

ELEMENTY METOD GRAFOWYCH I KOMBINATORYCZNYCH W BADANIACH FRAZEOLÓGICZNYCH

Piotr Pęzik¹ Bartosz Zieliński²

¹Pracownia Językoznawstwa Korpusowego i Komputerowego,
Instytut Anglistyki,
Uniwersytet Łódzki

²Katedra Informatyki,
Wydział Fizyki i Informatyki Stosowanej,
Uniwersytet Łódzki

ZIL @ IPI PAN 2014/05/12

- „Nie ma jednej frazeologii” (Lewicki, 1976).
- Frazeologia dystrybucyjna, neofirthiańska, *strong tea* vs. *powerful car*, tendencja idiomatyczna (Sinclair, 1996), *lexicality*
- (Korpusowa) “frazologia nadawcy” (Lewicki, 1976), frazematyka (Chlebda, 2003)
- Reprodukacja podstawowym kryterium formułiczności
- Por. niekompozycyjność w tradycyjnym ujęciu (czysta i figuratywna idiomatyczność)

FORMULICZNOŚĆ I JEJ PSYCHOLINGWISTYCZNE KRYTERIA IDENTYFIKACJI

- Czas produkcji, płynność, tempo, dystrybucja pauz w języku mówionym (Pawley & Syder, 1983)
- Akwizycja, utrata języka (Wray, 2002)
- Śledzenie ruchu gałek ocznych (Siyanova-Chanturia, Conklin, & Schmitt, 2011)
- Obrazowanie mózgu (Tremblay & Bayeen, 2010)

HASK, CZYLI SŁOWNIKI KOMBINACJI WYRAZOWYCH

- Dla NKJP: 106 723 “hasel” zawierających łącznie 5 636 857 kombinacji.
- Niektóre reguły są symetryczne -> kombinacje częściowo się powtarzają.
- Szczegóły: (Pęzik, 2013).

Node	F	Type	Positions	Colls	Coll	F	T	JD
mocny	13018	adj-noun	[-1 1 -2 2]	829	strona	891	27.12	0.92
mocny	13018	adj-noun	[-1 1 -2 2]	829	pozycja	401	19.25	0.89
mocny	13018	adj-noun	[-1 1 -2 2]	829	słowo	325	15.79	0.93
mocny	13018	adj-noun	[-1 1 -2 2]	829	wrażenie	274	15.76	0.91
mocny	13018	adj-noun	[-1 1 -2 2]	829	uderzenie	199	13.88	0.9
mocny	13018	adj-noun	[-1 1 -2 2]	829	kawa	189	13.4	0.89
mocny	13018	adj-noun	[-1 1 -2 2]	829	trunek	159	12.55	0.89
mocny	13018	adj-noun	[-1 1 -2 2]	829	argument	146	11.46	0.9
mocny	13018	adj-noun	[-1 1 -2 2]	829	herbata	127	10.93	0.86
mocny	13018	adj-noun	[-1 1 -2 2]	829	coś	210	10	0.89

TABLE: Zob. pelcra.pl/hask_pl, pelcra.pl/hask_en.

- Pilotażowe eksperymenty psycholingwistyczne

przymiotnik {
wiatr
pozycja
wpływ
skład
wola
emocja
ból
konkurencja
więź
strzał

przymiotnik {
uderzenie
akcent
strzał
uścisk
postanowienie
głos
argument
punkt
cios
makijaż

POZIOMY AKTYWACJI W ROZPOZNAWANIU SŁÓW

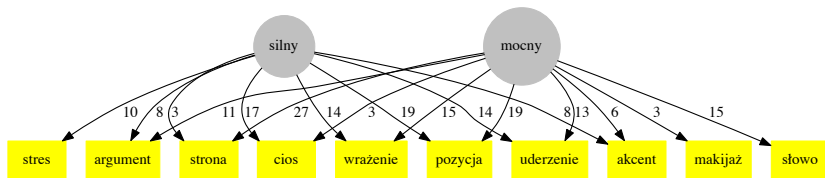
- 38/39 badanych wymienia **silny** dla 1. zbioru
- 39/39 badanych wymienia **mocny** dla 2. zbioru

siln(y/a/e) {
wiatr
pozycja
wpływ
skład
wola
emocja
ból
konkurencja
więź
strzał

mocn(y/a/e) {
uderzenie
akcent
strzał
uścisk
postanowienie
głos
argument
punkt
cios
makijaż

