



Using language corpora and wordnets in the translator's work: pros and cons

Łukasz Grabowski
University of Opole (Poland)

Translating Europe Workshop
23.05.2018, Warsaw

Plan of my talk

- Introduction
 - Language corpora
 - Wordnets
- Using language corpora in the translator's work
 - Parallel corpora
 - Monolingual reference corpora
- Using wordnets in the translator's work
 - Polish wordnet (*Słownosieć*) mapped onto Princeton WordNet of English
- Conclusions

What are language corpora?

- Custom-designed **collections of authentic linguistic data (texts)**, either **written texts** or **transcriptions of recorded spoken data**
- Compiled according to a number of criteria
 - Size
 - Text availability / copyright
 - Representativeness
 - Balance
 - Computer-readable format

Language corpora for translators

- Language corpora specifically designed for translators may include:
 - **originals and translations**
 - **Parallel corpora**
 - proved to be of great value for translators (Google Translate, TAUS Data Association, Linguee, Reverso Context etc. are all built on parallel corpora)
 - **Professional translators are more familiar with parallel search** (or parallel concordancing) **feature of translation memory systems (TMs)**
 - **non-translations** (i.e. native texts written originally in a given language)
 - **monolingual reference corpora**
 - **texts in different languages**
 - **originals and translations in the same language, or**
 - **only native texts originally produced in different languages)**
 - **bilingual and monolingual comparable corpora**
 - primarily used for research purposes (Descriptive Translation Studies)

Why should translators bother with language corpora?

- **Corpus as a source of linguistic information** helps to answer various questions related to how we use language
 - across genres, text types and registers
 - **adherence to discourse norms and conventions of text production in the target language** and culture is essential in translation
 - this information is implicitly recorded in texts collected in language corpora
 - particularly important when translators deal with **creative and stylistically varied texts** (rather than cliched and highly-repetitive ones)
- Corpora enable the **translators to go beyond their linguistic intuition**
 - **linguistic preferences of translators**, in particular those rendering texts into a target language which is not their native one, **may not always coincide with the native speakers' preferences**
- Corpora may help **activate the translators' memory** of relevant contexts of use of words and expressions (Pezik 2017)
 - Lexical priming theory (Hoey 2005)

Monolingual reference corpora

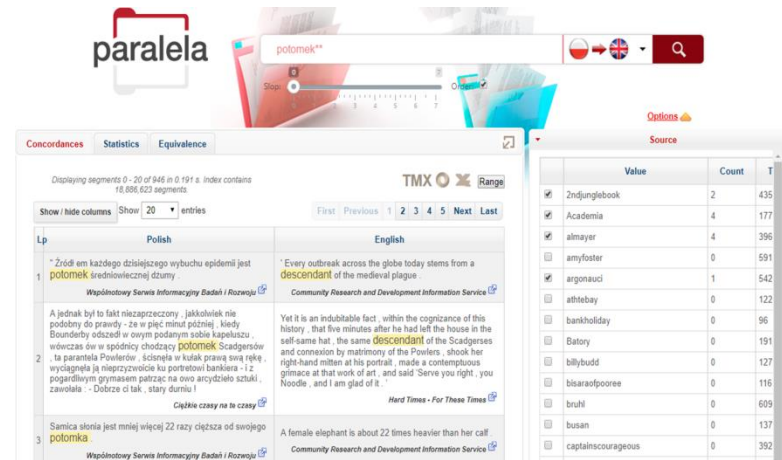
Examples

- *British National Corpus*
- *Corpus of Contemporary American English*
- *Narodowy Korpus Języka Polskiego*
- *Национальный корпус русского языка*
- *English Web 2015 (enTenTen15)*
- *Polish Web 2012 (plTenTen12)*
- *Russian Web 2011 (ruTenTen11)*

- Monco (<http://monitorcorpus.com/>)
 - a collection of English texts updated on a regular basis (circa 11 billion words, as of 9 May 2018)
- Frazeo (<http://frazeo.pl>)
 - an innovative search engine of press articles published on Polish news portals

Paralela

(Pęzik 2016)

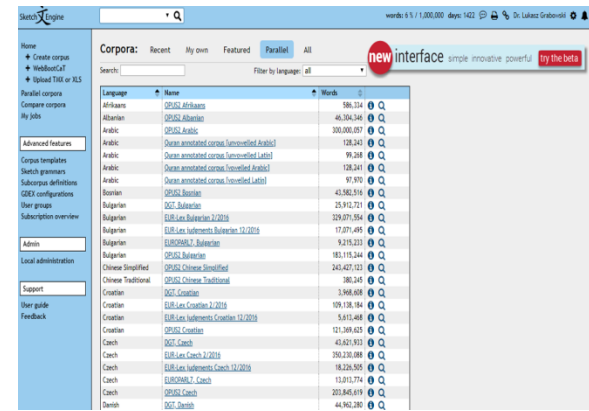


- **An English-Polish and Polish-English parallel corpus**
 - It includes **262 million words in 10,877,000 translation segments** found predominantly in legal texts
 - EU legislation, proceedings of the European Parliament etc.), press releases, medical texts (provided by the EMA) as well as film subtitles
- Translation segments are aligned at the sentence level (Pęzik 2016: 70), with **5.3% of the segments aligned manually.**

Everything in one place

SketchEngine (Kilgarriff et al. 2014)

<https://www.sketchengine.co.uk/>



The screenshot shows the SketchEngine web interface. The main content area displays a table of corpora with columns for Language, Name, and Words. The table lists various corpora for different languages, including Arabic, Bulgarian, Chinese, Croatian, Czech, and Dutch. The interface also includes a search bar, a filter for language, and a sidebar with navigation options like Home, Create corpus, and Upload TSV or XLS.

Language	Name	Words
Afrikaans	SPS2_Afrikaans	586,314
Arabic	SPS2_Arabic	46,346,546
Arabic	SPS2_Arabic	300,000,007
Arabic	Duqa annotated corpus (unsorted Arabic)	128,243
Arabic	Duqa annotated corpus (unsorted Latin)	99,268
Arabic	Duqa annotated corpus (unsorted Arabic)	128,243
Arabic	Duqa annotated corpus (unsorted Latin)	99,268
Bosnian	SPS2_Bosnian	43,562,516
Bulgarian	DGZ_Bulgarian	25,912,721
Bulgarian	EUR-Lex Bulgarian 1 (2019)	329,071,054
Bulgarian	EUR-Lex Instruments Bulgarian 12 (2019)	17,071,495
Bulgarian	EURPARAC-Bulgarian	9,219,233
Bulgarian	SPS2_Bulgarian	183,115,244
Chinese Simplified	SPS2_Chinese_Simplified	243,427,123
Chinese Traditional	SPS2_Chinese_Traditional	386,249
Croatian	DGZ_Croatian	3,948,608
Croatian	EUR-Lex Croatian 1 (2019)	109,138,104
Croatian	EUR-Lex Instruments Croatian 12 (2019)	5,615,488
Croatian	SPS2_Croatian	121,369,635
Czech	DGZ_Czech	43,625,933
Czech	EUR-Lex Czech 1 (2019)	20,326,088
Czech	EUR-Lex Instruments Czech 12 (2019)	18,228,205
Czech	EURPARAC-Czech	13,013,774
Czech	SPS2_Czech	203,846,619
Dutch	DGZ_Dutch	44,940,330

