

Граматичний потенціал Польсько-українського паралельного корпусу (PoIUKR)

Наталя Коциба

Інститут Міждисциплінарних Досліджень, Варшавський університет,
20 листопада 2009 р.

План презентації

- історія створення PoUKRy, можливості його застосування в лексикографії:
<http://www.domeczek.pl/~polukr>
- *два типи тагсетів в PoUKR: розширення тагсету Корпусу IPI PAN і MULTEXT-East*
- демонстрація програми УГТаг і пошуку по корпусу (Андрій Микуляк, Інститут Ядерної Фізики, Варшава)
- етапи роботи над корпусом: від паперових книжок до формату XML; демонстрація програми PLUczek

Морфосинтактичне тагування в PoIUKR

- джерела морфосинтактичної інформації в PoIUKR
- види тагсетів залежно від способу запису граматичної інформації (форма представлення)
- вибрані „проблематичні” категорії і як їх трактують існуючі тагсети (зміст)

Джерела морфосинтактичної інформації

- **Польська мова**: пакет ТаКІРІ, таґер Корпусу ІПІ (ПАН) + Вроцлавська політехніка, біля 200 тис. слів
- лематизатор, морфологічний аналізатор, уоднозначнювач
- **Українська мова**: граматичний словник УМІФ НАНУ, 250 тис. слів
- лематизатор, морфологічний аналізатор, немає уоднозначнення



“Словники України” on-line

Словозміна Синонімія Фразеологія

бути

- Реєстр**
- [бутербродик](#)
 - [бутербродний](#)
 - [Бутефліка](#)
 - [Бути](#)
 - [бути](#)
 - [Бутивля](#)
 - [бутил](#)
 - [бутилацетат](#)
 - [бутилен](#)
 - [бутилені](#)
 - [бутилкаучук](#)
 - [бутиловий](#)
 - [бутиль](#)
 - [Бутин](#)
 - [бутина](#)
 - [Бутини](#)
 - [бутинський](#)
 - [Бутир](#)
 - [бутираг](#)
 - [бутирометр](#)
 - [бутирофенон](#)
 - [бутири](#)
 - [бутифос](#)
 - [Бутівка](#)
 - [Бутівське](#)

бути – дієслово недоконаного виду

Інфінітив	бути	
	однина	множина
Наказовий спосіб		
1 особа		будьмо
2 особа	будь	будьте
МАЙБУТНІЙ ЧАС		
1 особа	буду	будемо, будем
2 особа	будеш	будете
3 особа	буде	будуть
ТЕПЕРІШНІЙ ЧАС		
1 особа	є	є
2 особа	є, еси	є
3 особа	є	суть, є
Активний дієприкметник		
Дієприслівник		
будучи		
МИНУЛИЙ ЧАС		
чол. р.	був	були
жін. р.	була	
сер. р.	було	
Активний дієприкметник		
Пасивний дієприкметник		
Безособова форма		
Дієприслівник		
бувши		

Український нетагований текст

- Львів розташований на етнічних українських землях і є одним з головних нервових вузлів українського народу, найважливішим клапаном його серця, вічним збудником честолюбства, гордості й потягу до волі.

Українська: вхідні таґи з лемами

Львів<JDJAJIJK><Львів 0|Львів 0|Лев 1|Лев 1|>
розташований<BDBAV?><розташований 0|розташований 0|розташувати 0|>
на<NONOZOPF><на 4|на 3|на 2|на 1|> **етнічних**<AVATAХ><етнічний 0|етнічний 0|етнічний 0|> **українських**<AVATAХJIGIJKGKJMGM><український 0|український 0|український 0|Український 0|Український 0|Український 0|>
землях<FM><земля 2|> **і**<SSSCN0Z0><і 1|і 3|і 2|> **є**<UPUOUNUKUMUL><бути 0|бути 0|бути 0|бути 0|бути 0|бути 0|> **одним**<HUNQHERQRERU><один 0|один 0|один 0|один 0|один 0|один 0|> **з**<PE><з 0|> **головних**<AVATAХ><головний 0|головний 0|головний 0|> **нервових**<AVATAХ><нервовий 0|нервовий 0|нервовий 0|> **вузлів**<MIMI><вузол 2|вузол 1|>
українського<ANADABJDJBKV><український 0|український 0|український 0|Український 0|Український 0|Українське 0|> **народу**<MBMCMVMC><народ 0|народ 0|нарід 0|нарід 0|>, **найважливішим**<AQAEAU><найважливіший 0|найважливіший 0|найважливіший 0|> **клапаном**<ME><клапан 0|>
його<FGODOBODOB><його 0|воно 0|воно 0|він 0|він 0|>
серця<NKNHNBNN><серце 0|серце 0|серце 0|серце 0|>,
вічним<AQAEAU><вічний 0|вічний 0|вічний 0|> **збудником**<MEME><збудник 1|збудник 2|> **честолюбства**<NB><честолюбство 0|>,
гордості<FCFBFF><гордість 0|гордість 0|гордість 0|> **й**<SSSCZ0><й 1|й 2|>
потягу<MFMCMGMBMCMFMGFDGD><потяг 2|потяг 2|потяг 2|потяг 1|потяг 1|потяг 1|потяг 1|потяга 0|Потяга 0|>
до<NGNFNENDNCNBNAHNNINJKNLNMNNPB><до 2|до 2|до 2|до 2|до 2|до 2|до 2|до 2|до 2|до 1|>
волі<UOFCFBFFGCGBGFGFGCGBGHGNFMF><воліти 0|воля 0|воля 0|воля 0|Воля 2|Воля 2|Воля 2|Воля 1|Воля 1|Воля 1|Воля 1|Воля 1|воло 0|віл 0|>.

