

UDC [811.161.2+811.162.1+811.111]:81'33

DOI <https://doi.org/10.32447/2663-340X-2022-11.23>

CORPUS-BASED APPROACH IN THE STUDY OF VERBAL PREDICATES

Nazarchuk Roksolana Zinoviivna

*Candidate of Philological Sciences,
Senior Lecturer at the Department of Applied Linguistics
Lviv Polytechnic National University
12, Bandery Str., Lviv, Ukraine*

Karamysheva Iryna Damirivna

*Candidate of Philological Sciences,
Associate Professor at the Department of Applied Linguistics
Lviv Polytechnic National University
12, Bandery Str., Lviv, Ukraine*

The article highlights the features of the corpus approach to the analysis of verbal predicates of Ukrainian, Polish and English languages. The problems of corpus linguistics are outlined, and the corpora of texts are reviewed, namely the General Regionally Annotated Corpus of Ukrainian (GRAC), the National Polish Corpus (NKJP), and the British National Corpus (BNC). These resources, meeting the requirements of authenticity, representativeness (balance), and sufficiency in volume (S. Buk, O. Demska, Ch. Fellbaum, S. Gries, A. Stefanowitsch M. Shvedova, V. Shyrovok and others), allow to optimize and to objectify the interpretation of language material as well as to obtain high reliability of results. The scientific novelty of the study is a comprehensive comparative analysis of the functional capabilities of the corpora of texts in Ukrainian, Polish, and English to determine their features, as well as outlining the prospects of applying a corpus-based approach to identify object connections of verbal predicates. The basic concepts of corpus linguistics are introduced: lemma, word form, token, label (tag), co-usage (collocation), concordance, frequency, and CQL query; they are illustrated by the implementations of the corresponding referentially specialized verbal predicates and their objects based on GRAC, NKJP, and BNC. In particular, it is shown that the simplest use of text corpora is to certify the use of a unit (concordance) and its frequency (i.e. the number of repetitions in the corpus). For languages rich in word change, it is important to be able to identify different forms of one word with its basic form – the lemma. The presence of a tag (POS, part of speech) allows acquiring the co-usage (collocations) of two words. It has been illustrated how the studied corpora allow further grouping, frequency ordering, and isolation of specific word combinations, and the markup of parts of speech facilitates the search for all collocations of verbal predicates with potential object-nouns at a distance of 4 tokens. The function of searching for co-usage (collocations) of verbs, grouped by lemmas or grammatical characteristics and ordered by frequency, is described. The advantages of semantic markup in the GRAC corpus for the automatic search of language units and their combinations are pointed out. All the described features of the text corpora are accompanied by illustrations from GRAC, NKJP, and BNC.

Key words: *corpus linguistics, corpus, GRAC, NKJP, BNC, verb, object, lemma, collocation*

Introduction. Linguistic studies using corpora, i.e. large collections of original texts, are gaining popularity along with the development of information technology that allows automatic search and analysis of language units and their combinations.

Characterizing the peculiarities of relations between a verbal predicate and the object, we consider it appropriate to use a corpus-based approach, which, according to researchers, ensures the objectivity of results and prevents introspective interpretation: “the standard procedures for accessing corpora (*concordances, collocate lists, frequency lists*) are a natural step towards identifying the relevant distributions in the first place” (Stefanowitsch, 2020, p. 59).

Problem setting. Following the works (Buk, 2001, p. 62–65; Demska, 2011, p. 83–89; Fellbaum, 2015, 2019; Gries, 2006; Stefanowitsch, 2020; Shvedova, 2020; Shyrovok, 2005, p. 12–13), text corpus means a set of language or speech data described linguistically competently, presented in the electronic form and fitted with appropriate specialized software, intended for a variety of studies, that meets the requirements of authenticity, representativeness (balance), sufficiency in volume. As O. Demska rightly points out, corpus “is formed from real fragments of written or spoken speech, without providing for the modification of speech reality, which turns it into an empirical category and allows considering the actual corpus

material as an empirical basis for linguistic study” (Demska, 2011, p. 41).

The authenticity provides for the involvement of texts of artistic, scientific, popular scientific works, periodicals, transcription of radio and television programs, etc. The Representativeness (or balance) of the collection means the reproduction of quantitative and qualitative diversity of areas of the real use of a particular language; for example, preserving the relative share of artistic and scientific texts. How big the corpus should depend on the language; the existing corpora of the main European languages range from hundreds of millions to over a billion tokens (along with a lemma as an initial form of a word, corpus linguistics uses the broader concept of a token as the smallest unit into which the corpus is divided, with any quantitative characteristic in a text; tokens are any sequence of characters between spaces or other separators: word form, number, punctuation mark, symbol (smiley, mathematical symbol, etc.). It should be noted that the mechanical increase in the size of corpora today does not require extraordinary efforts; a much bigger challenge is to maintain representativeness.