- A software platform for text analysis and text mining applications
- Access to **aligned parallel corpora** of various text types and genres
- Access to many **large monolingual corpora** with texts extracted from the web
 - **Term extraction** (from subject-specific texts uploaded by users) and **bilingual term extraction** (performed on TMs uploaded by users)
- A possibility to create **subject-specific corpora**
 - By **uploading materials** provided by users
 - By **providing keywords** and **extracting relevant texts automatically** from the web

Example 1

- **Multi-word terms** (nominalizations) extracted from a custom-designed corpus of **car reviews** (Top Gear)
 - 100 texts (60,372 words)

Multi-word	Score	F	RefF
<input type="checkbox"/> inside layout	W 1,385.85	<u>100</u>	<u>0</u>
<input type="checkbox"/> all-wheel drive	W 230.43	<u>17</u>	<u>3</u>
<input type="checkbox"/> hot hatch	W 206.09	<u>16</u>	<u>9</u>
<input type="checkbox"/> driving position	W 167.12	<u>15</u>	<u>28</u>
<input type="checkbox"/> four-wheel drive	W 155.95	<u>17</u>	<u>58</u>
<input type="checkbox"/> six-speed manual	W 153.33	<u>11</u>	<u>0</u>
<input type="checkbox"/> city car	W 149.45	<u>11</u>	<u>3</u>
<input type="checkbox"/> cruise control	W 139.78	<u>11</u>	<u>11</u>
<input type="checkbox"/> boot space	W 131.34	<u>10</u>	<u>7</u>
<input type="checkbox"/> climate control	W 128.08	<u>10</u>	<u>10</u>
<input type="checkbox"/> turbo engine	W 119.31	<u>9</u>	<u>6</u>
<input type="checkbox"/> automatic gearbox	W 117.61	<u>10</u>	<u>21</u>
<input type="checkbox"/> front-wheel drive	W 116.33	<u>9</u>	<u>9</u>
<input type="checkbox"/> 6-litre diesel	W 111.79	<u>8</u>	<u>0</u>
<input type="checkbox"/> turbo petrol	W 111.79	<u>8</u>	<u>0</u>
<input type="checkbox"/> petrol engine	W 108.42	<u>12</u>	<u>61</u>
<input type="checkbox"/> fuel economy	W 104.16	<u>12</u>	<u>68</u>
<input type="checkbox"/> steering wheel	W 102.21	<u>19</u>	<u>178</u>
<input type="checkbox"/> body roll	W 101.90	<u>8</u>	<u>11</u>
<input type="checkbox"/> wind noise	W 100.26	<u>8</u>	<u>13</u>
<input type="checkbox"/> infotainment system	W 97.94	<u>7</u>	<u>0</u>
<input type="checkbox"/> 4-litre petrol	W 97.16	<u>7</u>	<u>1</u>
<input type="checkbox"/> driving experience	W 97.12	<u>8</u>	<u>17</u>
<input type="checkbox"/> engine range	W 96.30	<u>7</u>	<u>2</u>
<input type="checkbox"/> manual gearbox	W 93.81	<u>7</u>	<u>5</u>
<input type="checkbox"/> 2-litre diesel	W 84.09	<u>6</u>	<u>0</u>
<input type="checkbox"/> infotainment screen	W 84.09	<u>6</u>	<u>0</u>
<input type="checkbox"/> trim level	W 83.42	<u>6</u>	<u>1</u>
<input type="checkbox"/> material quality	W 83.42	<u>6</u>	<u>1</u>
<input type="checkbox"/> centre console	W 79.86	<u>6</u>	<u>6</u>
<input type="checkbox"/> base model	W 79.86	<u>6</u>	<u>6</u>
<input type="checkbox"/> torsion beam	W 70.24	<u>5</u>	<u>0</u>
<input type="checkbox"/> km co2	W 70.24	<u>5</u>	<u>0</u>
<input type="checkbox"/> adaptive cruise	W 70.24	<u>5</u>	<u>0</u>
<input type="checkbox"/> radar cruise	W 70.24	<u>5</u>	<u>0</u>

<input type="checkbox"/> three-cylinder engine	W 70.24	<u>5</u>	<u>0</u>
<input type="checkbox"/> boot floor	W 70.24	<u>5</u>	<u>0</u>
<input type="checkbox"/> rear headroom	W 70.24	<u>5</u>	<u>0</u>
<input type="checkbox"/> plug-in hybrid	W 70.24	<u>5</u>	<u>0</u>
<input type="checkbox"/> blind-spot warning	W 70.24	<u>5</u>	<u>0</u>
<input type="checkbox"/> head-up display	W 69.68	<u>5</u>	<u>1</u>
<input type="checkbox"/> front end	W 69.26	<u>10</u>	<u>114</u>
<input type="checkbox"/> 5-litre diesel	W 69.07	<u>5</u>	<u>2</u>
<input type="checkbox"/> six-speed gearbox	W 67.87	<u>5</u>	<u>4</u>
<input type="checkbox"/> company car	W 67.49	<u>11</u>	<u>143</u>
<input type="checkbox"/> ride quality	W 67.28	<u>5</u>	<u>5</u>
<input type="checkbox"/> standard car	W 66.14	<u>5</u>	<u>7</u>
<input type="checkbox"/> company car tax	W 65.59	<u>5</u>	<u>8</u>
<input type="checkbox"/> sweet spot	W 64.50	<u>5</u>	<u>10</u>
<input type="checkbox"/> list price	W 63.85	<u>7</u>	<u>60</u>
<input type="checkbox"/> rear axle	W 58.34	<u>5</u>	<u>23</u>
<input type="checkbox"/> old car	W 58.12	<u>7</u>	<u>77</u>
<input type="checkbox"/> electric power	W 57.87	<u>6</u>	<u>51</u>
<input type="checkbox"/> cross-traffic alert	W 56.39	<u>4</u>	<u>0</u>
<input type="checkbox"/> eight-speed auto	W 56.39	<u>4</u>	<u>0</u>
<input type="checkbox"/> air con	W 56.39	<u>4</u>	<u>0</u>
<input type="checkbox"/> safety kit	W 56.39	<u>4</u>	<u>0</u>
<input type="checkbox"/> collision mitigation	W 56.39	<u>4</u>	<u>0</u>
<input type="checkbox"/> three-cylinder petrol	W 56.39	<u>4</u>	<u>0</u>
<input type="checkbox"/> electric power steering	W 56.39	<u>4</u>	<u>0</u>
<input type="checkbox"/> rear suspension	W 56.24	<u>5</u>	<u>28</u>
<input type="checkbox"/> standard kit	W 55.95	<u>4</u>	<u>1</u>
<input type="checkbox"/> supple ride	W 55.95	<u>4</u>	<u>1</u>
<input type="checkbox"/> driving environment	W 55.45	<u>4</u>	<u>2</u>
<input type="checkbox"/> 5-litre turbo	W 55.45	<u>4</u>	<u>2</u>
<input type="checkbox"/> handling balance	W 54.96	<u>4</u>	<u>3</u>
<input type="checkbox"/> rear-wheel drive	W 54.96	<u>4</u>	<u>3</u>
<input type="checkbox"/> 6-litre engine	W 54.96	<u>4</u>	<u>3</u>
<input type="checkbox"/> active safety	W 54.49	<u>4</u>	<u>4</u>