Приклади граматичних кодів, що використовуються в УМІФ (384)

Грамматичне значення	Морфкод	Приклад
Дієслово, інфінітив, доконаний вид, активний стан	VA	<i>прочитати</i>
Дієприкметник, чоловічий рід, однина, називний відмінок, доконаний вид, минулий час, активний	BA	<i>зрослий</i>
Невідмінюваний прикметник	AZ	<i>ультра</i>
Іменник загальний, жіночий рід, однина, давальний відмінок	FC	<i>квітці</i>
Предикатив (присудкове слово)	X0	<i>слід</i>

Польський нетаґований текст

- W dzisiejszym posiedzeniu komisji uczestniczy ekspert komisji pan profesor Jan Gajewski.



ulif2.dtd x ulif1.dtd x morphSh.xml x

Current Element

```
1 <?xml:version="1.0" encoding="UTF-8"?>
2 <!DOCTYPE cesAna SYSTEM "xcesAnalPI.dtd">
3 <cesAna xmlns:xlink="http://www.w3.org/1999/xlink" type="pre_morph" version="PI-1.2">
4 <chunkList xml:base="text.xml">
5 <chunk type="p" xlink:href="#dv1hd1">
6 <chunk type="s">
7 <group id="a1" rule="Uncertain: Dobre-PrepNG (z postmodyfikatorem dopełniaczowym) na koncu zdania lub nawiasu, lub przed czasownikiem" synh="a1" semh="a:
8 <tok id="a1">
9 <orth>W</orth>
10 <lex disamb_sh="0"><base>w</base><ctag>prep.acc.nwok</ctag></lex>
11 <lex disamb="1"><base>w</base><ctag>prep.loc.nwok</ctag></lex>
12 </tok>
13 <tok id="a2">
14 <orth>dzisiejszym</orth>
15 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.sg.inst.m1.pos</ctag></lex>
16 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.sg.inst.m2.pos</ctag></lex>
17 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.sg.inst.m3.pos</ctag></lex>
18 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.sg.inst.n.pos</ctag></lex>
19 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.sg.loc.m1.pos</ctag></lex>
20 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.sg.loc.m2.pos</ctag></lex>
21 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.sg.loc.m3.pos</ctag></lex>
22 <lex disamb="1"><base>dzisiejszy</base><ctag>adj.sg.loc.n.pos</ctag></lex>
23 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.pl.dat.m1.pos</ctag></lex>
24 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.pl.dat.m2.pos</ctag></lex>
25 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.pl.dat.m3.pos</ctag></lex>
26 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.pl.dat.f.pos</ctag></lex>
27 <lex disamb_sh="0"><base>dzisiejszy</base><ctag>adj.pl.dat.n.pos</ctag></lex>
28 </tok>
29 <tok id="a3">
30 <orth>posiedzeniu</orth>
31 <lex disamb_sh="0"><base>posiedzenie</base><ctag>subst.sg.dat.n</ctag></lex>
32 <lex disamb="1"><base>posiedzenie</base><ctag>subst.sg.loc.n</ctag></lex>
33 <lex disamb_sh="0"><base>posiedzieć</base><ctag>ger.sg.dat.n.perf.aff</ctag></lex>
34 <lex disamb="1"><base>posiedzieć</base><ctag>ger.sg.loc.n.perf.aff</ctag></lex>
35 </tok>
36 <tok id="a4">
37 <orth>komisji</orth>
38 <lex disamb="1"><base>komisja</base><ctag>subst.sg.gen.f</ctag></lex>
39 <lex disamb_sh="0"><base>komisja</base><ctag>subst.sg.dat.f</ctag></lex>
40 <lex disamb_sh="0"><base>komisja</base><ctag>subst.sg.loc.f</ctag></lex>
41 <lex><base>komisja</base><ctag>subst.pl.gen.f</ctag></lex>
42 </tok>
```

