

Leksykografia elektroniczna

Praktyki, metody, narzędzia

Krzysztof Nowak

Pracownia Łaciny Średniowiecznej

Instytut Języka Polskiego PAN

Plan wykładu

- **korpus w pracy leksykografia**

- po co korpus?
- jaki korpus?

- **narzędzia korpusowe**

- prezentacja
- przeszukiwanie

- **teoria i metody**

- **laboratorium 1**

- **narzędzia leksykograficzne**

- od surowego korpusu do hasła (SketchEngine etc.)
- systemy redakcji słowników

- **laboratorium 2**

Co zmieniły korpusy w opisie językoznawczym *tout court*?

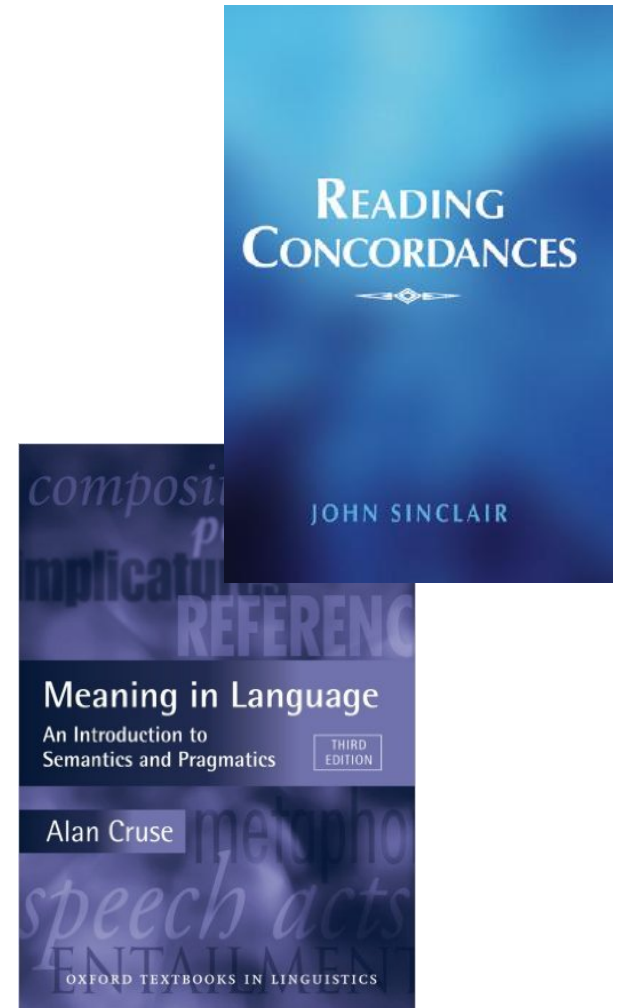
Po co korpus?

Zalety danych korpusowych

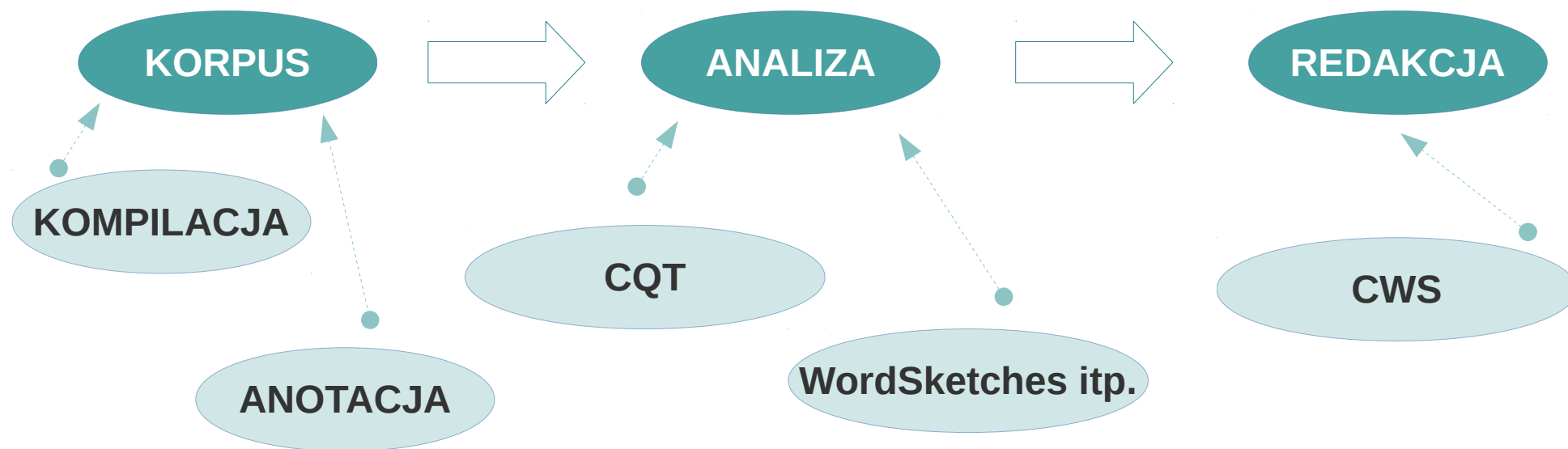
- obiektywizacja obserwacji (frekwencja wyrazów, znaczeń, zjawisk)
- śledzenie wzorców składniowych, pragmatycznych
- lepsze ujęcie frazeologii

