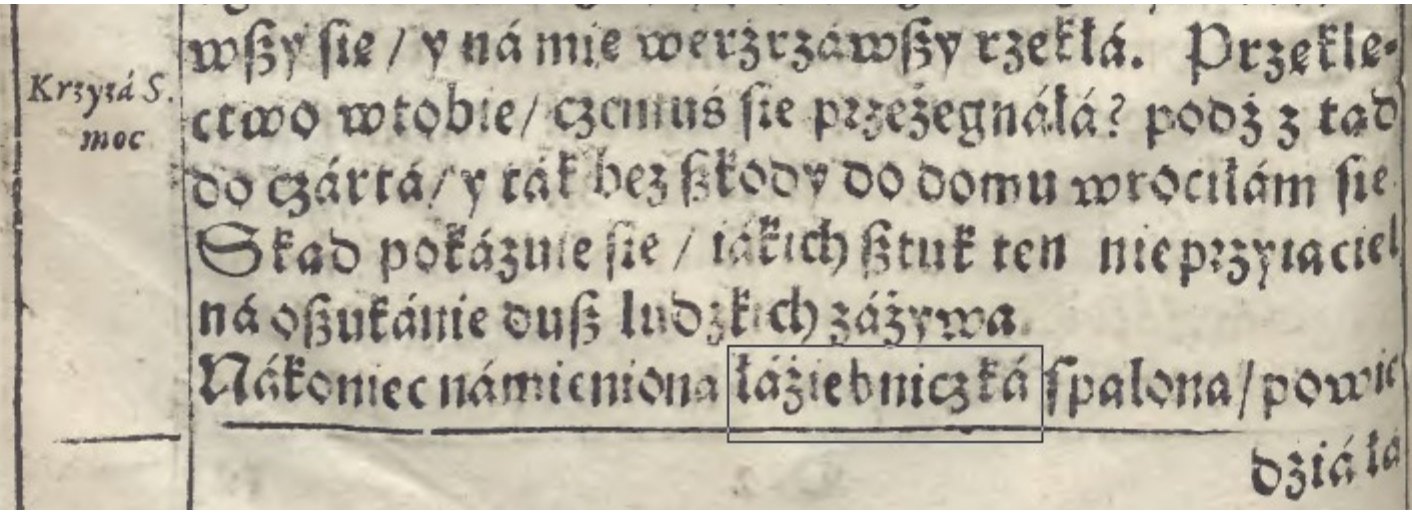


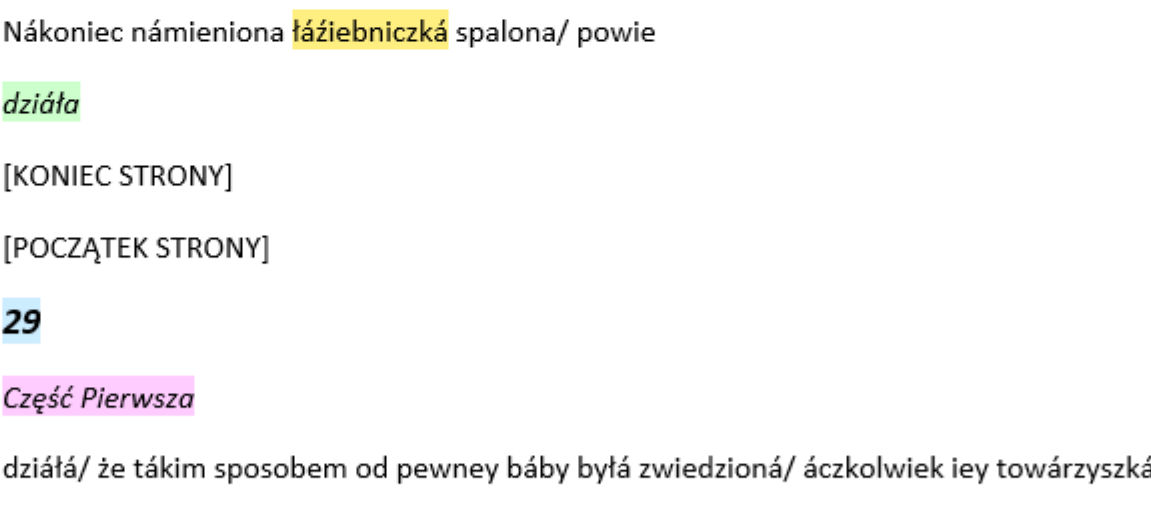
NOTUM PER IGNOTUM THE ENRICHMENT OF LEXICAL INFORMATION AND CORPUS RESOURCES BY USING THE RESULTS OF MORPHOLOGICAL ANALYSIS OF HISTORICAL TEXTS