Insert Element

Insert Sibling

Insert Entity

amp
apos
gt
lt
quot

Види таґсетів залежно від
способу запису граматичної
інформації (форма презентації)

Види таґсетів залежно від рівня кодування (словоформа чи значення притаманних їй категорій)

СИМВОЛЬНІ: граматична характеристика **словоформи**
передається в одному коді **British National Corpus**, корпус УМІФ
вимагає мало машинної пам'яті, але занадто багато людської

ЛАНЦЮЖКОВІ: кожна категорія, її атрибут і вартості цього
атрибута мають свій унікальний код; морфосинтактична
характеристика словоформи представлена **послідовністю кодів**
вартостей атрибутів **KIPi, CzNK, НКРЯ**

можуть бути економніші при пошуку, ніж символи, напр. коли
шукаємо допасування на рівні вартості якогось атрибуту, що
властивий кільком частинам мови (відмінок, рід для іменника і
прикметника)

Види таґсетів відносно позиціювання граматичної інформації

- **ПОЗИЦІЙНІ** CzNC, MULTEXT

Кожна категорія (виражена її вартістю в тазі) має фіксовану позицію у ланцюжку, вартості різних категорій можуть мати ту саму нотацію (t,n,y,p, 1, 2, 3)

ЧНК (біля 4 тис.)

volen: VsYS---XX-AP--- verb, passive participle, masculine, singular, any person, any tense, positive, passive

hraniční: AAIS4----1A---- standard adjective, masc. inanimate, singular, accusative, positive

Щоб задати пошук по одному з атрибутів, треба записати цілий ланцюжок

- ***УНІВЕРСАЛЬНІ** KIPi „*gen.*”, „case=„gen””
- (символьних кодів ця проблема не торкається)

Види таґсетів залежно від принципів та глибини категоризації граматичної інформації

- **ГНІЗДОВІ** (CzNC : 1.POS - 10 і 2.subPOS - 75)
- **ФЛЕКСЕМНІ** корпус ІПІ ПАН (польська)
Флексеми :: ЧМ (Януш Бень) множини слів зі спільними наборами змінних граматичних атрибутів
- Частини мови або їх підмножини об'єднуються у спільні класи флексем, якщо мають однакову граматичну характеристику (теоретична підстава до розбиття займенників на групи, що приєднуються до класів слів, які вони “заміняють”; виокремлення займенника *sie, siebie, winien*) **29**
- Принцип аглютинації для економності запису

Граматичні категорії польських флексем і їх вартості

- **number:** sg, pl;
- **case:** nom, acc, gen, dat, inst, loc, voc;
- **gender:** masculine personal m1 (*facet*), masculine animate m2 (*koń*), masculine inanimate m3 (*stół*), feminine f (*kobieta, żyrafa, książka*), два середні роди: n1 (*dziecko*), n2 (*okno*), і три *plurale tantum* p1 (*wujostwo*), p2 (*drzwi*), p3 (*okulary*);
- **person:** pri, sec, ter;
- **degree:** pos, comp, sup;
- **aspect:** imperf, perf;
- **negation:** aff, neg;
- **accentability** (Pol.: *akcentowość*): akc, nakc;
- **post-prepositionality** (Pol.: *poprzyimkowość*): praep, npraep;
- **accommodability** (Pol.: *akomodacyjność*): congr, rec;
- **agglutination** (Pol.: *aglutynacyjność*): nagl, agl;
- **vocability** (Pol.: *wokaliczność*): wok, nwok.

Механізм аліасів

n n1 n2

P p1 p2 p3

masc m1 m2 m3

noun subst depr ger xxs ppron12 ppron3 PPRON GNOUN
PROPNOUN

pron ppron12 ppron3 siebie PPRON

verb fin praet aglt będzie inf imps impt **past ppas pcon pant** ger
winien PART(PPAST) FUT PRES

Інтерпретація вибраних категорій в граматичних словниках і корпусах

Дієприслівники

- польська (вид вписано на рівні словника і визначає синхронність чи попередність, дієприслівники часу не мають, час дієслів визначається за допомогою виду або аналітично)

Флексеми:

adv.pres.prtcp.

adv.anter.prtcp.

