

# PolEval 2020

Task 2. Morphosyntactic tagging  
of Middle, New and Modern Polish

**Witold Kieraś, Marcin Woliński**  
(prezentuje Łukasz Kobyliński)

# Znakowanie morfosyntaktyczne

Znakowanie słów w tekście odpowiednimi tagami (znacznikami) morfosyntaktycznymi.

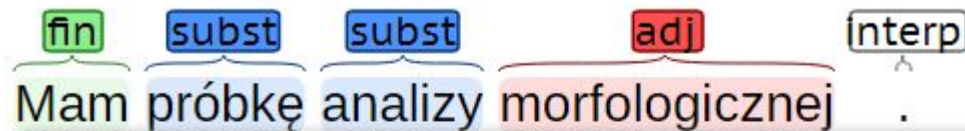
W języku polskim często: dezambiguacja morfosyntaktyczna (wskazywanie właściwego tagu ze zbioru możliwych).

Potrzebna definicja tagsetu:

- w języku angielskim: 36 - 200 tagów, np. the Penn Treebank:
  - NN - noun (singular); NNS - noun (plural);
  - VB - verb, base form; VBD - verb, past tense
- w języku polskim: 4 000 teoretycznie możliwych tagów!

Number	Tag	Description
1.	CC	Coordinating conjunction
2.	CD	Cardinal number
3.	DT	Determiner
4.	EX	Existential <i>there</i>
5.	FW	Foreign word
6.	IN	Preposition or subordinating conjunction
7.	JJ	Adjective
8.	JJR	Adjective, comparative
9.	JJS	Adjective, superlative
10.	LS	List item marker
11.	MD	Modal
12.	NN	Noun, singular or mass
13.	NNS	Noun, plural
14.	NNP	Proper noun, singular
15.	NNPS	Proper noun, plural
16.	PDT	Predeterminer
17.	POS	Possessive ending
18.	PRP	Personal pronoun
19.	PRPS	Possessive pronoun
20.	RB	Adverb
21.	RBR	Adverb, comparative
22.	RBS	Adverb, superlative
23.	RP	Particle
24.	SYM	Symbol
25.	TO	<i>to</i>
26.	UH	Interjection
27.	VB	Verb, base form
28.	VBD	Verb, past tense
29.	VBG	Verb, gerund or present participle
30.	VCN	Verb, past participle
31.	VBP	Verb, non-3rd person singular present
32.	VBZ	Verb, 3rd person singular present
33.	WDT	Wh-determiner
34.	WP	Wh-pronoun
35.	WPS	Possessive wh-pronoun
36.	WRB	Wh-adverb

# Tagset w języku polskim



mama:subst:pl:gen:f  
mamić:impt:sg:sec:imperf  
**mieć:fin:sg:pri:imperf**

Tagset jest pozycyjny

- mam - **mieć:fin:sg:pri:imperf**
  - mieć - lemat
  - fin - klasa gramatyczna (35 możliwości)
  - sg, pri, imperf - wymagane atrybuty gramatyczne (specyficzne dla każdej z klas)
  - dla wielu klas gramatycznych - również atrybuty opcjonalne
- ... daje to łącznie ok. 4000 możliwych tagów (ok. 1 000 pojawia się w rzeczywistym korpusie).

# Tagset w języku polskim (cd.)

Skąd tak wiele klas gramatycznych i atrybutów?

Czy *czytanie* to rzeczownik, który odmienia się przez przypadek, czy czasownik (ma aspekt: *rozpoczęto czytanie* ale nie *rozpoczęto przeczytanie*)?

Rozwiązanie: wydzielono osobną klasę gramatyczną dla form typu *czytanie*.

Klasy zostały wyodrębnione

- na podstawie cech fleksyjnych form (przez co się odmieniają, z czym uzgadniają),
- na podstawie cech składniowych (dystrybucyjnych),
- na podstawie cech semantycznych i pragmatycznych.

