



TOMASZ SZUTKOWSKI

 ORCID <http://orcid.org/0000-0003-3490-1707>

Uniwersytet Szczeciński

W POSZUKIWANIU KOLOKACJI I ICH ODPOWIEDNIKÓW PRZEKŁADOWYCH Z WYKORZYSTANIEM NARZĘDZI ELEKTRONICZNYCH I DANYCH KORPUSOWYCH

IN SEARCH FOR COLLOCATIONS AND THEIR TRANSLATIONAL EQUIVALENTS
WITH THE USE OF ELECTRONIC TOOLS AND CORPUS DATA

This paper presents a problem of stabilized word combinations (collocations) extraction and determining their translational equivalents in Polish and Russian languages. Multi-word language units have been the subject of linguistic research in various aspects of language for many years. However, still relatively few linguists and translators use in their work electronic corpora data, that in the 21st century should be considered as an important source of linguistic information, including real and potential translational equivalents. Author discusses the possibilities of using software as one of the methods of collocation extraction from text material and analyses selected units in a parallel corpus.

Przełom XX i XXI wieku naznaczony był gwałtownym rozwojem technologicznym w niemal wszystkich sferach ludzkiej działalności. Nie tylko nauki ścisłe i przyrodnicze zaczęły na wielką skalę wykorzystywać coraz większe możliwości obliczeniowe komputerów w celu gromadzenia i przetwarzania ogromnych ilości danych. Świat elektroniki nie pozostał jednak na peryferiach szeroko pojętych nauk humanistycznych, w tym także językoznawstwa, które wkroczyło w nową fazę rozwoju zwłaszcza pod szyldem lingwistyki komputerowej i korpusowej¹.

Na marginesie niniejszych rozważań pozostawiam kwestię genezy i aktualnego stanu rozwoju tej językoznawczej subdyscypliny².

¹ Zob. K. Church, M. Liberman, *The Future of Computational Linguistics: On Beyond Alchemy*, „Frontiers in Artificial Intelligence” 2021, vol. 4, s. 1–18.

² Zob. M. Świdziński, *Lingwistyka korpusowa w Polsce – źródła, stan, perspektywy*, „LingVaria” 2006, nr 1, s. 23–34.

W jej kontekście pragnę skupić się natomiast na kilku istotnych zagadnieniach, jakimi są: status wielowyrzowych jednostek języka (nazywanych niekiedy kolokacjami), możliwości ich wyodrębniania w konkretnych tekstach oraz wykorzystanie danych korpusowych w celu ustalenia optymalnych odpowiedników przekładowych. Z moich kilkunastoletnich obserwacji wynika, że w środowisku wielu językoznawców — szczególnie slawistów i rusycystów — zastosowanie narzędzi cyfrowych w badaniach nad językiem spotyka się z dość stanowczym oporem. Naturalnie, wciąż kompetencja specjalisty jest niezbędna w badaniach naukowych; bez niej choćby ogromna ilość wyekscerpowanych wartościowych danych pozostanie jedynie nieuporządkowanym i nieopracowanym materiałem. Warto jednak zwrócić uwagę na fakt, że dzięki różnorodnym programom, aplikacjom i korpusom ów materiał można nie tylko dużo szybciej pozyskać, ale przede wszystkim już na wstępnym etapie dokonać jego selekcji, segregacji czy klasyfikacji. Sposób wykorzystania technologii komputerowej w badaniach językoznawczych może mieć szeroką skalę i zróżnicowany charakter. Z pewnością nie we wszystkich tego typu działaniach niezbędna będzie profesjonalna kompetencja informatyczna, jak np. znajomość języków programowania, tworzenia złożonych baz danych, filtrów ekscerpcyjnych itp. W wielu sytuacjach wystarczy nawet podstawowa umiejętność obsługi komputera i zapoznanie się z nieskomplikowanymi narzędziami informatycznymi, aby móc skutecznie pozyskiwać oraz interpretować istotne dla badań naukowych dane. Celem niniejszego artykułu jest właśnie krótka prezentacja zautomatyzowanej metody ekscerpcji wielowyrzowych jednostek języka z dowolnych tekstów oraz ustalania ich odpowiedników na materiale polsko-rosyjskich relacji przekładowych.

Problem statusu wielowyrzowych jednostek języka nie jest w światowej lingwistyce ani nowy, ani tym bardziej zaskakujący. Można powiedzieć, że co najmniej od czasów Charlesa Bally'ego i jego koncepcji stylistyki funkcjonalnej, czyli od początku XX wieku, kwestia ustabilizowanych struktur leksykalnych na stałe wpisała się w dość wyraźnie zarysowany nurt badawczy, który później — dzięki m.in. pracom Wiktora Winogradowa — ukształtował współczesną frazeologię³. W skrajnej postaci przyjęła ona formę idiomatyki, jednak to nie tylko idiomy (np. „nie przebierać w słowach”, „być panem sytuacji”, „słomiana wdowa”, „ręka rękę myje”, „wyjść na światło dzienne”) wy-

³ Zob. M.A. Алексеенко, *Ещё раз о нерешённых проблемах фразеологии*, „Slavica Stetinensia” 1998, nr 8, s. 83–103.

kazują specyficzną spoiłość semantyczno-strukturalną. Między nimi a tzw. swobodnymi połączeniami wyrazowymi należy dostrzec bardzo obszerny i zróżnicowany jakościowo zbiór jednostek synkretycznych, peryferyjnych, częściowo zleksykalizowanych, znajdujących się niejako w pośrednim stadium procesu idiomatyzacji⁴ (np. „salon fryzjerski”, „świeżo malowane”, „szczyt kampanii”, „szczebel wojewódzki”, „wspominać nawiasowo”). Proces ten natomiast jest determinowany szczególnego rodzaju zależnościami lub inaczej – jak to ujął jeszcze w XIX wieku Jan Karłowicz – ustalonymi zwyczajem powinowactwami danej jednostki leksykalnej z innymi wyrazami⁵. Badania owych zależności czy powinowactw ma fundamentalne znaczenie dla badań językoznawczych, na co zwrócił uwagę Andrzej Bogusławski – „praktycznie zadaniem głównym”, „centralną kwestią lingwistyczną”, „pierwszoplanowym problemem badawczym” jest poszukiwanie i wyodrębnianie wielowyrzowych jednostek języka⁶, ponieważ ich liczba przybiera paradoksalnie milionowe wartości, co w porównaniu z jednowyrzowymi jednostkami jest wielkością niewspółmiernie większą⁷. Aspekt onomazjologiczny badań nad wielowyrzowymi jednostkami podkreśla szczególnie Wojciech Chlebda:

Jeżeli nazywanie jest aktem wieńczącym proces poznawania tego właśnie, a nie innego, wycinka, obszaru czy aspektu świata, to, chcąc mieć orientację w strukturze świata i naszych skomplikowanych z nim relacji, musimy w nieskończonym procesie poznawczym ustalać, jakie pojęcia o fragmentach tego świata jakim werbalizatorom odpowiadają. Mówiąc konkretniej: ustalmy granice werbalizatora, kształtowy zasięg i składowe danej formy słownej, a dojdziemy do ustalenia zasady odrębności i granic odpowiadającego mu pojęcia – tego właśnie, a nie innego, pojęcia o tym i tylko o tym wycinku, obszarze czy aspekcie rzeczywistości⁸.

⁴ T. Szutkowski, *Współczesna paremiografia rosyjska. Stan. Problemy. Perspektywy*, Volumina.pl, Szczecin 2015, s. 46.

⁵ J. Karłowicz, *Przyczynki do projektu wielkiego słownika języka polskiego*, „Rozprawy i sprawozdania Wydziału Filologicznego Akademii Umiejętności” 1876, t. IV, s. XXVIII–XXX.

⁶ A. Bogusławski, *Uwagi o pracy nad frazeologią*, w: Z. Saloni (red.), *Studia z polskiej leksykografii współczesnej*, t. 3, Wyd. Filii UW w Białymstoku, Białystok 1989, s. 18–19.