Analysis Procedure. The General Regionally Annotated Corpus of Ukrainian (GRAC) is publicly available and contains over 600 million words. The National Corpus of Polish (NKJP) comprises about 1.5 billion words, and the British National Corpus (BNC) contains about 100 million words. These corpora were created by research institutions of the respective countries, and they meet the already mentioned requirements.

An important point for the efficient use of the corpus of any language in linguistic studies is the presence of marking. Usually, corpus compilers keep relevant bibliographic data (e.g. author’s name, title, edition, year, page, etc.) for each fragment of a text, and add some grammatical information to each word (e.g., belonging to a part of speech, gender, case, etc.). The branching of the marking affects the qualitative opportunities of studies based on the corpus data. Thus, the semantic marking in the GRAC corpus of Ukrainian makes it possible to distinguish between abstract and common nouns denoting living beings or non-living beings, etc. The simplest use of language corpora is to certify the use

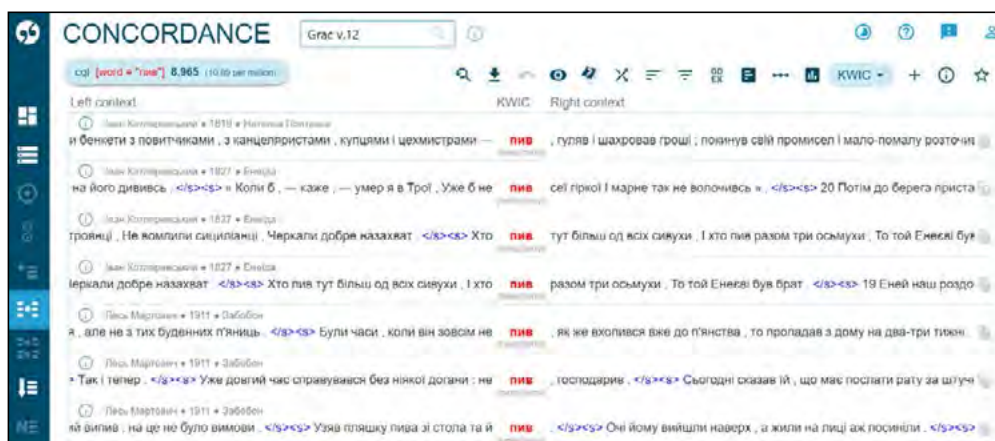


Figure 1. Examples of the use of the word form *пив* in the GRAC corpus of Ukrainian

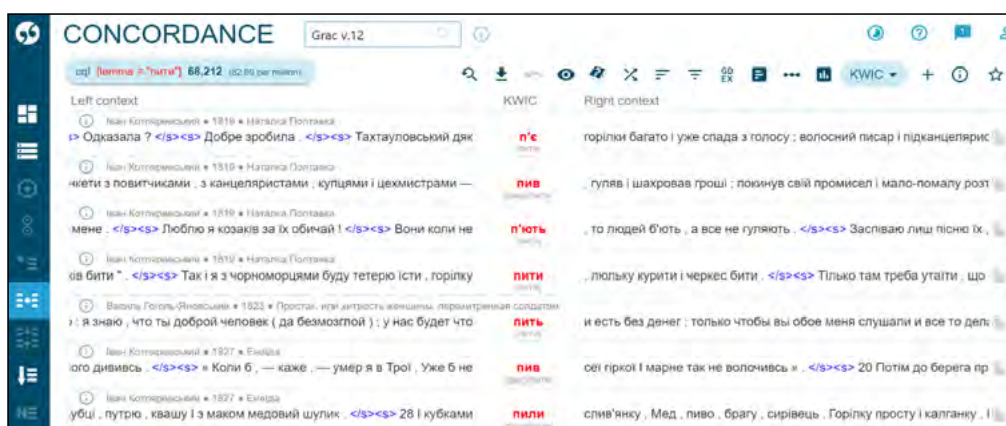


Figure 2. Examples of the use of the verb *пиво* in all its forms in the GRAC corpus

of a unit (*concordance*) and its frequency (i.e., the number of repetitions in the corpus). For a balanced corpus, the relative frequency $v_w = n_w / N$ of the appearance of a particular word w in it (i.e., the ratio of the frequency n_w to the total number of words N) corresponds to its relative frequency

in real language use; thus, a higher frequency of a unit w_1 than w_2 is an indicator that the unit w_1 is more commonly used in real language than the unit w_2 . Thus, the query [word = "пив"] finds in the GRAC corpus of Ukrainian 8,965 cases of the use of the word form *пив* (see Fig. 1). Given

