

Sustainable Development of the Practices of Digitization in National Library “Ivan Vazov” – Plovdiv

Ivan Kratchanov

National Library “Ivan Vazov”, 17 Avksentii Veleshki Str., Plovdiv 4000, Bulgaria
ivankra@gmail.com

Abstract. National Library “Ivan Vazov” in Plovdiv is the second largest library in Bulgaria. It is established as the second national depository of Bulgarian printed output and has contributed significantly to the preservation of the national cultural and historical heritage. This article offers an overview of the library’s history and current developments in the field of automation and digitization.

Keywords: Digitization, Plovdiv, Cultural Heritage, Digital Library.

1 Introduction to National Library “Ivan Vazov” – History, Holdings and Traditions

National Library “Ivan Vazov” in Plovdiv is the second national repository of Bulgarian textual heritage. The library has played an important role in the preservation of Bulgarian culture and history. It is the first cultural institution in Southern Bulgaria, established as a Regional Library and Museum of Eastern Rumelia in 1879.

After the Unification, the established legislature regulated equal rights and obligations for the two national libraries - in Sofia and Plovdiv. That is why the National Library in Plovdiv has developed as an archive of Bulgarian books and periodicals, a Bulgarian historical archive, a rich repository of manuscripts and Revival literature, unique collections of rare and valuable publications.

Today, National Library “Ivan Vazov” in Plovdiv is a cultural institute that continues to dynamically enrich and develop the traditions of its prominent founders. Patrons annually realize over 120,000 visits and loan 300,000 library documents. The library’s holdings are comprehensive and amount to over 1,900,000 library units – scientific, fiction, manuscripts, old-printed, rare and valuable publications, Bulgarian and foreign periodicals, photographs, maps, audiovisual and electronic documents, original works of art, personal libraries.

2 Traditions in the Field of Library Automation and the Digital Display of Holdings

The library traditionally follows and implements the latest trends in automation. 1979 is the year of the introduction of the library's first computer system IZOT-0310 and information search devices UPDML 9002-02 and IPU IZOT-0320. The first subscription to the AGRIS (FAO) database, stored on magnetic tapes, was done in 1980. A local library-information network was established in 1994.

The Digitization Centre was founded in 2008. Since then the library has participated in various projects, on national and international level, concerning the digitization and online display of its valuable holdings. Some noteworthy projects are Europeana Photography, EMBARK, BG08 "Digital Cultural and Historical Heritage of Plovdiv Municipality" and others. Digital copies of 95 Slavonic manuscripts from the 12th to the 18th Century may be accessed at the Manuscript collection of NALIS Repository, a prestigious academic digital library, founded by the Central Library of the Bulgarian Academy of Science, Sofia University "St. Kliment Ohridski" and the American University in Bulgaria.

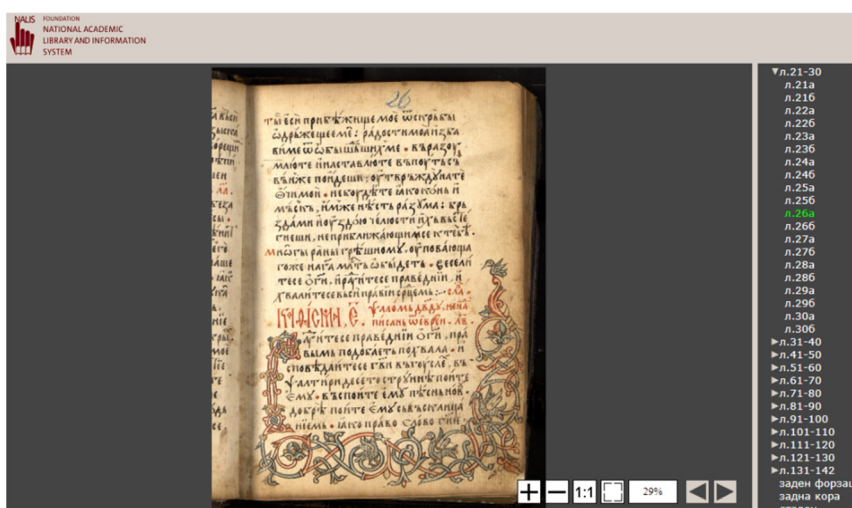


Fig. 1. An example of a page from Markovski psalter, 1638 (No. 5(207)) displayed in NALIS Repository