- українська (вид і час – морфологічна характеристика)
 1. Дієслово, дієприслівник, доконаний вид, *минулий час, активний стан **VW** *прочитавши* (*піду, прочитавши*)
 2. Дієслово, дієприслівник, недоконаний вид, минулий час, активний стан **UW** *читавши*.
 3. Дієслово, дієприслівник, недоконаний вид, *теперішній час, активний стан **UQ** *читаючи* (*робив/робитиме читаючи*)
- PolUKR

adverbial and adjectival participles characterised by aspect and *tense:

verb:part:perf , verb:part:imperf, verb:part:imperf:praet

Займенники

- проблема слов'янських займенників: 296 таїв для 309 займенників (Elena Paskaleva)
- польська: поділ на 1-2 ос., 3 ос і *siebie* (ów, jak?)
- українська: займенник-іменник, займенник-прикметник
- російська: також займенник-предикатив і займенник-прислівник
- чеська: 18 підкатегорій на рівні підЧМ
- PoUKR: підхід УМІФ, додатково поділ ІПІ ПАН на 1-2 і 3 особові

Windows Internet Explorer

http://korpus.pl/poliqarp/poliqarp.php?query=%5Bpos%3D%22pred%22%5D&corpus=1&showMatch=1&showContext=3&leftContext=5&rightContext=5&wideContext=50&hil

Google

korps.pl

0 blocked

Check

AutoLink

AutoFill

Send to

korps

pl

Settings

Strona

Narzędzia

WYSZUKIWANIE

Szukaj: [pos="pred"] w próbcie Korpusu IPI PAN (2. wydanie; 30M segmentów)

Szukaj

Pokaż opcje wyszukiwania »

WYNIKI

Zapytanie: [pos="pred"]

Znaleziono 1000 wyników
Wyświetlanie wyników 1 - 20

Następne 20

urządzeniem, za pomocą którego	można [można:pred]	było, bez konieczności wysłania
królom, chanom i emirom	wolno [wolno:pred]	dotykać relikwii. Ostatnim razem
co robił i myślał,	można [można:pred]	się było dowiedzieć tylko z
kierują się obcym liberalizmem.	To [to:pred]	hipokryci. Budują nie państwo
nie gorsi ani lepsi.	To [to:pred]	prawdziwa potęga. Nie dziwie
mieli nigdy łatwego życia.	Można [można:pred]	powiedzieć, że obaj jesteście
któremu choćby z racji wieku	wolno [wolno:pred]	bezkarnie wytykać błędy i śmieszności
nazywaliśmy zresztą Rakietą.	Szkoda [szkoda:pred]	, że zginął. Po
i przyświątynną medresą. Była	to [to:pred]	jedyna zagraniczna podróż przyszłego emira
biesiadach, czekaniu na nie	wiadomo [wiadomo:pred]	co. Omar odzywał się
co o tym sądzę.	To [to:pred]	był wrzesień dziewięćdziesiątego dziewiątego i
, mamrocząc pod nosem:	To [to:pred]	jest nasz dom, tygrysów
tygrysów dom i lwów,	to [to:pred]	kraj potężnych gór, zielonych
któremu zabronić wjazdu. Zabierało	to [to:pred]	wiele czasu i komplikowało życie
jest wrogów. Nie pierwszy	to [to:pred]	zamach i nie mnie jednego
udzielnego księstwa, z którego	można [można:pred]	by wyruszyć po władzę w

Start

warsztaty_ratysbon...

Словники України о...

Gmail - FW: konfere...

Poliqarp online - Win...

Microsoft Excel użyt...

analize

Internet

100%

PL

17:40

Предикативи польські

польська 26, українська 176, російська >1200,
чеська 0

To książka (*This* <is a> book) → займенник-
іменник

Модальні слова: *można, trzeba, wolno, wiadomo, trza, niepodobna, podobna, dość, dosyć*
→ модальні прислівники (диференціація на семантичному рівні).

Віддієслівні: *słyszać, widać, stać, czuć, znać* і
відіменникові: *szkoda, potrzeba, żal, wstyd, strach, pora, czas, brak, śmiech* → прислівники

Предикативи укр.