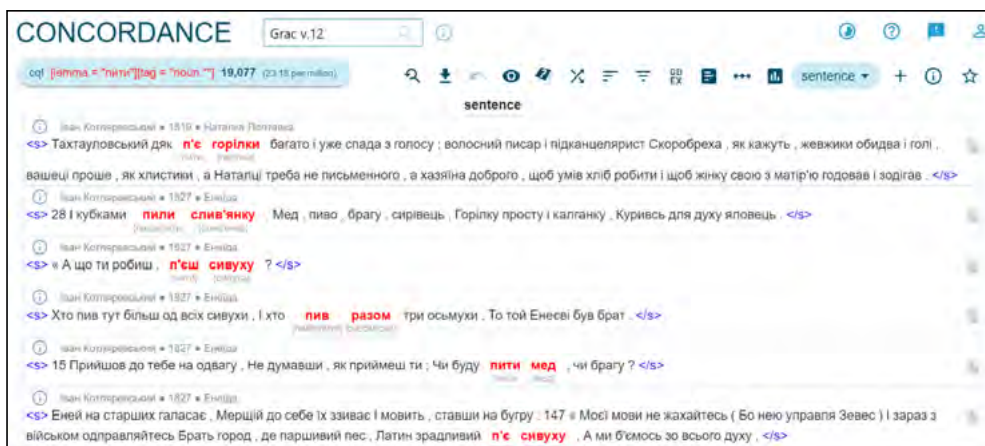


Figure 3. Search for combinations of nouns with the verb *пити* in the GRAC corpus

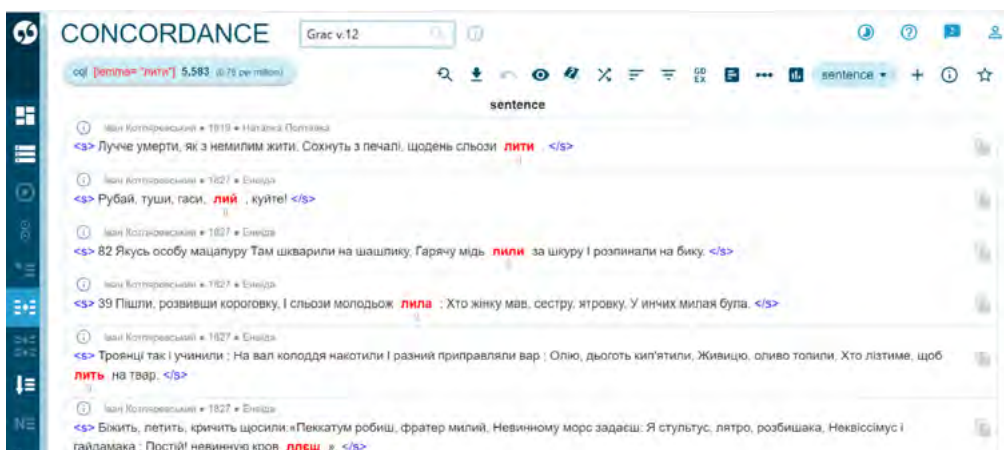


Figure 4. Examples of the use of the lemma *лити* in the GRAC corpus using the SketchEngine search server

Frequency CHANGE CRITERIA BACK TO CONCORDANCE

Show frequency per million

Word	Frequency	Frequency per million
лив	1,047	1.27
лліс	941	1.14
лити	739	0.90
Лліс	515	0.63
лліють	487	0.59
лили	431	0.52
лило	217	0.26
лий	193	0.23
лила	173	0.21
лїть	81	0.10
Ллий	79	0.10
лліш	72	0.09
ллю	55	0.07
лїйте	55	0.07
Ллив	53	0.06
Ллїльо	51	0.06
Лливо	49	0.06
лліте	46	0.06

Figure 5. Frequency of the use of different forms of the lemma *лити*

Lemma	Cooccurrences	Candidates	T-score	MI	LogDice
[лити]	178	3,363	13.34	12.93	9.35 ...
[дош]	1,089	61,183	32.99	11.36	9.06 ...
[відро]	289	13,120	16.99	11.66	8.98 ...
[лляти лити]	120	2,194	10.95	12.98	8.98 ...
[млині]	180	11,180	13.41	11.21	8.46 ...
[крокодилячий]	61	568	7.81	13.95	8.34 ...
[сльоза]	551	61,385	23.46	10.37	8.07 ...
[бруд]	93	9,520	9.64	10.49	7.66 ...
[безперестанку]	31	1,901	5.57	11.23	7.08 ...
[помил]	25	1,003	5.00	11.84	6.96 ...
[крокодил крокодилів]	20	107	4.47	14.75	6.85 ...
[цебро цебер]	21	698	4.58	12.11	6.78 ...
[надворі]	55	13,885	7.40	9.19	6.53 ...
[літ]	76	21,383	8.70	9.03	6.53 ...
[вода]	710	304,073	26.57	8.43	6.23 ...
[розтоплення]	15	1,310	3.87	10.72	6.16 ...

