

# Enhanced digital editions: retrieving POS tags from pre-digital word indexes

Joanna Byszuk, Magdalena Król, Maciej Eder



EADH 2018 Galway



**Why?**

# **Biblioteka Pisarzy Polskich – A Collection of Polish Writers: context**

# About Biblioteka Pisarzy Polskich

Two series: A and B:

34 volumes in B series, covering literary heritage of most important Old and Middle Polish writers

Linguistic edition by best scholars, printed for over 70 years, now forgotten or really unique (thus, highly priced)

Consequent framework supported by years of textual analysis

# About Biblioteka Pisarzy Polskich – construction

- images of the original sources
  - a philological transcription and a commentary
  - comprehensive indexes of all the words used in the edited texts:
    - their grammatical tags,
    - the information about inflection (if applicable),
    - precise location of the form in question in the original text
- > complete sets of POS tags identified manually by cutting-edge scholars in Polish philology

## Posłow Graeckich.

Jesli namniej przewinyl/ być mu w okowách.  
Wiec ja podobno z mniejszym niebezpieczeństwem  
Grzesze: bo sam sie trące swym wszeteczeństwem.  
Przełożonych występy miastá zgubiły/  
I szerokie do gruntu Cárstwa zniszczyły.

POSEŁ. HELENA.

PO. Dobra nowinę pániey swojej niosę/  
Rozumiem temu/ że już dawno tego  
Poselstwa czeka/ serce swe troskami  
I płączęm trapiąc: Ale oto prawie  
Ná czas wychodzi z domu. O Krolowa/  
Wdzięczney nowiny posła masz przed soba.

Hel. Day Boże/byś co przyniosł pociesznego.

PO. Posłowie twoi jáko przyjecháli/  
Ták odjeżdżają / á ty przedsie známi.

Hel. Byłeś sam w rádzie / czyś słyszał od kogo?

PO. Byłem przy wszystkim: y prosto mi stámtąd  
Jść ALEKsánder do ciebie rozkazał.

Hel. Jeszczeć niewidzę/ szczegobych sie prawie  
Ucieszyć miátá/wszakże powiedz przedsie  
Jáko co byto. PO. Powiem/ieno słuchay.

Skoro w rádzie zásiedli pánowie/Krol naprzód  
Te rzecz do nich uczynił. Niezwyktem nic nigdy  
Bez rády wászey czynić: ábych też zwykł kiedy/  
(czego w pámięci niemam) w tej sprawie koniecznie  
Syná swego bych niechciał: áby mie ojcowska

B iij

Miłość



180 Jesli namniej przewinł – być mu w okowách.  
Więc ja podobno z mniejszym niebezpieczeństwem Grzeszę, bo  
sam się trąc swym wszeteczeństwem. Przełożonych występy  
miastá zgubiły  
I szerokie do gruntu cárstwa zniszczyły. 5

POSEŁ – HELENA

PO.: Dobrą nowinę pániey swojej niosę.  
Rozumiem temu, że już dawno tego Poselstwa czeka, serce swe  
troskami

I płączęm trapiąc. Ale oto prawie 10

185 Ná czas wychodzi z domu. O, krolowa, Wdzięcznej  
nowiny posła masz przed sobą.

HEL.: Daj boże, byś co przyniosł pociesznego.

PO.: Posłowie twoi jáko przyjecháli,  
Ták odjeżdżają, á ty przedsie z námi! 15

190 HEL.: Byłeś sam w rádzie, czyś słyszał od kogo?

PO.: Byłem przy wszystkim i prosto mi stámtąd Iść Aleksánder  
do ciebie rozkazał.

HEL.: Jeszczeć nie widzę, z czego bych się prawie Ucieszyć miátá,  
wszakże powiedz przedsie 20

195 Jáko co byto. PO.: Powiem, jeno słuchaj.

Skoro w rádzie zásiedli pánowie, krol naprzód

Tę rzecz do nich uczynił: „Nie zwykłem nic nigdy Bez rády wászej  
czynić; á bych też zwykł kiedy

(Czego w pámięci nie mam), w tej sprawie koniecznie 25

200 Syná swego bych nie chciał, áby mie ojcowska

B3 Miłość

# Why is grammar tagging of old texts hard?

Difficult disambiguation of word sense:

Big number of word novelties:

- atypical word formation,
- improper inflection adjusted to rhyme and melody of a poem,
- fluency with morphology – free switching between categories

# Why is grammar tagging of old texts hard?

Non-trivial grammatical problems:

- anaphors,
- elipsis,
- non-neutral word order (used for different purposes)

Non-obvious cultural contexts in proper names, unknown properties of common nouns, aspects of verbs etc.



