



Embedding transcription and transliteration layers in the Digital Library of Polish and Poland-related News Pamphlets

Maciej Ogrodniczuk, Włodzimierz Gruszczyński

Institute of Computer Science / Institute of Polish Language
Polish Academy of Sciences

The 23rd International Conference on Asia-Pacific Digital Libraries (ICADL 2021)
December 2, 2021

Our experiment

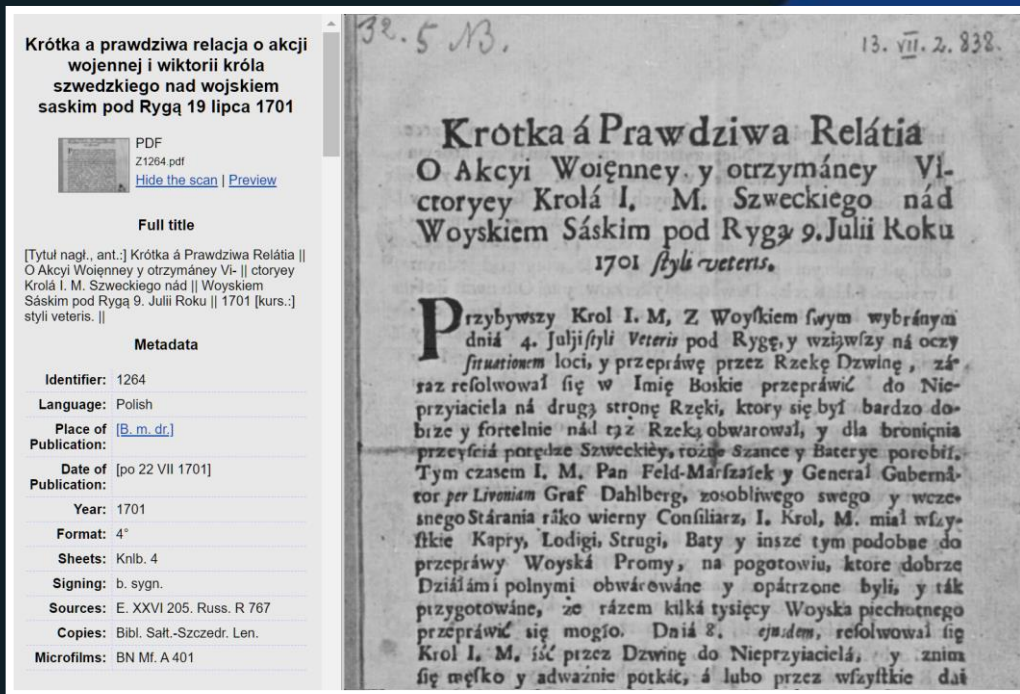
In a nutshell:

- old Polish prints are made available in an existing digital library as graphical PDFs (scanned images of prints)
- an electronic corpus of old Polish texts contains texts from the prints in two variants: transcription and transliteration
- a corpus search engine is available but we need to be able to search also in the digital library
- we can easily search in metadata
- but we would also like to search in both variants inside graphical PDFs

The Digital Library of Polish Prints from the 16th–18th Centuries

The source of images:

- a thematic DL with over 2 000 prepress documents
- managed by Eprints
- soon 200+ items to be extended with texts
- available at <https://cbdu.ijp.pan.pl>



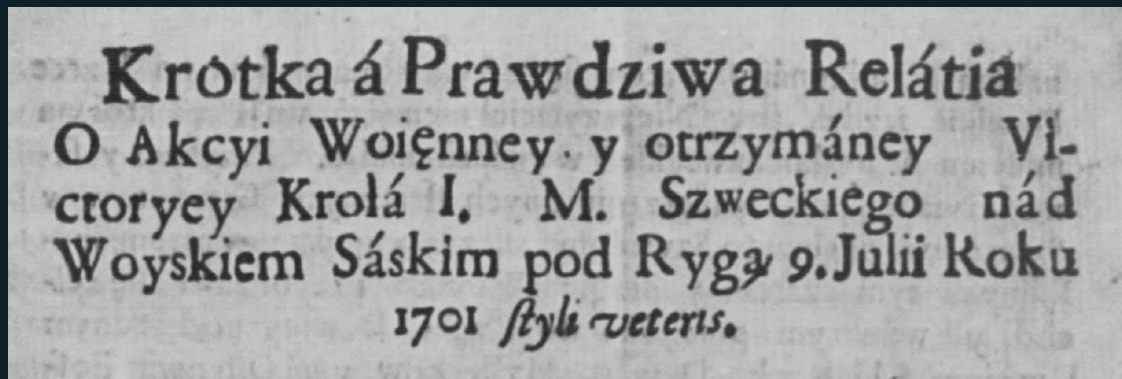
The Electronic Corpus of the 17th and 18th Century Polish Texts