⁷ Skalę i wagę tego zjawiska bardzo dobrze oddaje nowatorskie opracowanie leksykograficzne zespołu pod kierunkiem Profesora Wojciecha Chlebdy: W. Chlebda (red.), *Polsko-rosyjski słownik par przekładowych. Tom zbiorczy Podręcznego idiomatykonu polsko-rosyjskiego (z. 1–5)*, Wyd. Uniwersytetu Opolskiego, Opole 2014.

⁸ W. Chlebda, *Reprodukty na warsztacie*, w: W. Chlebda (red.), *Na tropach reproduktów. W poszukiwaniu wielowyrzowych jednostek języka*, Wyd. Uniwersytetu Opolskiego, Opole 2010, s. 9–10.

Wraz z rozwojem badań nad semantyką, strukturą i funkcją związków wyrazowych wykazujących pewien stopień leksykalizacji pojawiły się różne koncepcje dotyczące ich nazewnictwa. Trudno byłoby w granicach jednego artykułu przedstawić choćby zarys tego problemu. Można jednak wymienić co najmniej kilkanaście terminów funkcjonujących we współczesnym językoznawstwie: „frazologizm”, „frazem”, „kolokacja”, „odtworzalny wielowyrazowiec”/„reprodukt”⁹, „peryfraz”, „analizm”, „konstrukcja opisowa”, „multiwerbizm”, „połączenie konwencjonalne”, „połączenie frazeologiczne”¹⁰ i in. Za każdą z tych propozycji terminologicznych stoją mniej lub bardziej przekonujące racje i stanowiska. Specyficzne rozmycie niektórych z tych pojęć wymaga od badacza niemal każdorazowego zdefiniowania używanych terminów. Taka sytuacja ma miejsce przede wszystkim w przypadku frazeologizmu, który w zależności od tzw. wąskiego lub szerokiego rozumienia frazeologii¹¹ może oznaczać diametralnie różne jakościowo byty językowe. Nie wdając się jednak w szczegółowe analizy, należy zauważyć, że wspólnym mianownikiem jednostek języka kryjących się za wskazanymi terminami jest przede wszystkim cecha odtwarzalności oraz ich co najmniej binarna struktura. I to właśnie konglomerat tych dwóch parametrów optymalnie oddaje termin „kolokacja” (z łac. „collocatio” – ‘ułożenie, układ, umieszczenie, ustawienie’), który zaistniał w środowisku brytyjskich językoznawców w latach pięćdziesiątych XX wieku¹². Badacze ci podkreślali głównie występowanie określonych schematów strukturalnych oraz specyficzne predyspozycje w zakresie łączliwości jednego leksemu z innymi leksemami. Istotę kolokacji trafnie oddaje definicja Tona van der Woudena: „Collocation is a term that refers to the mechanism, or fact, that certain words are regularly found in the company of other words”¹³. Można zatem stwierdzić, że każdy leksem niejako wytwarza szczególnego rodzaju sieć wiązań z innymi jednostkami leksykalnymi, preferuje współwystępowanie jednego elementu lub kilku, a z in-

⁹ W. Chlebda, *Reprodukty na warsztacie...*, s. 10.

¹⁰ E. Białek, *Kolokacja w przekładzie. Studium rosyjsko-polskie*, UMCS, Lublin 2009, s. 8.

¹¹ T. Szutkowski, *Współczesna paremiografia rosyjska...*, s. 37.

¹² Zob. J.R. Firth, *Papers in Linguistics 1934–1951*, Oxford University Press, London 1957. Stan badań w zakresie kolokacji w językoznawstwie rosyjskim, polskim i zachodnioeuropejskim omawia Ewa Białek: E. Białek, *Kolokacja w przekładzie...*, s. 8–14.

¹³ T. van der Wouden, *Negative Contexts. Collocation, Polarity and Multiple Negation*, Routledge, London 1997, s. 5.

nymi nie wchodzi w językową koegzystencję (np. „odnieść sukces”, „ponieść porażkę”, ale nie „*ponieść sukces”, „*odnieść porażkę”). W zjawisku tym uczestniczy oddziaływanie — z jednej strony — uzusu, a z drugiej — normy językowej¹⁴. Nie zawsze jednak ich wektory są zwrócone w tym samym kierunku¹⁵. Kolokacjom towarzyszy również zjawisko wariantywności. Ponadto, nie wszystkie muszą odznaczać się wysoką frekwencją — niektóre z nich to okazjonalizmy, językowe efemerydy.

Ustalenie wspomnianych przez Chlebdę granic jednostek językowych jest w zasadzie prostą procedurą jedynie w odniesieniu do pojedynczych leksemów. O wiele trudniej wyodrębnić jednostki nieciągłe, wieloelementowe czy — używając nomenklatury lingwistyki komputerowej — N-gramowe¹⁶ (dwu i więcej-gramowe: bigramy, trigramy), czyli kolokacji. Zadaniem badacza jest więc ustalenie stopnia ich spoistości semantyczno-strukturalnej, a jednym z jej głównych wskaźników będzie zjawisko reprodukowalności w mowie potwierdzone nie tylko frekwencją, ale także wartością tzw. siły kolokacji. Cel ten można osiągnąć, odwołując się do konkretnych danych tekstowych.

Językoznawca staje zatem w pierwszej kolejności przed wcale niełatwym zadaniem ekscerpacji kolokacji. Może je pozyskać na kilka sposobów. Jednym z nich jest wciąż metoda tradycyjna, „ręczna”, czy też — inaczej mówiąc — fiszkowa. Z pewnością bardzo pracochłonna i czasochłonna, gdyż wymaga uważnej lektury oraz analizy mniej lub bardziej obszernego materiału tekstowego. Przyczyny, które utrudniają ekscerpację kolokacji (wielowyrazowego reproduktu), wymienia i omawia szczegółowo Chlebda. Są to: 1) nieobecność wyraźnie zakreślonych granic początku i końca jednostki; 2) występowanie w strukturze jednostki elementów „widocznych” i „niewidocznych”; 3) wariantywność (wymienność komponentów, możliwość poszerzania lub skracania struktury wyjściowej); 4) możliwość istnienia jednostki rozrzuconej w tekście; 5) polisemia (szczególnie

¹⁴ Zob. M. Bańko, *Słownik dobrego stylu, czyli wyrazy, które się lubią*, PWN, Warszawa 2006; E. Geller, A. Dąbrówka, *Słownik stylistyczny języka polskiego*, Świat Książki, Warszawa 2007.

¹⁵ Wpisując w wyszukiwarkę www.google.pl niepoprawne warianty „*ponieść sukces”, „*odnieść porażkę”, możemy odnaleźć co najmniej kilkadziesiąt przykładów ich użycia w tekstach.

¹⁶ Zob. D. Jurafsky, J. Martin, *Speech and Language Processing*, Prentice-Hall, New Jersey 2008.

neosemantyzacja jednostki już w języku funkcjonującej)¹⁷. Ten niełatwy jak widać proces można nieco przyspieszyć, wykorzystując odpowiednie narzędzia informatyczne, które jednak absolutnie nie zwalniają badacza z niezbędnej czujności, wykorzystania własnej kompetencji oraz lingwistycznej intuicji. Metoda automatycznej ekscerpacji dostarczy określonego językowego „surowca”, wymagającego mimo wszystko dalszej starannej i rzetelnej obróbki w zależności od ustalonego celu badawczego.