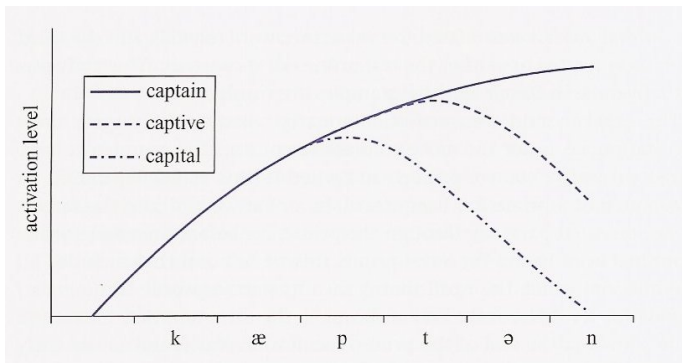
SILNY VS. MOCNY

- Wyniki są interesujące, zważywszy na stopniowalność rozróżnienia między frazemem a syntagmą.
- *Mocne słowa* : *silne słowa* 284:5 w zrównoważonej części NKJP
- *Silny stres* : *mocny stres* 101:1, *mocny makijaż* : *silny makijaż* 39:1



- Niektóre podzbiory kolokatów są unikalne (globalnie w całym grafie słownika) na przyjętym poziomie utrwalenia dla danych ośrodków kolokacji.
- Przykładem szczególnych przypadków unikalnych, jednoelementowych zbiorów są składniki tzw. “cranberry collocations”, czyli frazemy z **wyrazami uwięzionymi** (por. frazeologia.pl, Piotr Müldner-Nieckowski), np. *put the **kibosh** on, zbić kogoś z **pantałyku***
- Czy i w jaki sposób różne kombinacje (docelowo permutacje) kolokatów (primes) odpowiadają różnym poziomom aktywacji w procesie rozpoznawania ośrodków kolokacji (targets)?

POZIOMY AKTYWACJI W ROZPOZNAWANIU SŁÓW



(Warren, 2013)

- Teoria kohorty
- Punkty unikalności (jedyne pasujące słowo), izolacji (wstępne rozpoznanie słowa), rozpoznania (80% pewność), całkowitej akceptacji (100% pewność) (Field, 2004).

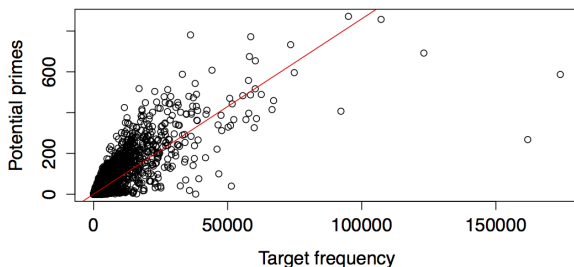
- Badania efektu torowania, również z wykorzystaniem korpusowych kryteriów dystrybucyjnych.
- Również w NLP – warsztat CogALex-IV@COLING 2014

whiskey {
gin
drink
bottle
soda

- Zawężenie badań do zaświadczonych w korpusach kontekstów frazeologicznych.

whiskey {
Irish
strong
Scotch
old

- Badania podstawowe: kompozycyjność/formuliczność języków.
- W jakim stopniu punkty unikalności określone na podstawie korpusów referencyjnych znajdują potwierdzenie w eksperymentach psycholingwistycznych?
- Weryfikacja korpusowych metod identyfikacji jednostek frazeologicznych. Które miary, lub które kombinacje miar skutkują wyłonieniem wyrazów najsilniej naprowadzających?
- Co decyduje o aktywacji? Częstość i rozproszenie czy uwięzienie?



