



Connecting Data for Digital Libraries: The Library, the Dictionary and the Corpus

Maciej Ogrodniczuk

Linguistic Engineering Group
Institute of Computer Science
Polish Academy of Sciences

Włodzimierz Gruszczyński

Institute of Polish Language
Polish Academy of Sciences

ICADL 2019

6 November 2019

Interconnecting related language resources

Repositories of related historical language data usually exist independently:

- digital libraries store source objects
- language corpora maintain their textual content
- electronic dictionaries provide linguistic context of individual words/phrases

But how about:

- inter-linking the repositories
- using the content from the corpus and the electronic dictionary to enhance digital library objects

The resources used in the study

The Digital Library:

- CBDU: The Digital Library of Polish and Poland-Related Ephemeral Prints from the 16th, 17th and 18th Centuries
- a thematic digital library of approx. 2,000 Polish and Poland-related pre-press documents
- dated between 1501 and 1729
- managed by EPrints
- planned to contain rich metadata (historical comments, glossaries of foreign interjections, explanations of background details and relations between library objects: translations, adaptations, alterations of the base text, their sources etc.)
- prints available only in image form (PDF files containing scanned originals)

<https://cbdu.ijp.pan.pl>

CBDU: Basic metadata

Relacja potrzeby i wiktoryi
otrzymanej pod Wiedniem
przez wojska króla Jana III



Full title

[Tytuł nagł., ant.:] RELACYA
Potrzeby która trwała godzin
14. y Wiktoryey otrzymaney [kurs.:] Die 12.
Septēbris. [ant.:] nād Nieprzyjacielem pod
Wiedniem, przez Woyskǎ Naiásnięzszego
y Niezwyciężonego KROLA Jego
Mości Polskiego Wielkiego Monar-
chy JANA Trzeciego, z pod Na-
miotow Węzyrskich z Obozu wysła-
na, tudzież Excerpta z Listu tegoż
Naiásnięzszego króla Jego Mości
do Krolowey Iey Mości piśanego,
[kurs.:] sub Die 13. Septēbris. Anno
Domini, 1683.

Metadata

Identifier: 1036

Language: Polish

Connections: War: B. Jak war. A z
wyjątkiem różnicy w
tytuł; w. 8 od góry
zamiast JANA-IDNA. Por.
poś 1095. 1096. 1097.
1689.

Place of
Publication: [B. m. dr.]

Date of
Publication: [po 13 IX 1683]

Year: 1683

Variant: A

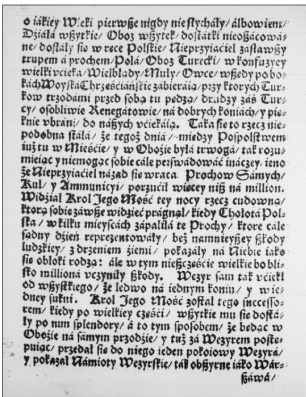
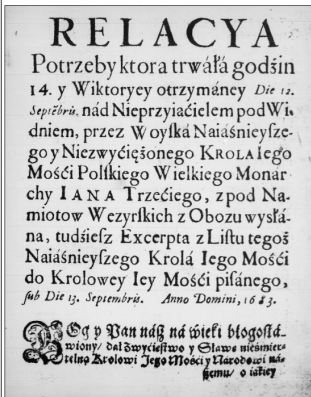
Format: 4^o

Sheets: Kniha. 4

Signing: Sygn. J⁴

Sources: E. XXV/1 204, XXV/III 354.
GK. 582

Copies: B1 24730 1, BN XVII. 3.



CBDU: Extended metadata

Historical Commentary

Opis dotyczy kampanii króla Szwecji Karola XII w Inflantach we wstępnej fazie tzw. Wojny Północnej (1700-1721) między Szwecją a koalicją Rosji, Danii, Saksonii, Prus i Hanoweru oraz od 1704 Rzeczypospolitej. Opisywane zwycięstwo szwedzkie nad połączonymi siłami sasko-rosyjskimi było ukoronowaniem tej kampanii, rozpoczętej nieudaną próbą zdobycia Rygi przez wojska saskie IX 1700. Wspomniana w tekście bitwa pod Narwą miała miejsce 20 XII 1700 i zakończyła się spektakularnym zwycięstwem Szwedów. Dokument ma niewątpliwie charakter propagandowy i skierowany był do zwolenników Karola XII w Rzeczypospolitej.