<input type="checkbox"/> five-year warranty	W 54.49	<u>4</u>	<u>4</u>
<input type="checkbox"/> rear legroom	W 54.02	<u>4</u>	<u>5</u>
<input type="checkbox"/> low speed	W 53.95	<u>5</u>	<u>34</u>
<input type="checkbox"/> standard equipment	W 51.50	<u>5</u>	<u>41</u>
<input type="checkbox"/> electric motor	W 51.50	<u>5</u>	<u>41</u>
<input type="checkbox"/> turning circle	W 50.17	<u>4</u>	<u>14</u>
<input type="checkbox"/> drive system	W 50.17	<u>4</u>	<u>14</u>
<input type="checkbox"/> car tax	W 49.64	<u>6</u>	<u>78</u>
<input type="checkbox"/> big car	W 49.22	<u>5</u>	<u>48</u>
<input type="checkbox"/> automatic transmission	W 45.81	<u>4</u>	<u>26</u>
<input type="checkbox"/> big news	W 45.15	<u>4</u>	<u>28</u>
<input type="checkbox"/> small car	W 44.51	<u>5</u>	<u>65</u>
<input type="checkbox"/> diesel engine	W 43.89	<u>6</u>	<u>103</u>
<input type="checkbox"/> front passenger	W 43.61	<u>4</u>	<u>33</u>
<input type="checkbox"/> driver assist	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> reversing camera	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> practical car	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> dual-zone climate	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> instrument cluster	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> adaptive cruise control	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> dual-zone climate control	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> radar cruise control	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> small crossover	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> lane departure	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> autonomous emergency	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> five-door hatch	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> central infotainment	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> crossover market	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> seven-seat mpv	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> 5-litre petrol	W 42.55	<u>3</u>	<u>0</u>
<input type="checkbox"/> supple suspension	W 42.55	<u>10</u>	<u>0</u>

Example 2

- **Bilingual wordsketches**

- **matched headwords and collocations** extracted from parallel corpora

- *Declaration/deklaracja* preceded by modifiers in EUROPARL7, English and Polish

declaration (noun) EUROPARL7, English freq = 5,532 (91.07 per million) **deklaracja** (noun) EUROPARL7, Polish freq = 1,620 (106.77 per million)

Use another candidate translation: [pisemny oświadczenie](#) [patrzeć niepodległość](#) [wpisać protokół](#) [składać protokół](#) [rejestr](#)
 Click on collocates to access reciprocal bilingual search or find [translated collocations](#)

modifier		52.62	a modifier		68.58
written	236	10.46	powszechny	29	10.27
written declaration			Powszechnej Deklaracji Praw Człowieka		
Written	76	9.62	paryski	21	9.18
: see Minutes			deklaracji paryskiej z		
unilateral	77	9.16	jednostronny	26	9.12
unilateral declaration of independence			jednostronnej deklaracji niepodległości Kosowa		
joint	258	9.08	pisemny	47	9.09
a joint declaration			pisemną deklarację		
solemn	53	9.08	głosotłowny	13	8.56
solemn declarations			głosotłowne deklaracje		
Lapsed	45	8.96	berliński	11	8.16
: see Minutes			deklaracji berlińskiej		
Laeken	35	8.22	końcowy	14	7.63
the Laeken declaration			deklaracji końcowej		
mere	36	8.01	pusty	8	7.56
mere declarations of			krajowy	58	7.47
grand	19	7.55	krajowych deklaracji w sprawie zarządzania		
the grand declarations			manhattański	6	7.45
final	79	7.50	wspólny	96	7.42
the final declaration			wspólnej deklaracji		
open	66	7.07	praski	6	7.42
the open declaration			jednoznaczny	8	7.29
formal	23	7.06	otwarty	12	7.25
a formal declaration			„otwartej deklaracji”		
Doha	17	6.98	samoświadczy	5	7.16
of the Doha declaration			uroczysty	5	7.07
non-binding	12	6.97	boloński	5	7.07
non-binding declarations			polityczny	60	6.85
fine	19	6.85	falszywy	5	6.82
fine declarations of			dauhański	5	6.79
ministerial	12	6.62	jasny	14	6.78
the ministerial declaration			milenijny	6	6.78
secret	13	6.55	wyraźny	11	6.66
secret declarations			ministerialny	4	6.66

- We can see which collocations for English *declaration* do not have their Polish equivalent with *deklaracja*

Example 2.1

- If we click on collocates, we obtain access to translated collocations

written (*adjective*) Alternative PoS: [write verb](#) (freq: 4,279)
EUROPARL7, English freq = [2,431](#) (40.02 per million)

pisemny (*adjective*) EUROPARL7, Polish freq = [1,700](#) (112.05 per million)

Use another candidate translation: [142 oświadczenie regulamin 149 art wyjaśnienie](#)) ([patrzeć](#)
Click on collocates to access reciprocal bilingual search or find [translated collocations](#)