W 2010 roku w jednym z artykułów wspólnie z Jurijem Fiedoruskowem¹⁸ omówiliśmy sposób ekscerpacji związków atrybutywnych z wykorzystaniem wyrażen regularnych (*regular expressions*)¹⁹. Metoda ta wymaga użycia specjalnych kodów, nazywanych inaczej ankodami (ros. аношкинский код), w postaci dwuliterowych identyfikatorów odpowiadających kombinacji parametrów semantyczno-gramatycznych²⁰. Procedura pozyskania zadanych modeli strukturalnych zakłada cztery etapy: 1) przygotowanie bazy ekscerpcyjnej w postaci tekstu lub tekstów (korpusu), 2) wprowadzenie tekstu do modułu analizatora morfologicznego w celu otrzymania odpowiednich danych w postaci tagów (znaczników) – do każdej jednostki leksykalnej zostaje przyłączony odpowiadający jej znacznik gramatyczny; 3) modyfikacja tekstu umożliwiająca zastosowanie filtrów ekscerpcyjnych składających się z ankodów i wyrażen regularnych; przykładowy tekst w postaci linearnej:

Беседа главного редактора газеты „Тамбовские известия” А.К. Клименко с одним из лидеров Тамбовского городского отделения партии пенсионеров [...]

przyjmuje postać następującego segmentu (wyraz + ankod):

Беседа га
главного йб

¹⁷ W. Chlebda, *Nieautomatyczne drogi dochodzenia do reproduktów wielowyrazowych*, w: W. Chlebda (red.), *Na tropach reproduktów...*, s. 19.

¹⁸ Ю. Федорущков, Т. Шутковски, *Лексико-грамматическая сочетаемость атрибутивных словосочетаний русского языка в контексте методов компьютерной экскерпции*, w: R. G. Tirado, L. Sokolova, I. Votyakova (red.), *La lengua y literatura rusa en el espacio educativo internacional: estado actual y perspectivas*, Rubiños, Granada 2010, s. 1564–1569.

¹⁹ Zob. J. Hopcroft, R. Motwani, J. Ullman, *Wprowadzenie do teorii automatów, języków i obliczeń*, przekł. B. Konikowska, PWN, Warszawa 2003.

²⁰ www.aot.ru (10.10.2021).

W POSZUKIWANIU KOŁOKACJI...

редактора аб
газеты гб
Тамбовские йт
известия еб;

4) zastosowanie zadanego ekscerptora, np. *.*йб\$ \n.*абv*, który wyselekcjonuje z danej próbki tekstu związki wyrazowe o strukturze ‘przymiotnik + rzeczownik’ w dopełniaczu liczby pojedynczej (np. „транспортного налога”, „прожиточного минимума”, „русского языка”)²¹.

Strukturę przykładowych filtrów ekscerptyjnych przedstawia poniższa tabela²²:

| Filtr ekscerptyjny | Wyekscerpowany związek wyrazowy |
|---------------------------|--|
| <i>^.*aa\$ \n.*\тйа</i> | красивый стол |
| <i>^.*аб\$ \n.*\тйб</i> | красивого стола |
| <i>^.*ав\$ \n.*\тйв</i> | красивому столу |
| <i>^.*аг\$ \n.*\тйг</i> | красивый стол |
| <i>^.*ад\$ \n.*\тйд</i> | красивым столом |
| <i>^.*ае\$ \n.*\тйе</i> | красивом столе |
| <i>^.*аж\$ \n.*\тйт</i> | красивые столы |
| <i>^.*аз\$ \n.*\тйу</i> | красивых столов |
| <i>^.*аи\$ \n.*\тйф</i> | красивым столам |
| <i>^.*ай\$ \n.*\тйх</i> | красивые столы |
| <i>^.*ак\$ \n.*\тйц</i> | красивыми столами |
| <i>^.*ал\$ \n.*\тйч</i> | красивых столах |

Omówiona pokrótce metoda umożliwia ekscerpowanie związków wyrazowych (w tym także kolokacji) o z góry założonej strukturze bez względu na stopień leksykalizacji ich komponentów składowych, co może być z punktu widzenia językoznawcy istotną wadą, gdyż wymaga skrupulatnej analizy całej listy otrzymanych wielowyrazowych struktur języka (selekcję danych może jednak znacząco usprawnić oprogramowanie statystyczne²³). Jej zaletą natomiast jest z pewnością możliwość szybkiego pozyskania danych do badań nad modela-

²¹ Ю. Федорущков, Т. Шутковски, *Лексико-грамматическая сочетаемость...*, s. 1567–1568.

²² Tamże, s. 1568.

²³ Zob. P. Wierzchoń, *Пять bardzo skuteczных (справданных) способов на массовое выделение многовыразовых сегментов подозранных о фразе-матичность (czyli reпродуктов)*, w: W. Chlebda (red.), *Na tropach reпродуктов...*, s. 106–107.

mi strukturalnymi grup wyrazowych (werbo-nominalnych, werbo-adwerbialnych, atrybutywnych itd.)²⁴.

Poszukiwanie kolokacji ułatwiają obecnie coraz lepiej rozwinięte narzędzia korpusowe. Dla badacza kluczową kwestią jest możliwość wyekscerpowania związków wyrazowych, które charakteryzują się już pewnym stopniem leksykalizacji, a więc nie kwalifikują się do zbioru tzw. swobodnych połączeń wyrazowych. Narodowy Korpus Języka Polskiego udostępnia użytkownikom specjalnie opracowany moduł ekscerpcyjny, nazywany kolokatorem²⁵. Narzędzie to umożliwia wyszukiwanie na podstawie kryterium leksykalnego ośrodka kolokacji. Po wpisaniu w okno kolokatora leksemu „alarm” otrzymujemy następujące wyniki (poniżej zamieszczam jedynie początkowe rekordy z listy rankingowej, zob. ilustrację na następnej stronie).

Kolokator Narodowego Korpusu Języka Polskiego podaje bardzo istotne dla badacza dane: liczbę wystąpień ośrodka kolokacji, pasujące współwystąpienia wraz z frekwencją, liczbę kolokacji spełniających zadane kryteria. Lista rankingowa jest jednak porządkowana według współczynnika χ^2 , dzięki któremu oblicza się wskaźnik siły kolokacji²⁶. Współczynnik ten, w przeciwieństwie do rozkładu częstości, pozwala na wyciągnięcie poprawnych wniosków dotyczących identyfikacji rzeczywistych, a nie pozornych wielowyrazowych jednostek języka²⁷, gdyż informuje o statystycznej istotności korelacji między dwiema określonymi zmiennymi, którymi w tym wypadku są komponenty danej kolokacji. Naturalnie, nie można uzyskanych danych traktować bezkrytycznie. Wyekscerpowane jednostki wymagają i tym razem od badacza analizy oraz weryfikacji. Struktura powyższej przykładowej listy rankingowej daje już pewien obraz łączliwości wybranego leksemu. Mamy zatem następujące związki wyrazowe: „fałszywy alarm”, „alarm powodziowy”, „alarm szalupowy”, „alarm przeciwpowodziowy”, „alarm bombowy”, „alarm pożarowy”, „wszczać alarm”, „alarm wibracyjny”, „alarm przeciwpowodziowy”, „alarm antyterrorystyczny”, „alarm przeciwlotniczy”,

²⁴ Zob. J. Fiedoruszkow, *Metody automatyzacji ekscerpcji konstrukcji atrybutywnych języka rosyjskiego*, w: W. Chlebda (red.), *Na tropach reproduktów...*, s. 59–85.

²⁵ www.nkjp.uni.lodz.pl (10.10.2021).

²⁶ Zob. E. Babbie, *Badania społeczne w praktyce*, PWN, Warszawa 2007.

²⁷ M. Hebal-Jezińska, *Podstawowe zasady korzystania z korpusów przy badaniu języka*, w: W. Chlebda (red.), *Na tropach korpusów...*, s. 29.

W POSZUKIWANIU KOLOKACJI...