- семельфактиви: *зирк, круть-верть* → вигуки
- відприслівникові: *зимно, безвітряно, безсніжно* → прислівники
- демінутивні дієслова: *спуцькати, спатки, ходитоньки, їсточки, їстки* → інфінітивні форми дієслова

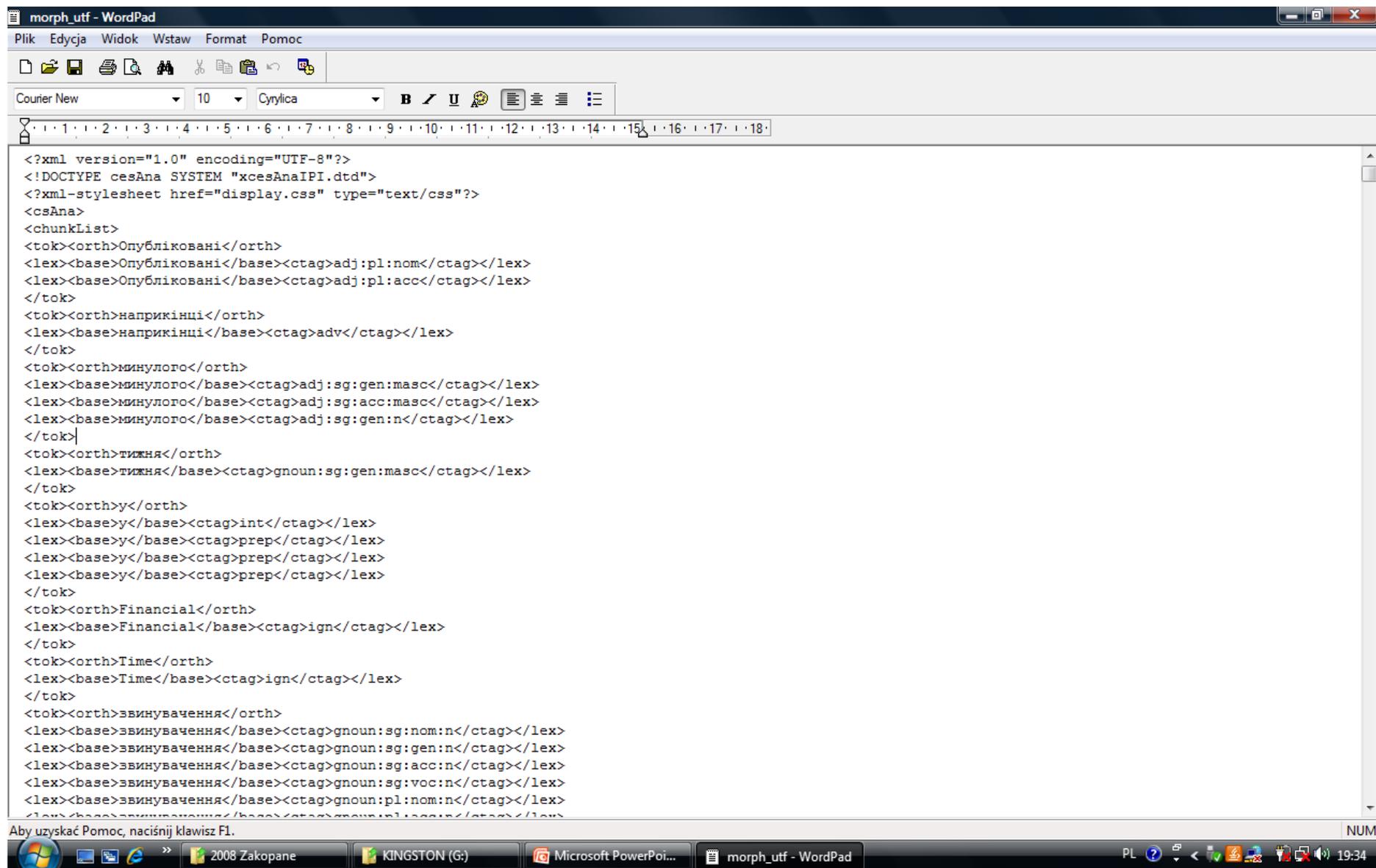
Прикметники

- УГС УМІФ: прикметники вищого і найвищого ступенів подано як окремі леми, ступінь як категорія не числиться взагалі
- IPI PAN:po-**polsku**;(adj vs **по-польськи** adv)
polsko-niemiecki ???
- **НКРЯ** вищий ступінь прикметників і прислівників включено до предикативів
- PoIUKR: розроблено правила визначення ступеня прикметника і його релематизації, реалізація на рівні конвертації тагів;
післяприйменникові прикметники трактуються як частини прислівників

Таґсет PoIUKR

- позиційний
- гніздовий (ЧМ і підЧМ) для користувача
- альяси підключаються на рівні пошуку