modifies		98.03		modifies		99.76		and/or		6.50		
statement	1,041	10.10	oświadczenie	1,285	12.95	verbal	6	8.96	oral	29	7.89	
. Written statements (Rule 149		. Oświadczenia pisemne (art .				written and oral						
explanation	154	9.23	wyjaśnienie	123	10.21	express	4	7.86	prior	3	6.22	
. Written explanations of vote in		. Pisemne wyjaśnienia dotyczące sposobu głosowania				audiovisual	3	6.08	formal	4	5.50	
declaration	236	9.17	deklaracja	47	8.79	detailed	4	4.70	parliamentary	7	4.58	
written declaration		pisemną deklarację				original	3	4.49	annual	3	3.84	
answer	171	8.25	tłumaczenie	19	8.26	comprehensive	3	3.83	urgent	3	3.71	
will receive written answers		tłumaczeń pisemnych				s	3	2.36				
reply	74	8.09	zapytanie	9	7.32							
a written reply		odpowiedź		37	7.20							
STATEMENT	11	7.23	pisemnej odpowiedzi									
. WRITTEN STATEMENT (RULE 120		pytanie		37	6.60							
press	33	7.01	pytanie pisemne									
the written press ,		forma		14	6.35							
Statement	7	6.48	w formie pisemnej									
consent	11	6.24	wersja	6	6.00							
written consent		kwestionariusz		3	5.82							
questionnaire	6	6.15	potwierdzenie	3	5.70							
confirmation	7	6.05	zapis	4	5.62							
translation	8	5.90	zgoda	4	5.15							
question	150	5.71	opinia	4	4.14							
a written question		dokument		3	3.71							
Declaration	9	5.63	poparcie	3	3.68							
constitution	12	5.53	rezolucja	4	3.62							
written constitution		wniosek		7	3.53							
notification	5	5.52	procedura	3	3.51							
exam	3	5.29	sprawozdanie	5	2.30							
version	9	5.11										
assurance	5	4.95										
record	7	4.86										
submission	3	4.82										
summary	3	4.79										
authorisation	4	4.71										
form	27	4.68										

When do translators typically use language corpora? Or should use?

- **Monolingual reference corpora** (e.g. *Narodowy Korpus Języka Polskiego*)
 - When we are worried that **the translation does not fit the norms and conventions (grammatical, stylistic etc.) of the target language** (Polish language)
 - Excessive **lexical or syntactic calques** and **idiosyncratic collocations** (*massage baths – wanny masażowe, computer key – guzik komputerowy* etc.), **overuse or underuse of certain grammatical structures** (e.g. passive voice in Polish) or **prefabricated formulas** etc.
- **Parallel corpora**
 - When we want to see **how translators dealt with similar translational problems in the past** (Zanettin 2003), and
 - to verify translational choices / equivalents (textual)
 - to verify dictionary equivalents

Deleveraging – delewarowanie?

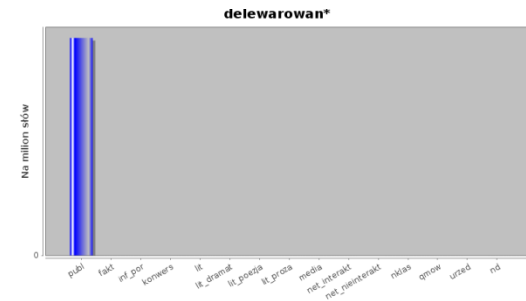
- *Most notably, in the wake of the Lehman collapse, the pattern of **deleveraging** and forced sales has been so intense that traditional price relationships have completely broken down, sending trading models haywire.*
- *W szczególności, po upadku Lehmanów, **delewarowanie** i wymuszone sprzedaże tak się nasiliły, że doszło do całkowitego załamania się tradycyjnych stosunków cenowych i rozchwiania modeli prowadzenia interesów.*
 - Financial world is stumbling blindly through a cognitive fog, Gillian Tett, published 6 March 2009. At: <http://www.ft.com/cms/s/0/da26af7c-09ee-11de-add8-0000779fd2ac.html>
 - Świat finansów potyka się i błądzi. At: <http://ft.onet.pl>

Deleveraging – delewarowanie?

- Data from Paralela
 - 12 occurrences of *deleveraging* (10 in RAPID subcorpus (EU Press Releases), 1 in EPP proceedings, 1 in NBP subcorpus)
 - Polish equivalents:
 - *oddłużanie*
 - *ograniczanie zadłużenia*
 - *pozbywanie się pasywów*
 - *eliminacja dźwigni finansowej*
 - ***delewarowanie (5 occurrences)*** ✓
 - *zrównoważanie (budżetów domowych) - deleveraging (of households ' balance sheets)*

delewarowanie

- Data from NKJP
 - 10 occurrences of *delewarowanie*



Przeszukiwany zbiór zawiera 240,192,461 słów. Znalaziono 10 akapitów pasujących do zapytania w 0.003s. Bieżąca strona zawiera 13 przykłady z 9 różnych tekstów.

1.	gdy nadmuchany balon pękł, rozpoczął się odwrotny proces -	delewarowania - czyli zmniejszania się rynku aktywów finansowych, głównie	Polityka	+
2.	minionych kilku lat prowadzą dziś intensywny proces	delewarowania , czyli spisywania na straty części aktywów, przy	Polityka	+
3.	przebiegało w trzech fazach: fazie stabilizacji, fazie	delewarowania oraz długotrwałej fazie zredukowanego zadłużenia i	Gazeta Ubezpieczeni...	+
4.	wzrost produktu krajowego brutto (PKB) w miarę postępującego	delewarowania portfeli konsumentów, wzrostu bezrobocia, stagnacji cen	Gazeta Ubezpieczeni...	+
5.	starań rządów i banków centralnych, biorąc pod uwagę procesy	delewarowania w bankach, firmach i gospodarstwach domowych, nie ma zbyt	Gazeta Ubezpieczeni...	+
6.	finansowej z sektora prywatnego do publicznego. Skala	delewarowania potwierdza, że potencjalny wzrost w większości krajów	Gazeta Ubezpieczeni...	+
7.	pełni dojrzeć do wzrostów. Sądzimy zatem, że trwający proces	delewarowania i wysoki stopień wzajemnych zależności ekonomicznych w skali	Gazeta Ubezpieczeni...	+
8.	azjatyckich w połączeniu z trwającym w tej chwili procesem	delewarowania w USA mogą w przyszłych latach wyprowadzić Azję na ścieżkę	Gazeta Ubezpieczeni...	+
9.	minęła pierwsza panika, w wieżowcach na Manhattanie ruszyło	delewarowanie . W czasach hossy ulubionym narzędziem bankierów z Wall	Polityka	+
10.	wskutek utraty zaufania do instytucji finansowych, z drugiej	delewarowanie , czyli pozbywanie się nadmiernego zadłużenia przez banki i	Polityka	+
11.	ale póki trwa, podgrzewa koniunkturę w gospodarce.	Delewarowanie jest niezbędne, ale procesowi temu towarzyszyć muszą	Polityka	+
12.	rządy na ratowanie sektora bankowego bilionów dolarów - bo	delewarowanie to proces, którego skala nie idzie w biliony, ale grube	Polityka	+
13.	o wysokiej rentowności i kredytów bankowych. Błyskawiczne	delewarowanie systemu finansowego także miało swój negatywny wpływ na	Gazeta Ubezpieczeni...	+

Delewarowanie

- Data from plTenTen12
 - 60 occurrences of *delewarowanie*

Polish Web 2012 (plTenTen12, RFTagger)