Kryteria ośrodka kolokacji:

Maksymalny odstęp: Zachowaj szyk:

Kryteria kolokatu:
 Część mowy: Kontekst z lewej: Kontekst z prawej:
 Wielkość próbki: Min. współwystąpienie:

Zaawansowane

| | | |
|-----|--|--|
| 3. | szałupowy | alarm_szałupowy(2). |
| 4. | przeciwpowodziowy | alarm_przeciwpowodziowy(24). |
| 5. | bombowy | alarm_bombowy(2), bombowy_alarm(2). |
| 6. | pożarowy | alarm_pożarowy(8). |
| 7. | wszczać | wszczał_alarm(2), wszczęła_alarm(10), wszczął_alarm(10), wszczęto_alarm(9), wszczęli_alarm(6), alarm_wszczęty(2), wszczęły_alarm(2), alarm_wszczął(2), alarm_wszczęły(2), wszczęłyby_alarm(1), alarm_wszczęli(1), alarm_wszczęto(1), wszczęłyby_alarm(1). |
| 8. | wibracyjny | alarm_wibracyjny(2). |
| 9. | przeciwpożarowy | alarm_przeciwpożarowy(1), przeciwpożarowa_alarm(1). |
| 10. | antyterrorystyczny | alarm_antyterrorystyczny(2), antyterrorystyczny_alarm(2). |
| 11. | przeciwlotniczy | alarm_przeciwlotniczy(10). |
| 12. | próbny | próbny_alarm(24). |
| | podnieśli_alarm(2), podniosła_alarm(6), podniósł_alarm(2), podniesiono_alarm(6), podnieść_alarm(6), podniosły_alarm(6), alarm_podnieśli(2), podnieście_alarm(1), alarm_podniesiono(2), alarm_podnieśli(2), podnieścili_alarm(2), alarm_podnieść(2), podniosłem_alarm(1), podniesiesz_alarm(1), podniosło_alarm(1), podniosłam_alarm(1), alarm_podnieście(1), podnieśliśmy_alarm(1), podniosłyby_alarm(1). | |
| 14. | wszczynać | wszczynają_alarm(2), wszczynano_alarm(2), wszczynać_alarm(2), wszczynają_alarm(2), alarm_wszczynany(1), wszczynali_alarm(1), wszczynując_alarm(1). |
| 15. | podniosły | podniosła_alarm(6), podniosły_alarm(6), alarm_podniosła(2), podniosła_alarm(1), ogłoszono_alarm(2), ogłosił_alarm(2), ogłoszony_alarm(2), ogłosić_alarm(2), alarm_ogłoszono(4), ogłosili_alarm(2), ogłoszą_alarm(2), ogłosi_alarm(2), ogłoszeń_alarm(2), ogłoszo_alarm(1), ogłosio_alarm(1), ogłosily_alarm(1), ogłosiliśmy_alarm(1), alarm_ogłosily(1), alarm_ogłosili(1), alarm_ogłoszono(1). |
| 17. | zarządzić | zarządził_alarm(2), zarządził_alarm(2), zarządzić_alarm(2), zarządzi_alarm(2), zarządziłem_alarm(2), zarządziłi_alarm(1), zarządziła_alarm(1), zarządzić_alarm(1). |
| 18. | spłoszyć | spłoszył_alarm(2), alarm_spłoszył(1). |
| 19. | ogłaszać | ogłaszam_alarm(2), ogłasza_alarm(2), ogłaszamy_alarm(2), ogłaszając_alarm(1), ogłaszającego_alarm(1), alarm_ogłasza(1), alarm_ogłaszając(1). |
| 20. | włączyć | włączył_alarm(2), włączyła_alarm(2), włączył_alarm(2), włączy_alarm(2), alarm_włączył(2), włączony_alarm(2), włącznie_alarm(1), włączyć_alarm(1), włączyliśmy_alarm(1). |
| 21. | bojowy | alarm_bojowy(2), bojowej_alarm(1). |
| 22. | odwołać | alarm_odwołano(10), alarm_odwołany(2), odwołano_alarm(2), odwołał_alarm(2), odwołać_alarm(2), odwołany_alarm(2), odwołał_alarm(1), odwołalem_alarm(1), odwołaliśmy_alarm(1), alarm_odwoła(1). |

próbny alarm”, „podnieść alarm”, „wszczynać alarm” itd. Dane korpusowe odzwierciedla m.in. słownik Mirosława Bańki, który dla celów poprawnościowych odnotowuje najczęstsze kolokacje²⁸. Jak wiadać, wiele z nich pokrywa się z wynikami próbnej ekscerpcji:

alarm (stan alarmowy) • **falszywy** ~ • ~ **bojowy**, **bombowy**, **lotniczy** • ~ **przeciwpożarowy**, **przeciwpowodziowy** • **ogłosić** ~, **podnieść**, **wszczać**, **zarządzić** • **bić na** ~, **uderzyć na** ~, **narobić** ~u *pot.*

alarm (urządzenie) • **włączyć** ~, **wyłączyć**

album • ~ **fotograficzny**, **pamiątkowy**, **rodzinny**

aleja • **główna** ~a, **boczna** • **szeroka** ~a • **cienista** ~a • **brzozowa** ~a, **kasztanowa**, **lipowa** i in. • ~a **parkowa**, **spacerowa** • **iść** ~ą • ~a **prowadzi dokądś**, **wiedzie**

²⁸ M. Bańko, *Słownik dobrego stylu...*, s. 3.

Poza narzędziami korpusowymi językoznawca ma także do dyspozycji inne sposoby pozyskania potrzebnego materiału badawczego, jak na przykład odrębne oprogramowanie, w którym analizie możemy poddać dowolny tekst (w korpusach zasób tekstów jest w danym momencie ograniczony i zamknięty). Mowa o programie kfNgram opracowanym przez Williama Fletchera²⁹ w celu generowania listy n-gramów w plikach tekstowych *txt lub HTML. Wartość „n” przybiera postać dowolnej liczby całkowitej i stanowi określoną sekwencję słów wedle zadanych parametrów. W pierwszym etapie należy sprawdzić, czy plik z przeznaczonym do badań tekstem ma nadany właściwy format (*txt lub HTML). Prosty sposób zmiany formatu jest wklejenie tekstu do znajdującego się w akcesoriach środowiska MS Windows notatnika. Tak sformatowany materiał językowy można poddać analizie. Interfejs programu jest dość prosty w układzie i strukturze. Wybrany plik tekstowy należy wprowadzić do systemu. Przy ustawieniu minimalnego, jednokrotnego wystąpienia, wszystkie składniki tekstu (łącznie ze znakami interpunkcyjnymi) zostaną potraktowane jak oddzielne segmenty. Następnie podzielony na n-gramy plik można poddać analizie frekwencyjnej według ustalonego schematu strukturalnego, a więc np. na bigramy, trigramy, czyli związki wyrazowe 2-, 3-składnikowe. Wyniki analizy są prezentowane w formie listy rankingowej z podaną frekwencją oraz liczbą potencjalnych wariantów (zob. ilustrację na następnej stronie).

Na podstawie analizy fragmentu tekstu o tematyce językoznawczej uzyskujemy bigramy, z których część zakwalifikujemy do zbioru kolokacji, jak np. „łączliwość leksykalna”, „łączliwość składniowa”, „łączliwość wyrazów”, „współwystępowanie wyrazów”, „jednostka wielowyrazowa”, „właściwości składniowe” itd. Należy oczywiście zwrócić szczególną uwagę na to, że program kfNgram podaje wyłącznie współczynnik częstości, który absolutnie nie decyduje o tym, czy dany bigram lub trigram uzyskuje status kolokacji. Na tym etapieznów konieczna jest kompetencja badacza-językoznawcy, który ów status oceni i zweryfikuje (z pewnością kolokacjami nie będą związki: „podobny wyraz”, „jeden język”, „typ wyrazu”). Warto zauważyć, że omawiane narzędzie informatyczne może mieć daleko szersze zastosowanie, znacznie wykraczające poza wyłączone poszukiwanie wielowyrazowych jednostek języka (np. w analizie stylometrycznej).

²⁹ www.kwicfinder.com/kfNgram (12.10.2021).