- Od czego zacząć? Słownik zawiera miliony kombinacji...
- Minimalne zbiory skutkujące punktami unikalności wydają się najprostsze do analizy (np. wyrazy uwięzione).
- Przewidywana konieczność generowania permutacji na późniejszym etapie.
- Problem: jak skutecznie znaleźć możliwie niewielkie zbiory?

ELEMENTY METOD GRAFOWYCH W BADANIACH FRAZEologiczNYCH

- Inne zastosowania modeli grafowych w problemach frazeologicznych.
- Sieci kolokacyjne: (Mehler, 2007)
- Inne problemy: weryfikacja poziomu synonimii metodami grafowymi
- Spójność synsetów w tezaursach i słowosieciach

PROBLEM "HITTING SET"

Niech $\{S_i\}_{i \in I}$ będzie skończoną rodziną zbiorów. Należy znaleźć najmniejszy zbiór C taki że $C \cap S_i \neq \emptyset$ dla każdego $i \in I$.

PROBLEM POKRYCIA ZBIORU

Niech $\{S_i\}_{i \in I}$ będzie skończoną rodziną zbiorów. Należy znaleźć najmniejszy podzbiór $J \subseteq I$ taki że $\bigcup_{i \in I} S_i = \bigcup_{j \in J} S_j$.

Niech $(U, V, E \subseteq U \times V)$ będzie grafem dwudzielnym. Znaleźć minimalny podzbiór $C \subseteq U$ taki że dla każdego $v \in V$ istnieje $u \in C$ takie że $(u, v) \in E$.

PROBLEM "HITTING SET"

- $U := \bigcup_{i \in I} S_i, V := I,$
- $(s, i) \in E$ wtedy i tylko wtedy gdy $s \in S_i$.

SET COVER

- $U := I, V := \bigcup_{i \in I} S_i,$
- $(i, s) \in E$ wtedy i tylko wtedy gdy $s \in S_i$.

Niech $(U, V, E \subseteq U \times V)$ będzie grafem dwudzielnym.

WPROWADZAMY NASTĘPUJĄCĄ NOTACJĘ:

- $E(u) := \{v \in V \mid (u, v) \in E\}$ dla każdego $u \in U$.
Dla $A \subseteq U$ definiujemy $E(A) := \bigcup_{u \in A} E(u)$.
- $E^{-1}(v) := \{u \in U \mid (u, v) \in E\}$ dla każdego $v \in V$.
Dla $B \subseteq V$ definiujemy $E^{-1}(B) := \bigcup_{v \in B} E^{-1}(v)$.

PROBLEM PRZEFORMUŁOWANY

Znaleźć dla grafu $(U, V, E \subseteq U \times V)$ minimalny $C \subseteq U$ t.ż. $E(C) = V$.

ALGORYTM ZACHŁANNY

- Zaczynamy od $C := \emptyset$.
- Dopóki $E(C) \neq V$ zwiększamy $C := C \cup \{u\}$ (u maksymalizuje $|E(u) \setminus E(C)|$).

GWARANCJE ALGORYTMU ZACHŁANNEGO

Niech C^* będzie rozwiązaniem optymalnym, a C znalezionym przez algorytm zachłanny.

- $|C| \leq H_k |C^*|$, gdzie $H_k = \sum_{i=1}^k 1/i \simeq \ln k$, $k := \max\{|E(u)| \mid u \in U\}$.
- $|C|/|C^*| \leq \ln n - \ln \ln n + O(1)$, gdzie $n := |V|$.

Istnieją lepsze algorytmy przybliżone niż zachłanny, ale ...

TWIERDZENIE

Minimalnego pokrycia zbioru nie można przybliżyć w czasie wielomianowym ze stosunkiem $|C^*|/|C|$ lepszym niż $(1 - o(1)) \ln n$.