# Tagset w języku polskim (cd.)

Leksem	Fleksem	Symbol	Przykład
Rzeczownik	rzeczownik	subst	<i>kot, profesorowie</i>
	forma deprecjatywna	depr	<i>profesory</i>
Liczebnik	liczebnik główny	num	<i>sześć, dużo</i>
	liczebnik zbiorowy	numcol	<i>sześcioro, trojga</i>
Przymiotnik	przymiotnik	adj	<i>polski</i>
	przymiotnik	adja	<i>polsko</i>
	przyprzymiotnikowy	adjp	<i>polsku</i>
	przymiotnik poprzyminkowy	adjc	<i>wesół</i>
Przysłówek		adv	<i>bardziej, kiedy</i>
Zaimek	nietrzecioosobowy	ppron12	<i>ja, tobie</i>
	trzecioosobowy	ppron3	<i>on, jemu</i>
	SIEBIE	siebie	<i>sobą</i>
Czasownik	forma nieprzeszła	fin	<i>jadam</i>
	forma przyszła	bedzie	<i>będę</i>
	czasownika być		
	aglutynant czasownika być	aglt	<i>-śmy</i>
	pseudoimiesłów	praet	<i>jadał</i>
	rozkaźnik	impt	<i>jadaj</i>
	bezosobnik	imps	<i>jadano</i>
	bezokolicznik	inf	<i>jadać</i>
	imiesłów przysłówkowy współczesny	pcon	<i>jadając</i>
	imiesłów przysłówkowy uprzedni	pant	<i>zjadłszy</i>
	odslownik	ger	<i>jadanie</i>
	imiesłów przymiotnikowy czynny	pact	<i>jadający</i>
	imiesłów przymiotnikowy bierny	ppas	<i>jadany</i>

Czasownik typu WINIEN (forma teraźniejsza)		winien	<i>winna, powinni</i>
Predykatyw		pred	<i>trzeba, słychać</i>
Przyimek		prep	<i>pod, we</i>
spójnik	współrzędny	conj	<i>oraz, lub</i>
	podrzędny	comp	<i>że, aby</i>
Wykrzyknik		interj	<i>ach, psiakrew</i>
Burkinostka		burk	<i>omacku, trochu</i>
Kublik		qub	<i>nie, -ż, również</i>
Skrót		brev	<i>dr, np</i>
Ciało obce		xxx	<i>errare, humanum</i>
Interpunkcja		interp	<i>;, -, , ]</i>

# Zadanie historycznofleksyjne PolEval

Na danych reprezentujących 400 lat rozwoju języka polskiego obejmujących:

- XVII i XVIII wiek — okres średniopolski — <https://korba.edu.pl/>,
- XIX wiek — okres nowopolski — <http://korpus19.nlp.ipipan.waw.pl/>,
- okres współczesny — <http://nkjp.pl/>.

Opis fleksji historycznej wymagał dostosowania tagsetu — w każdym z korpusów stosowany jest nieco inny opis, zostały one ujednolicone na potrzeby tego zadania.

Celem zadania jest ujednoznacznienie fleksyjne zaprezentowanych danych i dodanie interpretacji dla słów nieznanym analizatorowi, a więc takie samo zadanie, jakie wykonują tagery tekstów współczesnych.

# Zadanie historycznofleksyjne PolEval

Wydaje się nam, że interesującym elementem jest to, że dane nie są jednorodne, opisują trzy ściśle związane ze sobą języki.

**Udostępnione dane zawierają informację o czasie powstania poszczególnych tekstów (niekiedy przybliżonym).**

Spodziewamy się, że najlepsze tagery będą korzystały z tej informacji.

# Tagset historyczny (wybrane różnice)

- Trzy wartości liczby: pojedyncza, podwójna, mnoga, np.
  - *Dwie*<sub>DUAL</sub> *żabie*<sub>DUAL</sub> *upragnione*<sub>DUAL</sub> *po polach biegały*
- Uwzględniono tzw. formy krótkie lub niezłożone przymiotników paralelne do używanych współcześnie form „długich”, czyli „złożonych”:
  - *równy* adj:sg:nom:m:pos — *równien* adjb:sg:nom:m:pos
  - *dawnego* adj:sg:gen:m.n:pos — *dawna* adjb:sg:gen:m.n:pos
  - *polskiemu* adj:sg:dat:m:pos — *polsku* adjb:sg:dat:m.n:pos
  - *piękną* adj:sg:acc:f:pos — *pięknę* adjb:sg:acc:f:pos

W języku współczesnym ich pozostałością są formacje typu *z dawna* i *po polsku*.



# Tagset historyczny (wybrane różnice)

- W okresie średniopolskim dopiero kształtuje się rodzaj męskoosobowy, dlatego w korpusach historycznych stosowane jest inne oznaczenia rodzajów męskich. W danych PolEval przyjęto dla wszystkich okresów oznaczenie m — męski, a bardziej szczegółowe oznaczenia manim1 i manim2 — tylko w kontekstach różnicujących odpowiednie formy.