The source of texts:

- nearly 13.5 million words
- TEI-encoded texts following the format of the National Corpus of Polish
- in both transcribed and transliterated forms
- containing 40 prints from the digital library of old prints – ready to be used by the digital library

```
<fs type="morph">  
  <f name="orth">  
    <string>nazajutrz</string>  
  </f>  
  <f name="translit">  
    <string>nazaiutrz</string>  
  </f>  
</fs>
```

Transcription vs. transliteration: the differences



Krótká á Prawdziwa Relátia
O Akcyi Woïenney y otrzymáney
Victoryey Krolá I. M. Szweckiego
nád Woyskiem Sáskim pod Rygá
9. Julii Roku 1701 styli veteris.

Krótká a prawdziwa relacja
o akcji wojennej i otrzymanej
wiktoriej Króla I. M. Szwedzkiego
nad wojskiem saskim pod Rygą
9. julii roku 1701 styli veteris.

Transcription vs. transliteration: corpus search

Query

[translit="nazaiutrz"]

Displayed layer

transliterated

á

é

Á

É

Left context

Result

Right context

Text ID

Date

M. Pan Woiewodá Krákowski
Hetman Polny cum eadem
apparentia

nazaiutrz
[\[nazajutrz:adv\]](#)

Die[...] ma praesentis wiachał.
Woysko do Soboty przeszły
stało

AwLwow 1693

I. W. I. M. Pan Woiewodá Krákowski Hetman Polny cum eadem apparentia **nazaiutrz** Die[...] ma praesentis wiachał.
Woysko do Soboty przeszły stało pod Báriszem, w Niedzielę miało się daley ruszyć zá lázłowiec ku Wasiłowu,
co ieśli się stało czekamy in momentá wiadomości.

Transcription vs. transliteration: where to store them in the DL?

Krótká a prawdziwa relacjá o akcji wojennej i wiktórii króla szwedzkiego nad wojskiem saskim pod Rygą 19 lipca 1701

Full title

[Tytuł nagł., ant.:] Krótká á Prawdziwa Relátia || O Akcyi Woiénney y otrzymaney Vi- || ctoryey Krolá I. M. Szweckiego nád || Woyskiem Sáskim pod Rygą 9. Julii Roku || 1701 [kurs.:] styli veteris. ||

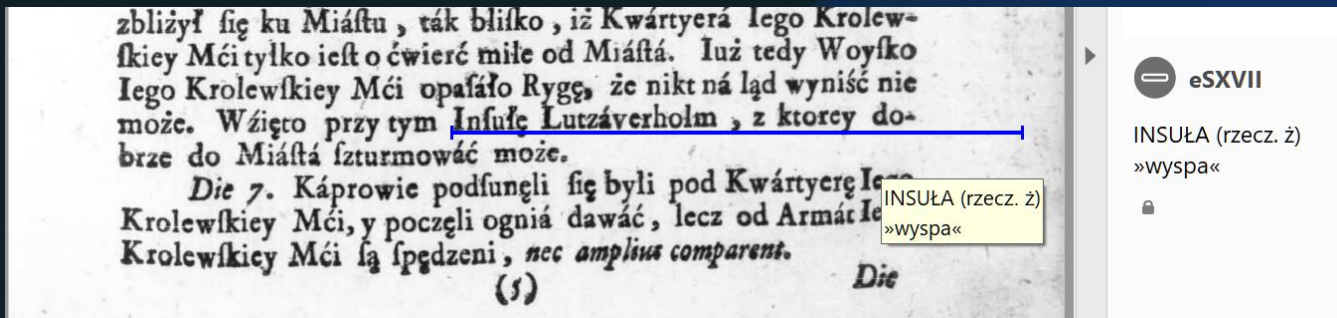
Translation into Contemporary Polish

Transkrypcja: Krótká a prawdziwa relacjá o akcji wojennej i otrzymanej wiktórii Króla I. M. Szwedzkiego nad wojskiem saskim pod Rygą 9. julii roku 1701 styli veteris. Przybywszy Król I. M. z wojskiem swym wybranym dnia 4 julii styli veteris pod Rygę i wzięwszy na oczy situationem loci i przeprawę przez rzekę Dźwinę, zaraz resolwował się w Imię Boskie przeprawić do Nieprzyjaciela na drugą stronę rzeki, który się był bardzo dobrze i fortelnie nad tąż rzeką obwarował y dla bronienia przejścia potędze szwedzkiej różne szańce i baterie porobił. Tymczasem I. M. Pan Feldmarszałek