Jaka teoria?

- rola kontekstu w opisie znaczenia (Kilgariff: „I don't believe in word senses”)
- Patrick **Hanks**, projekt *Corpus Pattern Analysis*: wyrazy aktywują potencjał znaczeniowy w kontekście
- John **Sinclair**, „Reading Concordances”
- Alan **Cruse**



Procedura



Jaki korpus dla leksykografii? (1)

jaki korpus dla leksykografii?

nie istnieje korpus do wszystkiego
(inny dla języka ogólnego, inny dla słownika dokumentacyjnego)

reprezentatywność

słownik ogólny

język

zróżnicowanie regionalne i
dialektalne

czas

synchronia ≠ diachronia

język mówiony czy pisany?

kanal

Sieć ≠ gazeta ≠ książka ≠
czat

domena

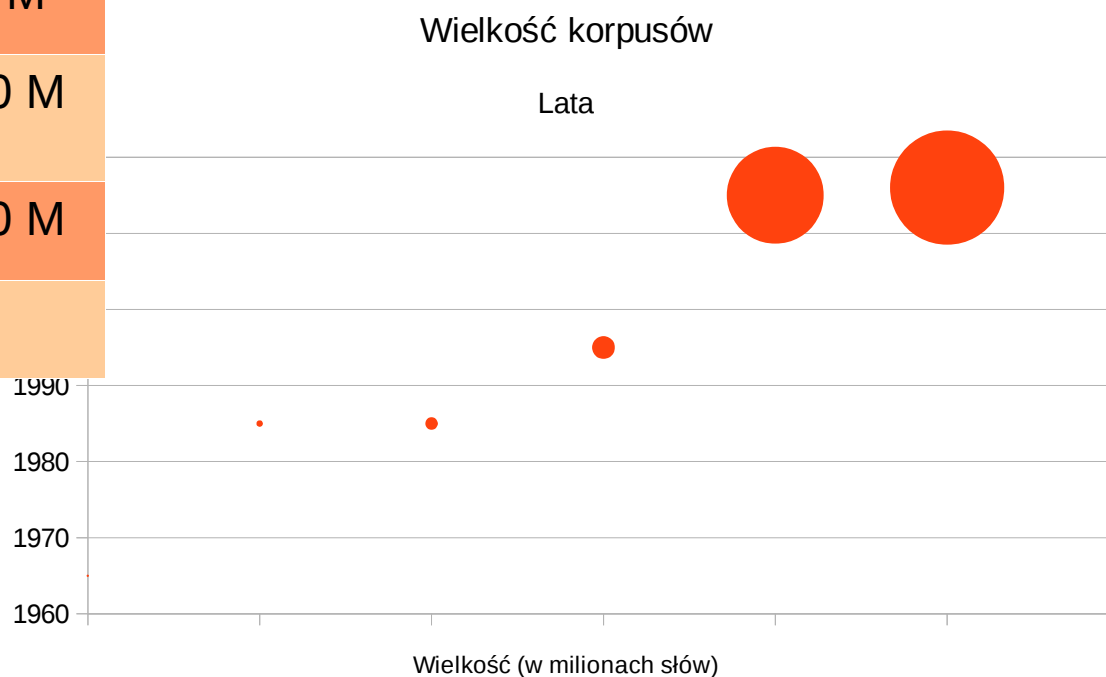
słownictwo specjalistyczne?