# Podjęcia do znakowania morfosyntaktycznego

Uczenie maszynowe, czy reguły lingwistyczne

- ręczne reguły tworzone przez lingwistów
- metody mieszane, np. tager TaKIPi (jęz. polski) - łączą podejście nadzorowane z regułami przygotowanymi przez lingwistów
- metody nienadzorowane - HMM (ukryte modele Markowa) i algorytm Expectation Minimization (EM)
- metody nadzorowane - oparte na algorytmach ML, np. tager Brilla (rule-based); oparte na Conditional Random Fields (CRF) - uwzględniają sekwencyjność języka; wykorzystujące sieci neuronowe

# Krótką historia tagerów dla języka polskiego

## Dotychczasowe podejścia:

- 2007: znakowanie oparte na regułach (TaKiPI), ok. 88% dokładności,
- 2010: tager Brilla dostosowany do j. polskiego (PANTERA), ok. 89% dokładności,
- 2012: uczenie pamięciowe (WMBT), ok. 90% dokładności,
- 2012: warunkowe pola losowe (WCRFT, Concraft), ok. 91% dokładności,

**Tymczasem dla języka angielskiego dokładność przekracza 97%**

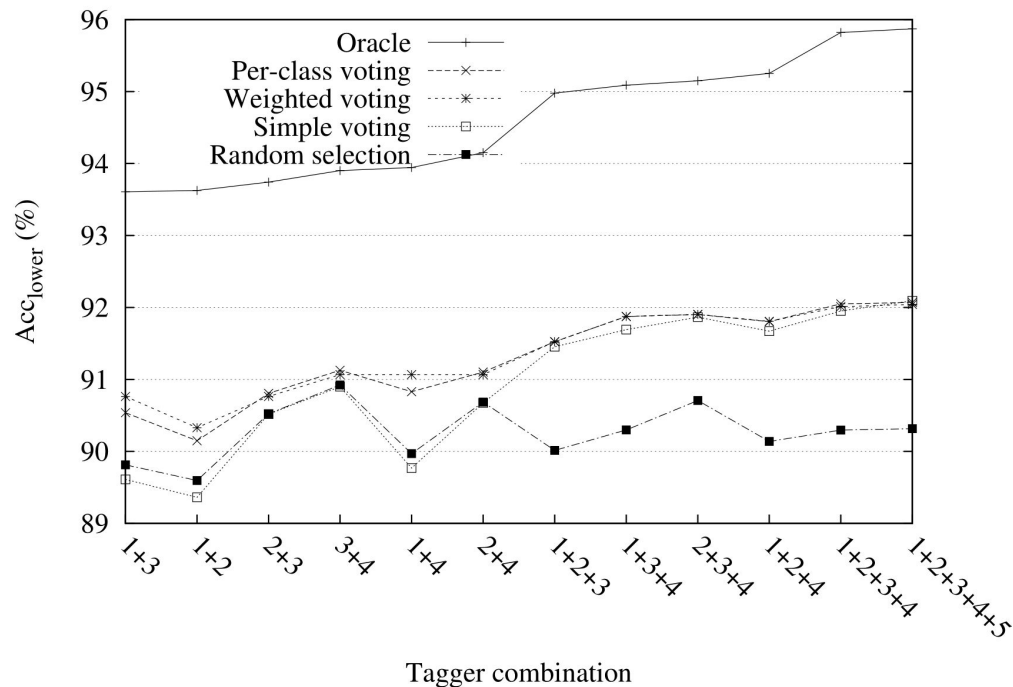
# Znakowanie a klasyfikacja

Skoro tagery są klasyfikatorami -  
możemy wykorzystać ensembling do  
poprawy jakości znakowania.

2014: PoliTa,

ok. 92% dokładności znakowania

8% błędnych znakowań to 80  
milionów błędów w korpusie o  
rozmiarze 1 miliarda słów!



# A może meta-uczenie?

Nauczmy nowy klasyfikator na podstawie wyników działania pojedynczych tagerów:

- który z klasyfikatorów wybrać w danym kontekście?