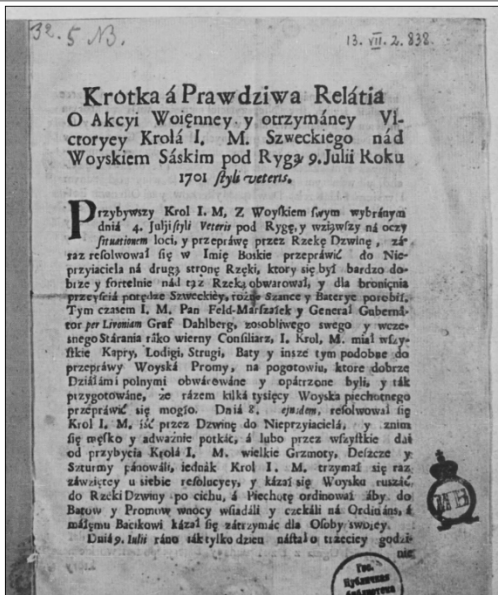
Dictionary

Entry	Explanation
wziąć na oczy	rozpatrzyć się, przyjrzeć się
resolwować się	zdecydować się, powziąć zamiar
fortelnie	rozważnie, sprytnie, pomysłowo
kapra	rodzaj łodzi

Translation into Contemporary Polish

Transkrypcja:

Krótką a prawdziwą relacją o akcji wojennej i otrzymanej wiktoryi Króla I. M. Szwedzkiego nad wojskiem saskim pod Rygą 9. Julii roku 1701 stylu veteris.



The resources used in the study

KORBA (Baroque) Corpus:

- The Electronic Corpus of the 17th and 18th Century Polish Texts (until 1772)
- over 13M tokens
- diachronic supplement to the National Corpus of Polish
- rich structural annotation (e.g., the location of the search phrase in the appropriate text page)
- text annotation (e.g., tagging words from foreign languages)
- 0.5M subcorpus manually annotated with morphosyntactic tags, then used as training data for a disambiguating tagger, then used for automatic annotation of the whole corpus

<https://korba.edu.pl>

ELEKTRONICZNY KORPUS TEKSTÓW POLSKICH Z XVII I XVIII W. (DO 1772 R.)

Korpus

Korba automatyczna ▾

Zapytanie

potrzeba

á

é

Ä

Ë

KONSTRUKTOR ZAPYTAŃ

ODRZUĆ OBCE SEGMENTY

METADANE ▾

Warstwa wyświetlania
uwspółcześniona ▾Liczba wyników na stronę
10 ▾

Wyszukaj

Znaleziono 500 wyników.

Lp	Lewy kontekst	Rezultat	Prawy kontekst	Skróć	Data
1	oñcjaliſtów dozór, rekomenduję i upominam. 16-to. Czeladzi	potrzeba [potrzeba:subst:sg:nom:f]	na dole ująć, żeby ich tak wiele nie zjeżdżało	InsGór_3	między 1706 a 1743
2	potrzebować nie będą, tedy i z nich rachunek oddać	potrzeba [potrzeba:subst:sg:nom:f]	. 13. Po dystrybucie zwyż mianowanych 4313 beczek na	InsGór_3	między 1706 a 1743
3	druzdy robotnicy, aby na dowołaniu byli, kędy największa	potrzeba [potrzeba:subst:sg:nom:f]	będzie, ile że jako obywatel w takim razie dopomóc	InsGór_3	między 1706 a 1743

The resources used in the study

The Electronic Dictionary:

- The Electronic Dictionary of the 17th–18th Century Polish
- started to be published on paper in 1999
- converted to electronic form in 2004 and further developed online
- currently 39k entries in various stages of development
- 77k grammatical forms in the entries (only confirmed inflectional paradigms)
- quotations transliterated from 1100 sources

<http://sxvii.pl>

Electronic dictionary



o słowniku kwerendy

ELEKTRONICZNY SŁOWNIK JĘZYKA POLSKIEGO XVII I XVIII WIEKU

znajdź hasła

zaczynające się od



szukaj



☐ przeszukaj także hasła w indeksie

☒ a fronte ☐ a tergo

A B C Ć D E F G H I J K L Ł M N O Ó P R S Ś T U V W X Y Z Ź Ż

POTRZEBA I, POTRZEB

rzecz.

ż

W TRAKCIE OPRACOWANIA



> Notowanie w słownikach

> Formy gramatyczne

Resource linking: From corpus to dictionary

KORBA

INFORMACJE - INSTRUKCJA LOGOWANIE

①

Zapytanie

[base="skrzydło"]

②

podniósł [skrzydła](#) [[skrzydło:subst:pl:acc:n](#)] namiotu



③

SKRZYDŁO

Część mowy: rzecz.