THE ELECTRONIC CORPUS OF THE 17TH AND 18TH C. POLISH TEXTS (UP TO 1772) - STEP BY STEP

STEP ONE collecting materials (from a scan to a ready-to-work TEI-tagged text)



the fragment of the Polish translation of „Mallus Malleficarum” by Kramer and Sprenger – 1614



manual transliteration and structural annotation (WORD format)

<note place="baŋia" xml:id="not_4.4.3.19-note"> Krzyża S. moc</note>
Przekleństwo wtobie! czemuś się przeznęła? podź z tą! do czarta! y tak bez szkody
do domu wrociłam sie. Stąd pokazuje się: iakich sztuk ten nieprzyjaciel na oszukanie
dusz ludzkich żążywa. Nakoniec namieniona **łaziebniczka** spalona/ powie
</note>
<note type="running-head" xml:id="not_4.4.3.20-head">Część Pierwsza</note>
działa/ że takim sposobem od pewney bąby były zwiędziona/ áczkolwiek iej towarzyszą
towarzyszą różny obyczajem/ to jest: iż szatana w osobie ludzkiej podkła w drodze
gdy wymysem płodzenia wszeteczności/ szła miłośnika swe/o/ nawiedzać: z ktor
obyczajem/ to jest: iż szatana w osobie ludzkiej podkła w drodze/ gdy wymysem
płodzenia wszeteczności/ szła miłośnika swe/o/ nawiedzać: z ktor szatanem
obyczajem ludzkim cielenie złączywszy się/ spytał iej szatan/ ieslihy go poznała/
odpowiedziała ona/ że zo by namniej nie zna: odpowiedział on. Iestem szatan/ y iesli

conversion to XML (compatible with TEI format)

The highlighted form is *łaziebniczka* 'a woman serving in baths'
– a word that haven't been recognised before. Let's see the
path that the string should go through to be recognised.

łaziennik m III przestarz. in. łaziebnik: Łaziennik wyciągał członki, szorował skórę, wycierał, a na ostatek namaszczał ciało olejami wonnymi z mirrą, nardem, cynamonem. Krasz. Kartki 669. // L

masculine form of „łaziebniczka” from Doroszewski's dictionary – 1958

STEP TWO from a transliterated text to a transcribed and morphosyntactically tagged text

AMEBA SUPERTOOL (from Pol. *ameba* ‘amoeba’, Gr. ἀμειβή ‘change’)

a tool converting texts from their transliterated to transcribed version

TRANSKRYBER

formulas of transcription for the converter originally created by Janusz Bień and his team for the purposes of the IMPACT project