Query (delewarowanie)-n 60 (0.01 per million) ⓘ

Page 1 of 3 Go Next | Last

blogbank.p... poprzednio) czy pozwolić na naturalne **delewarowanie** . Stymulowanie za wszelką cenę i jego skutki

blogbank.p... wczoraj w dół. Dane nieciekawe, może być dalsze **delewarowanie** na dolarze. Za dużo niewiadomych zatem AT i z

blogbank.p... wypłynęła na nim. Teraz mamy powszechne **delewarowanie** (czynnik działający deflacyjnie), a Grecja

blogbank.p... banki", bo była wcześniej, od lat przeżywa " **delewarowanie** "... Już myślałem, że ze swym postem się nieco

mises.pl , ze komunikat jest w porządku. Murphy: Czy **delewarowanie** ma negatywny wpływ na gospodarkę? Autor:

obserwator... instytucje finansowe będą prowadzić **delewarowanie** zbyt szybko i na zbyt dużą skalę. Na razie władze

obserwator... budzi nie tylko to, że prowadzone w Europie **delewarowanie** może negatywnie wpłynąć na notowania

bossa.pl oszczędności o około 3 p.p. PKB przy czym **delewarowanie** dotyczyło jedynie sektora przedsiębiorstw -

blogbank.p... co dzieje się ze złotem i srebrem bez dodruku = **delewarowanie**) a nie interwencjonizm w realną gospodarkę:

ifin24.pl ... przy zwiększającym się zadłużeniu (niestety, **delewarowanie** będzie długie i bolesne). "Zamożność" jest w

trystero.p... racjonalnego zarządzania, jak się rozpocznie **delewarowanie** bezrobocie skoczy do 25%, bo gdzie indziej w

obserwator... do 3-procentowego deficytu względem PKB oraz **delewarowanie** i rekapitalizacja sektora finansowego.

obserwator... . Jeszcze kilka lat restrukturyzacji Tak zwane **delewarowanie** - czyli zmniejszanie stosunku aktywów do

internacjo... kanałami. W języku Banku Światowego "globalne **delewarowanie** " (masowa wyprzedaż aktywów firm celem

trystero.p... w nieruchomościach a więc majątek w dół, dalsze **delewarowanie** , mniejsza konsumpcja. I tak praktycznie z

ibs.edu.pl... Obecnie wariant ekstremalny jest możliwy, bo **delewarowanie** objęło cały glob. W moim szkicu są dwa punkty

pkpplwiat... recesją u naszych partnerów eksportowych w UE " **delewarowanie** " w UE i potrzeba sanacji banków, co pociągnie za

gb.pl funduszy i ich jakości, jak też poprzez **delewarowanie** . Konieczność szybkiego podniesienia

wp.pl banków jest nadal pilna". Dodał też, że " **delewarowanie** jest potencjalnie bardziej niebezpieczne w

wp.pl i uchroni nas przed kryzysem . a teraz **delewarowanie** jest potencjalnie bardziej niebezpieczne w

Page 1 of 3 Go Next | Last

Sprovokować poprawę?

- Example from Paralela (JRC-Acquis sub-corpus)
 - *Another important issue is then discussed — the forced sale of mortgaged property . The GP notes that this area is characterised by a large number of procedures , timeframes and costs and that this would hinder cross-border lending activity. It therefore suggests a gradual approach to encourage improvements in forced sales procedures: (...)*
 - *Inną ważną kwestią omówioną w dalszej kolejności jest wymuszona sprzedaż nieruchomości obciążonej hipoteką . ZK zaznacza , że obszar ten charakteryzuje duża liczba procedur , okresów i kosztów i stanowiłoby to przeszkodę dla transgranicznej działalności kredytowej . Sugeruje zatem zastosowanie stopniowego podejścia do sprovokowania poprawy w zakresie procedur wymuszonej sprzedaży*

Sprovokować poprawę

- Data from NKJP
 - No occurrences of *sprovokować poprawę*
- Data from Frazeo
 - No occurrences of *sprovokować poprawę*
- Data from plTenTen12
 - 2 occurrences of *sprovokować poprawę* (in 7.7 billion words)
 - Jakość mediów – tych drukowanych i tych elektronicznych – jego zdaniem, dramatycznie spada. Na dodatek, nie widać żadnego mechanizmu w mediach prywatnych, który mógłby **sprovokować poprawę** sytuacji. (Source: <http://www.dziennikarzerp.pl/2010/12/kulawa-debata-publiczna/>)
 - Dodam jeszcze, że zamiast ubolewać, co jest nie tak, trzeba było samemu **sprovokować poprawę** - ale nie przez narzekanie, tylko przez rzucenie konkretnymi propozycjami- co i jak poprawić. (Source: <http://www.forum.fantazyzone.pl/printview.htm?t=51&start=0&sid=63ee71cac16f6f8ccdd3d53d34eb6104>)

- Data from pITenTen12



– *Przynieść poprawę, sprzyjać poprawie, (s)powodować poprawę, zaowocować poprawą, skutkować poprawą*

poprawa (noun) Polish Web 2012 (pITenTen12, RFTagger) freq = 568,711 (60.58 per million)