Figure 6. Collocations of the lemma *лито* at a distance of at most 4 tokens

query: ((lemma = "сльоза") [0,2]lemma="лити")|lemma="" 538 (0.65 per token)

sentence

- Іван Котляревський • 1819 • Наталка Полтавка
- <<> Лучче умерти , як з немилим жити , Сохнуть з печалі , щодень **сльози лита** </>
- Іван Котляревський • 1827 • Євста
- <<> 39 Пішли , розвивши корогову , і **сльози молодож лила** ; Хто жнку мав , сестру , ягровку , У инчих милая була . </>
- Владимир Винниченко • 1911 • Чортяк Пятра і Білий Володя .
- <<> О , Бідний , покинутий муж сидить , сумує , колихає дитинку й піркі **сльози лля** . </>
- Михайло Івченко • 1819 • Шурал воячки
- <<> А от як каже Проліс Ярину , то вона тужить , та **сльози лля** , та в тузі квітні шлює . </>
- Дмитрий Чубика • 1919 • Півний притупок
- <<> Дужче **сльози лйте** ! </>
- Клима Пондичук • 1920 • Червока мартов . Нариски й епіюдиони з часів революції
- <<> Опенка не **литиме сліз** , а шукатиме свого щастя на землі . Йі прекрасна ухмілка стане жорстокою і вимоги її серця будуть жадібними , бо краса люднини квадри жорстока ... </>
- Василь Король-Стерняк • 1929 • Іммерка

Figure 7. Collocations of the lemmas *лито* and *сльози* at a distance of at most 2 tokens

NARODOWY KORPUS JĘZYKA POLSKIEGO

Wyszukiwarka korpusowa PEŁCRA dla danych NKJP (Czytanie)

Wyniki: 100

lać**

Przeszukiwany zbiór zawiera 240.192.461 słów. Znaleziono 4,529 akapitów pasujących do zapytania w 0.257s. Bieżąca strona zawiera 101 przykładów z 29 różnych tekstów.

1.	Wszystko jedno na kim: Ciebie leją, to ty szukasz słabszego, żeby go	lać	-	Czarne okna
2.	-A jaka ma być? - powiedział Pędzina. - Chcesz ją	lać	, to teraz masz okazję.	Czarne okna
3.	Wiesz co? Strasznie chce mi się	lać	-	Kiler
4.	- Ależ mi się chce	lać	- wyszeptał Szarlej. - Nie zdoberzę...	Narrenturn
5.		lać	go! - rzywał Kozioł, uciepiony pleców woźnicy. - Lać go.	Mieć przeznaczenia
6.	- Lać go! - rzywał Kozioł, uciepiony pleców woźnicy.	lać	go, kumotry, gdzie popadnie i czym popadnie! Nie żalować! i kaedy doszli do Goicnica, byli zupełnie przemoczem.	Prowiek i inne sprawy...
7.	Zaczęło	lać	chwile później, Delikatne odczarki na Havelock były mczym w	Maria
8.	dwadzieścia rupii i rozwiął swój hamak w ogrodzie. Deszcz zaczyna	lać		
9.	spodobały się wujownicze wystąpienia pisarzy, którzy jeszcze niedawno	lać	lzy po Stalimie i pouczali Kisiela, co jest dobre dla socjalistycznej	Kisielewscy / Jan Au...
10.	gnali żydowscy klezmy, a żydowscy kielnerzy, niby szukanzo	lać	w chłopów wóde, ile tylko wiażło. Siedzieli tu znovu ci sami	Zycie ideologiczne, ...
11.	jak najdalej, bo zatrzymywali każdego, kto wyglądał na studenta, i	lać	każdego, przy kim znaleźli studencką legitymację. Na bezładnym placu	Zycie ideologiczne, ...
12.	rzeczy, lekkie, niepotrzebne i zapominane. Czas płynął też z nieba,	lać	się jak leniwa patoka, rozpryskiwał o melniczną łuskę bruka, lecz ani	Opowiesci galicyjsko...
13.	matymnym zbierniku. Opisał się na rękach i próbował ująć. Wodospad	lać	się po jego białej koszuli i niepiał go. Mock, analizując w myślach	Koniec świata w Bres...

Figure 8. The search result for the lemma *lać* in the NKJP corpus

the volume of the GRAC corpus, we obtain a frequency of 10.89 per million of tokens.

For languages rich in inflection, it is important to be able to identify different forms of a word with its basic form – *lemma* (for example, the word forms *пив, п'ють*, etc. with the verb *numu* in the infinitive form). The availability of appropriate marking makes it possible to take into account in one search query with a given lemma also all the uses of its various derived forms. Thus, the result of the query [lemma

= "ПИТИ"] in the GRAC will be 68,212 examples with the verb *numu* in all its forms (see Fig. 2).