– <https://bitbucket.org/jsbien/pol>

id_1768	1EM	‘za	mow	i’\$	mów
id_1769	1EM	‘za	mowi	T*’	mówi
id_1770	1EM	‘za	mowic	\$	mówic
id_1771	1EM	‘za	mowi	A*’	mówi
id_1772	1EM	‘za	mowic	\$	mówic
id_1773	1EM	‘u	mow	i’\$	mów
id_1774	1EM	‘u	mowi	T*’	mówi
id_1775	1EM	‘u	mowic	\$	mówic
id_1776	1EM	‘u	mowi	A*’	mówi
id_1777	1EM	‘u	mowic	\$	mówic
id_1778	1EM	‘u	mow	i’\$	mów
id_1779	1EM	‘u	mowi	T*’	mówi
id_1780	1EM	‘u	mowic	\$	mówic
id_1781	1EM	‘u	mowi	A*’	mówi
id_1782	1EM	‘u	mowic	\$	mówic
id_1783	1EM	‘o	mow	i’\$	mów
id_1784	1EM	‘o	mowi	T*’	mówi
id_1785	1EM	‘o	mowic	\$	mówic

At this point, every string should be converted to the modern form. The diacritical marks such as acutes or graves are erased (if not in use nowadays) or are added (if the opposite situation occurs). The orthography is modernised (contemporary rules of orthography are applied).

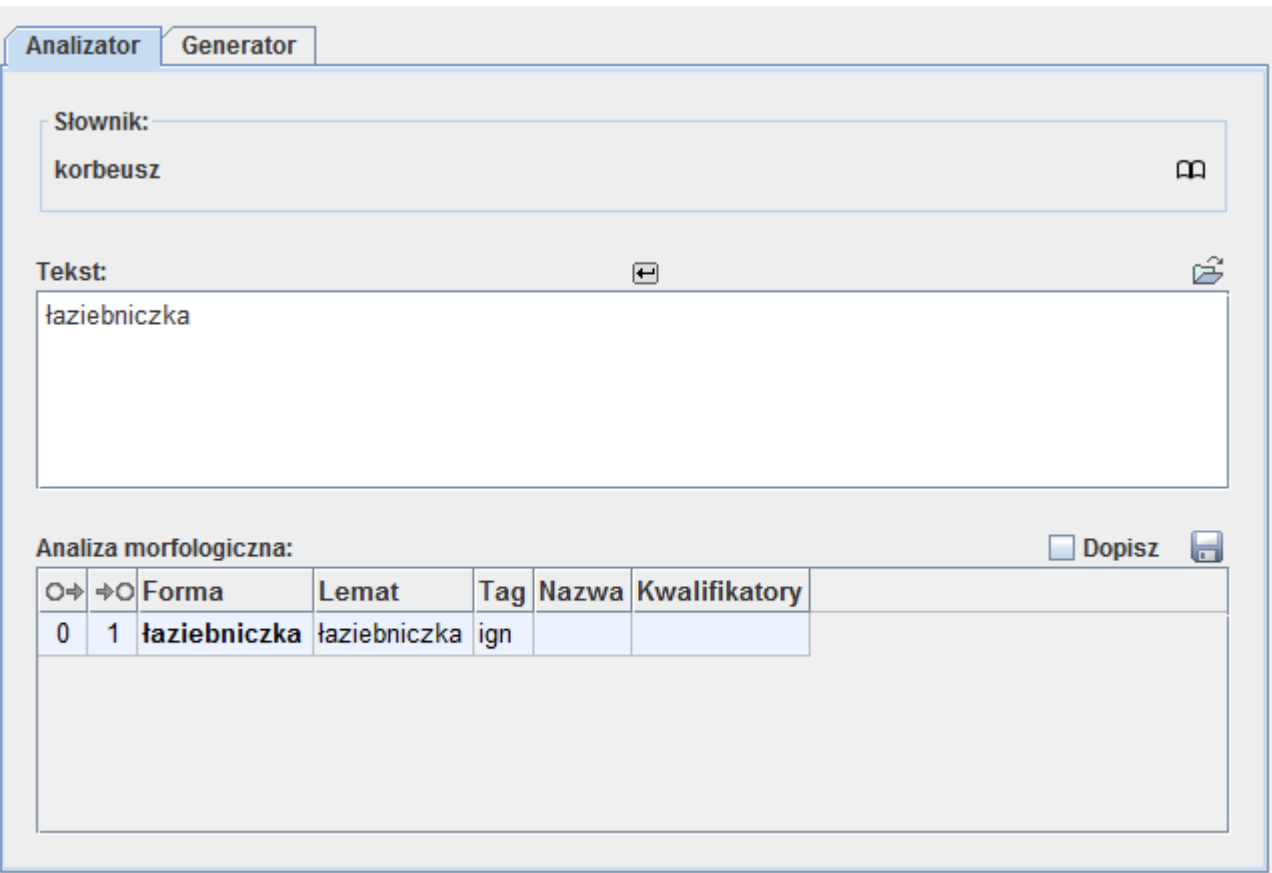
Special care should be taken here; if the formula is not restricted enough (especially by adding exceptions), it is extremely easy to deform a correctly written string.