Back in 2019

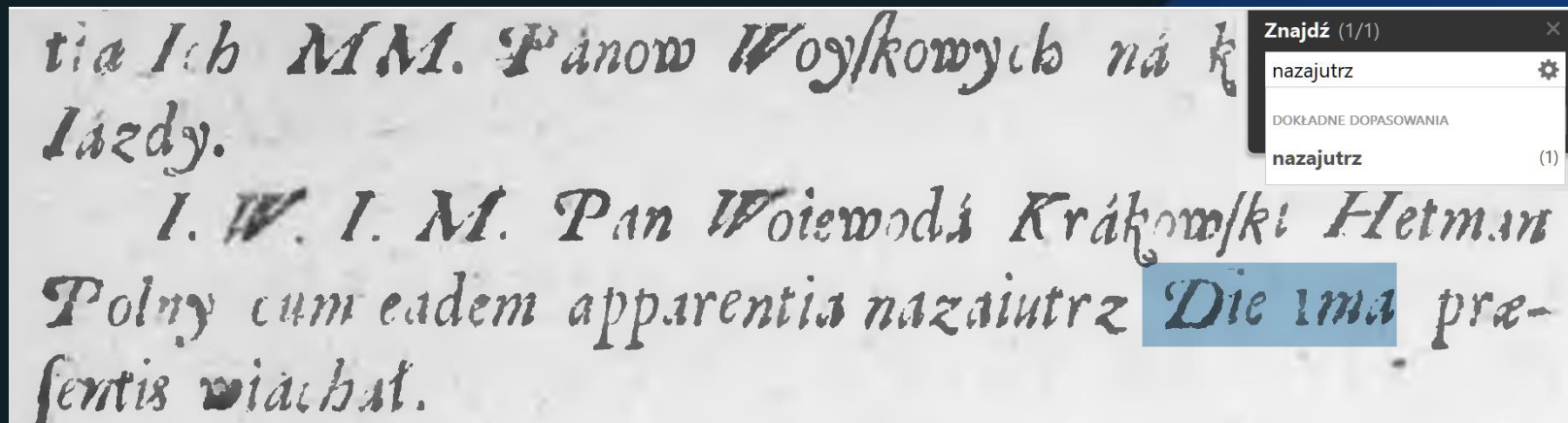
Adding texts from the corpus to PDF:

1. obtaining a dirty OCR of a print
2. retrieval of corpus transliterations
3. token alignment and fine-tuning
4. token replacement and creation of a clean version of the print
5. adding dictionary annotations



Experiments in adding word variants

1. adding word variants directly in the single hidden textual layer
→ wrong positioning of text



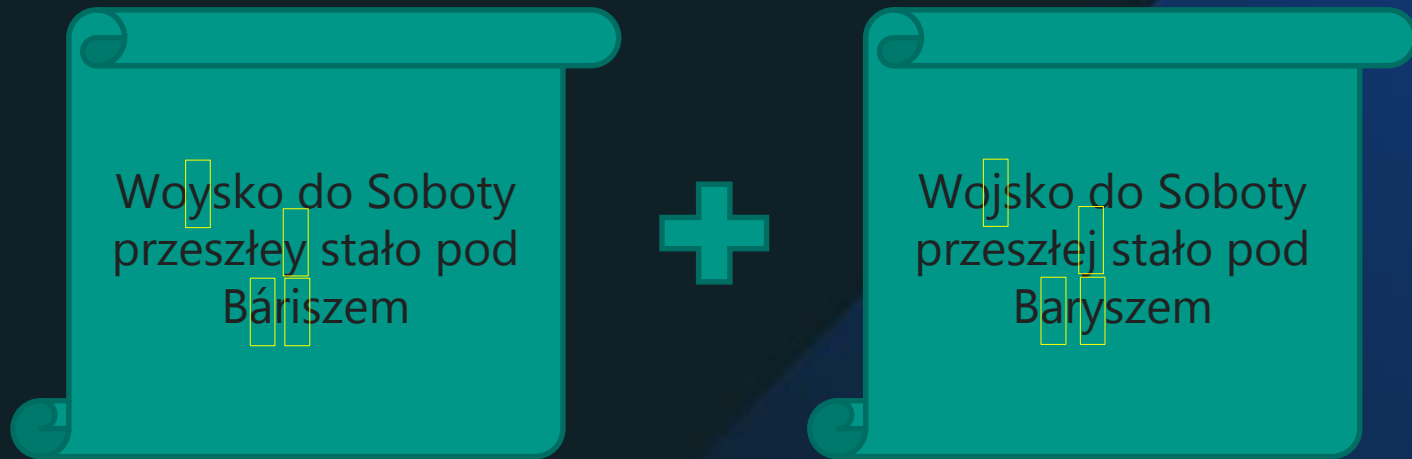
So what about...