W POSZUKIWANIU KOLOKACJI...

kfNgramPhraseFrames Browser beta 13 Oct 06

Open File Exit Find < Freqs Left Freqs Right > Help

33 phrase-frames with 100 variants in tekst.txt-02-Freq-phraseFrames.txt.
Select item to show phrase variants. Click on heading to re-sort.

| Rank | Phrase-Frame | Frequency | Variants |
|------|---------------|-----------|----------|
| 19 | * sie | 3 | 3 |
| 20 | * oraz | 3 | 3 |
| 21 | z * | 2 | 2 |
| 22 | wlasciwosci * | 2 | 2 |
| 23 | laczliwosc * | 2 | 2 |
| 24 | która * | 2 | 2 |
| 25 | * która | 2 | 2 |
| 26 | * laczliwosc | 2 | 2 |
| 27 | * ich | 2 | 2 |
| 28 | * jako | 2 | 2 |

laczliwosc *: rank 23, frequency 2, variants 2

Alpha Freq Copy p-frame + vars Copy all as sorted

laczliwosc skladniowa 1
laczliwosc leksykalna 1

Kolokacje w procesie translacji odgrywają fundamentalną rolę. Dotykamy w ten sposób kluczowego zagadnienia, jakim jest pojęcie i istota jednostki tłumaczenia. Białek wyróżnia dwa główne typy tych jednostek:

Przez leksykalną jednostkę przekładu rozumiem pojedyncze leksemy, przez ponadleksykalną zaś — związki wyrazowe, zdanie lub większe fragmenty oryginału, które tylko i wyłącznie jako całość są jednostką sensu, i tylko dla takiej całości może być dobrany odpowiednik w języku przekładu. Przy ponadleksykalnych jednostkach przekładu sens wypływa z integracji sensów wszystkich komponentów, a przekład z różnych względów na poziomie wyrazów nie jest możliwy lub nie jest wskazany³⁰.

W definicji tej badaczka podkreśla wagę integracji sensów komponentów wielowrazowej jednostki języka (a tym samym jednostki tłumaczenia), której nie można zignorować w procesie przekładu, gdyż w przeciwnym razie takie tłumaczenie może wprowadzić do tekstu w języku docelowym sens niezgodny lub całkowicie sprzeczny z intencją autora oryginału. Praca tłumacza to przede wszystkim ciągłe poszukiwanie adekwatnych środków wyrazu sensu, którego nośnikami są jednostki tekstu języka wyjściowego. W kontekście omawianej problematyki kolokacji należy mieć na względzie elementarny fakt — jak podkreśla Chlebda — że mówiący zazwyczaj posługują się w swoich wypowiedziach reprodukowanymi jednostkami ponad-

³⁰ E. Białek, *Kolokacja w przekładzie...*, s. 46.

leksykalnymi. Jednostki jednowyrazowe w zasobach języka tworzą zdecydowanie mniej liczny zbiór³¹. To jakże ważne spostrzeżenie nakazuje tłumaczowi — z jednej strony — dostrzeganie w tekście oryginału poszczególnych kolokacji jako jednostek będących nośnikiem sensu, z drugiej natomiast — dobieranie do nich ekwiwalentnych jednostek w języku przekładu. Obecnie w epoce intensywnego rozwoju technologii komputerowych sam słownik nie jest już wystarczającym źródłem ekwiwalencji przekładowej³². Korzystanie z korpusów i innych elektronicznych baz danych powoli staje się w codziennej pracy tłumaczy normą i nieodłącznym elementem praktyki zawodowej. Z moich obserwacji wynika, że zwłaszcza adepci niełatwego niewątpliwie rzemiosła translatorskiego, ale także sztuki translatorskiej, niezbyt chętnie korzystają z nowoczesnych narzędzi znacząco ułatwiających i przyspieszających realizację zadania doboru ekwiwalentnych jednostek tłumaczenia (wyjątkiem jest niewątpliwie translator Google).

Bardzo praktycznym narzędziem w ekscerpowaniu kolokacji oraz poszukiwaniu ich tekstowych ekwiwalentów jest czeski korpus elektroniczny InterCorp³³, przydatny nie tylko w eksploracjach języka czeskiego, gdyż dzięki innym podkorpusom umożliwia on analizę jednostek np. w polsko-rosyjskich relacjach przekładowych. Posługując się już wcześniejszym przykładem leksemu „alarm” jako ośrodka kolokacji, można zobrazować zasób potencjalnych wielowyrazowych jednostek języka (w tym przypadku — bigramów) wyekscerpowanych w podkorpusie języka polskiego (zob. ilustrację na następnej stronie)³⁴.

Z tabeli tej badacz może odczytać wiele bardzo wartościowych informacji, ponieważ korpus oblicza jednocześnie kilka istotnych w analizie kolokacji wskaźników statystycznych. Oprócz frekwencji są to: wartość MI3 i MI (prawdopodobieństwo wystąpienia dwóch słów jednocześnie), T-score (tzw. miara kontrastu — im wyższy wy-

³¹ W. Chlebda, *Nieautomatyczne drogi...*, s. 16.

³² O źródłach ekwiwalentów przekładowych pisze Wojciech Chlebda: W. Chlebda, *Korpusologia użytkowa dla początkujących i zaawansowanych*, w: W. Chlebda (red.), *Na tropach korpusów. W poszukiwaniu optymalnych zbiorów tekstów*, Wyd. Uniwersytetu Opolskiego, Opole 2013, s. 12–13.

³³ <https://www.korpus.cz> (15.10.2021).

³⁴ https://www.korpus.cz/kontext/collx?maincorp=intercorp_v13_pl&view-mode=kwic&pagesize=40&attrs=word&attr_vmode=visible-kwic&base_viewattr=word&refs=%3Ddoc.id&q=~OuekoW62Iuio&cattr=word&cfomw=-5&ctow=5&cminfreq=3&cminbgr=1&cbgrfns=3&cbgrfns=m&cbgrfns=t&cbgrfns=d&cbgrfns=p&cbgrfns=r&cbgrfns=s&cbgrfns=l&csortfn=d&collpage=1 (15.10.2021).

W POSZUKIWANIU KOLOKACJI...

| | Filtr | word | Freq | MI3 | MI | T-score | logDice ▼ | MI.log_f | relative freq. [%] | min_sensitivity | log likelihood |
|-----|-------|-----------------|------|--------|--------|---------|-----------|----------|--------------------|-----------------|----------------|
| 1. | p/n | Alarm | 74 | 25.915 | 13.496 | 8.602 | 10.012 | 58.267 | 20.274 | 0.037 | 1255.402 |
| 2. | p/n | falszywy | 78 | 25.546 | 12.975 | 8.831 | 9.977 | 56.693 | 14.130 | 0.039 | 1261.687 |
| 3. | p/n | Falszywy | 58 | 26.115 | 14.399 | 7.615 | 9.797 | 58.711 | 37.908 | 0.029 | 1068.886 |
| 4. | p/n | włączył | 60 | 23.968 | 12.154 | 7.744 | 9.490 | 49.964 | 8.000 | 0.030 | 897.740 |
| 5. | p/n | alarm | 52 | 22.796 | 11.395 | 7.208 | 9.110 | 45.242 | 4.727 | 0.026 | 721.352 |
| 6. | p/n | Czerwony | 34 | 21.341 | 11.166 | 5.828 | 8.622 | 39.699 | 4.033 | 0.017 | 460.304 |
| 7. | p/n | włączy | 25 | 22.580 | 13.292 | 5.000 | 8.590 | 43.307 | 17.606 | 0.013 | 415.674 |
| 8. | p/n | pożarowy | 22 | 24.421 | 15.503 | 4.690 | 8.486 | 48.608 | 81.481 | 0.011 | 456.184 |
| 9. | p/n | uruchomił | 19 | 20.728 | 12.232 | 4.358 | 8.139 | 36.644 | 8.444 | 0.010 | 286.030 |
| 10. | p/n | przeciwpożarowy | 17 | 23.016 | 14.841 | 4.123 | 8.109 | 42.896 | 51.515 | 0.009 | 326.740 |
| 11. | p/n | wyłączyć | 20 | 18.952 | 10.308 | 4.469 | 7.829 | 31.382 | 2.225 | 0.010 | 246.472 |
| 12. | p/n | biją | 16 | 19.240 | 11.240 | 3.998 | 7.795 | 31.844 | 4.244 | 0.008 | 218.132 |
| 13. | p/n | czerwony | 22 | 18.744 | 9.825 | 4.685 | 7.743 | 30.806 | 1.592 | 0.011 | 256.287 |
| 14. | p/n | wyłączony | 15 | 18.964 | 11.150 | 3.871 | 7.702 | 30.915 | 3.989 | 0.008 | 202.596 |
| 15. | p/n | falszywego | 14 | 19.170 | 11.555 | 3.740 | 7.672 | 31.293 | 5.283 | 0.007 | 197.131 |
| 16. | p/n | bije | 18 | 18.389 | 10.049 | 4.239 | 7.642 | 29.589 | 1.860 | 0.009 | 215.290 |
| 17. | p/n | Ogłaszam | 13 | 19.632 | 12.232 | 3.605 | 7.638 | 32.280 | 8.442 | 0.007 | 195.656 |
| 18. | p/n | bojowy | 13 | 19.507 | 12.106 | 3.605 | 7.629 | 31.949 | 7.738 | 0.007 | 193.298 |
| 19. | p/n | Ogłosić | 11 | 21.717 | 14.798 | 3.317 | 7.489 | 36.772 | 50.000 | 0.006 | 210.470 |
| 20. | p/n | włączony | 13 | 17.934 | 10.533 | 3.603 | 7.422 | 27.796 | 2.600 | 0.007 | 164.260 |