In 2017 its Digital Library becomes available online at <http://digital.plovdiv.bg/BG/Pages/LibIvanVazov.aspx>. It is a part of a web portal, which unites seven of the largest cultural institutions in Plovdiv. The library offers nine collections:

- BOOKS, which currently are 193, predominantly Statues of professional organizations from Plovdiv and the region, from the end of the 19th century and the beginning of the 20th century. They are of interest to researchers because

they provide insight into the way of life of the era. The plans are to continue uploading the most valuable and compelling holdings of the library.

- PERIODICAL PUBLICATIONS is a collection of newspapers and magazines currently containing 216 full-text titles, with approximately 20,000 separate issues. Particularly valuable are the newspapers and magazines from the East-Rumelian period: “Maritsa” newspaper - the first Bulgarian newspaper after the Liberation, “Narodniy glas”, “Nauka”, “Zora”, etc., as well as the Revival-period collection of periodical publications.
- MANUSCRIPTS - The rich collection includes Slavonic, Greek, Ottoman and Persian manuscripts on parchment and paper from the 11th to the 19th century. They are displayed with priority in NALIS Repository. The Digital Library currently displays 16 manuscripts, with the intention to upload more.
- GRAPHIC PUBLICATIONS is a collection of art prints, lithographs, etchings, engravings, posters, original paintings, etc. Especially interesting are the projects for monumental works – a total of 95 projects for large-scale wall murals, frescoes, ceramic tilework, most of them realized in many Bulgarian towns. The projects were created by renowned Bulgarian artist such as Dimitar Kirov, Yoan Leviev, Encho Pironkov.



Fig. 2. Yoan Leviev, Anna Grebenarova. Design for ceramic piece, 1966

- CARTOGRAPHIC PUBLICATIONS presents valuable possessions of the library such as the oldest map of Bulgaria, created by a Bulgarian and printed in

Bulgarian language - "Map of the present Bulgaria, Thrace, Macedonia and the adjacent lands".

- The PHOTOGRAPHS digital collection will be developed in the future. Currently 10 of them are displayed in the Digital Library. The library has a collection of approximately 4,000 photographs and postcards, portraits, events and sites of historical significance.



Fig. 3. Dimitri Ermakov. Plovdiv, Sahat Tepe, 1876

- The ARCHIVES DIGITAL COLLECTION will also be developed in the future. The Bulgarian historical archives in the National Library "Ivan Vazov" are of nationwide value, revealing key moments of the political, economic and cultural development of Bulgaria. The documents cover the period from the 12th to the 20th centuries, the most numerous being from the second half of the 19th century.
- In the future we will also develop our AUDIOVISUAL DOCUMENTS collection. The software platform of our Digital library allows us to publish audio and video file formats. The library has a rich collection of classic films on 35mm reels - masterpieces from the birth of cinema in 1895 to the early 40's of the XX century, including classic films by the Lumiere brothers, David Griffith, Charles Chaplin, Buster Keaton, Fritz Lang, John Ford, Dziga Vertov, Sergei Eisenstein, Luis Bunuel, Orson Wells and others. The library also has a large collection of recorded music with approximately 16,000 vinyl records.
- The library also has a rich collection of SHEET MUSIC PUBLICATIONS. Of particular interest is the manuscript "Bulgarska kitka", which may be seen in the Digital library. It was created in 1881 by the Czech composer Franz Schwestka especially for the needs of the newly established brass musical ensemble of Plovdiv.

