

Beyond Philology No. 17/3, 2020
ISSN 1732-1220, eISSN 2451-1498

<https://doi.org/10.26881/bp.2020.3.02>

Corpus analysis in applied linguistics: Selected aspects

JOANNA REDZIMSKA

*Received 27.04.2020,
received in revised form 16.09.2020,
accepted 1.10.2020.*

Abstract

Recently, teaching and learning processes have been significantly influenced by modern technologies. Thus, the teacher's position as the only authority in the classroom has been changed into playing the role of a guide or a facilitator who should possess the knowledge and skills to use modern technologies and to freely access data. This change is particularly visible in the field of teaching and learning languages with the application of various educational platforms and software. Since this situation has been widely discussed since the 1990s, for the sake of this article only selected aspects have been taken into account. The major focus of the present article is to present language corpus analysis as a method of activating teachers and students as participants in the Data-Driven Learning (DDL) process.

Keywords

corpus analysis, DDL, activation

Analiza korpusowa w językoznawstwie stosowanym: wybrane aspekty

Abstrakt

Rozwój technologii w znaczny sposób wpłynął na proces nauczania i uczenia się języka obcego. W konsekwencji, nauczyciel zmienił swoją pozycję z jedyne go autorytetu w klasie na rolę przewodnika oraz moderato ra, który powinien posiadać wiedzę i umiejętności pozwalające na wykorzystanie technologii i ogólnie dostępnych danych językowych. Widać to szczególnie w dziedzinie nauczania języków obcych, gdzie wykorzystywane są platformy i komputerowe programy edukacyjne. W związku z faktem, iż wpływ technologii na proces kształcenia opisywany jest w literaturze przedmiotu już od roku 1990, niniejszy artykuł omawia jedynie wybrane aspekty z tego zakresu. Główna uwaga poświęcona jest zagadnieniu analizy korpusowej jako metody aktywizacji nauczycieli i uczniów/ studentów poprzez proces uczenia się opartego na danych (Data-Driven Learning).

Słowa kluczowe

analiza korpusowa, DDL, aktywizacja

1. Introduction

The development of technology and the first computers paved the way for changes in all fields of research, including teaching and learning foreign languages. Thus, the traditional methods of introducing knowledge to students as well as the practice of various skills embraced the possibility of methods connected with computers, virtual reality, and free, easy language resources available for public use.

A language resource that is of core interest to this work is represented by the language corpus and teaching/learning method that is Data-Driven Learning (DDL). One of the most obvious applications of a language corpus is that it can function

as a source of knowledge about the target language's forms, use or statistics. Thus, in this respect language corpora constitute an alternative to a dictionary where the focus is mostly on meaning and possible examples where the form is used. One should also bear in mind that a language corpus as a whole always has a digital form, compared to dictionaries that traditionally have a printed form which is subsequently accompanied by a digital form. Yet, the aim of this work is to present how corpus analysis enhances language teaching and learning by offering methods and data that are not available elsewhere. However, bearing in mind the pace of the development of corpus linguistics as well as the abundance of publications connected with this field, for the sake of this article only selected aspects and corpora are further discussed. Thus, the following parts introduce a number of suggestions related to the practical application of language corpora and analysis on the basis of selected corpora for English and Polish.

2. Corpus linguistics

Although corpus linguistics has gained its position relatively recently, the origins of corpus linguistics, yet in a form different from the contemporary one, may be traced back to the 13th century (O'Keefe and McCarthy 2010). As O'Keefe and McCarthy point out, the need for preparing wordlists and the creation of concordances were methods of Bible exegesis where scholars (mostly monks) and their students indexed the Bible hoping to find divine authorship. Another example mentioned by O'Keefe and McCarthy with reference to religious texts is the work by Anthony of Padua who first listed concordances in the Vulgate Bible. Further developments in the methods of indexing texts for wordlists and concordances were expanded on other kinds of texts, for example Shakespeare's works were annotated for concordances until the late 18th century (O'Keefe and McCarthy 2010).

However, it is the 20th century with the advent of computers that brought about the most significant breakthrough in the corpus approach to language. The first attempts to create a machine-readable language corpus were made in the 1960s by Francis and Kučera (the Brown Corpus). Yet, with the generative approach to language at that time, their effort met with a significant amount of criticism. Generative grammar emphasizes the importance of a speaker's intuition and it concentrates on an explanatory adequacy, looking for universal language paradigms and principles. Corpus linguistics, by contrast, focuses on descriptive adequacy and examines the well-formedness and grammaticality of sentences (Meyer 2002). At the end of the 20th century, corpus linguistics gained its position and significance as a field of study and it has been acquiring greater importance ever since.