Warianty fonetyczne: SKRZYDŁO, KRZYDŁO

Znaczenia:

1. »jedna z dwóch kończyn ptaka służących do latania«
2. »boczne oddziały wojska ustawionego w szyku bojowym«

[Odnosnik do słownika](#)



Resource linking: From dictionary to corpus

①

SKRZYDŁO

rzecz.

n

Warianty fonetyczne: SKRZYDŁO, *KRZYDŁO

> Notowanie w słownikach

> Formy gramatyczne

> Znaczenia

Więcej cytatów w Korpusie Barokowym ■■



②

Zapytanie

(({base="skrzydło"&pos="subst"}))(({base="krzydło"&pos="subst"}))

Resource linking: From dictionary to corpus



o słowniku kwerendy

①

SKRZYDŁO

rzecz.

n

Warianty fonetyczne: SKRZYDŁO, *KRZYDŁO

▼ Formy gramatyczne

lp	M.	skrzydło	■ ■
	D.	skrzydła	■ ■
	C.	skrzydłu	■ ■
	B.	skrzydło	■ ■
	N.	skrzydłem	■ ■
	Ms.	skrzydle	■ ■
lm	M.	skrzydła	■ ■
	B.	skrzydła	■ ■
	Ms.	krzydłach	■ ■ skrzydłach ■ ■



②

Zapytanie

[orth="krzydłach"&pos="subst"&number="pl"&case="loc"]

Experiment 1: From corpus to digital library

Initially:

- CBDU prints available only as graphical PDFs
- unsuccessful attempts at dirty OCR

But:

- several prints were transliterated to be included in the corpus
- we could use this content to add the textual layer

Experiment 1: From corpus to digital library

Conversion workflow:

- 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

Experiment 1: From corpus to digital library

OCR quality for a sample page

	Token count
Actual number of textual tokens on the page	126
Number of tokens in dirty OCR	136
Number of properly OCR-ed tokens	22
Number of tokens in the corresponding corpus text	125
Number of correctly aligned tokens	124
Number of correctly transferred tokens	108

Experiment 1: From corpus to digital library

Transfer errors caused by:

- the quality of the source prints
 - noise in the source image
 - text bleeding through from the reverse page
- overlap of regions and tokens

Experiment 2: Dictionary annotation

The general idea:

- linking from metadata is easy but inconvenient to use
- manual linking requires an enormous amount of work
- dictionary entries could be used
- but some words are still understandable
- so how about filtering them with the contemporary dictionary?

Experiment 2: Dictionary annotation

Annotation workflow:

- 1 Filtering out modern vocabulary from e-SXVII using PoliMorf dictionary

Subset of the e-SXVII vocabulary	Entries
Meaningful dictionary entries	15 702
Lemmata present in the corpus	10 440
Lemmata older than 20th century	2 652
— with gerund correction	2 479

- 2 Creating dictionary annotations

Experiment 2: Dictionary annotation

fię Wojsko, Kawalerią y Intanterią która ielzcze po wieczce przeszley nie weszła była, wprowadził do Miastá; potym rzęsiłto z Dział strzelac zaczął, tak iż kilka koni pod temi zabito, którzy byli przy Krolu I. Mści.

Die 6. ejusdem. Ruszył Krol I. Mści z całym Wojskiem, y zbliżył się ku Miastu, tak blisko, iż Kwartyerá Iego Krolewskiej Mści tylko jest o ćwierć mile od Miastá. Iuż tedy Wojsko Iego Krolewskiej Mści opasało Rygę, że nikt ná ląd wynisć nie może. Wzięto przy tym Insulę Lutzáverholm, z ktorey do brze do Miastá szturmować może.

Die 7. Káprowie podsunęli się byli pod Kwartyerę Iego Krolewskiej Mści, y poczęli ognia dawać, lecz od Armáe Iego Krolewskiej Mści są spędzeni, *nec amplius comparent.*

(5)

Die

eSXVII

INSULÁ (rzech. ż)
»wyspa«



INSULÁ (rzech. ż)
»wyspa«

Experiment 2: Dictionary annotation

Observations for a sample print (29 entries in the manually created dictionary):

- 21 entries cannot be suggested:
 - 14 entries are absent from e-SXVII
 - 5 entries are multi-word
 - 2 entries are stubs
- 3 entries are unexplained in the manual dictionary
- 4 entries can be found in the modern dictionary
- one entry is valid to be suggested by our mechanism

Conclusions

And future work:

- the idea can be applied to digital libraries using other layered representation formats such as DjVu
- the proposed tools can be replaced
- linking individual content objects with external resources is independent from any technical solution used for digital library management
- was it worth it? yes! digital libraries play the central role in linguistic research!

Thank you!

And the funding institutions:

- The work was financed by a research grant from the Polish Ministry of Science and Higher Education under the National Programme for the Development of Humanities for the years 2019–2023 (grant 11H 18 0413 86)
- The conference fee was sponsored by the grant holder, the Institute of Polish Language, Polish Academy of Sciences
- The travel to ICADL 2019 was sponsored by the Institute of Computer Science, Polish Academy of Sciences