In 2019 a new functionality was introduced in the Digital Library of National Library “Ivan Vazov”– indexing and searching in PDF files’ text contents. The machine-encoded text is obtained through the use of optical character recognition (OCR) software. By including the OCR step, full-text indexing and internal document search capability can be applied, making it easier for users to discover and use the materials. Digital-born PDF files do not need to be processed.

The main activities within the scope of the new functionality are as follows:

- Development of the software platform to upload and display PDF files with the possibility to search the contents of the file (if OCR had been implemented).
- One-time migration service for all existing collections in the Digital Library in order to replace the existing images in the platform with corresponding PDFs.
- Purchasing of ABBYY Finereader 14 software product for OCR and processing of PDF files.

It was important to work closely and exchange ideas with the software developer, so that our intentions could be realized as fully as possible and to adapt to the limitations of Microsoft’s software platform SharePoint, on the basis of which our Digital Library was created.

3 Developments of the Digital Library Necessitated by the Introduction of the Ability to do Content Search

The updated capabilities of the software platform required changes in the user interface in order to accommodate the new functionalities. New search fields were added and content search is applied to both global and collection search. Searching is done in the indexed contents of the PDF files and the provided metadata, and is augmented by the boolean operators supported by SharePoint 2010.

A help box is added to the search fields, which displays static text with short search instructions, as well as a link, leading to an external page with comprehensive guidelines.

Especially significant are the changes in the collection "Periodical publications", because of the necessity to display a multitude of search results from many issues. So far, this was not necessary because the search was done only with respect to the metadata of the title.

A new gallery tool for viewing and navigating PDF files was implemented, which has the following capabilities:

- Search by keyword in the contents of the file, mark the matches found with distinctive colour, display the number of matches and provide controls to move to matches.
- Suitable navigation of the pages of the document, with instruments such as thumbnail view of individual pages and a blank field where a desired page number may be written.
- Fullscreen view.
- Help menu with additional tools - navigate to the first or last page of the file, change the orientation of the page and hand-tool.

The administrative part of the software platform did not require a major overhaul and the most important addition was a section for PDF file upload in the metadata entry form.

The contract with the software developer includes also a one-time migration service of PDF files for all collections in order to replace the existing images in the platform. The purpose is to decrease the time needed for the replacement. The personnel of the Digitization Center are currently working on the preparation of the complete batch of files for the replacement. The name of each PDF file must include the unique ID number of the corresponding record in the Digital Library system, which it will replace.

4 Specifics of OCR and PDF File Processing

The primary “master” files, stored on the library’s servers and created for the purposes of long-term digital storage will be used for the purposes of OCR.

Predominantly, the master files are 24-bit color images, scanned at 300 ppi from the original paper source. Scanning from the originals is generally acknowledged to produce higher quality master images (Klijn, 2008). In this way, the degree of recognition will be approaching its maximum.

The master files will be processed to a single PDF file for each unit of cultural heritage. For example, a single PDF file should be a monograph or a newspaper issue. The resultant file will have appropriately lowered image quality, with a size suitable for online display. The PDF will have its images aligned with hidden machine-readable text, product of the recognition. MRC compression method will not be used.

After a thorough review of the available software, which included sampling the experience and opinions of libraries from Bulgaria, Russia, Ukraine, Serbia and other countries, where the recognition of Cyrillic text is of special relevance, it was decided to purchase a licensed version of ABBYY FineReader 14, which was the most widely used software to perform OCR and in our tests was the best at recognizing Cyrillic text. Other options were considered as well (such as Adobe Acrobat XI Pro) but the results achieved were not as satisfactory.

Using ABBYY FineReader 14 on master files obtained from well-preserved originals, written in modern Bulgarian language, yields very high OCR success rate, most often above 99%. However, texts that we have to deal with, those of high cultural and historical value and with expired copyright, are predominantly from the period before the Orthographic Reform of 1945. The accuracy of OCR software is language dependent: alphabet; old letters without the coding tables; old grammar, obsolete words, phrases and idioms; dictionaries; multi-lingual documents (Andreev & Kirov, 2009). At the same time, despite the care that has been taken by the library to preserve the originals, old textual documents present challenges to OCR, such as the natural darkening of the paper, faded print, in-library binding in proximity to the text, etc.