Фрагмент українського тексту з тагами PoIUKR



```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE cesAna SYSTEM "xcesAnaIPI.dtd">
<?xml-stylesheet href="display.css" type="text/css"?>
<cesAna>
<chunkList>
<tok><orth>Опубліковані</orth>
<lex><base>Опубліковані</base><ctag>adj:pl:nom</ctag></lex>
<lex><base>Опубліковані</base><ctag>adj:pl:acc</ctag></lex>
</tok>
<tok><orth>наприкінці</orth>
<lex><base>наприкінці</base><ctag>adv</ctag></lex>
</tok>
<tok><orth>минулого</orth>
<lex><base>минулого</base><ctag>adj:sg:gen:masc</ctag></lex>
<lex><base>минулого</base><ctag>adj:sg:acc:masc</ctag></lex>
<lex><base>минулого</base><ctag>adj:sg:gen:n</ctag></lex>
</tok>
<tok><orth>тижня</orth>
<lex><base>тижня</base><ctag>gnoun:sg:gen:masc</ctag></lex>
</tok>
<tok><orth>у</orth>
<lex><base>у</base><ctag>int</ctag></lex>
<lex><base>у</base><ctag>prep</ctag></lex>
<lex><base>у</base><ctag>prep</ctag></lex>
<lex><base>у</base><ctag>prep</ctag></lex>
</tok>
<tok><orth>Financial</orth>
<lex><base>Financial</base><ctag>ign</ctag></lex>
</tok>
<tok><orth>Time</orth>
<lex><base>Time</base><ctag>ign</ctag></lex>
</tok>
<tok><orth>звинувачення</orth>
<lex><base>звинувачення</base><ctag>gnoun:sg:nom:n</ctag></lex>
<lex><base>звинувачення</base><ctag>gnoun:sg:gen:n</ctag></lex>
<lex><base>звинувачення</base><ctag>gnoun:sg:acc:n</ctag></lex>
<lex><base>звинувачення</base><ctag>gnoun:sg:voc:n</ctag></lex>
<lex><base>звинувачення</base><ctag>gnoun:pl:nom:n</ctag></lex>
<lex><base>звинувачення</base><ctag>gnoun:pl:acc:n</ctag></lex>
```

Тагсети МТЕ

- Міжнародний проект MULTEXT-East:
<http://nl.ijs.si/ME/>
- морфосинтаксичні специфікації для польської і української мов від версії МТЕ-4
- конвертер для польської мови з формату Корпусу ІРІ РАН (скор. КІПІ) до МТЕ-4(pl), доступний за адресою:
<http://domeczek.pl/~polukr/mte-conv>

Конверсія польських тагів

- Конверсію зроблено на основі формату тагів морфологічного аналізатора Морфеуш, які також використовуються для Корпусу ІРІ РАН. З розмічених за допомогою програми ТаКІРІ текстів корпусу витягнуто автоматично **1295** тагів.
- Потім таги були поділені на їх найменші значеннєві частини і записані до реляційної бази даних, де кожна вартість займала окрему колонку. Коди були замінені на їх відповідники в МТЕ
- Вигеновано нові коди з відповідним порядком атрибутів для кожної частини мови.
- Велика частина тагів мала однозначну проекцію, але багато з них також були поділені на групи.

Типи проєкцій тагів

- **Однозначна** проєкція на МТЕ (**1192** тагів): порівняльний і найвищий ступінь прикметників, дієслова, дієприкметники, герундії (відієслівні іменники), кількісні іменники, депреціативні іменники, особові і зворотні займенники, множина іменників, прийменники.
- Таги, котрі **поділено** на наступні категорії в МТЕ: насамперед “кублики” і неособові займенники.
- **Нові таги**: збірні числівники, кілька форм займенників.
- Таги, котрі були **зведені** до одного.

Розширення тагів КІПІ

- З **1298** оригінальних тагів **101** отримали більше, ніж одну проекцію в МТЕ:
- кожен з 60 тагів для прикметників основного ступеня порівняння був поділений на 13 більш детальних;
- 18 іменникових тагів “subst” отримали від 2 до 7 тагів в МТЕ;
- “кублики” були поділені на 7 категорій з 27 унікальними тагами
- предикативи поділено на 3 категорії з 4 тагами

Розроділ “кубликів” в проекції МТЕ

Категорія	Приклад	Таг МТЕ	Кількість лексем
C	alboż	1	11
I	hej	1	179
P	jakoś, się	16	85
Q	że	2	74
R	wczoraj	4	233
S	ponad	2	7
X	mocium	1	8