A tag (POS, part of speech) mark makes it possible to obtain collocations of two words. For example, to identify all nouns that occur after the verb ПИТИ, you can use the search [lemma = "ПИТИ"] [tag = "noun.*"] (the result is 19,077 examples of the use in the GRAC (see Fig. 3). The corpus allows further grouping, frequency ordering and allocation of specific phrases.

#	Kolokacja	Pasujące współwystąpienia (kliknij na frekwencję, aby wyświetlić przykłady)	Ogolem	Chi ²
1.	zary	zar_leje /50/, zar_lal /28/, leje_zar /15/, lal_zar /15/, zar_lejacy /10/, lejacy_zar /5/, zar_lać /1/, lejacy_zar /1/, lejacy_zar /1/, lalo_zary /1/, lejacy_zary /1/, lejacy_zary /1/, lejacy_zary /1/	110	676,091.36
2.	zar	zar_leje /50/, zar_lal /28/, leje_zar /15/, lal_zar /15/, zar_lejacy /10/, zary_lejacego /5/, lejacy_zar /5/, zarem_lejacy /1/, lejacy_zary /1/, zarze_lejacy /1/, zar_lać /1/, lejacy_zar /1/, lalo_zary /1/, lejacy_zary /1/, zarem_leje /1/	119	522,269.16
3.	deszcz	lal_deszcz /50/, deszcz_lal /27/, deszcz_leje /45/, leje_deszcz /39/, lać_deszcz /3/, deszcz_leje /3/, leja_deszcz /2/, deszcz_lalo /2/, deszcz_laly /2/, deszcz_lala /2/, deszcz_leja /2/, deszcz_lalo /2/, leja_deszcz /2/, lejacy_deszcz /2/, deszcz_laly /2/, deszcz_lać /2/, lejacego_deszcz /2/, laly_deszcz /2/, deszczach_leje /1/, lejacy_deszcz /1/, leja_deszczem /1/, deszcze_lal /1/, deszcz_lejacy /1/, deszcz_lejac /1/, deszcz_lanie /1/, lać_deszcz /1/, lejacy_deszczem /1/, deszcz_lejacy /1/, deszcz_lejac /1/, deszczami_lejacy /1/, deszcze_leje /1/	240	301,618.6
4.	strumień	leje_strumieniami /45/, lal_strumieniami /22/, lala_strumieniami /10/, lalo_strumieniami /13/, laly_strumienie /3/, lalo_strumieniem /3/, lala_strumieniem /3/, strumieniem_lalo /3/, strumieniami_leje /2/, lać_strumieniami /2/, lejacy_strumieniami /2/, laly_strumieniami /2/, strumieniami_lala /2/, leja_strumieniami /2/, lejacego_strumieniami /1/, strumienie_laly /1/, strumieniami_lal /1/, lejacych_strumieni /1/, strumieniami_lejacego /1/, leje_strumieniem /1/, strumienia_lejacej /1/, lal_strumień /1/, strumienie_leja /1/, lejacy_strumieniami /1/, lal_strumieniem /1/, strumien_lal /1/, laly_strumieniami /1/, leja_strumieniami /1/	139	300,931.86
5.	wosk	lanie_wosku /25/, lania_wosku /25/, lanego_wosku /2/, lala_wosku /2/, laly_wosk /2/, lać_wosk /2/, laniem_wosku /2/, lejacy_wosk /2/, wosku_lanego /2/, lania_wosku /1/, wosku_lania /1/, lalimiy_wosku /1/, leja_wosk /1/, laniu_wosku /1/, wosk_lać /1/, lejany_wosk /1/, wosk_lanie /1/, leja_wosk /1/, lala_wosku /1/, woda_leja /15/, woda_lala /13/, lanie_wody /11/, leja_woda /28/, lania_wody /23/, lać_wody /20/, lala_woda /20/, lal_wody /13/, lali_wody /13/, leje_wody /13/, leja_wody /12/, lać_wody /11/, lać_woda /13/, wody_lejacej /2/, wody_leja /6/, laniu_wody /6/, leja_wody /2/, lejacy_wody /5/, lejacy_wody /5/, woda_leja /5/, lejany_woda /5/, laniem_wody /5/, lejacy_woda /5/, wodg_lać /5/, leje_wody /5/, lanie_woda /5/, lali_woda /5/, leje_woda /5/, wody_leje /5/, wody_laly /5/, wode_lal /5/, lei_wode /5/, woda_lejaca /5/, lalimiy_wode /2/	57	282,072.75

Figure 9. The search results for collocations of the lemma *lać* with nouns



Figure 10. The search for collocations of the lemma *drink* with nouns at a distance of up to 4 tokens

HELP		FREQ
1	FOOD	244
2	WATER	140
3	CORFEE	88
4	TEA	80
5	WINE	76
6	PEOPLE	62
7	ALCOHOL	57
8	MILK	50
9	BEER	42
10	LOT	31
11	CUP	29
12	DRINK	26
13	DAY	25
14	CHAMPAGNE	23
15	TIME	21
16	GLASS	20
17	WHISKY	20