a modifier	15.59	prec_prep	30.38	prec_verb	14.71	prec_na	8.15
znaczący +	10,974 9.06	poprzez +	2,201 5.54	ulec +	10,917 9.75	wpłynąć +	5,439 9.89
znaczącą poprawę		poprzez poprawę		obiecować +	3,057 9.50	wpływać na poprawę	
znaczący +	4,513 8.73	celem +	253 5.23	obiecuję poprawę		nadzieja +	3,944 9.61
znaczącą poprawę		celem poprawy		nastąpić +	4,944 8.46	nadzieję na poprawę	
zdecydowany +	3,319 8.46	ku +	585 4.99	nastąpiła poprawa		wpływać +	4,891 8.90
zdecydowaną poprawę		ku poprawie		ulegać +	3,454 8.30	wpływa na poprawę	
wyraźny +	2,446 8.24	nad +	3,791 4.82	ulega poprawie		szansa +	5,773 8.37
wyraźną poprawę		nad poprawą		odczuć +	1,097 8.05	sposób +	3,641 8.35
widoczny +	1,834 7.83	pomimo +	504 4.64	przynieść +	4,155 7.97	sposób na poprawę	
widoczną poprawę		mimo +	922 4.59	następować +	1,648 7.69	perspektywa +	336 7.51
radykałny +	1,081 7.61	mimo poprawy		następuje poprawa		perspektyw na poprawę	
radykałnej poprawy		dla +	11,264 4.36	zauważyłam +	591 7.60	recepta +	403 7.37
stopniowy +	638 7.39	dla poprawy		zauważyłam poprawę		recepty na poprawę	
stopniowej poprawy		do +	48,123 4.33	obiecać +	742 7.59	liczyć +	1,859 7.33
systematyczny +	796 7.32	do poprawy		obiecал poprawę		liczyć na poprawę	
systematyczną poprawę		na +	70,355 4.19	zauważyć +	1,670 7.51	metoda +	338 7.25
zauważalny +	493 7.31	na poprawę		zauważać +	661 7.40	metoda na poprawę	
zauważalną poprawę		oprócz +	656 4.17	spodziewać +	1,350 7.37	wpływ +	2,586 7.09
chwilowy +	688 7.21	. Oprócz poprawy		przynosić +	1,742 7.18	wpływ na poprawę	
chwilową poprawę		zamiast +	336 4.17	odnotować +	708 7.18	pomysł +	1,057 6.87
zadnej +	589 7.15	zamiast poprawy		oczekiwać +	1,160 7.13	pomysł na poprawę	
zadnej poprawy .		dzięki +	1,131 3.84	służyć +	2,149 7.05	zależać +	447 6.75
odczuwalny +	406 7.07	dzięki poprawie		służyć poprawie		zależy na poprawie	
odczuwalną poprawę		odnośnie +	170 3.78	sprzyjać +	969 6.88	pozwolić +	789 6.74
natychmiastowy +	835 7.01	odnośnie poprawy		sprzyja poprawie		pozwoli na poprawę	
natychmiastowej poprawy		o +	12,387 3.48	odczuwać +	907 6.81	środek +	502 6.44
nieznaczący +	396 6.89	o poprawę		rokować +	307 6.79	środków na poprawę	
nieznaczącą poprawę		wskutek	56 2.59	nie rokuje poprawy		wskazywać +	667 6.22
kliniczny +	555 6.54	wskutek poprawy		zauważać +	535 6.75	nacisk +	269 6.08
poprawy klinicznej		bez +	1,270 2.48	spowodować +	1,502 6.71	nacisk na poprawę	
ogólny +	2,453 6.53	bez poprawy		zaowocować +	495 6.66	oczekiwać +	277 6.03
ogólną poprawę		obok	91 2.08	skutkować +	525 6.53	oczekiwaniu na poprawę	
trwały +	881 6.50	, obok poprawy		skutkuje poprawą		lek +	169 5.94
trwałej poprawy		prócz	36 2.08	postanawiać +	249 6.46	przeznaczyć +	199 5.83
ciągły +	1,178 6.48	prócz poprawy		postanawiam poprawę		przeznaczyć na poprawę	
ciągłej poprawy		poza +	357 1.96	obserwować +	792 6.33	oczekiwanie +	145 5.82
jednoczesny +	448 6.46	. Poza poprawą					

sprovokować

- Sprovokować *coś*
- Sprovokować *kogoś do czegoś*

has	obj4		
		4.00	
zajście	<u>34</u>	8.46	
sprowokował całe zajście			
dyskusja +	<u>225</u>	8.02	
zamieszki	<u>16</u>	7.89	
reakcja +	<u>126</u>	7.49	
bójka	<u>9</u>	7.23	
debata	<u>25</u>	7.15	
skandal	<u>12</u>	6.98	
awantura	<u>22</u>	6.92	
kłótnia	<u>11</u>	6.92	
konflikt	<u>31</u>	6.77	
odwet	<u>7</u>	6.68	
wybuch	<u>12</u>	6.63	
pojawić	<u>9</u>	6.48	
zamieszka	<u>5</u>	6.44	
incydent	<u>8</u>	6.40	
odzew	<u>7</u>	6.38	
interwencja	<u>9</u>	6.34	
burda	<u>5</u>	6.28	
zadyma	<u>5</u>	6.19	
katastrofa	<u>9</u>	6.13	
wojna	<u>50</u>	6.10	
korekta	<u>15</u>	6.09	
lawina	<u>6</u>	6.08	
inwazja	<u>5</u>	6.01	
refleksja	<u>13</u>	5.98	

» the verb *sprovokować*, 'provoke sth'
has a negative aura of meaning in Polish

post	do		
		4.91	
kłótnia	<u>66</u>	8.16	
się sprowokować do kłótni			
brać	<u>54</u>	8.08	
sprowokować do brania			
bójka	<u>34</u>	7.75	
sprowokować do bójki .			
pyskówka	<u>15</u>	7.63	
reakcja	<u>41</u>	7.58	
agresja	<u>28</u>	7.58	
sprowokować do agresji .			
zastanowić	<u>34</u>	7.36	
sprowokować do zastanowienia się nad			
dyskusja +	<u>325</u>	7.18	
sprowokować do dyskusji			
przemyslenie	<u>39</u>	7.13	
sprowokować do przemyśleń			
rozprawić	<u>9</u>	7.01	
prawić	<u>9</u>	7.01	
myślenie +	<u>120</u>	6.96	
sprowokować do myślenia			
odpowiedź	<u>41</u>	6.95	
sprowokować do odpowiedzi			
polemika	<u>16</u>	6.93	
zabrać	<u>17</u>	6.93	
sprowokować do zabrania głosu			
bałakać	<u>8</u>	6.85	
awantura	<u>19</u>	6.84	
sprowokować do awantury			
odezwać	<u>8</u>	6.84	
atak +	<u>116</u>	6.83	
sprowokować do ataku .			
zadać	<u>14</u>	6.83	
sprowokować do zadania			
wypowiedź	<u>70</u>	6.68	
sprowokować do wypowiedzi			
pleść	<u>7</u>	6.66	
trajkotać	<u>7</u>	6.66	
powiadać	<u>7</u>	6.66	
zachowanie	<u>35</u>	6.64	

Limitations of language corpora

from the perspective of translators

- **Lack of sense disambiguation**
 - When we translate, we search for sense rather than form
 - *e.g. balkonik* (see later example)
 - framework for walking, part of a building, an undergarment?
- **Limited size of some language corpora**
 - **Negative evidence**
 - the fact that a given word or expression does not occur in the corpus does not mean that it is not used
- **If not updated, corpora do not catch up with reality**
 - No occurrences of *fake news*, *lokowanie produktu*, *hejtować* in NKJP
- **Loads of data mean problems with their filtering and interpretation**
 - The need to learn corpus functionalities (corpus search engines, query syntax etc.)
 - **Translators are busy people**
 - They may not have time for learning all this
 - Resource allocation for specific roles in the translation project

Wordnets

- Large **relational lexico-semantic databases**
- **Princeton WordNet** (Fellbaum 1998) was built on psycholinguistic principles
 - <http://wordnetweb.princeton.edu/perl/webwn>
 - **intended to reflect the structure of human lexical memory** (cf. Miller 1998)
 - originated as an experiment on mapping lexical memory of children,
 - gradually evolved into a huge electronic resource covering a large part of the lexical system of English (Fellbaum, 1998).
- Words are grouped in the so-called **synsets**
 - nouns, verbs, adjectives and adverbs
 - **represent lexicalised concepts** (Miller 1998)
 - **(cognitive) sets of synonymous lexical units (LUs)**
 - Synonymy is a **conceptual relation established on the basis of a linguist's intuition and dictionary definitions**
- **Synsets** are linked via semantic relations
 - hyponymy, hypernymy, partial-synonymy etc.