# **Contemporary applications of linguistic resources from old editions**

# Contemporary applications of linguistic resources (not only) from old editions

Unflattening existing critical editions into digital ones

Diachronic investigations

# Contemporary applications of linguistic resources (not only) from old editions

- Possible training set for automatic tagging of historical texts (yields a large set of manually annotated data)
- More precise description of grammatical trends (instead of descriptions of linguistic phenomena) – quantitative input needed in historical data
- An enhancement for Universal Dependencies

# Why Kochanowski to start with?



# Jan Kochanowski

A writer with good ear for dialects and spoken language, using a sublime and sophisticated Polish at the same time:

- mixed registers,
- interesting word formation,
- playing with the sound and the sense of words

Wide range of topic covered in his works.

# Jan Kochanowski's heritage

A strong influence on the development of Polish literature.

Visible presence in modern culture and education.

Multiple resources are now prepared, including:

- dictionaries,
- critical editions,
- indexes,
- theoretical descriptions.

# Odprawa posłów greckich

- Manuscript created around 1577,
- Many editions preserved, until 1877 printed in collected works (at least 15 different versions)
- Written for the wedding celebrations of Jan Zamoyski and Krystyna Radziwiłłówna
- A typical example of Kochanowski's works: Mix of archaic and innovative vocabulary and structures

# Unflattening print critical editions



# Unflattening print critical editions – challenges

- Editions printed between 1889 (1952) and now
  - obtaining them
- Cracking the case(s) open
  - From a printed index to a usable structure

ROZBOJCA (f) *sb m – 1g N*: Rozbojca C4v/5. G: v rozbojce Dv/26. A: rozboycę C4v/4. entry  
ROZDRAPAC (f) *vb pf – inf*: Włcy mieli rozdrapać D2v/20. frequency  
ROZESLAĆ (f) *vb pf – imp 2 sg*: Ispiegł rozelli D/3. morphosyntax  
ROZEZNAWAĆ (f) *vb impf – praet 3 sg m*: Vznawca twarzy rozmawiał C4/15. example  
ROZGA (f) *sb f – 1g N*: rozgą wyrolił D2/7. location  
ROZGNIEWAĆ (f) *vb pf – praet 3 sg m*: rozgniewał C/7.  
**ROZKAZAĆ** (5) *vb pf – inf*: Rozkazać służyć A4/22. *praet 3 sg m*: Iść ... rofkazał B3/18. *plusq 1 sg m*: iam był rofkazał D2v/18. *imp 2 sg*: Rołkaż być pogotowiu D/2. *con praet 2 sg m*: był ... / ... rofkazał C4v/13.  
ROZKOSZ (z) *sb f – pl N*: rołkożyć nałze (f) / Niepewne Bv/20. G: oni by rołkoży trwałitych vitywali A4v/26.  
ROZLEGAC SIE (f) *vb pf – praet 3 sg m*: Isept ... / Rozlegal (f) sie po iak B4v/7.  
ROZMAITY (f) *ai – pl G*: lidibá chorob rozmaitych B2/11.  
ROZMYSL (z) *sb m – 1g N*: przemy rozmył D/10. G: w tym rozmyłu trzebá B3v/8.

# Challenge: getting the data

- Availability / Copyright

OCR problems:

- Various print
- Various state of preservation
- Typesetting – e.g. italics matter, delimitation hard to catch automatically

ROZBOJCA (f) *sb m – 1g N: Rozbojca C4v/5. G: v rozbojce Dv/26. A: rozboycę C4v/4.* entry  
ROZDRAPAC (f) *sb pf – inf: Włcy mieli rozdrapć D2v/20.* frequency  
ROZESLAĆ (f) *sb pf – imp 2 sg: łpięgi rozelli D/3.* morphosyntax  
ROZEEZNAWAĆ (f) *sb impf – praet 3 sg m: Vznawca twarzy rozmawiał C4/15.* example  
ROZGA (f) *sb f – 1g N: rozgą wyrolił D2/7.* location  
ROZGNIEWAĆ (f) *sb pf – praet 3 sg m: rozgniewał C/7.*  
ROZKAZAĆ (f) *sb pf – inf: Rozkazać flużyć A4/22. praet 3 sg m: Iść ... rozkazał B3/18. plusq 1 sg*  
*m: iam był rozkazał D2v/18. imp 2 sg: Rołkaż być pogotowiu D/2. con praet 2 sg m: był ... / ...*  
*rozkazał C4v/13.*  
ROZKOSZ (f) *sb f – pl N: rozkoiżyć naite (f) / Niepewne Bv/20. G: oni by rozkoiży trwiliżych*  
*stowali A4v/26.*

# Unflattening print critical editions – steps

1. OCR and cleaning
2. Identifying elements in each entry
3. Extracting elements
4. Verifying index of elements
5. Building tags
6. Tagging text
7. Verifying accuracy & comparing with automatic taggers

# Conclusions

# Conclusions

Possibilities:

- Easier verification of text tagging.
- Tagging independent of transliteration and transcription.
- With more ready editions we can unify grammar description within the whole corpus.

Thank you!