nik wskaźnika, tym większe prawdopodobieństwo, że dana kombinacja słów jest bardziej ustabilizowana, a więc może stanowić kolokację), logDice (tzw. logarytm Dice'a – pomiar częstotliwości bigramu ze średnią częstotliwością jego składowych), frekwencja względna oraz wskaźnik asocjacji (stosunek częstotliwości bigramu do częstotliwości jednego z jego składników)³⁵. W zamieszczonej powyżej liście rankingowej wybrano jako wskaźnik selekcji parametr logDice, którego obliczenia siły kolokacji nie zależą od wielkości korpusu, co z kolei bardzo ułatwia porównywanie danych między różnymi bazami tekstowymi. Wskaźnik logDice przyjmuje wartość od $-\infty$ do $+14$. Unifikując m.in. warianty graficzne (np. mała i duża litera) i gramatyczne, wysoki wskaźnik siły kolokacji wykazują następujące jednostki: „falszywy alarm”, „włączyć alarm”, „czerwony alarm”, „alarm pożarowy”, „uruchomić alarm”, „alarm przeciwpożarowy”, „wyłączyć alarm”, „bić na alarm”, „ogłaszać/ogłosić alarm”, „alarm bojowy”, „włączony alarm”. Kolokacje te przyjmują wartość logDice w przedziale (7,422, 9,977). W korpusie można zweryfikować kontekstowo każde wyekscerpowane wystąpienie analizowanej jednostki³⁶:

³⁵ Wszystkie wskaźniki wraz ze wzorami zostały dokładnie omówione pod adresem: https://wiki.korpus.cz/doku.php/pojmy:asociaelni_miry?redirect=1#mi-score_a_mi3 (15.10.2021).

³⁶ https://www.korpus.cz/kontext/view?maincorp=intercorp_v13_pl&view-mode=align&pagesize=40&attrs=word&attr_vmode=visible-kwic&base_viewattr=word&refs=%3Ddoc.id&q=~asWYgci066KI (15.10.2021).

Vjskytů: 7 | l.p.m.: 0,06 (vztaženo k celému korpusu) | ARF: 2,89 | Výsledek je setříděn 1 / 1

Vyber řádků: základní

InterCorp.v13 - Polish

| | | | |
|--------------------------|-------------|---|-----|
| <input type="checkbox"/> | ._SUBTITLES | Czerwony alarm . | ... |
| <input type="checkbox"/> | ._SUBTITLES | - Czerwony alarm w sercu . | ... |
| <input type="checkbox"/> | ._SUBTITLES | Ochrona ... tu jadalnia ... mamy czerwony alarm powtarzam czerwony alarm . | ... |
| <input type="checkbox"/> | ._SUBTITLES | Ochrona ... tu jadalnia ... mamy czerwony alarm powtarzam czerwony alarm . | ... |
| <input type="checkbox"/> | ._SUBTITLES | - Odwołaj czerwony alarm . | ... |
| <input type="checkbox"/> | ._SUBTITLES | Czerwony alarm ! | ... |
| <input type="checkbox"/> | ._SUBTITLES | Czerwony alarm ! | ... |

◀ silnie toksyczny . Kratery wskazują na bombardowania z orbity . - Jak dawno temu ? - Według tempa rozpadu promieniotwórczego , osiemset dziewięćdziesiąt dwa lata temu . Musiały tu żyć miliony . Wygląd . - Wyłączam silniki . - Odwołaj czerwony alarm . Przydziel duży nacisk . Tom , upewnij się , że B'Elanna ma wystarczającą pomoc w maszynowni . Chcę mieć z powrotem silniki warp . Tak jest . Nasi przyjaciele są nadal na orbicie . - Krążą jak sępy . ▶

1 / 1

Na podstawie analizy poszczególnych przykładów i kontekstów ich użycia wyłania się kolokacja „czerwony alarm” w znaczeniu ‘stan gotowości, mobilizacja’ (synonim kolokacyjny: „czerwony alert”). Podobne badanie można przeprowadzić w innym języku. Mając na względzie polsko-rosyjskie relacje przekładowe, wykorzystamy w tym celu podkorpus języka rosyjskiego. Po wpisaniu w kolokator leksemu „тревога” uzyskujemy następujące wyniki wyszukiwania:

| Filtr | word | Freq | MI3 | MI | T-score | logDice | log likelihood | min_sensitivity | MI_log_f | relative freq. [%] | |
|-------|-------|-----------|-----|--------|---------|---------|----------------|-----------------|----------|--------------------|--------|
| 1. | p / n | вызывает | 44 | 21.075 | 10.156 | 6.627 | 9.484 | 535.521 | 0.039 | 38.661 | 4.977 |
| 2. | p / n | сигнал | 38 | 20.855 | 10.360 | 6.160 | 9.440 | 473.300 | 0.034 | 37.953 | 5.732 |
| 3. | p / n | Ложная | 17 | 22.499 | 14.324 | 4.123 | 8.921 | 328.827 | 0.015 | 41.402 | 89.474 |
| 4. | p / n | вызывают | 15 | 17.834 | 10.021 | 3.869 | 8.394 | 179.292 | 0.013 | 27.783 | 4.532 |
| 5. | p / n | Тревога | 12 | 19.262 | 12.092 | 3.463 | 8.365 | 179.739 | 0.011 | 31.016 | 19.048 |
| 6. | p / n | тревога | 13 | 17.649 | 10.248 | 3.603 | 8.275 | 159.566 | 0.012 | 27.046 | 5.306 |
| 7. | p / n | ложная | 11 | 19.505 | 12.586 | 3.316 | 8.266 | 173.300 | 0.010 | 31.276 | 26.829 |
| 8. | p / n | поводу | 32 | 18.240 | 8.240 | 5.638 | 8.204 | 303.086 | 0.013 | 28.812 | 1.319 |
| 9. | p / n | Сигнал | 11 | 18.233 | 11.315 | 3.315 | 8.196 | 151.924 | 0.010 | 28.116 | 11.111 |
| 10. | p / n | воздушной | 9 | 18.136 | 11.796 | 2.999 | 7.956 | 130.731 | 0.008 | 27.162 | 15.517 |
| 11. | p / n | ложной | 9 | 17.709 | 11.369 | 2.999 | 7.932 | 125.006 | 0.008 | 26.178 | 11.538 |
| 12. | p / n | голосе | 10 | 15.633 | 8.989 | 3.156 | 7.695 | 104.971 | 0.009 | 21.556 | 2.217 |
| 13. | p / n | объявлена | 7 | 17.125 | 11.511 | 2.645 | 7.597 | 98.678 | 0.006 | 23.935 | 12.727 |
| 14. | p / n | чувство | 16 | 15.557 | 7.557 | 3.979 | 7.412 | 136.139 | 0.008 | 21.410 | 0.821 |
| 15. | p / n | тепловой | 6 | 16.917 | 11.748 | 2.449 | 7.393 | 86.697 | 0.005 | 22.860 | 15.000 |
| 16. | p / n | Внимание | 8 | 14.617 | 8.617 | 2.821 | 7.359 | 79.803 | 0.007 | 18.934 | 1.713 |
| 17. | p / n | сигналом | 6 | 15.748 | 10.578 | 2.448 | 7.332 | 76.431 | 0.005 | 20.583 | 6.667 |
| 18. | p / n | охватила | 6 | 15.080 | 9.910 | 2.447 | 7.271 | 70.725 | 0.005 | 19.283 | 4.196 |
| 19. | p / n | поднять | 8 | 14.151 | 8.151 | 2.818 | 7.206 | 74.613 | 0.007 | 17.910 | 1.240 |
| 20. | p / n | Повторяю | 6 | 14.458 | 9.288 | 2.446 | 7.186 | 65.472 | 0.005 | 18.074 | 2.727 |

