2. Historical Linguistic Studies in the area of:

- Morphology and Morphosyntax;
- Morphology;
- Phonetics;
- Lexicology;
- Etymology;
- Derivation;
- Phraseology;
- Textology;
- Orthography.

3. University education on all levels (bachelor, master, doctor) in the field of:

- Palaeoslavonic and Old Church Slavonic Studies;
- History of Bulgarian Language;
- History of Literary Bulgarian;
- Old Bulgarian Literature;
- Medieval History;
- Computer and Corpus based linguistics.

4. Preparing the editions (both traditional and electronic) of :

- Medieval texts;
- Dictionaries, Glossaries etc.;
- Textbooks, Handbooks, Manuals etc.

5. Presenting Bulgarian Cultural Heritage

Abstract. The article presents the results of the project “ICT Tools for Historical Linguistic Studies” funded by the European Social Fund, OP Human Resources. The main project goal was to elaborate electronic tools for creating a Historical Dictionary of Diachronic Type that should present the history of the Bulgarian words from their first written occurrence until today. By the end of the project the partnership (Faculty of Slavic Studies to Sofia University, Institute for Bulgarian Language, BAS and PAM Publishing Company, Sofia) had at their disposal a set of *Old Bulgarian Unicode fonts*, meant for publishing medieval texts and a *converter* that converts no Unicode documents into the new

standard. The convertor allowed the participants in a relatively short time to create a *Diachronic text corpus of Bulgarian medieval texts* that contains already more than 90 texts dated from 10 to 18 century. The corpus software enables editing the texts and turned out to be an excellent tool for preparing electronic editions of the Old Bulgarian (OCS) manuscripts. In addition to the corpus an *electronic dictionary of Old Bulgarian* is available, which contains the digitized version of *Старобългарски речник*, produced by IBL. Both tools are accessible on the project website at the address histdict.uni-sofia.bg.

The *Standard of the Historical Dictionary* took shape during the project course and respective software for elaborating new dictionary entries was designed and tested. The article displays also screenshots that demonstrate the functionalities of both the corpus and dictionary software.

Anna-Maria Totomanova
St. Clement of Ohrid University of Sofia
15 Tsar Osvoboditel blvd.
1000 Sofia, Bulgaria
atotomanova@abv.bg

Figures:

Fig. 1. Cyrillica Bulgarian 10 U.

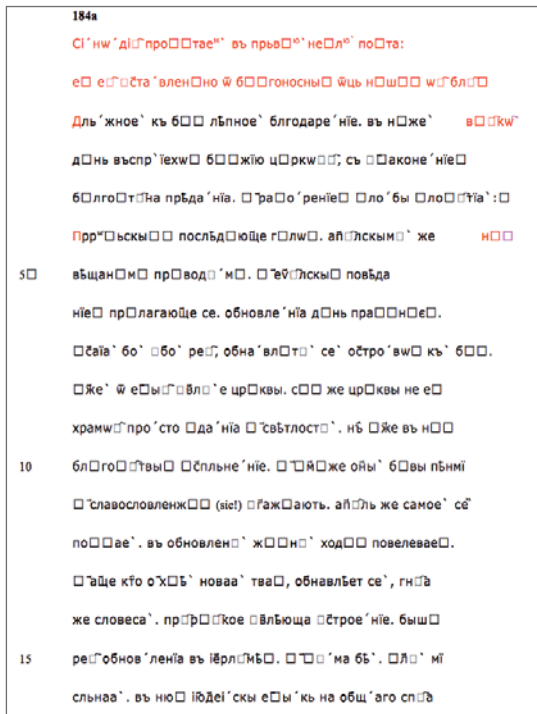


Fig. 2. Converter interface.

