# Tagset adaptation to language changing over time. The case of the Electronic Corpus of the 17th- and 18th-century Polish Texts

Aleksandra Wieczorek
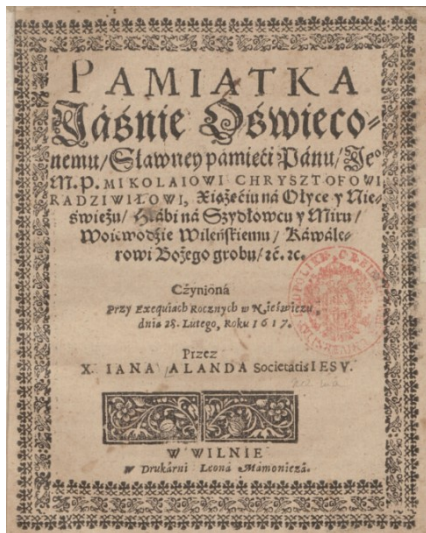
# THE ELECTRONIC CORPUS OF 17TH- AND 18TH-CENTURY POLISH TEXTS

- Cryptonym: KORBA (KORpus BARokowy 'baroque corpus')
- Principal investigator: Włodzimierz Gruszczyński
- Team: members of Institute of Polish Language and Institute of Computer Science, Polish Academy of Sciences

- Funding: Polish Ministry of Science and Higher Education, National Programme for the Development of Humanities grant (project number 0036/NPRH2/H11/81/2012 and 0413/NPRH7/H11/86/2018)
- Project duration: 2013-2018 and 2019-2023
- Coordinating body: Institute of Polish Language, Polish Academy of Sciences
- Cooperation: Institute of Computer Science, Polish Academy of Sciences
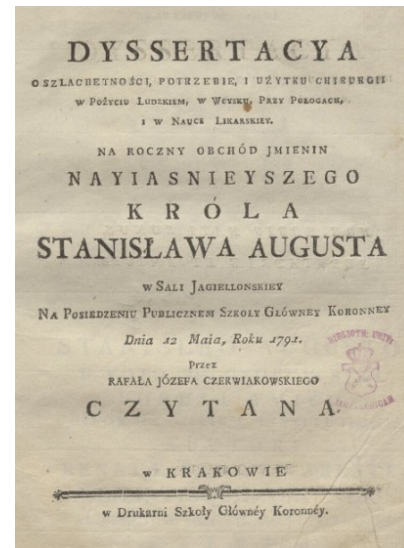
www.korba.edu.pl

PAN IJP Institute of Polish Language Polish Academy of Sciences

# THE ELECTRONIC CORPUS OF 17TH- AND 18TH-CENTURY POLISH TEXTS



1617

- 13.5M → 25M tokens
- richly annotated
- usage: research of grammar and lexis of 17th–18th centuries
- 1601–1800 (Baroque, Enlightenment)



1791

# Changes in Polish grammar during 17th and 18th c.

- Appearance of new grammatical categories or their values (e.g. of "masculine-personality" category)

- Disappearance of some grammatical categories or their values (e.g. disappearance of dual number)

- Changes of inflectional paradigms (e.g. of numerals)

- Changes of inflectional endings

- …

# 1. The new category: masculine-personality

# Grammatical cathegory of gender in Polish — nouns

## m (masculine)

pan 'gentleman'

lew 'lion'

dom 'house'

## f (feminine)

królowa 'queen'

żyrafa 'giraffe'

harfa 'harp'

## n (neuter)

dziecko 'child'

prosię 'piglet'

okno 'window'