Przyjmując ponownie za kryterium siły kolokacji wartość wskaźnika logDice, można wyodrębnić w analizowanych tekstach następujące kolokacje: „вызывать тревогу”, „сигнал тревоги”, „ложная тревога”, „тревога по поводу”, „воздушная тревога”, „тревога в голосе”, „объявить тревогу”, „чувство тревоги”, „тревога по тепловой защите”, „тревога охватила кого-либо”, „поднять тревогу” i in.

W POSZUKIWANIU KOLOKACJI...

Korpus InterCorp umożliwia również wyszukiwanie ekwiwalentów nie tylko na poziomie jednostek jednowyrazowych, ale także całych kolokacji. Analiza poszczególnych kontekstów daje tłumaczowi szerokie spektrum potencjalnych translatów, znacznie w wielu przypadkach wykraczających poza dane słowników dwujęzycznych. Po wpisaniu w okno wyszukiwarki kolokatora jednostki „czerwony alarm”, korpus wyświetla następujące dane³⁷:

The screenshot shows a search interface for the InterCorp v13 corpus. The search term is 'czerwony alarm'. The results are displayed in a table with two columns: 'InterCorp v13 - Polish' and 'InterCorp v13 - Russian'. The table contains several rows of results, each with a checkbox, a subtitle, and the corresponding text in both languages. A tooltip is visible over one of the results, providing a detailed explanation of the context.

| InterCorp v13 - Polish | InterCorp v13 - Russian |
|--|--|
| <input type="checkbox"/> .SUBTITLES Czerwony alarm . | <input type="checkbox"/> .SUBTITLES Потеше . |
| <input type="checkbox"/> .SUBTITLES - Czerwony alarm w sercu . | <input type="checkbox"/> .SUBTITLES Кардиологам приготовиться . |
| <input type="checkbox"/> .SUBTITLES Ochrona ... tu jadalnia ... mamy czerwony alarm ... powtarzam czerwony alarm . | <input type="checkbox"/> .SUBTITLES Центральная охрана, это столовая . У нас чрезвычайная ситуация . Повторю : Чрезвычайная ситуация . |
| <input type="checkbox"/> .SUBTITLES Ochrona ... tu jadalnia ... mamy czerwony alarm ... powtarzam czerwony alarm . | <input type="checkbox"/> .SUBTITLES Центральная охрана, это столовая . У нас чрезвычайная ситуация . Повторю : Чрезвычайная ситуация . |
| <input type="checkbox"/> .SUBTITLES - Odwołaj czerwony alarm . | <input type="checkbox"/> .SUBTITLES Я отработку за вас оставляю часть смены . |
| <input type="checkbox"/> .SUBTITLES Czerwony alarm ! | <input type="checkbox"/> .SUBTITLES Тревога ! |
| <input type="checkbox"/> .SUBTITLES Czerwony alarm ! | <input type="checkbox"/> .SUBTITLES Тревога ! |

Tooltip text: « не хотела причинить ему вред, но он не послушал меня и не остановился . Эти идиоты видели слишком много старых фильмов о тюрьме . 28-е марта - 16 дней до поединка Иди к телефону и объяви чрезвычайную ситуацию . Центральная охрана , это столовая . У нас чрезвычайная ситуация . Повторю : Чрезвычайная ситуация . Айсмен ! Айсмен . До того , как ты попал сюда , мы считали тебя героем . Но потом ты приехал и вел себя так . »

Wyłączając z analizy tłumaczenie opisowe lub opuszczenia, na podstawie wyekscerpowanych danych można ustalić co najmniej dwie pary przekładowe: „czerwony alarm” → „чрезвычайная ситуация”, „czerwony alarm” → „тревога”. W podkorpusie języka rosyjskiego możemy co prawda odnaleźć jednostkę „красная тревога”, jednak w korpusie paralelnym nie ma ani jednego przykładu polskiego translatu:

The screenshot shows a search interface for the InterCorp v13 corpus. The search term is 'красная тревога'. The results are displayed in a table with two columns: 'InterCorp v13 - Polish' and 'InterCorp v13 - Russian'. The table contains several rows of results, each with a checkbox, a subtitle, and the corresponding text in both languages.

| InterCorp v13 - Polish | InterCorp v13 - Russian |
|--|---|
| <input type="checkbox"/> .SUBTITLES пожар . - Пожар ! Ты сукин ... Красная тревога | <input type="checkbox"/> .SUBTITLES Повторю - красная тревога . - Пожар ! - |
| <input type="checkbox"/> .SUBTITLES красная тревога . - Пожар ! - Пожар ! Красная тревога | <input type="checkbox"/> .SUBTITLES ! Готовы ? Сам ! Торопитесь ! Торопитесь ! Пошли |
| <input type="checkbox"/> .SUBTITLES она это перенесет . Это еще что такое ? Красная тревога | <input type="checkbox"/> .SUBTITLES ! Красная тревога ! Положе , на заводе произошла авария |
| <input type="checkbox"/> .SUBTITLES . Это еще что такое ? Красная тревога ! Красная тревога | <input type="checkbox"/> .SUBTITLES ! Пожме , на заводе произошла авария . Назад ! |
| <input type="checkbox"/> .SUBTITLES покрову планеты ? - 78.6 % Всем станциям . Красная тревога | <input type="checkbox"/> .SUBTITLES . Переходим на планетарные запасы энергии . Минстер президент , |

³⁷ https://www.korpus.cz/kontext/view?maincorp=intercorp_v13_pl&viewmode=align&pagesize=40&attrs=word&attr_vmode=visible-kwic&base_viewattr=word&refs=%3Ddoc.id&q=~buIQek8Sk6o8 (15.10.2021).

Ekwiwalencję tych dwóch kolokacji należałoby zatem zbadać na podstawie kontekstów ich użycia odrębnie w języku polskim i języku rosyjskim. Warto zwrócić uwagę na to, że słowniki polsko-rosyjskie nie odnotowują tej kolokacji. W takich sytuacjach wartość danych korpusowych jest tym bardziej wymierna i cenna.