Polish wordnet (plWordNet), Słowosieć

- <http://plwordnet.pwr.wroc.pl/wordnet/>
- In plWordNet, **synsets are sets of synonymous LUs that enter the same constitutive relations**
 - e.g. hyponymy, hypernymy, meronymy, antonymy (Piasecki et al. 2009, Maziarz et al. 2012)
- **Each word is defined implicitly with reference to other words**
 - *samochód* 'car' is a kind of *pojazd drogowy* 'road vehicle'; it is a whole consisting of *silnik* 'engine', *spryskiwacz* 'windscreen washer', *podwozie* 'chassis' and so on; its close counterpart is the colloquial *fura* and *bryka* 'wheels, etc.

	No. of lemmas	No. of LUs	No. of synsets	Monosemous lemmas	Polysemous lemmas
Verbs	19984	41017	29892	11283	8701
Nouns	133845	176938	132628	109652	24193
Adverbs	8010	14040	11260	4692	3318
Adjectives	29228	54021	46705	16580	12648
All	191067	286016	220485	142207	48860

Mapping of wordnets (at synset level)

- **Linking plWordNet and Princeton WordNet synsets** corresponding in **meaning and position in wordnet structure**
 - Currently nouns
- Mapping direction: **plWordNet > Princeton WordNet**
 - using a set of interlingual relations
 - I-synonymy
 - I-hyponymy and hypernymy
 - I-meronymy and holonymy
 - I-partial synonymy
 - I-inter register synonymy
 - I interparadigmatic synonymy

Using wordnets in translation practice

- Using mapped wordnets can be summarized as follows:
 - **Determining the sense of a source language word**
 - **Finding possible translations (candidates)**
 - **Choosing a target language equivalent (a definitive translation) (Rudnicka & Piasecki 2013)**
- Examples

Situation 1

A known word in an unfamiliar context

- *After I finished studying for my courses one weekend, I called my mum. I told her that we were going to have a confirmation service during **chapel**, and I was not looking forward to attending. The topics chosen for **chapel** often seemed negative and left me feeling down.*
 - *Encountering Angels: True Stories of How They Touch Our Lives Every Day*, by Judith MacNutt (2016), Chosen Books*
 - Source:
https://books.google.pl/books?id=mgl_CgAAQBAJ&pg=PT23&dq=%22during+chapel%22&hl=pl&sa=X&ved=0ahUKEwi0iu6kpqrbAhUD2SwKHQYODd4Q6AEIMDAB#v=onepage&q=%22during%20chapel%22&f=false

* Thanks to E. Rudnicka and T. Piotrowski (2017) for suggesting this example

chapel

- Synsets in Princeton WordNet
 - **{chapel 1}** (artefact)
 - a place of worship that has its own altar
 - **{chapel 2; chapel service 1}** (act)
 - a service conducted in a place of worship that has its own altar

- Relevant synset in Princeton Wordnet
 - {**chapel service 1; chapel 2**} (act)
 - a service conducted in a place of worship that has its own altar;
 - he was late for chapel
 - Hypernymy to {service 3, religious service 1, divine service 1}
 - I-hypernymy {**nabożeństwo 1**} ✓
 - forma modlitewnego zgromadzenia wiernych danego wyznania

Situation 1.2

We don't know the equivalent

*At **Low Mass** with a congregation present, or if the church is large, incense may be blessed beforehand by the celebrant. At Missa Cantata, or **Low Mass** with an opening hymn, the **thurible** may be taken in procession.*

- Cooper, I. (2010). Ceremonies of The Young Rite, p. A-3

Source: <https://books.google.pl/books?id=hxX9AgAAQBAJ&pg=RA1-PA3&dq=%22low+mass%22+%2B+church&hl=en&sa=X&ved=0ahUKEwjlr5fQ-MPaAhUQmrQKHANhA544ChDoAQg8MAU#v=onepage&q=%22low%20mass%22%20%2B%20church&f=false>

- Relevant synsets in Princeton WordNet
 - **{low mass 1}**
 - a Mass recited without music
 - Hypernymy to {Mass 4}
 - (Roman Catholic Church and Protestant Churches) the celebration of the Eucharist
 - I-synonymy to **{cicha msza 1}** ✓
 - msza, podczas której nie ma muzyki, diakona ani subdiakona, a na ołtarzu palą się dwie świece; inaczej nazywana mszą czytaną
 - **{thurible 1; censer 1}**
 - a container for burning incense (especially one that is swung on a chain in a religious ritual)
 - Hypernymy to {vessel 3}
 - an object used as a container (especially for liquids)
 - I-synonymy to **{kadzielnica 1; trybularz 1}** ✓
 - naczynie kościelne, które służy do spalania kadzidła

Situation 2

We know the meaning of a SL word but we don't know its TL equivalent

*At Low Mass with a congregation present, or if the church is large, incense may be blessed beforehand by the **celebrant**. At Missa Cantata, or Low Mass with an opening hymn, the thurible may be taken in procession.*

- Cooper, I. (2010). Ceremonies of The Young Rite, p. A-3

- Synsets in Princeton Wordnet
 - {**celebrant 1**; celebrater 1; celebrator 1}
 - A person who is celebrating
 - {**celebrant 2**}
 - an officiating priest celebrating the Eucharist

- Relevant synset in Princeton WordNet
 - {celebrant 2}
 - an officiating priest celebrating the Eucharist
 - Hypernymy to {priest 1}
 - a clergyman in Christian churches who has the authority to perform or administer various religious rites
 - I-synonymy to {**celebrans 1; celebrant 1; oficjant 1**} ✓
 - duchowny odprawiający nabożeństwo

Situation 3

We are not sure about the meaning of a SL word and we don't know its TL equivalent

- *Kiedy przechodząc do kuchni pokonuje próg, **balkonik**, o który się opiera, wydaje cichy, ale wyraźny stukot.*
 - (Super Express, 2006; PWN_1302900002719)

Synsets in Słowosieć

- **{balkonik 1}** (art) ✓
 - specjalne urządzenie na kółkach, przy pomocy którego osoby starsze, niepełnosprawne i mające problem z chodzeniem, mogą się łatwiej poruszać
- **{balkonik 2}** (loc)
 - Markedness (diminutives) to {balkon 1}
 - element architektoniczny, przestrzeń wystająca poza mury budynku, często ogrodzona balustradą
- **{balkonik 3; balkonетка 1}** (art.)
 - rodzaj biustonosza z miseczkami skrojonymi z trzech części, który podnosi i zaokrągla piersi