1. checking whether it is possible to store more than one layer of the text in the PDF
2. testing which tools are capable of merging both layers
3. checking whether the most common working environment (for our users) can make use of that



Experiments in adding word variants

2. merging two hidden textual layers to make them concurrently searchable



How does it work?

The tools used for merging both layers:

- CPDF
- PDFTK

The tools used to search in both layers:

- standalone Acrobat Reader
- Acrobat Reader plugin for Edge, Firefox and Chrome

Coherent PDF Command Line Tools and C/Python API Community Release

Powerful, free tools to manipulate PDF files

[View on GitHub](#)



[Current version: 2.4 (21st June 2018)]

The Coherent PDF Command Line Tools can be used to manipulate PDF files in a variety of ways. For example:

- Merge PDF files together, or split them
- Encrypt and decrypt
- Scale, crop and rotate pages
- Read and set document information
- Copy, add or remove bookmarks
- Stamp logos, text, dates, page numbers
- Add or remove attachments
- Losslessly compress PDF files



PDFtk the pdf toolkit

PDFtk is a simple tool for doing everyday things with PDF documents. It comes in three flavors: *PDFtk Free*, *PDFtk Pro*, and our original command-line tool *PDFtk Server*.

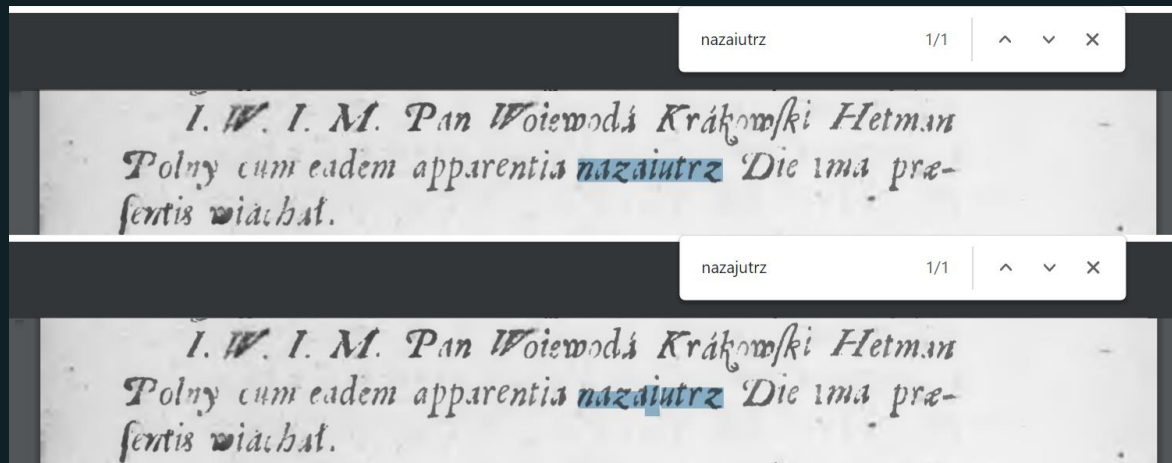
PDFtk *Free*

PDFtk Free is our friendly graphical tool for quickly **merging and splitting PDF documents and pages**. It is free to use for as long as you like.

Power Users: PDFtk Free comes with our command-line tool, *PDFtk Server*. So you get both the GUI and the command-line interface to PDFtk!

The result

1. transliterated and transcribed text is searchable in the browser and offline
2. position of the text is properly shown:



3. only one layer is available for copying text

Thank you!