Z powyższych rozważań wypływa kilka zasadniczych wniosków. Po pierwsze, w badaniach językoznawczych zdecydowanie większą uwagę należy poświęcić wielowyzrazowym jednostkom języka w różnych aspektach – formalnym, semantycznym i funkcjonalnym (pragmatycznym). W wielu ośrodkach naukowych widać od kilkunastu lat postęp tych badań i ich konkretne rezultaty. Przykładem tego jest przede wszystkim seria idiomatykonu, nad którą intensywnie pracuje opolski zespół filologów pod kierunkiem Profesora Wojciecha Chlebdy. Powyższy postulat nie wynika absolutnie z subiektywnych preferencji lub językoznawczej „mody”, lecz z obiektywnych przesłanek dotyczących ontologii języka naturalnego i jego specyfiki, o czym była mowa na wstępie niniejszego artykułu. Po drugie, mając na względzie fakt, że ilościowo kolokacje przeważają nad jednowyrazowymi jednostkami, proces ich wyodrębniania zwłaszcza dla celów leksykograficznych i przekładowych można znacząco przyspieszyć i ułatwić, stosując odpowiednie narzędzia elektroniczne. Jak starałem się to ukazać, nie są one nazbyt skomplikowane, co powinno – w moim przekonaniu – zachęcić zarówno językoznawców starszego, jak i młodszego pokolenia do korzystania z nich w codziennej pracy naukowej i dydaktycznej. Zaprezentowane dane statystyczne tym bardziej podnoszą jakość analizy, a wnioski wyprowadzone na jej podstawie stają się jeszcze bardziej obiektywne. Po trzecie, elektroniczne bazy danych w postaci korpusów (ale nie tylko) znajdują bardzo praktyczne zastosowanie w tłumaczeniach, a proces ustalania par przekładowych na podstawie realnych danych tekstowych (nie kwestionując absolutnie wartości tradycyjnych słowników dwujęzycznych) nie tylko może przyczynić się do podniesienia jakości tłumaczonych tekstów, ale także jest w stanie realnie wspomóc osiągnięcia współczesnej leksykografii przekładowej. Na koniec warto przywołać coraz bardziej aktualne i popularne w językoznawstwie pojęcie frazeotranslacji, które wydaje się swego rodzaju naukowym spoiwem powyższych konstatacji³⁸. Temu zagadnieniu należałoby jednak poświęcić odrębne opracowanie.

³⁸ Zob. J.-P. Colson, *Computational phraseology and translation studiem: from theoretical hypotheses to practical tools*, w: G. Corpas Pastor, Jean-Pierre Colson (red.), *Computational Phraseology*, John Benjamins, Amsterdam-Philadelphia 2020, s. 65–81.

REFERENCES

- Alekseyenko, Mikhail. "Yeshchë raz o nereshënykh problemakh frazeologii." *Slavica Stetinensia*, 1998, no. 8: 83–103 [Алексеенко, Михаил. "Ещё раз о нерешённых проблемах фразеологии." *Slavica Stetinensia* 1998, no. 8: 83–103].
- Babbie, Earl. *Badania społeczne w praktyce*. Warszawa: PWN, 2007.
- Bańko, Mirosław. *Słownik dobrego stylu, czyli wyrazy, które się lubią*. Warszawa: PWN, 2006.
- Białek, Ewa. *Kolokacja w przekładzie. Studium rosyjsko-polskie*. Lublin: UMCS, 2009.
- Bogusławski, Andrzej. "Uwagi o pracy nad frazeologią." *Studia z polskiej leksykografii współczesnej*. Tom 3. Saloni, Zygmunt. Ed. Białystok: Wyd. Filii UW w Białymstoku, 1989: 13–30.
- Chlebda, Wojciech. "Korpusologia użytkowa dla początkujących i zaawansowanych." *Na tropach korpusów. W poszukiwaniu optymalnych zbiorów tekstów*. Chlebda, Wojciech, Ed. Opole: Wyd. Uniwersytetu Opolskiego, 2013: 7–15.
- Chlebda, Wojciech. "Nieautomatyczne drogi dochodzenia do reproduktów wielowyrazowych." *Na tropach reproduktów. W poszukiwaniu wielowyrazowych jednostek języka*. Chlebda, Wojciech. Ed. Opole: Wyd. Uniwersytetu Opolskiego, 2010: 15–35.
- Chlebda, Wojciech. Ed. *Polsko-rosyjski słownik par przekładowych. Tom zbiorczy Podręcznego idiomatykonu polsko-rosyjskiego (z. 1–5)*. Opole: Wyd. Uniwersytetu Opolskiego, 2014.
- Chlebda, Wojciech. "Reprodukty na warsztacie." *Na tropach reproduktów. W poszukiwaniu wielowyrazowych jednostek języka*. Chlebda, Wojciech. Ed. Opole: Wyd. Uniwersytetu Opolskiego, 2010: 7–13.
- Church, Kenneth. Liberman, Mark. "The Future of Computational Linguistics: On Beyond Alchemy." *Frontiers in Artificial Intelligence*, 2021, vol. 4: 1–18.
- Colson, Jean-Pierre. "Computational phraseology and translation studies: from theoretical hypotheses to practical tools." *Computational Phraseology*. Corpas Pastor, Gloria. Colson, Jean-Pierre. Eds. Amsterdam-Philadelphia: John Benjamins Publishing Company, 2020: 65–81.
- Fedorushkov, Yuriy, and Shutkovski, Tomash. "Leksiko-grammaticheskaya sochetayemost' atributivnykh slovosochetaniy russkogo yazyka v kontekste metodov komp'yuternoy ekspterpsii." *La lengua y literatura rusas en el espacio educativo internacional: estado actual y perspectivas*. Tirado, Rafael Guzmán, and Sokolova, Larisa, and Votyakova, Irina. Eds. Granada: Rubiños, 2010. 1564–1569 [Федорушков, Юрий, and Шутковски, Томаш. "Лексико-грамматическая сочетаемость атрибутивных словосочетаний русского языка в контексте методов компьютерной экперпсии." *La lengua y literatura rusas en el espacio educativo internacional: estado actual y perspectivas*. Tirado, Rafael Guzmán, and Sokolova, Larisa Votyakova, and Irina. Eds. Granada: Rubiños, 2010: 1564–1569].
- Fiedoruszkow, Jurij. "Metody automatyzacji ekscerpcji konstrukcji atrybutywnych języka rosyjskiego." *Na tropach reproduktów. W poszukiwaniu wielowyrazowych jednostek języka*. Chlebda, Wojciech. Ed. Opole: Wyd. Uniwersytetu Opolskiego, 2010: 59–85.
- Firth, John Rupert. *Papers in Linguistics 1934–1951*, London: Oxford University Press, 1957.

- Geller, Ewa. Dąbrowka, Andrzej. *Słownik stylistyczny języka polskiego*. Warszawa: Świat Książki, 2007.
- Hebal-Jeziarska, Milena. "Podstawowe zasady korzystania z korpusów przy badaniu języka." *Na tropach korpusów. W poszukiwaniu optymalnych zbiorów tekstów*. Chlebda, Wojciech. Ed. Opole: Wyd. Uniwersytetu Opolskiego, 2013: 17–30.
- Hopcroft, John. Motwani, Rajeev. Ullman, Jeffrey. *Wprowadzenie do teorii automatów, języków i obliczeń*. Warszawa: PWN, 2003.
- Karłowicz, Jan. "Przyczynki do projektu wielkiego słownika języka polskiego." *Rozprawy i sprawozdania Wydziału Filologicznego Akademii Umiejętności 1876*, tom 4: 28–30.
- Jurafsky, Dan. Martin, James. *Speech and Language Processing*. New Jersey: Prentice-Hall, 2008.
- Szutkowski, Tomasz. "Subkompetencja frazeologiczna studentów-neofilologów jako składnik kompetencji językowej tłumacza. Badanie pilotażowe." *Roczniki Humanistyczne*, 2019, t. LXVII, z. 7: 93–109.
- Szutkowski, Tomasz. *Współczesna paremiografia rosyjska. Stan. Problemy. Perspektywy*. Szczecin: Volumina.pl, 2015.
- Świdziński, Marek. "Lingwistyka korpusowa w Polsce — źródła, stan, perspektywy." *LingVaria*, 2006, no. 1: 23–34.
- Wierzchoń, Piotr. "Pięć bardzo skutecznych (sprawdzonych) sposobów na masowe wyodrębnianie wielowyrzowych segmentów podejrzanych o frazematyczność (czyli reproduktów)." *Na tropach reproduktów. W poszukiwaniu wielowyrzowych jednostek języka*. Chlebda, Wojciech. Ed. Opole: Wyd. Uniwersytetu Opolskiego, 2010: 87–125.
- Wouden, Ton van der. *Negative Contexts. Collocation, Polarity and Multiple Negation*. London: Routledge, 1997.