korpus równoległy?

leksykografia dwujęzyczna

Jaki korpus dla leksykografii? (2)

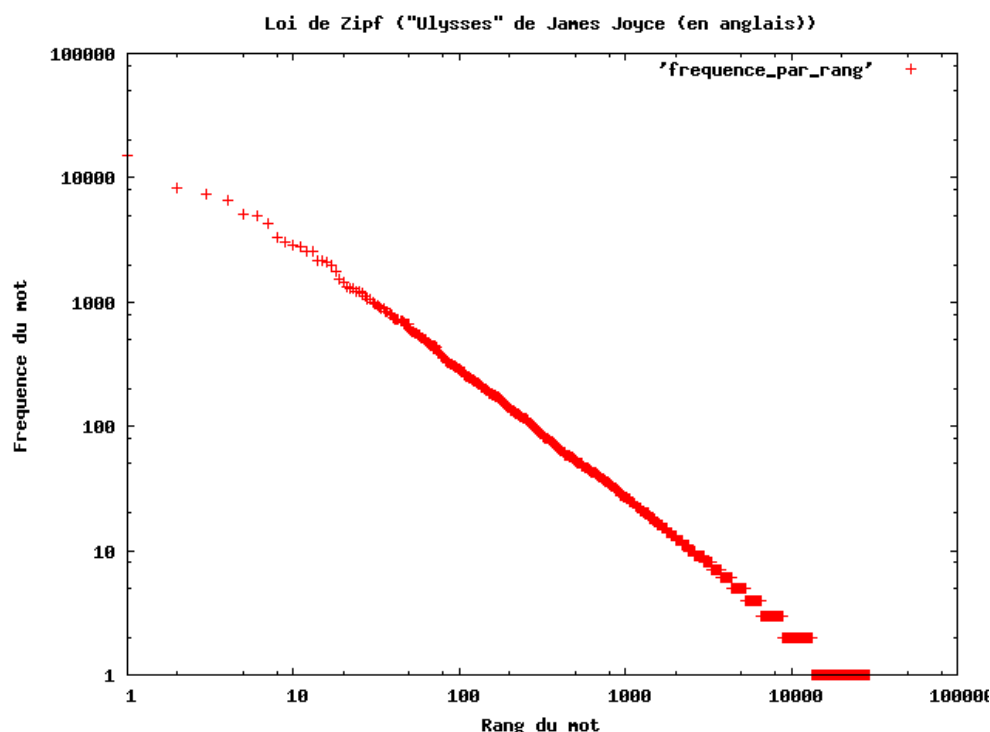
Nazwa korpusu	Wielkość
Brown / LOB Corpus	1 M
Birmingham	8 M
British National Corpus	100 M
Narodowy Korpus Języka Polskiego	1 800 M
Oxford English Corpus	2 500 M
Web as Corpus	?



Jaki korpus dla leksykografii? (3)

Prawo Zipfa...

$$\text{frekwencja} \approx \frac{\text{constans}}{\text{ranking}}$$



... i jego konsekwencje

- mały korpus daje względne pojęcie o słowach o dużej frekwencji
- opis np. 100 k najpopularniejszych słów wymaga już bardzo dużego korpusu
- podobnie jak wyrazy (nawet o dużej frekwencji) zachowują się ich rzadsze znaczenia i użycia

Przygotowanie korpusu

kodowanie znaków

ąęćłźżóśń

metadane

```
<tytul>Washington  
Post</tytul>  
<rok>1960</rok>  
<kanal>prasa<kanal>
```

typografia

Lorem
ipsum dolor sit
amet.

tokenizacja

Litwo|ojczyzno|
moja|!|

lemmatyzacja

Litwa|ojczyzna|
mój|!|

PoS

PROP|NOUN|
PRON|INT|

inne

Narzędzia korpusowe

Klasyfikacja (Kilgariff & Kosem)

- lokalne ≠ sieciowe
- związane z określonym korpusem ≠ niezależne
- gotowy korpus ≠ Web as corpus
- proste ≠ zaawansowane