Every move can throw the baby out with the bathwater.

KORBEUSZ

morphological analyser customised to analyse 17th and 18th C. texts, based on Morfeusz – an analyser for contemporary Polish.

<http://sgjp.pl/morfeusz/>



working on igns – seeking the balance between analysis paralysis and corpus utopia

the ign-lists (a set of strings unrecognised by Korbeusz and tagged as *ign* – from Lat. *ignotum* ‘unknown’)

A	B	C
1 text	transit	orth
2 DoBiKontr	MATYJASZ	MATYJASZ
3 JUNI	JUNI	JUNI
4 DoBiKontr	JUNI	JUNI
5 DoBiKontr	MATYJASZA	MATYJASZA
6 DoBiKontr	SUPRA	SUPRA
7 DoBiKontr	PENDENTA	PENDENTA
8 DoBiKontr	JUNI	JUNI
9 DoBiKontr	wielobnoici	wielobnoici
10 DoBiKontr	senjorów	senjorów

currently: ca. 0.5M segments

A	B
1 orth	liczba - orth
2 jn	4172
4 imp	2898
5 Zaczym	2699
6 dat	1595
7 Czeczta	1558
8 Czeczta	1345
9 Wm	1205
10 lichm	1001

ca. 0.2M unique segments (not lemmas)

Evaluating the igns is the key to enrich our resources successfully. The main factor is frequency. More frequent strings should be added to dictionary first to decrease the amount of unrecognised strings. It is also a good time to check out if the transcription is applied correctly.