- *Kiedy przechodząc do kuchni pokonuje próg, **balkonik**, o który się opiera, wydaje cichy, ale wyraźny stukot.*
(Super Express, 2006; PWN_1302900002719)
- Relevant PlWordNet synset
 - {**balkonik 1**} (domain: artefact)
 - specjalne urządzenie na kółkach, przy pomocy którego osoby starsze, niepełnosprawne i mające problem z chodzeniem, mogą się łatwiej poruszać
 - Hypernymy to {chodzik 1} {podpora 1; podpórka 1; podparcie 2}
 - I-synonimy to {**walker 5; Zimmer frame 1; Zimmer 1**}
 - a light enclosing framework (trade name Zimmer) with rubber castors or wheels and handles ✓


Situation 4

A known word in an unfamiliar context

- Skoro jest wilgoć, to i pojawiły się wreszcie moje **kołpaki** (...) Moje dwie gromady, umiejscowione w odległych od siebie miejscach, wysypują się zawsze w październiku. Chociaż podobno powinny to czynić już od maja. Grzyby te trochę mi się kojarzą z czubajkami kaniami- także mają gładki jasnobrązowy wierzchołek kapelusza i odstające łuski – u młodych białawe, u starszych brązowiejące.

– Source: <http://puszcza.net.pl/ciesza-sie-sojki-i-slimaki>

- Czernidłak kołpakowaty

- **Identifying the sense of the source language word**
 - {**kołpak 1**; dekiel 1}
 - osłona felgi samochodowej
 - Kołpak może być wykonany z metalu lub tworzyw sztucznych.
 - {**kołpak 2**}
 - nakrycie głowy pochodzenia tureckiego, popularne w XIV-XVII wieku w Europie Wschodniej; była to czapka bez daszka, obrzeżona futrem wilczym, rysim lub lisim z aksamitną główką, niekiedy przybierała kształt cylindryczny, zakończony spiczasto
 - Tatarskie kołpaki przypominała za to Czapka Monomacha, nakrycie głowy moskiewskich carów, wskazując też na powiązania kulturowe państw ruskich stanowiących lenno Złotej Ordy.
 - {**kołpak 3**; bełdka kołpak 1; sowa 2; kołpaczek 2; czubajka 2} 
 - inna nazwa czernidłaka kołpakowatego
 - {**kołpak 4**; dzwoniak 1; delikatka 1; rumieniaczek 1; bełdka 2; rumieniak 2}
 - ludowa nazwa niektórych grzybów z rodzaju dzwonków
 - {**kołpak 5**; podbołotuszka; rozeta pomarszczona 1; płachetka 1; bydlarka 1; niemka 2; płachta 2; turek 3}
 - inna nazwa płachetki zwyczajnej - grzyba jadalnego z rodziny zasłonakowatych

- **Finding possible translation candidates**
 - {**kołpak 3**; bełdka kołpak 1; sowa 2; kołpaczek 2; czubajka 2}
 - inna nazwa czernidłaka kołpakowatego
 - Inter-register synonyms
 - {czernidłak kołpakowaty 1}
 - **I-inter-register synonyms**
 - {shaggymane 1; Coprinus comatus 1; shaggy cap 1; shaggymane mushroom 1}
- **Choosing a target language equivalent**
 - {**shaggymane 1; Coprinus comatus 1; shaggy cap 1; shaggymane mushroom 1**} ✓

Limitations of wordnets

- **Polish WordNet has been mapped onto Princeton WordNet only partially**, which translates into **limited coverage**
 - and at the level of synsets (abstract units)
- **Translators do not deal with synsets but with lexical units** (form-meaning mappings)
 - Mapping at the level of lexical units (nouns) is in the pipeline (Rudnicka et al. 2017)
 - It considers formal, semantic and **translational criteria**
- **There may be gaps and mismatches between PWN and pIWN** (Rudnicka et al. 2016), which are due to:
 - Wordnet-specific structural and methodological differences
 - Specificity of inter-wordnet mapping procedure
 - Systemic differences between English and Polish
 - Differences in terms of lexicalization of concepts
 - often due to cross-cultural differences
- **Wordnets for different languages (other than English) are created using different methods**
 - Transfer, transfer-and-merge (automatic), transfer-and-merge (semi-automatic and corpus-based) - pIWN

Conclusions

- Using parallel and reference corpora as well as wordnets is **a realistic scenario in the translator's work**
 - It offers **exposure to authentic linguistic data** 😊
 - It may help **improve the textual fit of translations** 😊
 - But it can be **time-consuming and labour-intensive** and it may require learning new skills related to the use of those resources 😞
- **Language corpora and wordnets may complement (e-)dictionaries, CATs and TMs**
 - Notably when translating creative texts

References (selected)

- Grabowski, Ł. (2018). “Stance bundles in English-to-Polish translation: a corpus-informed study”. *Russian Journal of Linguistics*, 22 (2), 404-422
- Kilgarriff, A., Baisa, V., Bušta, J., Jakubíček, M., Kovář, V., Michelfeit, J., Rychlý, P. & Suchomel, V. (2014). “The Sketch Engine: ten years on”. *Lexicography* 1(1): 7–36.
- Krueger, R. (2017). Working with corpora in the translation classroom
- Peźik, P. (2017). Korpusy referencyjne, korpusy równoległe, ekwiwalencja frazeologiczna. Presentation delivered at CLARIN-PL workshops (3-4 Feb 2017) in Łódź
- Podhajecka, M. (2013). Calques in Polish translations of English press articles: Linguistic innovations or mistranslations?” In B. Lewandowska-Tomaszczyk and M. Thelen (eds.), *Translation and Meaning Part 10*. Maastricht: Zuyd University of Applied Sciences. 165–179.
- Rudnicka, E. & Piotrowski, T. (2015). Dwujęzyczna SłowoSieć – możliwości wykorzystania w pracy tłumacza. Presentation delivered at CLARIN-PL workshops in IBL in Warsaw.
- Rudnicka, E., Witkowski, W. & Grabowski, Ł. (2016). “Towards a methodology for filtering out gaps and mismatches across wordnets: the case of noun synsets in plWordNet and Princeton WordNet”. In: B. Barbu Mititelu, C. Forascu, Ch. Fellbaum, P. Vossen (Eds.), *Proceedings of the Eighth International Global WordNet Conference 2016*, 27-30 Jan 2016, Bucharest, Romania, pp. 344-351
- Rudnicka, E., Bond, F., Grabowski, Ł., Piasecki, M. & Piotrowski, T. (2017). “Towards equivalence links between senses in plWordNet and Princeton WordNet”. *Lodz Papers in Pragmatics*, 13 (1), 3-24.
- Zanettin, F. (2002). Corpora in Translation Practice. In: *Proceedings of the LREC Workshop, Language Resources for Translation Work and Research*. 2002. 10-14.

Thank you

- for your attention
 - lukasz@uni.opole.pl
 - ResearchGate: LukaszGrabowski2