Figure 11. The list of noun collocations of the lemma *drink* sorted by frequency

SEARCH	FREQUENCY	CONTEXT	OVERVIEW
51 GW5 W_fict_prose	1	and offered coffee. He said: 'No thank you, I don't drink coffee. But might I trouble you for a cup of hot water?'	
52 H85 W_fict_prose	1	was a sparse but efficient notetaker) to a quiet table by the window, drink some coffee and study what she had written down. If she understood it then	
53 H81 W_fict_prose	1	and she knew Dana was deliberately waiting until she left the flat before appearing to drink the coffee left ready for her. Dana didn't come to the showroom either	
54 HA5 W_fict_prose	1	I'll phone the airport now for you and see what's available while you drink your coffee. Do you have a preference between Gatwick and Heathrow?'	
55 HA6 W_fict_prose	1	use it? He smiled! What a sweet girl you are. Drink your coffee. Lissa, he ordered. 'You'll feel like a	
56 HGD W_fict_prose	1	moonlight. "How convenient," he laughed, sitting back down to drink his coffee. There is one flaw in the idea, though. Tonight	
57 HGF W_fict_prose	1	Lisa accepted coffee, pleased that it was afternoon when the English don't generally drink coffee. She was delighted that the coffee was real and that I had used	
58 HGF W_fict_prose	1	to the Polish cafes and Hungarian tea-rooms where she could talk with other emigres and drink Viennese coffee in fluted glasses. They argued! They talked! They	
59 HH0 W_fict_prose	1	to drink it. There is no sugar in here. I don't drink coffee without sugar. I am surprised and slightly reassured by this thread of	
60 HR9 W_fict_prose	1	to be very careful when prowling along the shelves. A First Spiritualist doesn't drink coffee, or eat white bread or cheese (apart from Gorgonzola – the Good	
61 HTR W_fict_prose	1	doubt if life is like that. Now, just give me five minutes to drink some coffee and I'll ring Geoffrey. Loretta was rinsing the cups when:	
62 J13 W_fict_prose	1	But I stop. Suzie waits. Her green eyes watching me. I drink some coffee. Albie never had a car. Ludo says.'	
63 JY0 W_fict_prose	1	will always be my only mother. And Dad, of course. Drink your coffee, or you won't be ready when David comes.'	
64 JY3 W_fict_prose	1	the only way to keep the shareholders happy! Guy Sterne leaned forward to drink some more coffee, his wide mouth suddenly twisting humorously. At least Sar	
65 JY8 W_fict_prose	1	that gauzy dress, but it was too late, hard seen. Drink your coffee. Damn it, I don't want any coffee!	
66 KBV W_fict_prose	1	poor little bezzars." I'm not likely to forget. Gabby. Drink up your coffee," said Charley calmly. 'I know where I rate	

Figure 12. The list of collocations of the lemmas *drink* and *coffee*

A. Stefanowitsch analyses the case (Stefanowitsch, 2020, p. 47–48) when the study of the use of the adjective *implacable* in the 450 million corpus of contemporary American COCA proves its appearance mainly in a specific context, which implies rivalry between people or even their enmity. This fact is not reflected in the illustrative materials of the Merriam-Webster dictionary, which leads to misinterpretation of the connotations of this adjective.

The greatest opportunities for researchers are offered by corpora with flexible queries using the CQL (Corpus Query Language). For example, you can search for all potential objects of the verb *ламату* in the GRAC corpus using the query `([tag = «noun.*»][{0,3}[lemma= «ламати»]) ([lemma = «ламати»][{0,3}[tag = «noun.*»]) within <s/>`. The result is a list of uses of the verb *ламату* in all its forms (lemma= «ламати») together with a noun ([tag = «noun.*»]) in one sentence (within <s/>), between which there may be up to three other words ([{0,3}). For example: В світлому отворі дверей миготіли постаті, і деякі, переступаючи в захваті поріг, ламали тут священний затишок різким човганням взуття.

The detailed characteristics of the analysed text corpora are presented below.

The GRAC is the largest corpus of the Ukrainian language, publicly available and convenient for conducting both general and specialized linguistic studies. The GRAC contains rich grammatical and semantic marking, is integrated into the SketchEngine [https://parasol.vmguest.uni-jena.de/grac_crystal/#dashboard?corpname=grac12] search system, which allows advanced searching with full use of marking using the CQL.

The total volume of version 12 of the corpus published in 2021 is 822,959,896 tokens, 640,932,211 words (7,591,461 unique words total), among which there are 2,693,775 lemmas; the

GRAC contains 50,803,520 sentences collected from 97,245 documents and is the largest digitized base of the Ukrainian language.