# Grammatical cathegory of gender in Polish — adjectives

**m**
dobr**y** 'good':

**f**
dobr**a**:

**n**
dobr**e**:

pan

lew

dom

królowa

żyrafa

harfa

dziecko

prosię

okno

# 3 masculine subgenders in modern Polish

- personal = m1

- animate impersonal = m2

- inanimate = m3

# 3 masculine subgenders in modern Polish

sg. N pan m1

… 

A pan-a m1

…

sg. N lew m2

…

A lw-a m2

…

sg. N dom m3

…

A dom m3

…

pl. N pan-owie m1
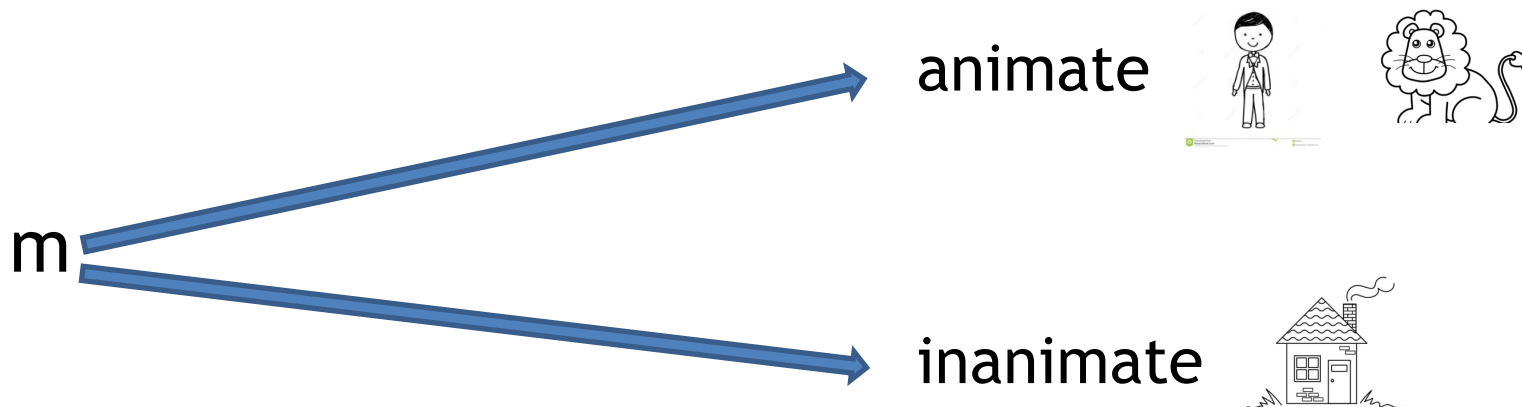
…

A pan-ów m1

…

pl. N lw-y m2

…

A lw-y m2

…

pl. N dom-y m3

…

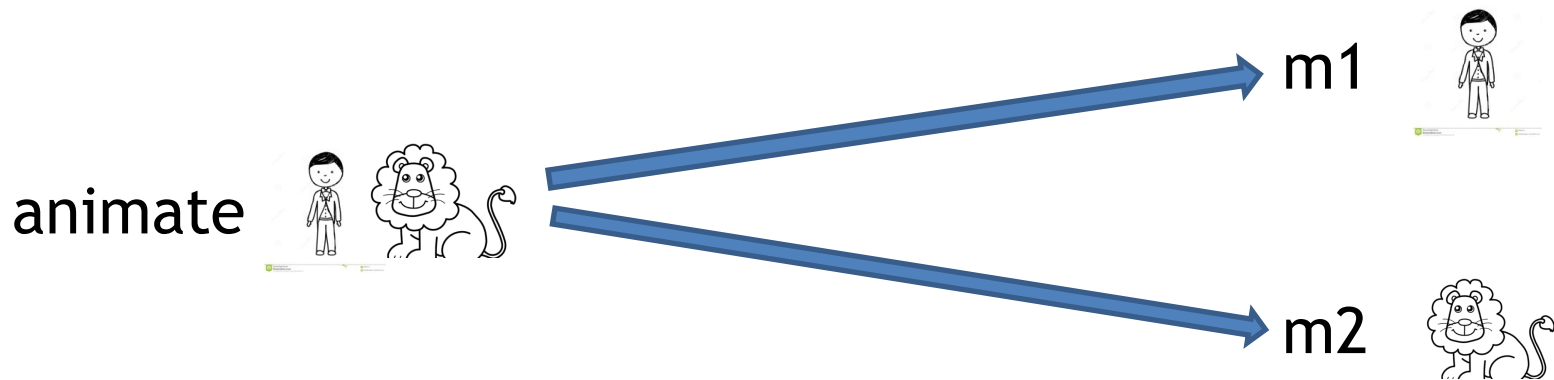A dom-y m3

…

# 3 masculine subgenders in modern Polish

sg. N pan m1

… 

A pan-a m1

…

pl. N pan-(owie) m1

…

A pan-ów m1

…

sg. N lew m2

…

A lw-a m2

…

pl. N lw-y m2

…

A lw-y m2

…

sg. N dom m3

…

A dom m3

…

pl. N dom-y m3

…

A dom-y m3

…

# Changes of the inflection of masculine nouns in Old Polish (up to 15th c.)
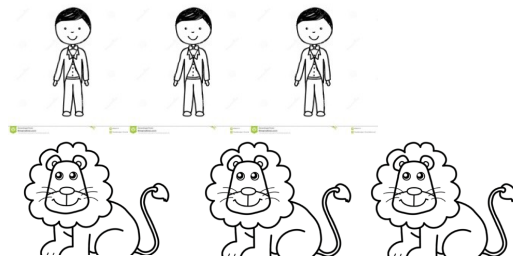
m

animate

inanimate

# Changes of the inflection of masculine nouns in 17th—18th centuries

# Changes of the inflection of masculine nouns in 17th—18th centuries

- N. pl pan-y / pan-owie

- N. pl lw-y / lw-owie

# Tagset

A set of markers signifying part of speech and morpho-syntactic features specific for this part of speech.

harfa [subst:sg:nom:f]

pos number case gender

# Tagging the gender in modern Polish

The gender is assigned to the particular noun, e.g.

pan [subst:sg:nom:m1]            lew [subst:sg:nom:m2]

pan-owie 'gentlemen'            lw-y 'lions'

   [subst:pl:nom:m1]               [subst:pl:nom:m2]