ŹRÓDŁO

Feige, Uriel. "A threshold of $\ln n$ for approximating set cover."
JACM 45(4) (1998), 634-652.

PROBLEM KORPUSOWYCH PUNKTÓW MINIMALNOŚCI DLA KOLOKACJI LEKSYKALNYCH

GRAF KOLOKACJI

- U — zbiór ośrodków kolokacji.
- V — zbiór kolokatów.
- $(u, v) \in E$ — słowo u pojawia się w korpusowym słowniku kolokacji (z częstością powyżej granicznej) wraz z kolokatem v .

PROBLEM

Dla każdego ośrodka kolokacji $u \in U$ znaleźć (o ile istnieje) najmniejszy zbiór kolokatów $C_u \subseteq V$ jednoznacznie wyznaczający u , tzn. taki że

$$\forall u' \in U . C_u \subseteq E(u') \Rightarrow u = u'.$$

Niech $u \in U$. Zdefiniujmy

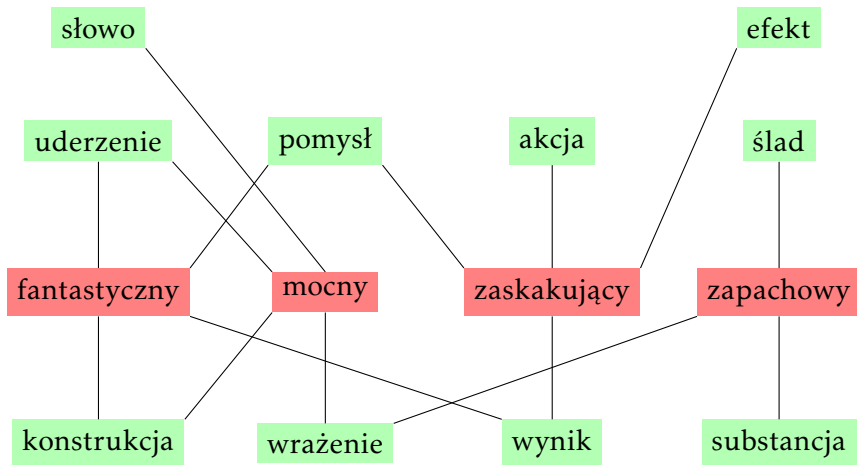
GRAF DWUDZIELNY $\mathcal{G}'_u := (U_u, V_u, E_u)$, GDZIE

- $U_u := E(u)$, $V_u := E^{-1}(E(u)) \setminus \{u\}$,
- $E_u := \{(v, u) \in U_u \times V_u \mid (u, v) \notin E\}$.

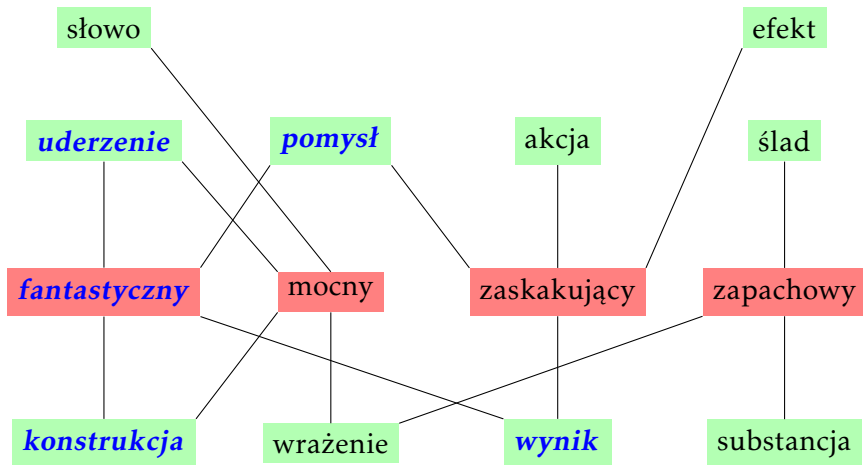
OBSERWACJA

Dla ośrodka kolokacji $u \in U$ znalezienie najmniejszego unikalnego zbioru kolokatów jest tożsame z rozwiązaniem grafowego uogólnienia problemu “Set cover” dla grafu \mathcal{G}'_u .