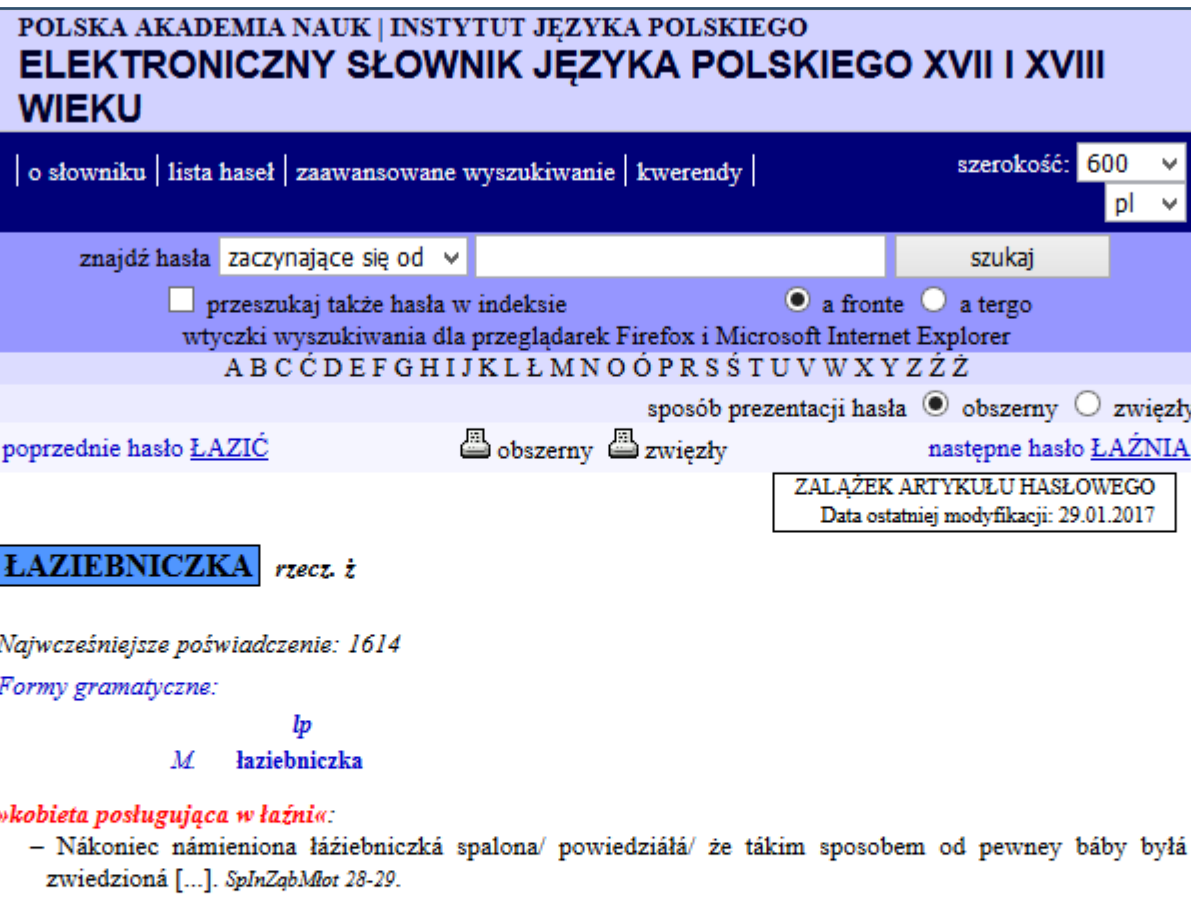
246466 SpInZabMlot	zwiędziony	zwiędziony
246467 SpInZabMlot	zle	zle
246468 SpInZabMlot	zebymci	zebymci
246469 SpInZabMlot	sercecznie	sercecznie
246470 SpInZabMlot	Bazyliu	Bazyliu
246471 SpInZabMlot	nowu	nowu
246472 SpInZabMlot	świątobliwaze	świątobliwaze
246473 SpInZabMlot	poscigać	poscigać
246474 SpInZabMlot	Rävenspurgu	Rävenspurgu
246475 SpInZabMlot	záchosua	záchosua
246476 SpInZabMlot	łaziebniczka	łaziebniczka
246477 SpInZabMlot	wazyła	wazyła
246478 SpInZabMlot	obywateli	obywateli
246479 SpInZabMlot	gorniego	gorniego

The fragment of the ign-list below shows the discussed string.

STEP THREE, FOUR, FIVE ...

enrichment of lexical information and corpus resources

1. adding new entries or new inflectional forms to the e-SXVIII



2. creating new formulas for Transkryber

To make the string recognizable by Korbeusz it is necessary to modernise the orthography. In our case, at least two formulas should be used:

~ if „z” stands before „i”, it should be replaced with „z”
~ every „d” should be replaced with „a”

So it goes like the following:

id_of_formula	on/off	initials_of_author	right_context	replacee	left_context	replacer
id_1	1EM	..*	z	i*	..*	z
id_2	1EM	..*	z	i*	..*	z