# How to tag the gender in Middle Polish?

pan [subst:sg:nom:m1]

pan-owie / pan-y
  [subst:pl:nom:m?]

lew [subst:sg:nom:m2]

lw-owie / lw-y
  [subst:pl:nom:m?]



*lwowie*

Institute of Polish Language
Polish Academy of Sciences
PAN IJP

# Tagging the gender in KorBa

pan [subst:sg:nom:m]

pan-a [subst:sg:gen:m]

pan-owie
   [subst:pl:nom:manim1]

pan-y [subst:pl:nom:m]

lew [subst:sg:nom:m]
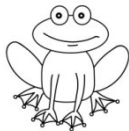
lw-a [subst:sg:gen:m]

lw-owie
   [subst:pl:nom:manim1]

lw-y [subst:pl:nom:m]

# 2. The disappearance of dual number
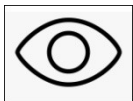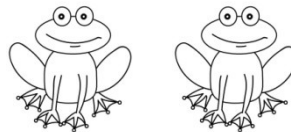
# Category of number in Old Polish

| **sg** | **du** | **pl** |
|---|---|---|

żaba   żabie   żaby 
'frog'

oko   oczy   oka 
'eye'

# The disappearance of dual during 17ᵗʰ—18ᵗʰ c.

**sg**            **du**            **pl**

żaba            żabie            żaby

'frog'

oko            oczy            oczy

'eye'            oka

# How to annotate forms like "oczy" in KorBa?

**Up to 1740**

oczy → du
oka → pl

**After 1740**

oczy → pl
oka → pl

# Conclusions

- The development of the grammatical system should be reflected in the morphosyntactic annotation of the corpus.
- There is no universal solution — for each grammatical feature, different issues should be taken into account:
  - substantive issues, e.g. the degree of prevalence of a given linguistic phenomenon; the timeframe for linguistic change
  - practical issues, e.g. the possibility of creating useful corpus queries and conducting linguistic research in future.

# Thank you!

www.korba.edu.pl