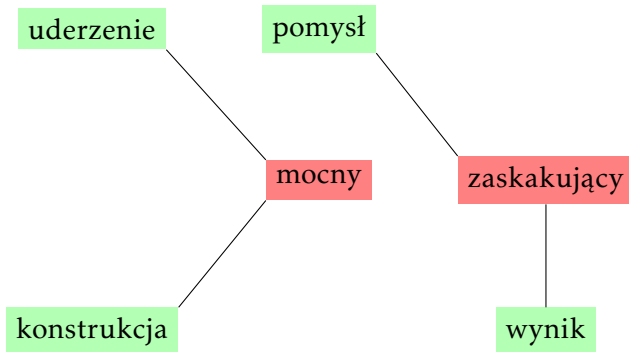
PRZYKŁAD GRAFU KOLOKACJI

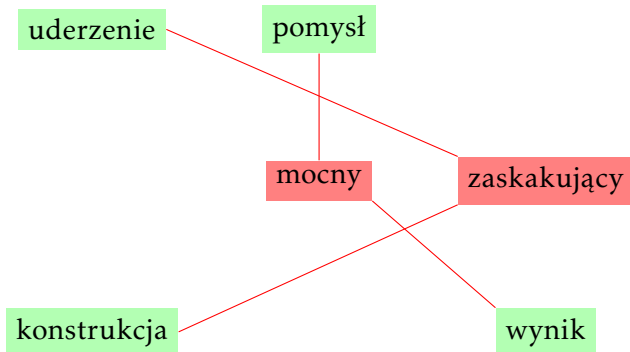


PRZYKŁAD GRAFU KOLOKACJI

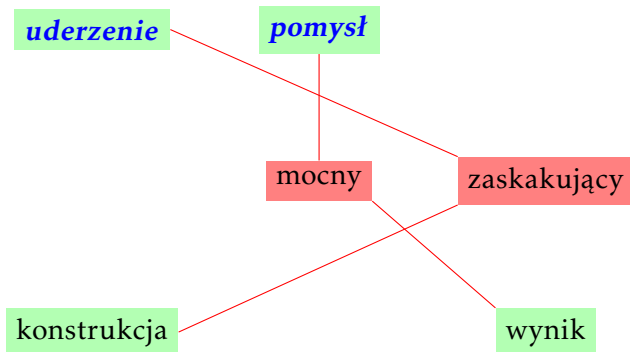


PODGRAF DLA $u = \text{fantastyczny}$





GRAF $\mathcal{G}'_{\text{fantastyczny}}$, NP. $C_{\text{fantastyczny}} = \{\text{uderzenie, pomysł}\}$



SZKIELET ALGORYTMU POSZUKIWANIA KORPUSOWYCH PUNKTÓW MINIMALNOŚCI

U — ośrodki kolokacji, V — kolokaty, $E \subseteq U \times V$ — kolokacje.

```
for all  $u \in U$  do  
  if  $\text{IsUNIQUE}(E(u), u)$  then  
    PRINT(Min. unik. podzb. kolokatów dla  $u$ :  $\text{MINCOL}(u)$ )  
  else  
    PRINT(Brak unikalnych podzbiorów kolokatów dla  $u$ )  
  end if  
end for
```

```
function IsUNIQUE( $C, u$ )  
  for all  $u' \in U \setminus \{u\}$  do  
    if  $C \subseteq E(u')$  then  
      return false  
    end if  
  end for  
  return true  
end function
```

Require: u ma unikalny zbiór kolokatów

function MINCOL(u)