Нові таги

Таг КІПІ	Таг МТЕ	Опис тагу МТЕ	Слововживання	Приклад
ppron3:sg:gen:f:ter:nakc:praep	Pp-3f--sgy-n	Pronoun Type=personal Person=third Gender=feminine Number=singular Case=genitive Clitic=yes Syntactic_Type=nominal	44	<i>niej</i>
ppron3:sg:gen:f:ter:nakc:praep	Pp-3f--sgasn	Pronoun Type=personal Person=third Gender=feminine Number=singular Case=genitive Clitic=agglutinant Definiteness=short-art Syntactic_Type=nominal		<i>ń</i>
ppron3:sg:acc:f:ter:nakc:praep	Pp-3f--say-n	Pronoun Type=personal Person=third Gender=feminine Number=singular Case=accusative Clitic=yes Syntactic_Type=nominal	11	<i>nią</i>
ppron3:sg:acc:f:ter:nakc:praep	Pp-3f--saasn	Pronoun Type=personal Person=third Gender=feminine Number=singular Case=accusative Clitic=agglutinant Definiteness=short-art Syntactic_Type=nominal		<i>ń</i>

Скорочення тагів КІПІ, статистика

- особові займенники 3-тьої особи (флексема *ppron3* у КІПІ) передбачають **287** різних тагів КІПІ для опису **5** лем з їх **23** формами → **65** тагів МТЕ.
- 1-ша і 2-га особа (флексема *ppron12*); **146** оригінальних тагів КІПІ → **30** тагів МТЕ.
- **433** таги для **42** форм особових займенників у КІПІ зведено до **95** тагів у версії МТЕ
- таги на словоформу: починаючи від форми *nim* з **53** інтерпретаціями в КІПІ, **33** для *nich* і **25** для *nimi* (16 форм з 10 чи більше інтерпретаціями) до *tu*, *jetu*, *ja* з **3** чи **4** інтерпретаціями.

Таги особових займенників жін. роду третьої особи однини у знахідному відмінку

Таг КІПІ	Таг МТЕ	Форма слова
ppron3:sg:acc:f:ter:акс:npraep	Pp-3f--san-n	<i>ja</i>
ppron3:sg:acc:f:ter:акс:praep	Pp-3f--say-n	<i>nia</i>
ppron3:sg:acc:f:ter:nакс:npraep	Pp-3f--san-n	<i>ja</i>
ppron3:sg:acc:f:ter:nакс:praep	Pp-3f--say-n	<i>nia</i>
ppron3:sg:acc:f:ter:npraep	Pp-3f--san-n	<i>ja</i>
ppron3:sg:acc:f:ter:praep	Pp-3f--say-n	<i>nia</i>

Інтерпретація:

Pp-3f--san-n: Pronoun Type=personal Person=third Gender=feminine Number=singular Case=accusative Clitic=no Syntactic_Type=nominal

Pp-3f--say-n: Pronoun Type=personal Person=third Gender=feminine Number=singular Case=accusative Clitic=yes Syntactic_Type=nominal

Поділ на слова у реченнях: *moglibyście*

KIPi

- `<orth>mogli</orth><lex
disamb="1"><base>móc</base><ctag>praet:pl:m1:imperf</ctag></lex><ns/>`
- `<orth>by</orth><lex
disamb="1"><base>by</base><ctag>qub</ctag></lex><ns/>`
- `<orth>ście</orth><lex
disamb="1"><base>być</base><ctag>aglt:pl:sec:imperf:nwok</ctag></lex>`

MTE

- `<w lemma="móc" ana="Vmpis-pmy">mogli</w>
<w lemma="by" ana="Q">by</w>
<w lemma="być" ana="Vapip2p--sa">ście</w>`

Переглянутий поділ на слова:

- `<w lemma="móc" ana="Vmpis2pmy-y">mogliście</w>`
- `<w lemma="móc" ana="Vmpcp3pmy-y">mogliby</w>`
- `<w lemma="móc" ana="Vmpcp2pmy-y">moglibyście</w>`

Фрагмент індексу MSD

Tag MTE	Опис тагу	Лексеми	Приклад
Vmeis2sf--y	Verb Type=main Aspect=perfective VForm=indicative Tense=past Person=second Number=singular Gender=feminine Clitic=yes	85	<i>powiedziałaś/powiedzieć, zrobiłaś/zrobić, przyszłaś/przyjść</i>
Vmeis2sm--y	Verb Type=main Aspect=perfective VForm=indicative Tense=past Person=second Number=singular Gender=male Clitic=yes	274	<i>przyszedeś/przyjść, powiedziadeś/powiedzieć, zrobiadeś/zrobić,</i>
Vmeis2sn--y	Verb Type=main Aspect=perfective VForm=indicative Tense=past Person=second Number=singular Gender=neuter Clitic=yes	1	<i>pozostałoś/pozostać, przeszłoś/przejsć</i>
Vmeis-pf	Verb Type=main Aspect=perfective VForm=indicative Tense=past Person=second Number=plural	619	<i>odbyły/odbyć, rozpoczęły/rozpocząć, zaszły/zajść</i>