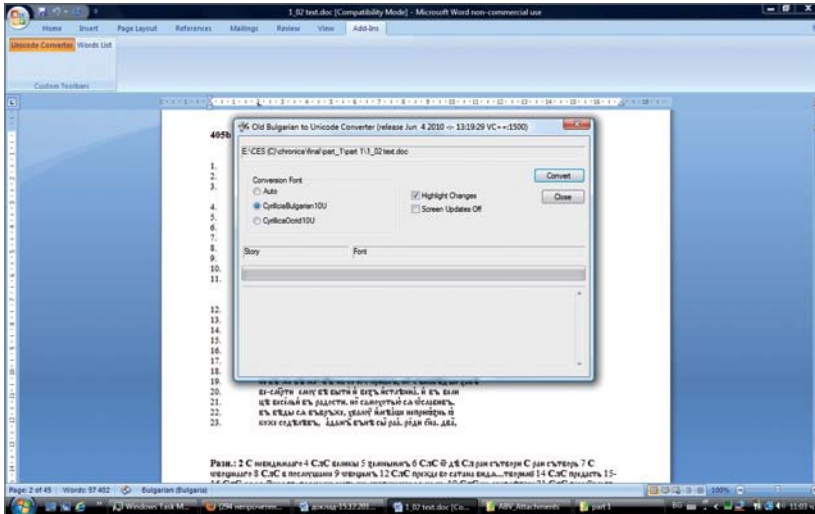


Fig. 3. Digitalized Старобългарски речник Interface (Lexeme search).

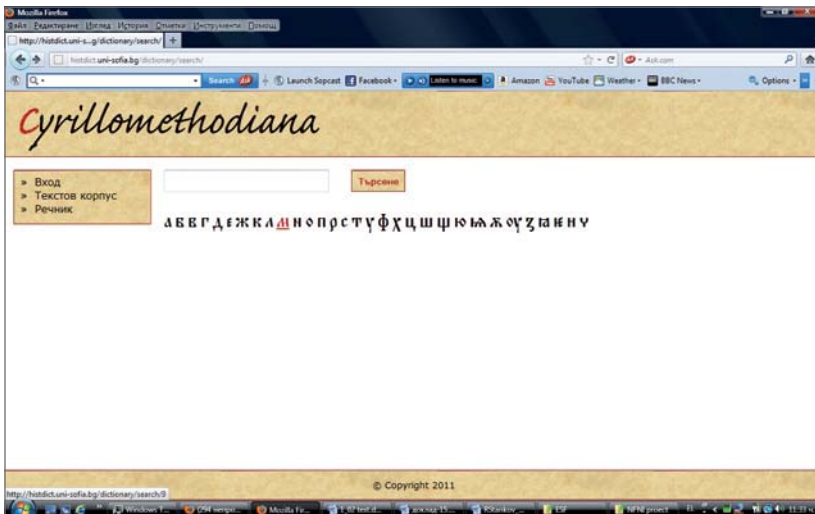


Fig. 4. Digitalized Старобългарски речник Interface (Dictionary entries).

АЛЕКСАНДРОВЪ

АЛЕКСАНДРОВЪ прил притеж от ЛИ

1. Александров, на Александър — синът на Симон Киринеец ꙗзѣдѣша мною ходоаштоу едномуу снмоуу кувѣннннѣ ... отцоу александроуу. н рѣфоуу *М Мк 15.21 З А СК*

2. Александров, на Александър — презвитер в Сид [Памфилия], умр. мъченически при имп. Аврелиан [270—275 г.] съприниманнѣ бѣдн надълекаштнн съмрътн александровѣ *С 161.2*

Изч *М З А СК С Гр* [тоу] 'Аλεξάνδρου **АЛЕКСАНДРОВЪ** **АЛЕКСАНДРОВЪ** **АЛЕΞΑΝΔΡΟΒЪ** Нвб александров *ОА ВА СРв* Александров *ФИ СТИл, РЛФИ* Александрово *ср МИПК, Пр.* в им