$n \leftarrow |E(u)|, \mathcal{C} \leftarrow \{E(u)\}$

for all $C \in 2^{E(u)} \setminus \{E(u)\}$ **do**

if ISUNIQUE(C, u) **then**

if $|C| < n$ **then**

$n \leftarrow |C|, \mathcal{C} \leftarrow \{C\}$

else if $|C| = n$ **then**

$\mathcal{C} \leftarrow \mathcal{C} \cup \{C\}$

end if

end if

end for

return \mathcal{C}

end function

Pętla wykonywana $2^{|E(u)|} - 1$ razy (to bardzo źle), ale:

- W korpusowym słowniku kolokacji zbiory kolokatów dla większości ośrodków kolokacji są niewielkie (< 10).
- Można te rozmiary regulować modyfikując wartości progowe częstości.
- Można też przyjąć że interesują nas tylko niewielkie unikalne zbiory kolokatów i tylko takie rozważać w algorytmie.
- W ostateczności można użyć algorytmów przybliżonych.

ULEPSZENIA ALGORYTMU: REPREZENTACJA ZBIORÓW KOLOKATÓW

- $\text{IsUNIQUE}(C,U)$, gdzie dla każdego $u' \in U \setminus \{u\}$ sprawdzamy czy $C \subseteq E(u')$, jest wąskim gardłem.
- Kluczowa jest reprezentacja zbiorów $E(u)$ zapewniająca:
 - efektywne wykonywanie operacji na zbiorach
 - i efektywną serializację/deserializację (zbiory danych mogą nie mieścić się w pamięci).

POMYSŁ: KOMPRESOWANE BITMAPY

- Reprezentujemy $E(u)$ jako skompresowaną bitmapę $(B_{uv})_{v \in V}$,
gdzie $B_{uv} := \begin{cases} 1 & \text{if } (u,v) \in E \\ 0 & \text{if } (u,v) \notin E. \end{cases}$
- Istnieją homomorficzne algorytmy kompresji.
- Istnieją implementacje dla Javy: JavaEWAH
(Lemire et al., Data & Knowledge Engineering 69(1), 2010)

- Podzbiory należy efektywnie generować w kolejności rozmiaru (1-elementowe, 2-elementowe, ...).
- Można zatrzymać się na podzbiorach K -elementowych dla arbitralnie ustalonego K .
- Dla ustalonego k , k -elementowe podzbiory powinny różnić się minimalnie.

REVOLVING DOOR ALGORITHM

Knuth, *The Art of Computer Programming*, Vol 4.

```

function MINCOL( $u$ )
  for  $n \in 1.. \min\{K, |E(u)|\}$  do      ▷  $K$  — wybrana wartość graniczna
     $\mathcal{C} \leftarrow \emptyset$ 
    for  $C \in \text{REVDOOR}(n, E(u))$  do
      if ISUNIQUE( $C, u$ ) then
         $\mathcal{C} \leftarrow \mathcal{C} \cup \{C\}$ 
      end if
    end for
    if  $\mathcal{C} \neq \emptyset$  then
      return  $\mathcal{C}$ 
    end if
  end for
  return  $\emptyset$ 
end function

```


INNE PROBLEMY:

ANALIZA SPÓJNOŚCI GRUP OŚRODKÓW KOLOKACJI

- Problem: Zmierzyć jak spójne są synsety — stopień w jakim synonimy mają również podobne zbiory kolokatów.
- Pomysł: Użyć miar spójności w ważonych grafach dwudzielnych.
- Korzystamy z miar lokalnych — miara dla każdego synonimu.
- Jeśli potrzebna jest miara globalna synsetu, miary lokalne można uśrednić

OZNACZENIA: GRAF DWUDZIELNY WAŻONY (U, V, a)

- U — synset, V — kolokaty synsetu,
- $a : U \times V \rightarrow \mathbb{R}_+$ — miara współwystępowania ośrodka kolokacji z kolokatem w korpusie językowym.