The SketchEngine search server allows searching for the number of uses of a word form or lemma in the corpus, as well as finding their collocations (common uses) with other units. For example, the search for the lemma *лumu* gives 5,583 occurrences (see Fig. 4). If necessary, you can view the frequency of different forms of the basic lemma (see Fig. 5).

For further analysis, the server has the function of finding collocations, grouped by lemmas or other grammatical characteristics and ordered by frequency (see Fig. 6).

To identify special common uses, the SketchEngine allows complex searches using the CQL; for example, the query `([lemma = «сльоза»][{0,2}[lemma=«лити»]) ([lemma=«лити»][{0,2}[lemma=«сльоза»]) within <s/>` starts a search for collocations of the verb *лumu* and the noun *сльоза* in their various forms; both units are present in the same sentence (within <s/>), and between them there may be up to two other words (see Fig. 7).

The NKJP is the corpus of Polish created in cooperation with the Institute of Fundamentals of Informatics of the PAS (the Polish Academy of Sciences), the Institute of the Polish Language of the PAS, the PWN Scientific Publishing House of the PAS, and the University of Lodz with the support of the Ministry of Science and Education of Poland.

To work with the corpus, the PELCRA search engine was used, which shows the frequency of uses of a certain word form or lemma with subsequent analysis of collocations. The NKJP contains a balanced part (about 300 million words) as well as a complete part (about 1.5 billion words). For example, the search result for the lemma *lać* is presented in Fig. 8.

The PELCRA system makes it possible to search for collocations (for example, Fig. 9 shows the search results for all the collocations of the lemma *lać* with nouns immediately before or after it).

The BNC is the result of the work of a specially formed consortium (which united the Universities of Oxford and Lancaster, several commercial publishing houses, and the British Library); contains almost 100 million words of British English of the late twentieth century. The current version of the corpus was completed in 1994 (the corpus is non-dynamic, like the NKJP); 90% of all textual information are written sources of various genres and 10% are conversational sources. The BNC was created specifically for linguistic studies and, despite its relatively small size, the corpus is well balanced, and therefore the results of searches in it are representative, i.e. give an undistorted picture of the functioning of contemporary British English (according to C. Fellbaum, experts consider the BNC to be "a reliable source for English researchers" (Fellbaum, 2019, p. 749).

The publicly available program allows making basic searches for exact matches, all forms of a given lemma, collocations, etc. The corpus has a

partial morphological marking, but shows a lack of semantic one.

For example, a basic search for the lemma *drink* as a verb can be done by entering the query *drink_v** in the field; the search result will be general information about the frequency of this lemma; in particular, it was found 2,981 times. The marking of parts of speech makes it possible to search for all collocations of specialized reference verbs with potential noun objects at a distance of at most 4 tokens (see Fig. 10). The search result will be the list of all relevant nouns (see Fig. 11); the examples of the use can be viewed by highlighting the desired one (see Fig. 12).

Conclusion. Linguistic studies using text corpora allow automatic search and analysis of language units and their combinations. The text corpora of Ukrainian, Polish, and English (GRAC, NKJP, BNC) meet the requirements of authenticity, representativeness (balance), and sufficiency in volume. The introduction of research text corpora into linguistic use makes it possible to optimize and objectify the analysis of language material and provides for highly reliable results.

ABBREVIATIONS

- BNC – British National Corpus
GRAC – The General Regionally Annotated Corpus of Ukrainian
NKJP – Narodowy Korpus Języka Polskiego

BIBLIOGRAPHY

1. Бук С. Н. Велика проза Івана Франка: електронний корпус, частотні словники та інші міждисциплінарні контексти. Львів: ЛНУ імені Івана Франка, 2021. 424 с.
2. Демська О. Текстовий корпус: ідея іншої форми. Київ: ВПЦ НАУКМА, 2011. 282 с.
3. Широков В., Бугаков О., Грязнухіна Т. та ін. Корпусна лінгвістика. Київ: Довіра, 2005. 471 с.
4. Fellbaum Ch. How flexible are idioms? A corpus-based study// *Linguistics*. 2019. Vol. 57. № 4. P. 735–767.
5. Fellbaum Ch. The treatment of multi-word units// *Oxford Handbook of Lexicography*/ ed. by P. Durkin. Oxford: Oxford U-ty Press, 2015. P. 411–425.
6. Gries S.Th. Corpus-based methods and cognitive semantics: The many senses of to run // *Corpora in Cognitive Linguistics: Corpus-based Approaches to Syntax & Lexis*/ ed. by S.Th. Gries & A. Stefanowitsch. Berlin; New York: Mouton de Gruyter, 2006. P. 57–100.
7. Shvedova M. The General Regionally Annotated Corpus of Ukrainian (GRAC, uacorp.org): Architecture and Functionality // *Proceedings of the 4th International Conference on Computational Linguistics and Intelligent Systems (COLINS 2020)*. 2020. Vol. I. P. 489–506.
8. Stefanowitsch A. *Corpus Linguistics: A Guide to the Methodology*. Berlin: Language Science Press, 2020. 508 p.

