

# Conjugate, decline and spell like years ago

## A corpus-based morphological analyzer of 19<sup>th</sup> century Polish

Joanna Bilińska   Witold Kieraś   Magdalena Derwojedowa

Chair of Formal Linguistics

Institute of Polish Language  
University of Warsaw



Slavic Diachronic Corpus Linguistics Conference  
Tromsø, April, 21-22, 2015

# The project

*Automatic morphological analysis of Polish texts from 1830-1918 period with respect to evolution of inflection and spelling, grant awarded by Polish National Science Centre (DEC-2012/07/B/HS2/00570)*

2013-2016

<http://www.f19.uw.edu.pl>

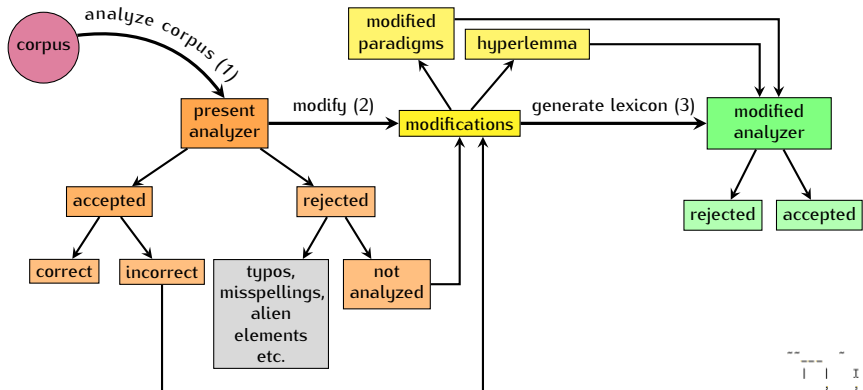


# Motivation

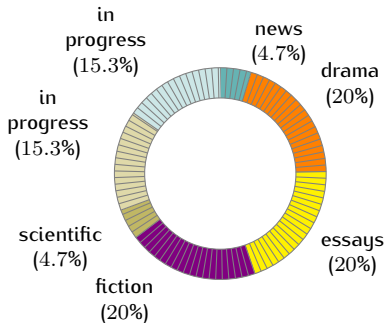
- growing interest in NLP of old(er) Polish
- availability of adaptable tools
- extensive, detailed inflectional description of Polish
- computer-readable text from the period



# A diachronic morphological analyzer in three steps



# The corpus



Samples collected per style

- 5 stylistic subcorpora, 200,000 segments each
- 5-13 samples/year (at least one sample of each style per year)
- more than 500 authors (inc. 119 anonymous), 480 books, 135 newspapers & magazines, 54 digital libraries
- repository ownCloud
- downloadable at <http://www.f19.uw.edu.pl/>



## The good, the bad and the ugly: excerpting samples

- stroke of luck: a text in original spelling and inflection
- time machine: aging a text from modern edition with original spelling and forms (a trick for poor quality scans)
- everyday humdrum: OCR of the digital scans (Fine Reader, tesseract)
- bad luck: typing in a text (good text, no contemporary re-editions, poor scan)

Jako bezwarunkową tezę stawiam:  
„Ropa w Galicyi, gdziekolwiek takowa występuje, znajduje się na powtórnem złożysku” (sekundäre Lagerstätte) i jest skutkiem suchej destylacji z wielkiej ilości nagromadzonych szczątków organicznych, zawartych w starszych warstwach, jak są dotychczas w Karpatach znane pokłady łupku bitumicznego utworu neokomskiego.

wielu innych występów ochrania! Ależ gdy z drugiej strony, swąmy chciwość wielu dzierzawców, i niektórych dziedziców, mianowicie spekulantów, kupujących dobra bez pieniędzy, którzy w folwarkach, gdzie w porównaniu pól dworskich, pańszczyzna jest zamala, nie chcąc wydawać na najem, a pragnąc przez wyciągnięcie intrat wyjść z interesów, nieraz

Good and poor quality scans.



## A sample and metadata

Jako bezwarunkową tezę stawiam: „Ropa w Galicyi, gdziekolwiek takowa występuje, znajduje się na powtórnem złożysku” (sekundäre Lagerstatte) i jest skutkiem suchej destylacyi z wielkiej ilości nagromadzonych szczątków organicznych, zawartych w starszych warstwach, jak są dotychczas w Karpatach znane pokłady łupku bitumicznego utworu neokomskiego.

autor: Olszewski, Stanisław  
tytuł: Przyczynek do teorii pochodzenia i występowania nafty w Galicyi  
data wydania: 1881  
miejsce wydania: Lwów  
redaktor: Radziszewski, Br.  
tytuł książki:  
tytuł gazety, czasopisma, serii wydawniczej: Kosmos. Czasopismo Polskiego Towarzystwa Przyrodników im. Kopernika  
nr: R. 6  
wydawnictwo: nakł. Polskiego Towarzystwa Przyrodników im. Kopernika  
numery stron: 524-527  
styl: popularnonaukowy  
źródło: Śląska Biblioteka Cyfrowa  
link: <http://www.sbc.org.pl/dlibra/docmetadata?id=20199>  
uwagi: nie ma w HINT, chemia?