АГНЬЦЪ

АГНЬЦЪ *а м*

1. Агнец, агне ꙗзѣте се азѣ послымаѣ въи. ꙗко агньца по срѣдѣ влѣкѣ *М Лк 10.3 СРв. С 534.26* горѣ възиграша сѣ ꙗко овьнн. ꙗ хлѣмн ꙗко агньци овьнн *СП 113.4 СРв. СП 113.6 СЕ Зб 10—11* мо на агньцель на пасѣж. призьрн ꙗ нехѣ. на си брашѣна твоѣ. ꙗ на агньцъ сѣ. ꙗ стн н. ꙗкоже ститн нзволн агньцъ. ꙗже приведе авелѣ во вѣсѣжарбелма *СЕ 16b 2, 4, 5* юдѣи же събѣзаште агньцъ закалахѣ. а ꙗже отѣ поганѣ. въ пѣвѣтѣ бѣ *К 13b 14 СРв. С 450.21* акѣ овьнѣ на заколнннн вѣденѣ вѣстѣ. н акѣ агньцъ *С 434.25—26 СРв. С 437.2* Образно. ꙗгда же овѣдоваша. ꙗла снмоновн. петроу нс. снмоно новннѣ лювншн лн ма павѣ (снхѣ). ꙗла емоу ен ꙗн. тѣ вѣсн ꙗко люваѣ тѣ. ꙗла емоу пасн агньца моѣ *М Йо 21.15 З А*

2. В христианството — название на Исус Христос, който е принесен в жертва като изкупление за греха на човечеството ꙗ ꙗже нашѣ. прѣдѣложен сѣ салѣ. агньцъ нпороуинѣ. за жнкоѣтѣ вѣснго анра. призьрн на нѣи. ꙗ на (х)лѣкѣ сѣ. ꙗ на чашнѣ снѣж. ꙗ съ(т)ворн ꙗк прѣвстѣе тѣло твоѣ. хѣ *СС 1b 4* бѣко ꙗ ꙗже нашѣ. вѣсѣдожнтельо. истнннѣнѣ агньцѣ. въземѣн грѣхѣи вѣснго анра. не прѣзьрн дѣшѣ молаѣцѣ сѣ тѣвѣ *СЕ 15a 4* стоаши на крѣстѣ агньцъ. н два вѣкѣ *С 437.15*

АГНЬЦЪ БОЖЪН

ὁ ἀμνὸς τοῦ θεοῦ Агнец божий — изкупителна жертва [за Исус Христос]

въ оутрѣн дѣнь вндѣ нсѣ градѣшта кѣ свѣѣ. ꙗ ꙗла се агньцъ бѣжн. въземѣн грѣхѣи анра вѣснго *М Йо 1.29 З А СК Б* ꙗ оузьрѣ нсѣ градѣшта. ꙗла се агньцъ бѣжн *М Йо 1.36 З А* сн во вѣса вѣшѣ. да отѣнемѣн грѣхѣи анроу. агньцъ н снѣ божнн. воѣлѣ на съпаснѣжѣ стрѣсть сѣ вѣмн прндѣтѣ. н на проданнн станѣтѣ *С 331.25*

М З А СК Б СП СС СЕ К С Гр ἀρνὶν ἀμνὸν ἀμνὸς **АГНЬЦЪ** **АГНЬЦЪ** **АГНЬЦЪ** **АГНЬЦЪ** Нвб агнец *остар ОА ВА НТ НГр* ЕтБАН ЕтМл МлБТР *АР РБЕ РРОДД*

Fig. 5. Dictionary Entry Editing Tool Interface.

ВѢДРЪНЪ

Лема Добави нов контейнер

Форма:

Граматична група Добави нов контейнер

Разширение:

Част на речта:

Значение Добави нов контейнер

Дефиниция:

Пример Добави нов контейнер

Пример на старобългарски:

Издание Добави нов контейнер

Съкращение:

Страница:

Fig. 6. Corpus Interface (Text search)