Фрагмент лексикону

- absurdami absurd N-mnnp1 17
- absurdem absurd N-mnns1 307
- absurdom absurd N-mnnpd 6
- absurdowi absurd N-mnnsd 4
- absurdu absurd N-mnnsг 578
- absurdy absurd N-mnnpa 59
- absurdy absurd N-mnnpn 58
- absurdzie absurd N-mnns1 17
- absurdów absurd N-mnnpг 163
- aby aby C 201168
- ac ac X 1099
- ach ach I 1170

15 тисяч найбільш вживаних лем витягнуто за допомогою програми Poliqarp

У лексиконі налічується 175848 форм цих слів (приблизно 11.72 форм на лему).

Фрагмент тагованого тексту польського перекладу Джорджа Оруела “1984”

- <p id="Opl.5">
- <s id="Opl.5.1">
- <w lemma="być" ana="Vmpis-sm">Był</w>
- <w lemma="jasny" ana="A-pm--sn">jasny</w>
- <c>,</c>
- <w lemma="zimny" ana="A-pm--sn">zimny</w>
- <w lemma="dzień" ana="N-mnnsa">dzień</w>
- <w lemma="kwietniowy" ana="A-pmn-sa">kwietniowy</w>
- <w lemma="i" ana="C">i</w>
- <w lemma="zegar" ana="N-mnnpn">zegary</w>
- <w lemma="bić" ana="Vmpis-pmn">biły</w>
- <w lemma="trzynasty" ana="Mlof--si">trzynastą</w>
- <c>.</c>
- </s>

Бібліографія

- INTERA unified tagset project www.elda.org/intera
- Hanna Dalewska-Greń *Języki Słowiańskie*, Warszawa, PWN 1997.
- Tomas Erjavec et al. *Multext-East specifications for Slavic languages*, Budapest, 2003.
- Jan Hajič. Positional Tags: Quick Reference (Czech „HM” Morphology), 2000.
- Elena Paskaleva. *Balkan South-East Corpora Aligned to English*. In: *The Proceedings of the Workshop on Common Natural Language Processing Paradigm for Balkan Languages*, EACL 2007
- Adam Przepiórkowski and Marcin Woliński. [A Flexemic Tagset for Polish](#). In: *The Proceedings of the Workshop on Morphological Processing of Slavic Languages*, EACL 2003.
<http://nlp.ipipan.waw.pl/~adamp/Papers/2003-eacl-ws12/ws12.pdf>
- Широков В.А et al. *Корпусна лінгвістика*. Київ: Довіра, 2005.

Документація до PoIUKR

статті і слайди доступні зі сторінки <http://www.domeczek.pl/~natko/>

[The Current State of Work on the Polish-Ukrainian Parallel Corpus \(PoIUKR\).](#)

Proceedings of the International Workshop within MONDILEX project "Problems of Slavic Lexicography" Kyiv, 2-4 February 2009.

[Морфосинтаксичне тагування польсько-українського паралельного корпусу](#)

[\(PoIUKR\).](#) Proceedings of the International Conference “MegaLing'2008. Horizons of Applied Linguistics and Linguistic Technologies, Parthenit – Crimea, Ukraine, 20-27 September 2008”.

Ivan Derzhanski and Natalia Kotsyba. [Towards a Consistent Morphological Tagset for Slavic Languages: Extending MULTEXT-East for Polish, Ukrainian and Belarusian.](#) Metalanguage and Encoding Scheme Design for Digital Lexicography: MONDILEX Third Open Workshop, Bratislava, 15–16 April 2009.

Natalia Kotsyba, Andriy Mykulyak, Igor V. Shevchenko. [UGTag: morphological analyzer and tagger for Ukrainian language.](#) Proceedings of the international conference [Practical Applications in Language and Computers \(PALC 2009\)](#), Łódź, 6-8 April 2009.

Natalia Kotsyba, Adam Radziszewski and Ivan Derzhanski. [Integrating the Polish language into the MULTEXT-East family: morphosyntactic specifications, converter, lexicon and corpus. \(presentation slides\).](#) Proceedings of [Research Infrastructure for Digital Lexicography: MONDILEX Fifth Open Workshop](#), October 14, 2009, Ljubljana, Slovenia, within [INFORMATION SOCIETY 2009](#) 12th International multiconference.