

Corpus of frequency dictionary of contemporary Polish (1963-67) – new annotation

Łukasz Szałkiewicz

Institute of Computer Science, Polish Academy of Sciences

lukasz.szalkiewicz@ipipan.waw.pl

1. About Corpus

The original purpose of the corpus was to create a general frequency dictionary of contemporary Polish. The work started in 1967. Partial results were published between 1972 and 1977, the completed dictionary in 1990. The corpus was later augmented in various respects, both by manual editing and auto-

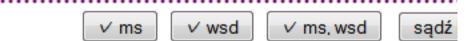
3. Screen from Anotatornia (tool for corpus annotation)

ukasz.szalkiewicz (superanotator) Anotatornia NKJP (ркодиксузма (z WSD), 8003) Strona główna Zweryfikowane Bieżąca transza Sł. sensów Zmiana hasła Wyloguj

Poziomy anotacji segmentacja granice zdań morfoskładnia sensy słów

412 (tr. 83, at. 818). Zachmurzenie duże z lokalnymi większymi przejaśnieniami i miejscami niewielkie opady deszczu ze śniegiem . ■ Od północnego zachodu ochłodzenie . ■ Temperatura w nocy od około , minus trzech stopni w części północno - zachodniej do około zera w centrum i około plus trzech stopni na południowym wschodzie . ■ Temperatura maksymalna odpowiednio od około minus trzech do około plus dwóch i około plus pięciu stopni . ■

seg:**zwer** | sn:**zwer** | ms:**doos** | wsd:**doos** (kf, 2011-04-10 19:52:08) (s.zurowski, tr. 84, at. 828



mated procedures.

Corpus data contain 10,000 samples divided into 5 parts:

•essays,

• news,

• scientific texts,

fiction

• plays.

Every sample is approximately 50 words long, they all come from texts published between 1963 and 1967 and contain bibliographic description of its source. Each word is tagged with its base form and some morphological properties. Sentence boundaries are also marked.

In 2001 corpus authors agreed to publish the data in the Internet under GNU licence.

2. New linguistic annotation

• ZACHMURZENIE subst:sg:non		zachmurzenie zachmurzenie subst:sg:nom:n <u>wybierz</u> <u>dodaj</u> 🛛
 duże duży adj:sg:nom:n:pos wyl 	<u>bierz</u> dodaj 🛛	duże puży adj:sg:nom:n:pos <mark>wybierz</mark> dodaj 🛛
 z z prep:inst:nwok wybierz de 	odaj 🛛	z z prep:inst:nwok <mark>wybierz</mark> dodaj ⊠
 lokalnymi LOKALNY adj:pl:inst:n:pos 	<u>vybierz</u> dodaj ⊠	lokalnymi LOKALNY adj:pl:inst:n:pos <u>wybierz</u> dodaj 🛛
 większymi wielki adj:pl:inst:n:com w 	<u>ybierz dodaj ⊠</u>	większymi <mark>duży adj:pl:inst:n:com <mark>wybierz</mark> dodaj</mark> 🛛
 przejaśnieniami przejaśnienie subst:pl:inst 	:n <u>wybierz</u> <u>dodaj</u> 🛛	przejaśnieniami przejaśnienie subst:pl:inst:n <u>wybierz</u> <u>dodaj</u> 🛙
 <i>i</i> I conj wybierz dodaj ☑ 		i I conj <mark>wybierz</mark> dodaj ⊠
 miejscami мтезясе subst:pl:inst:n wyb 	<u>ierz dodaj</u> 🛛	miejscami мтезясе subst:pl:inst:n <mark>wybierz</mark> dodaj 🛛
 niewielkie NIEWIELKI adj:pl:nom:m3:p 	os <u>wybierz</u> dodaj 🛛	niewielkie NIEWIELKI adj:pl:nom:m3:pos <u>wybierz</u> dodaj 🛛
 opady opad subst:pl:nom:m3 wyl 	bierz dodaj 🛛	<i>opady</i> орад subst:pl:nom:m3 <u>wybierz</u> <u>dodaj</u> 🛛
• doczezu		doczczu

4. Using original annotation – tagset changes

Manual annotation is costly task in the development of a corpus. One way to reduce the cost was to use data from original annotation as outcome of work imaginary annotator, called 'KF' (in Polish Korpus Frekwencyjny). So original annotation was treated as if it were made by the first annotator. But these data were congruent with IPI PAN tagset not with NKJP tagset. Some automatical changes were needed to prevent repeated discrepancies between KF and second, real annotator. There were created **tagsets overview of the** differences and it was transformed in a three-part list: what can be done automatically, semi-automatically and manually. Examples:

NKJP tagset		IPI PAN tagset		
word	tag	tag	automatic change for all occurrences	
że	comp (or qub, but much less)	conj	comp	
stąd	adv	qub	adv	
tu	adv	qub	adv	
dlaczego	adv	qub	adv	
we	prep:wok	prep	prep:wok	
ponad	prep or qub	prep	prep, if next is num	

5. Additional works

• there were no ready **sentence segmentation** in corpus data, so it was added auto-

New annotation – according to rules established in annotation of National Corpus of Polish (NKJP).

- 1. All corpus data has been manually verified by two independent annotators (but see next paragraph), connecting with the server via a web interface.
- 2. Any conflicts have been resolved by a referee (so-called 'superannotator').
- 3. The linguistic annotation was carried out at three levels: word-level segmentation, sentence-level segmentation and morphosyntax.
- 4. There was used the same purpose-built tool, Anotatornia.
- 5. NKJP tagset and the same guidelines.
- NKJP tagset differentiates between coordinate conjunctions (Pol. spójniki równorzędne; conj), e.g. *i, lub, oraz,* and subordinate conjunctions (Pol. spójniki podrzędne), sometimes called complementisers (comp), e.g. *że, aby, bowiem*.
- There where the class of particle-adverbs in IPI PAN tagset, separate from the class of adverbs (adv) which consists of only de-

matically according to Marcin Milkowski's rules

- some minor changes in the texts of the corpus , such as the removal of double dots in the text of the drama
- methods of changing the bases to begin with a capital letter were developed, because in the original annotation all saved lemmas were always lower case, resulting in many discrepancies

6. Bibliography

You can see the most current data, corpus documentation and selected bibliography on http://clip.ipipan.waw.pl/PL196x

7. People

Anna Andrzejczuk (annotator, tagsets conversion)

adjectival or gradable adverbs. The class of adverbs is larger in NKJP tagset than in IPI PAN, also includes traditional adverbs which are neither de-adjectival nor gradable.

Michał Lenart (computer engineer)
Łukasz Szałkiewicz (superannotator, tagsets conversion)
Sebastian Żurowski (main annotator)