REFERENCES

1. Buk S. N. (2021). Velyka proza Ivana Franka: elektronnyi korpus, chastotni slovnyky ta inshi mizhdystyplinarni konteksty [Ivan Franko's great prose: electronic corpus, frequency dictionaries, and other interdisciplinary contexts]. Lviv: LNU imeni Ivana Franka [in Ukrainian].
2. Demska O. (2011). *Tekstovyi korpus: ideia inshoi formy* [Text corpus: the idea of another form]. Kyiv: VPTs NaUKMA [in Ukrainian].
3. Fellbaum Ch. (2019). How flexible are idioms? A corpus-based study. *Linguistics*, 57 (4), 735–767.
4. Fellbaum Ch. (2015). The treatment of multi-word units. In P. Durkin (Ed.), *Oxford Handbook of Lexicography* (pp. 411–425). Oxford: Oxford University Press.

5. Gries S.Th. (2006). Corpus-based methods and cognitive semantics: The many senses of to run // In S.Th. Gries & A. Stefanowitsch (Eds.), *Corpora in Cognitive Linguistics: Corpus-based Approaches to Syntax & Lexis* (pp. 57–100). Berlin; New York: Mouton de Gruyter.
6. Shvedova M. (2020). The General Regionally Annotated Corpus of Ukrainian (GRAC, uacorporus.org): Architecture and Functionality. *Proceedings of the 4th International Conference on Computational Linguistics and Intelligent Systems (COLINS 2020)*, I, 489–506.
7. Shyrokov V., Buhakov O., Hriaznukhina T. a.o. (2005). *Korpusna linhvistyka [Corpus linguistics]*. Kyiv: Dovira [in Ukrainian].
8. Stefanowitsch A. (2020). *Corpus Linguistics: A Guide to the Methodology*. Berlin: Language Science Press.

КОРПУСНИЙ ПІДХІД У ДОСЛІДЖЕННІ ДІЄСЛІВНИХ ПРЕДИКАТІВ

Назарчук Роксолана Зіновіївна

*кандидат філологічних наук,
старший викладач кафедри прикладної лінгвістики
Національного університету «Львівська політехніка»
вул. С. Бандери, 12, Львів, Україна*

Карамишева Ірина Дамірівна

*кандидат філологічних наук,
доцент кафедри прикладної лінгвістики
Національного університету «Львівська політехніка»
вул. С. Бандери, 12, Львів, Україна*

У статті висвітлено особливості корпусного підходу до аналізу дієслівних предикатів української, польської, англійської мов. Окреслено проблеми корпусної лінгвістики, здійснено огляд корпусів текстів, а саме Генерального регіонально анотованого корпусу української мови (GRAC), Національного корпусу польської мови (NKJP) і Британського національного корпусу (BNC). Згадані ресурси, відповідаючи вимогам автентичності, репрезентативності (збалансованості), достатності за обсягом (С. Бук, О. Демська, К. Фельбавм, А. Стефановіч, В. Широков, М. Шведова та ін.), дають змогу оптимізувати й об'єктивізувати тлумачення мовного матеріалу, одержати високу достовірність результатів. Наукова новизна дослідження полягає у комплексному зіставному аналізі функційних можливостей корпусів текстів української, польської, англійської мов з метою визначення їхніх особливостей, а також окреслення перспектив застосування корпуснобазованого підходу для вияву об'єктних зв'язків дієслівних предикатів. Уведено основні поняття корпусної лінгвістики: лема, словоформа, токен, мітка (тег), співживання (колокація), конкорданс, частотність, запит CQL; їх проілюстровано реалізаціями відповідних референтно спеціалізованих дієслівних предикатів і їхніх об'єктів на базі GRAC, NKJP і BNC. Зокрема, показано, що найпростіше використання корпусів текстів полягає в засвідченні вживання тої чи іншої одиниці (concordance) та її частотності (тобто кількості повторень у корпусі). Для мов, багатих на словозміну, важлива можливість ідентифікації різних форм одного слова з його базовою формою – лемою. Проілюстровано, як досліджені корпуси уможливають подальше групування, упорядкування за частотністю та виокремлення конкретних словосполук, а розмітка частин мови полегшує пошук усіх колокацій дієслівних предикатів з потенційними об'єктами-іменниками на відстані 4 токенів. Описано функцію пошуку співживань дієслів, позрупованих за лемами чи граматичними характеристиками та впорядкованих за частотністю. Вказано на переваги семантичної розмітки у корпусі GRAC для автоматичного пошуку одиниць мови та їхніх поєднань. Усі описані особливості корпусів текстів супроводжено ілюстраціями з GRAC, NKJP і BNC.

Ключові слова: корпусна лінгвістика, корпус, GRAC, NKJP, BNC, дієслово, об'єкт, лема, колокація