A sample: in yellow — 19 c. forms,  
underlined — alien elements.

A metadata file.



# The dictionary

Hasło	Część mowy	Rodzaj/aspekt
podbieg	subst	m3
podbiegać	v	ndk
podbiegający	pact	ndk
podbieganie	ger	ndk
podbiegły	adj	
podbiegnąć	v	dk
podbiegnięcie	subst	n2
podbiegnięcie	ger	dk
podbiegunowość	osc	f
<b>podbiegunowy</b>	<b>adj</b>	
podbieleć	v	ndk
podbiełający	pact	ndk
podbiełanie	ger	ndk
podbiełany	ppas	ndk
podbiełenie	ger	dk
podbiełić	v	dk
podbiełony	ppas	dk
podbiełecz	subst	m1
podbiełec	v	ndk

Wszytkie formy    Formy bazowe

przymiotnik [SJPDor.]  
P4

	l. p.					l. m.		
	m1	m2	m3	n1,n2	ż	p1	m1	pozostale
M./(W.)	podbiegunowy		podbiegunowe		podbiegunowa	podbiegunowi		podbiegunowe
D.			podbiegunowego		podbiegunowej			podbiegunowych
C.			podbiegunowemu		podbiegunowej			podbiegunowym
B.	podbiegunowego	podbiegunowy	podbiegunowe		podbiegunową	podbiegunowych		podbiegunowe
N.			podbiegunowym		podbiegunową			podbiegunowymi
Ms.			podbiegunowym		podbiegunowej			podbiegunowych
Złoż.								podbiegunowo+

**Odysłacz**

nazwa cechy    [podbiegunowość](#)

przymiotnik „zanegowany”    [niepodbiegunowy](#)

PODBIEGUNOWY 'polar' in *Grammatical dictionary of Polish*  
(„Słownik gramatyczny języka polskiego”, Saloni et al. 2012)

<http://www.sgjp.pl/kuznia/slownik/>





# The analyzer for modern Polish

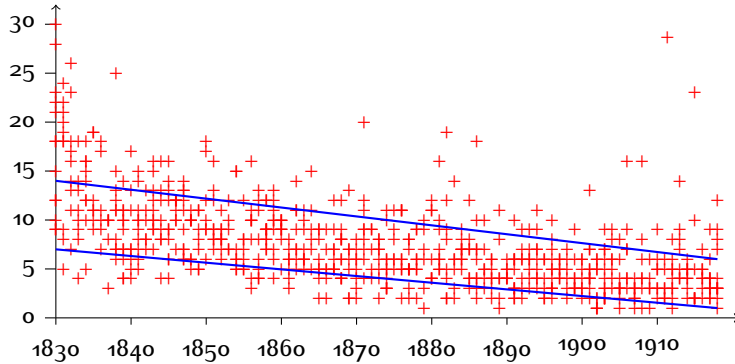
```
Morfeusz analyzer, version: 1.9.0
Setting dictionary search path to: .
Using dictionary: sgjp (default)
Utwór iłów solnych odznacza się mniejszemi ilościami zazwyczaj cięższėj
[0,1,Utwór,utwór,subst:sg:acc:m3,pospolita,_
 0,1,Utwór,utwór,subst:sg:nom:m3,pospolita,_]
[1,2,iłów,ił:s1,subst:pl:gen:m3,pospolita,geol.
 1,2,iłów,ił:s2,subst:pl:gen:m2,pospolita,_
 1,2,iłów,ił:s2,subst:pl:gen:m3,pospolita,_]
[2,3,solnych,solny,adj:pl:acc:m1.p1:pos,_,_
 2,3,solnych,solny,adj:pl:gen:m1.m2.m3.f.n1.n2.p1.p2.p3:pos,_,_
 2,3,solnych,solny,adj:pl:loc:m1.m2.m3.f.n1.n2.p1.p2.p3:pos,_,_]
[3,4,odznacza,odznaczać,fin:sg:ter:imperf,_,_]
[4,5,się,się,qub,_,_]
[5,6,mniejszemi,mniejszemi,ign,_,_]
[6,7,ilościami,ilość,subst:pl:inst:f,pospolita,_]
[7,8,zazwyczaj,zazwyczaj,adv,_,_]
[8,9,cięższėj,cięższėj,ign,_,_]

```

A string from a sample analyzed by *Morfeusz 2.0* (Woliński, 2014)  
with SGJP-based lexicon. Marked not analyzed 19c. forms.