There are a number of ways to improve the accuracy of the recognized texts. A straightforward solution is to include post-OCR manual corrections as another stage to the digitization process. However, with the human efforts needed to correct OCR errors, it becomes quite a tedious job (Andreev & Kirov, 2009). Considering the large amount

of time necessary, to achieve high levels of accuracy (around 98%), the labour-intensive cleaning required to remove OCR errors means the two-step process may be no more efficient than manually inputting texts from scratch, a procedure that suits small- to medium-scale projects (Strange, McNamara, Wodak, & Wood, 2014). A way to mitigate this issue, while still using manual labour, is to design the digital library software platform in such a way that it involves the users of the resources and allows them to correct OCR mistakes. This solution was first implemented by the National Library of Australia in their newspaper digital collection. It is considered a successful practice (Holley, 2009) and has been incorporated by other institutions.

Another method of improving the quality of OCR, relevant to the software ABBYY FineReader, is the option to train the program to patterns, the interpretation of which the software deems uncertain. This is useful in cases of non-standard fonts and is especially important for Cyrillic texts, where the training to recognize specific letter symbols is essential. Such are the letters Ъ, Ѫ, ІѪ, А, ІА, etc., which were gradually removed from the modern written language, eventually reducing the number of letters in the alphabet to the current 30.

When reviewing the practices of libraries in Bulgaria, which offer PDF files with recognized text in their digital libraries, it is evident that there is no uniform standard for the ways in which OCR of texts before the Orthographic Reform of 1945 is performed. For instance, some of the reviewed libraries were tempted to replace the archaic letter symbols with their modern equivalents, which is done in order to aid the search of the contents of older texts, so that users would not have to write the required search expression twice – in the old and new spelling. However, this cannot be considered a good practice, because the spelling conversion is not fixed in the sense that old letter symbols are often replaced by more than one modern letter symbol. For instance, the letter “Ъ” is replaced by modern “Е” or “Я” and “Ѫ” is replaced by modern “Б” or “А”.

There have been collaborative work with our partners from the project CLaDa -BG to develop the best methodology for optical character recognition and the consequent methods of searching in the text. The library’s role would be in providing texts from its rich holdings, from different periods, with different fonts, formats, etc., and also to test and apply the developed resources, such as thesauri, search-engine-complementing instruments and others.

5 Conclusion

The new file format will lead to changes in the work cycle of the Digitization Centre, which in terms of web presentation focused exclusively on the .jpg image format. An important part of the future digitization work will be the mastery and long-term establishment of OCR-related activities, aiming to ensure the highest possible level of resource usability. Therefore, it is important to share digitization experience with other partners, working in the same field, with the aim to form a comprehensive strategy and collaborative solutions.

Involvement in developing tools to automate the submission of data to Europeana is of top priority as well. Stronger presence there is mandatory, but the great amount of manual work disrupts our potential to contribute in a meaningful and visible way.

References

- Andreev, A., & Kirov, N. (2009). Hausdorff distances for searching in binary text images. *Serdica : Journal of Computing*, 3 (1), 24.
- Holley, R. (2009). Many Hands Make Light Work: Public Collaborative OCR Text Correction. *Australian Historic Newspapers*.
- Klijn, E. (2008). The current state-of-art in newspaper digitization : a market perspective. *D-Lib Magazine* 14(1/2).
- Strange, C., McNamara, D., Wodak, J., & Wood, I. (2014). Mining for the Meanings of a Murder: The Impact of OCR Quality on the Use of Digitized Historical Newspapers. *Digital Humanities Quarterly*, 8 (1).

Received: July 01, 2019

Reviewed: July 10, 2019

Finally Accepted: July 23, 2019