PROSTA MIARA SPÓJNOŚCI W GRAFIE DWUDZIELNYM WAŻONYM

WAGA WIERZCHOŁKA $u \in U$:

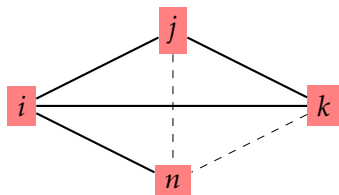
$$d_u := \sum_{v \in V} a_{uv}$$

PROSTA LOKALNA MIARA SPÓJNOŚCI DLA $u \in U$:

$$S_u := \frac{1}{d_i(|U| - 1)} \sum_{v \in V} \sum_{u' \in U \mid u' \neq u} a_{uv} a_{u'v}$$

Watts, Duncan J., and Steven H. Strogatz.
Collective dynamics of small-world networks.
Nature 393.6684 (1998), 440–442.

- Stosunek ilości trójkątów z wierzchołkiem w danym wierzchołku do ilości możliwych.
- Istnieją niezliczone wersje dla grafów ważonych.
- W grafie dwudzielnym trójkątów nie ma.
- Należy zatem przerobić graf dwudzielny na jednodzielny.



Trójkątów dla i : 1

Możliwych trójkątów dla i : 3

PRZYKŁAD LOKALNEJ MIARY OPARTEJ NA TRÓJKĄTACH

PRZERABIAMY (U, V, a) NA (U, w) GDZIE

$$w : U \times U \rightarrow \mathbb{R}_+, \quad w_{uu'} := \frac{1}{\sqrt{d_u d_{u'}}} \sum_{v \in V} a_{uv} a_{u'v}.$$

LOKALNA MIARA SPÓJNOŚCI DLA $i \in U$

$$C_i := \frac{\sum_{j, k \in U \mid j \neq k \neq i} w_{ij} w_{ik} w_{kj}}{\sum_{j, k \in U \mid j \neq k \neq i} w_{ij} w_{ik}}.$$

- Chlebda, W. (2003). *Elementy frazematyki, wprowadzenie do frazeologii nadawcy* (2nd ed.). Leksem.
- Field, J. (2004). *Psycholinguistics: the key concepts*. London ; New York: Routledge.
- Lewicki, A. M. (1976). *Wprowadzenie do frazeologii syntaktycznej: teoria zwrotu frazeologicznego*. Uniwersytet Śląski.
- Mehler, A. (2007). Large text networks as an object of corpus linguistic studies. *Corpus Linguistics. An International Handbook of the Science of Language and Society*. de Gruyter, Berlin/New York.
- Pawley, A., & Syder, F. H. (1983). Two puzzles for linguistic theory: Nativelike selection and nativelike fluency. *Language and communication*, 191, 225.

- Pęzik, P. (2013). Paradygmat dystrybucyjny w badaniach frazeologicznych. powtarzalność, reprodukcja i idiomatyzacja. In P. Stalmaszczyk (Ed.), *Metodologie językoznawstwa. ewolucja języka, ewolucja teorii językoznawczych*. Wydawnictwo Uniwersytetu Łódzkiego.
- Sinclair, J. (1996). The search for units of meaning. *Textus*, 9(1), 75–106.
- Siyanova-Chanturia, A., Conklin, K., & Schmitt, N. (2011, April). Adding more fuel to the fire: An eye-tracking study of idiom processing by native and non-native speakers. *Second Language Research*, 27(2), 251–272. Available from <http://slr.sagepub.com/cgi/doi/10.1177/0267658310382068>

- Tremblay, A., & Bayeen, H. (2010). Holistic processing of regular four-word sequences: A behavioural and ERP study of the effects of structure, frequency and probability on immediate free recall. In D. Wood (Ed.), *Perspectives on formulaic language: acquisition and communication*. London ; New York: Continuum.
- Warren, P. (2013). *Introducing psycholinguistics*. Cambridge: Cambridge University Press.
- Wray, A. (2002). *Formulaic language and the lexicon*. Cambridge (GB); New York: Cambridge university press.