Statystyki wyboru tagerów w zbiorze testowym:

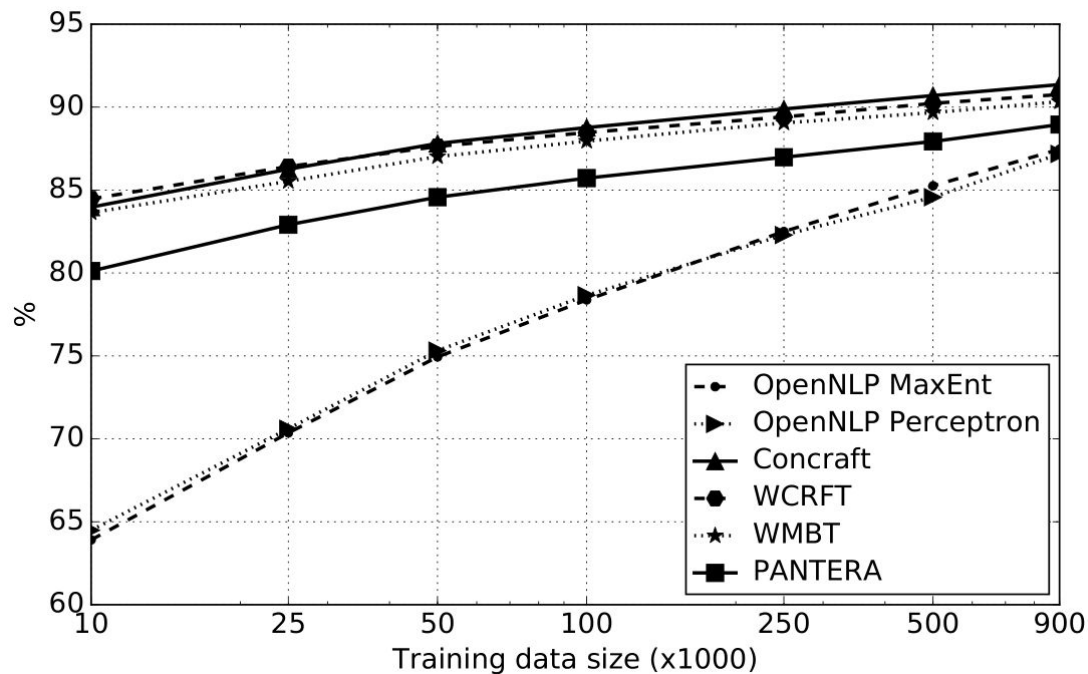
	pantera	wmbt	wcrft	concraft
Total	0.235	4.527	20.491	74.747
Known words	0.294	3.885	22.585	73.235
Unknown words	0.0	7.062	12.193	80.745

Tagger	5-2 split		
	$Acc_{lower}$	$Acc_{lower}^K$	$Acc_{lower}^U$
pantera	88.3646	91.9884	7.0421
wmbt	90.1567	92.0960	46.6339
wcrft	90.6354	92.7668	42.8029
concraft	91.1535	93.1072	47.3090
polita	91.7264	93.7206	46.9726
Meta libsvm	92.0797	93.9697	49.6637
Meta xgboost	92.1040	93.9842	49.8268

# Wpływ rozmiaru danych uczących

Rozmiar (zaanotowanego) zbioru uczącego ma bardzo duży wpływ na jakość znakowania.

... ale tworzenie takiego zbioru jest bardzo kosztowne.



# Stan dla języka polskiego - wyniki PolEval 2017

<b>System name</b>	<b>Acc (%)</b>	<b>deep network</b>	<b>hand-crafted features</b>	<b>character-level embeddings</b>	<b>word-level embeddings</b>
<i>Toygger</i>	94.6343	<i>yes</i>	<i>yes</i>	<i>no</i>	<i>yes</i>
<i>KRNNT_AB</i>	93.8083	<i>yes</i>	<i>yes</i>	<i>no</i>	<i>no</i>
<i>NeuroParser</i>	93.6109	<i>yes</i>	<i>no</i>	<i>yes</i>	<i>no</i>
<i>AvgPer_Forced</i>	90.9134	<i>no</i>	<i>yes</i>	<i>no</i>	<i>no</i>
<i>Concraft</i>	91.6115	<i>no</i>	<i>yes</i>	<i>no</i>	<i>no</i>
<i>WCRFT</i>	91.1693	<i>no</i>	<i>yes</i>	<i>no</i>	<i>no</i>
<i>WMBT</i>	90.6722	<i>no</i>	<i>yes</i>	<i>no</i>	<i>no</i>