Wybrane korpusy wyposażone w interfejs sieciowy

- NKJP
- Cobuild's Bank of English
- korpusy Institut für Deutsche Sprache (Mannheim)

Podstawowe funkcje narzędzi korpusowych

- możliwość ładowania własnego korpusu
- możliwość współdzielenia
- indeksacja
- obsługa anotacji
- obsługa kodowania

Podstawowe *modi*

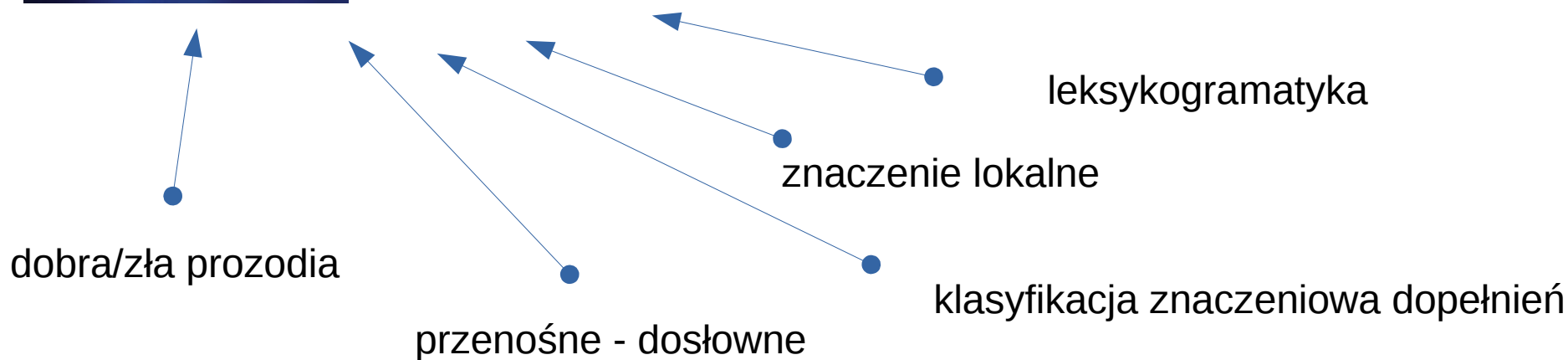
- listy frekwencyjne
 - kandydaci na wyrazy hasłowe
 - filtrowanie PoS
 - filtrowanie wg metadanych (np. domen)
- listy kolokacji (wg miar asocjacji):
 - znaczenie (Firth)
 - stabilne połączenia
- konkordancja
 - KWIC
 - wyszukiwanie *wildcards*
 - dopasowanie wyświetlania
 - sortowanie wg kontekstu, typu tekstu, czasu powstania itd.
- randomizacja

TXM / CQPWeb

Jak pracować z konkordancją?



- 1. Sprawdzić kontekst: ... 2L, 1L, 1R, 2R ...**
- 2. Zredukować pojedyncze wystąpienie do powtarzających się wzorców**
- 3. Postawić hipotezę**
- 4. Analizować kolejne wystąpienia**
- 5. Sprawdzić, co wnoszą do ujęcia**



Laboratorium 1: szkic hasła słownikowego na podstawie NKJP

1. Wybór słowa
2. Opracowanie szkicu hasła na podstawie wyszukiwarki PELCRA:
 - a) kolokator,
 - b) konkordancje.
3. Zwrócić uwagę na:
 - a) częste połączenia i inne elementy kontekstu,
 - b) gatunki tekstu,
 - c) nacechowanie.
4. Skonfrontować ze znaczeniem odnotowanym w wybranym słowniku: czy analiza korpusowa wniosła cokolwiek?
5. Czy hasło powinno znaleźć się np. w słowniku ogólnym? Słowniku kieszonkowym?
6. Jakie są wady poszczególnych narzędzi analizy korpusu?