input: łaziebniczka

output: łaziebniczka

Sprawy dawne z duchownemi o dziesięciny aby były kasowane, chociażby były osądzone i na egzekucy były, a teraz znowu dla zgody z nimi, jabyim **rozmał** in genere wszystkie, począć dawać według postanowienia; rozumiałbym, żebyśmy tym tylko płacił, którzy w naszych parafiach służą i robą i dawać; jen tak, co by sie dobrze wychowali, służąc Panu Bogu, a jeżeliby co postronnym dalekim przyszło dawać, tedy tak dawać z łanu, jako pobór **dawamy**, i tylko w ten czas, kiedy poboru niemasz, gdyż mają wsi i wielkie opatrzenia. A w każdej sprawie in **caus** simplicibus mechy był każdy bliższy do odwoła tak z duchownych, jako świeckich. Bo teraz duchowni świeckich każą kłopotować, co by nie miało być: **ciatusz** **propor** być ma in **caus** simplicibus ad **passonem**. Gdy granice archybispa, biskup, opat i każdy z duchownych z którym służebicem odprawuje, niech też sam

„dawamy” is an old form of the lexeme DAWAĆ ‘to give’ (plural form of the first person), was added to the SGJP as a variant of modern form „dajemy”

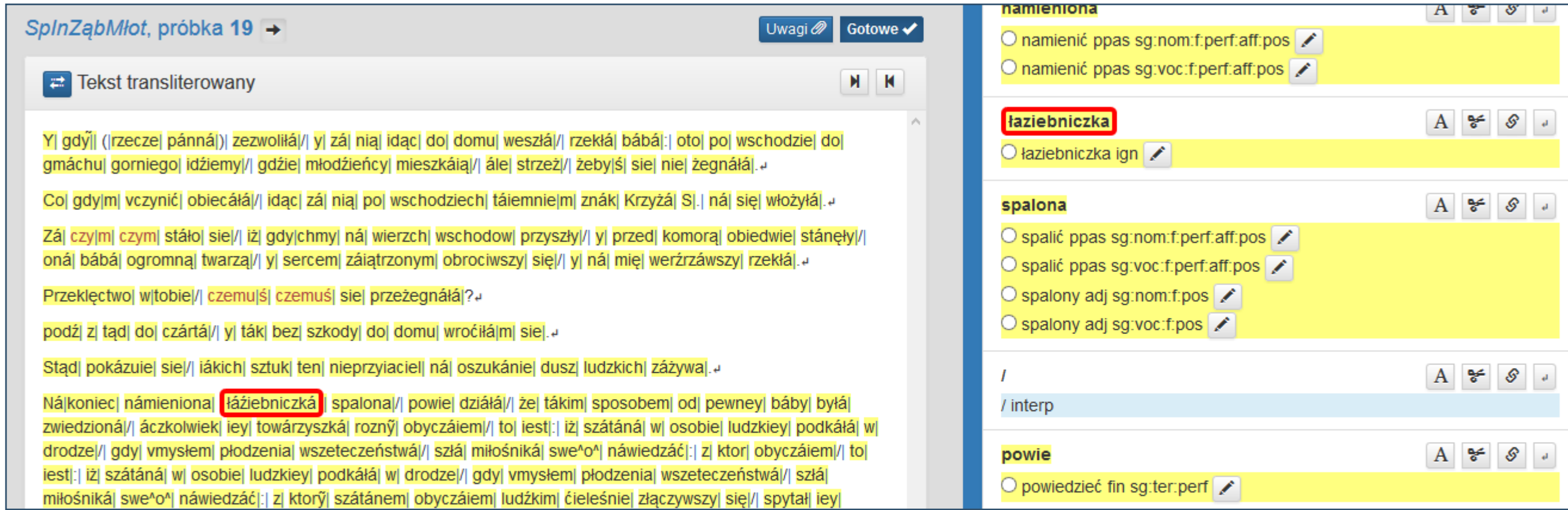
and then, retrieving information by Korbeusz

Analiza morfologiczna:	Forma	Lemat	Tag	Nazwa	Kwalifikatory
0 1	łaziebniczka	łaziebniczka	subst sg nom f	nazwa pospolita	