# The performance of non-modified analyzer on 1830-1918 corpus



Percent of unrecognized unique segments per file.

# Modifications

The screenshot displays the 'Kuźnia (Smithy) editor' interface with 11 modification rules arranged in two columns. Each rule consists of a number, a text input field, a dropdown menu, and a plus sign. Rules 2, 4, 6, 7, and 11 are highlighted with red boxes. Rule 9 has a loading spinner icon.

Rule	Text	Dropdown	Plus
1	slab y	Wybierz	+
2	slab ego	Wybierz	+
2	slab égo	SGJP: daw.	+
3	slab emu	Wybierz	+
3+	slab u	Wybierz	+
4	slab ym	Wybierz	+
4	slab ém	SGJP: daw.	+
4	slab em	SGJP: daw.	+
5	slab e	Wybierz	+
6	slab a	Wybierz	+
7	slab ej	Wybierz	+
7	slab éj	SGJP: daw.	+
8	slab a	Wybierz	+
9	slab l	Wybierz	+
10	slab ych	Wybierz	+
11	slab emi	SGJP: daw.	+
11	slab émi	SGJP: daw.	+
11	slab ymi	Wybierz	+

*Kuźnia* (Smithy) editor

# Days of future past: modified lexicon

1	podbiegunow <b>y</b>
2	podbiegunow <b>ego</b>
3	podbiegunow <b>emu</b>
4	podbiegunow <b>ym</b>
5	podbiegunow <b>e</b>
6	podbiegunow <b>a</b>
7	podbiegunow <b>ej</b>
8	podbiegunow <b>ą</b>
9	podbiegunow <b>i</b>
10	podbiegunow <b>ych</b>
11	podbiegunow <b>ymi</b>
12	podbiegunow <b>o+</b>

2 G/A masc sg

4 I masc sg

7 D/C/L/V fem sg

11 I pl

1	podbiegunow <b>y</b>
2	podbiegunow <b>ego</b> podbiegunow <b>ého</b> <i>daw.</i>
3	podbiegunow <b>emu</b>
4	podbiegunow <b>ym</b> podbiegunow <b>ém</b> <i>daw.</i> podbiegunow <b>em</b> <i>daw.</i>
5	podbiegunow <b>e</b>
6	podbiegunow <b>a</b>
7	podbiegunow <b>ej</b> podbiegunow <b>éj</b> <i>daw.</i>
8	podbiegunow <b>ą</b>
9	podbiegunow <b>i</b>
10	podbiegunow <b>ych</b>
11	podbiegunow <b>ymi</b> <i>daw.</i> podbiegunow <b>emi</b> <i>daw.</i> podbiegunow <b>émi</b> <i>daw.</i>
12	podbiegunow <b>o+</b>



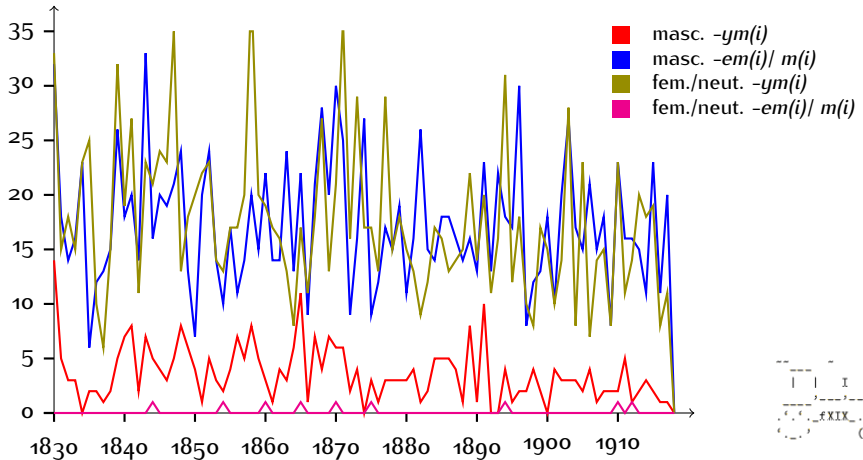
# Wavering instrumental

Kopczyński (1780) differentiate between instrumentalis of masculine and neuter:

	<i>sg</i>	<i>pl</i>
<i>masc</i>	<b>ym</b>	<b>ymi</b>
<i>neut, fem</i>	<b>em</b>	<b>emi</b>



# Distribution of instrumental endings across genders

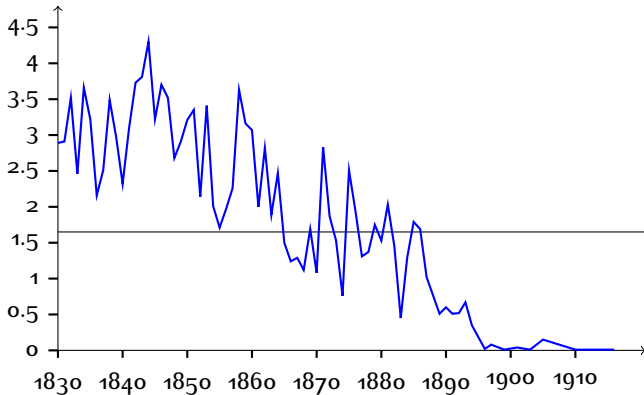


# Applications

- part-of-speech tagging,
- tracing dynamics of language change (spelling and inflection)
- following adaptation stages of loanwords
- marking neologisms/disused/obsolete words with date



## Example: dynamics of é

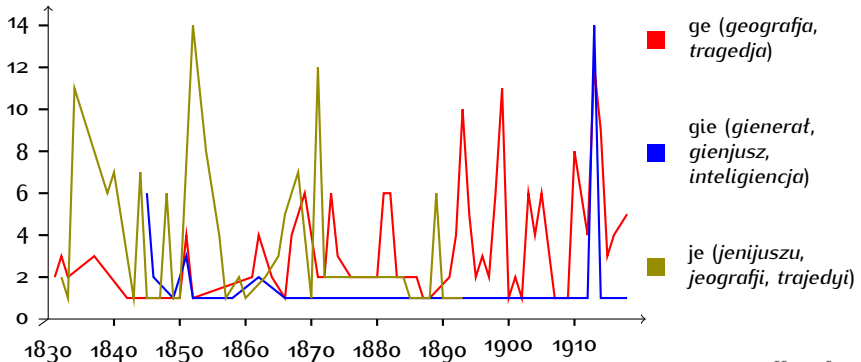


Percent of unique segments with é per year.





## Example: *ge:gie:je* alternations



Distribution of *je:gie:ge* spelling



## References

- DERWOJEDOWA M., KIERAŚ W., SKOWROŃSKA D. i WOŁOSZ R., *Współczesne narzędzia leksykograficzne a analiza tekstów dawniejszych*, „Polonica”, t. XXXIV s. 21--27 2014a.
- DERWOJEDOWA M., KIERAŚ W., SKOWROŃSKA D. i WOŁOSZ R., *Zasób leksykalny polszczyzny II poł. XIX wieku a możliwość automatycznej analizy morfologicznej tekstów z tego okresu*, [w:] *Leksyka języków słowiańskich w badaniach synchronicznych i diachronicznych*, red. M. GĘBKĄ-WOLAK, J. KAMPER-WAREJKO i A. MOROZ, s. 183--195, 2014b.
- DERWOJEDOWA M., SKWROŃSKA D. i KIERAŚ W., *Korpus polszczyzny XIX wieku – od mikrokorpusu do korpusu średniej wielkości*, „Prace Filologiczne”, s. 249--254 2014c.
- KOPCZYŃSKI O., *Grammatyka dla szkół narodowych*, I wyd., Warszawa 1780.
- SALONI Z., WOLIŃSKI M., WOŁOSZ R., GRUSZCZYŃSKI W. i SKOWROŃSKA D., *Słownik gramatyczny języka polskiego*, II wyd., Warszawa 2012, CD.
- WOLIŃSKI M., *Morfeusz Reloaded*, [w:] *Proceedings of the Ninth International Conference on Language Resources and Evaluation, LREC 2014*, red. N. CALZOLARI, K. K. CHOUKRI, T. DECLERCK, H. LOFTSSON, B. MAEGAARD, J. MARIANI, A. MORENO, J. ODIJK i S. PIPERIDIS, s. 1106–1111, ELRA, Reykjavík, Iceland 2014, ISBN 978-2-9517408-8-4, URL <http://www.lrec-conf.org/proceedings/lrec2014/index.html>.



<http://www.f19.uw.edu.pl/>

funded by the Polish National Science Centre  
grant DEC-2012/07/B/HS2/00570