SketchEngine

Sumaryzacja danych: WordSketches

- kolokacje + gramatyka
- „jednostronicowe podsumowanie gramatycznego i kolokacyjnego zachowania wyrazu” (Kosem & Kilgariff 2012)

Sumaryzacja danych: SketchDiffs

- porównanie kolokacji między wyrazami hasłowymi

goal *(noun)*
ukWaC freq = **168,184** (107.82 per million)

<u>object of</u>	<u>59,154</u>	<u>3.20</u>	<u>subject of</u>	<u>25,630</u>	<u>2.00</u>	<u>adj subject of</u>	<u>2,159</u>	<u>1.40</u>	<u>modifies</u>
score	<u>8,555</u>	11.03	score	<u>1,029</u>	8.39	galore	<u>27</u>	7.30	ultimate
achieve	<u>9,504</u>	9.69	disallow	<u>270</u>	8.28	achievable	<u>43</u>	6.95	winning
concede	<u>1,418</u>	9.31	concede	<u>223</u>	7.52	attainable	<u>15</u>	6.77	long-term
accomplish	<u>588</u>	7.88	gape	<u>76</u>	6.48	unrealistic	<u>14</u>	5.59	league
reach	<u>1,937</u>	7.42	kick	<u>93</u>	5.45	intact	<u>19</u>	5.04	primary
net	<u>351</u>	7.42	orientate	<u>37</u>	5.12	worthy	<u>36</u>	4.82	second

clever/intelligent ukWaC freqs = **20,593** | **26,128**

clever 6.0 4.0 2.0 0 -2.0 -4.0 -6.0 intelligent

<u>and/or</u>	<u>4,831</u>	<u>9,801</u>	<u>2.10</u>	<u>3.40</u>	<u>modifier</u>	<u>4,933</u>	<u>3,168</u>	<u>0.90</u>	<u>0.50</u>	<u>modifies</u>	<u>11,415</u>	<u>16,772</u>
funny	<u>231</u>	<u>103</u>	6.9	5.7	fiendishly	<u>45</u>	0	8.1	--	ploy	<u>72</u>	0
inventive	<u>22</u>	<u>24</u>	5.9	5.5	devilishly	<u>17</u>	0	6.8	--	clog	<u>51</u>	0
witty	<u>133</u>	<u>166</u>	8.3	8.2	awfully	<u>15</u>	0	6.1	--	trick	<u>166</u>	0
resourceful	<u>12</u>	<u>29</u>	5.8	6.3	terribly	<u>25</u>	0	6.1	--	twist	<u>94</u>	0
insightful	<u>11</u>	<u>30</u>	5.2	6.1	dead	<u>14</u>	0	5.9	--	eh	<u>40</u>	0
affectionate	<u>6</u>	<u>32</u>	4.5	6.2	diabolically	<u>9</u>	0	5.9	--	chap	<u>48</u>	0
thoughtful	<u>14</u>	<u>121</u>	5.0	7.7	amazingly	<u>17</u>	<u>7</u>	5.9	5.0	wordplay	<u>21</u>	0
compassionate	0	<u>27</u>	--	5.9	exceptionally	<u>28</u>	<u>25</u>	5.9	5.9	ruse	<u>17</u>	0
literate	0	<u>26</u>	--	5.9	averagely	0	<u>7</u>	--	6.1	Huh	<u>18</u>	0
well-informed	0	<u>25</u>	--	6.0	supposedly	0	<u>28</u>	--	6.2	lyric	<u>83</u>	<u>80</u>
adaptive	0	<u>39</u>	--	6.1	ferociously	0	<u>8</u>	--	6.2	creature	<u>11</u>	<u>137</u>
informed	0	<u>59</u>	--	6.1	highly	0	<u>572</u>	--	6.9	agent	<u>9</u>	<u>465</u>
thought-provoking	0	<u>20</u>	--	6.2	farquely	0	<u>28</u>	--	6.0	guess	0	<u>25</u>

Systemy redakcji słowników

Po co?

- kontrola formy hasła
- kontrola struktury hasła
 - następstwo i występowanie części
 - skróty
- kontrola języka opisu i definicji
- połączenie z narzędziami korpusowymi (*TickBox Lexicography*) i leksykalnymi bazami danych
- wyeliminowanie niespójności

Komponenty



- interfejs edycji



- baza danych



- zarządzanie użytkownikami

Następne zajęcia

- **data**

- piątek, 15 IV 2016, godz. 18

- **temat**

- Tworzenie i zastosowania słowników komputerowych (Wordnet, Framenet, Babelnet i spółka)

- **wymagania**

- założenie konta trial na:

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

<https://cqpweb.lancs.ac.uk/>

http://ucts.uniba.sk/aranea_about/