Cyrillomethodiana

» Вход
» Текстов корпус
» Речник

Текстов корпус

Заглавие:	Псалтирь на св. Кирилъ
Заглавие на латински:	Psalterium Scripturae divinitus inspiratae sancti patris Antiochi
Жанр:	Слова
Дата на ръкописа:	Първата половина на XI век
Правопис:	Руски
Заглавие:	Државна конституция (ЗКО и ЛЗ) - Официален прѣпис
Заглавие на латински:	CODEX
Жанр:	Юридически
Дата на ръкописа:	XIV век
Правопис:	Ресавски
Заглавие:	ОБЛОБО ОУПЪТНИ СЪВЕЩАНИЯ - ВЪ ТЪМЪ КЪТЪ НА РЕПЪСАНИИ СЪСТАВИ СЪСТАВИ
Заглавие на латински:	...
Жанр:	историко-апокалиптично съчинение
Дата на ръкописа:	трета четвърт на XIII век
Правопис:	Руски
Заглавие:	Описание на Псалтиря - Официален прѣпис

© Copyright 2011

Fig. 7. Corpus functionalities (Metadata editing)

Редактиране на текст	
Заглавие:	<input type="text" value="Слеса на св. Доротеј"/>
Заглавие на латински:	<input type="text" value="S. Dorotheus archimandrita - Doctrinae diversae"/>
Жанр:	<input type="text" value="поучителни слова"/>
Автор:	<input type="text" value="св. Доротеј"/>
Превод?	<input type="checkbox"/>
Дата на ръкописа:	<input type="text" value="средата на XIV век"/>
Дата на превода:	<input type="text"/>
Дата на преписа:	<input type="text" value="средата на XIV век"/>
Правопис:	<input type="text" value="Среднобългарски / Търновски"/>
Име на ръкописа:	<input type="text" value="ръкопис № 1054 от свирката на М. П. Пегодин"/>
Хранилище на ръкописа:	<input type="text" value="Руска национална библиотека - Санкт Петербург"/>
Сигнатура на ръкописа:	<input type="text" value="Пог. 1054"/>
Страници:	<input type="text" value="226 (л. 191а-л. 304б)"/>
Нормализиран текст?	<input type="checkbox"/>

Fig. 8. Corpus Interface (Entering/editing texts)

Редактиране на текст	
Заглавие:	<input type="text" value="Пандекте на Аветик"/>
Заглавие на латински:	<input type="text" value="Pandectes Scripturae divinus inspiratae sancti patris Aetio"/>
Жанр:	<input type="text" value="Слова"/>
Автор:	<input type="text" value="Превод на текст"/>
Превод?	<input type="checkbox"/>
Дата на ръкописа:	<input type="text" value="Първата половина на XI век"/>
Дата на превода:	<input type="text" value="Първата половина на XI век"/>
Дата на преписа:	<input type="text" value="Първата половина на XI век"/>
Правопис:	<input type="text" value="Руски"/>
Име на ръкописа:	<input type="text" value="Пандекте Аветика"/>
Хранилище на ръкописа:	<input type="text" value="Москва, ГИМ"/>
Сигнатура на ръкописа:	<input type="text" value="ГИМ, Воскр. 30р"/>
Страници:	<input type="text" value="615 (от 2а до 309б)"/>
Нормализиран текст?	<input type="checkbox"/>
Текст:	<input type="text" value="За Присъмяния антиноха уърноуица мвръгы"/>

Fig. 9. Corpus functionalities (Footnote)

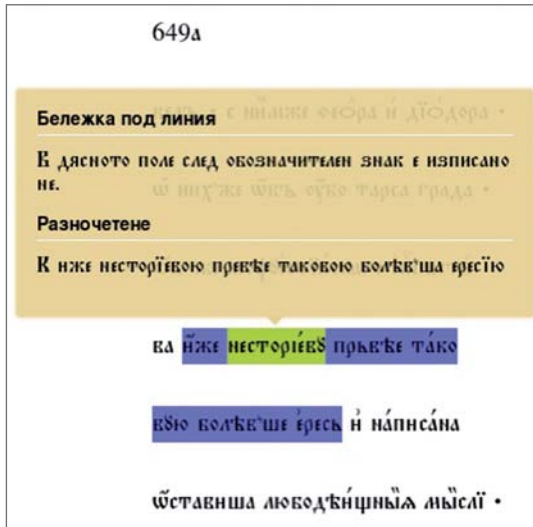


Fig. 10. Corpus functionalities (Variant readings)



Fig. 11. Corpus functionalities (Red letters)