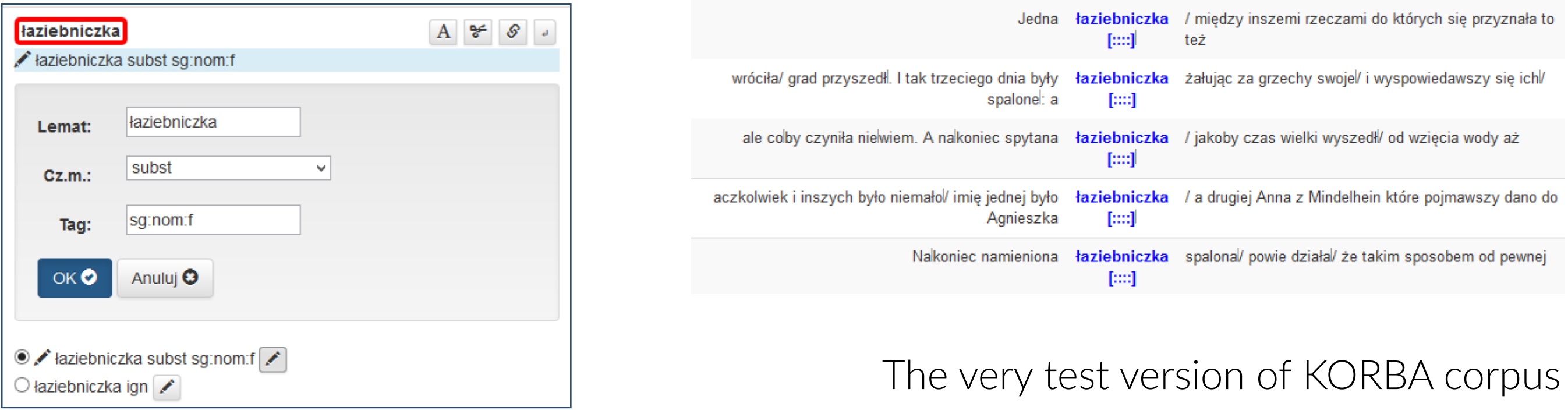
Analiza morfologiczna:	Forma	Lemat	Tag	Nazwa	Kwalifikatory
0 1	dawamy	dawać	fin pl pri imperf		

ANOTATORNIA

a tool supporting manual morphosyntactic annotation of texts; the annotators disambiguate interpretation of the segment by selecting the one provided beforehand by Korbeusz or adding their own manually (as it is shown below).



The discussed word (*łaziebniczka*) is still unrecognised by Korbeusz due to the lack of harvesting. In such case, the adequate values of grammatical categories can be chosen manually by typing a lemma and selecting corresponding POS and tags.



The very test version of KORBA corpus

ABOUT OUR PROJECT:

title:	Electronic corpus of 17th and 18th century Polish texts (up to 1772), acronym: KorBa (KORpus BARokowy)
the aim:	the creation of a corpus of 17th and 18th century Polish texts (another stage of development of the Polish National Corpus (Narodowy Korpus Języka Polskiego, NKJP, see: http://nkjp.pl/) and tools for its processing (searching, filtering, summarizing statistical data, etc.)
corpus size:	around 12 million tokens
duration:	2013 – 2018
coordinator:	Institute of Polish Language, Polish Academy of Sciences
cooperation:	Linguistic Engineering Group, Institute of Computer Science, Polish Academy of Sciences
main investigator:	Włodzimierz Gruszczyński
funding:	A Ministry of Science and Higher Education National Programme for the Development of Humanities grant 0036/NPRH2/H11/81/2012

CONTACT US:

Renata Bronikowska r.bronikowska@wp.pl	Emanuel Modrzejewski modrzejewski.emanuel@gmail.com
--	--

Institute of Polish Language, Polish Academy of Sciences

www.ijp-pan.krakow.pl/en

The Electronic Corpus of the 17th and 18th c. Polish Texts (up to 1772)

korba@ijp-pan.krakow.pl