As far as the applicability of corpus linguistics is concerned, McEnery and Wilson (2011) highlight that corpus linguistics is a useful tool for identifying and characterising particular aspects of language use as well as researching these aspects from a linguistic perspective. Further the two authors (McEnery and Wilson) point out that multiple areas of linguistics derive from corpus linguistics, yet each area requires different methodology to analyse language, which has its consequence in the distinction between corpus-based and non-corpus based studies. Since corpus linguistics accounts for the complexity of language as a communicative tool with the application of interfering data (a corpus-based analysis), it stands in opposition to the generative approach whose major task is to study context-independent and most of all universal rules of language (non-corpus based studies) (Meyer 2002).

Consequently, the above-mentioned aspects raise the question of the reasons for creating different kinds of corpora. According to Renouf (2007), the three main arguments for the creation of corpora centre around the issue of science (the scientific drive for the observation and the analysis of data to test various scientific hypotheses), a pragmatic need (defined in practical

categories of the availability of data, funding and formal and technological solutions that are required for such research) and 'a fluke' (understood as an opportunity to start a new initiative that meets certain research or market demands). Moreover, Renouf (2007) mentions that the above factors highly influence both the size and the possible applications of a corpus with the tendency for small and specialised corpora, e.g. Freiburg-LOB Corpus of British English (FLOB) or the Freiburg-Brown Corpus of American English (FROWN) to compare relatively modern corpora with earlier corpora.

Thus, the application of language corpora is the most significant aspect motivated by the need for the investigation of language use in context, where the research data that is collected from a vast array of language users is the greatest benefit to the analysis (Meyer 2002). The usability of a given corpus is partially defined by its size as Meyer (2002) states that large corpora are particularly necessary for inferring details connected with grammatical constructions, forms, frequency, context or communicative power, whereas smaller corpora also possess scientific potential as long as they contain a collection of particular constructions. Undoubtedly, these are lexicographers that benefit from the use of corpus analyses by inferring information about lexical units, their range, morphological realisations and possible meanings; additionally, most of the lexicographic analysis is a largely automatic process (performed by means of computer programmes that provide data such as frequencies of words, lemmas, key words in context, tagged parts of speech) (Meyer 2002). Furthermore, the above method, as Mayer (2002) claims, is also widely applied to studying meanings and the actual uses of words which, without a corpus, are difficult to identify.

Additionally, language corpora are a way of registering language variations of different kinds, such as sociolinguistic characteristics (gender, age, ethnicity) that are represented in metadata. Following Meyer (2002), there is a choice of software

that can be used for the above purpose, an example of which is SARA (available at natcorp.ox.ac.uk/archive/SARA/index.xml).

Historical linguistics can also profit from corpus linguistics and corpus analysis. Two examples of this kind of corpora are the LOB and FLOB corpora (two parallel synchronic corpora) where one can compare language changes as well as variation in grammar and lexis (Renouf 2007). However, as Renouf (2007) points out, diachronic corpora are very often based on chronologically ordered texts or corpora that offer a selection of consequent texts (RDULES unit of the AVIATOR project available at rdues.bcu.ac.uk/aviator.shtml), which allows for the analysis of productive and creative aspects of language, collocation changes as well as word sense or meaning modifications.

Still other fields like translation studies or contrastive analysis develop due to the use of parallel corpora which (according to Meyer 2002) provide information about syntax, morphology or pragmatic aspects of translated text that can be further contrasted and compared. Parallel corpora, based on bilingual dictionaries created for this purpose, can be used for training translators and although it is a demanding task, there is software like Paraconc (paraconc.com) that facilitates the above mentioned procedures (Meyer 2002).

2.1. Examples of corpora

Corpus linguistics has gained its popularity recently, which has as its consequence the fact that a growing number of scholars and businesses are interested in projects which allow for the creation of corpora and making such corpora publicly available. As Lee (2010) points out these are not only English language corpora that are commonly used for corpus analysis but also public corpora for other languages which find their application in language study and research. The access to corpora is offered by distribution agencies and archives sites, with International Computer Archive of Modern and Medieval English (ICAME) (icame.uib.no), Linguistic Data Consortium (LDC) (ldc.upenn).

edu), CLARIN-PL (Common Language Resources and Technology Infrastructure available at clarin-pl.eu/) for Polish, and the Oxford Text Archive (OTA) (ota.bodleian.ox.ac.uk/repository/xmlui) to name a few, but as Lee (2010) highlights, access may be restricted due to the copyright or funding of these corpora.

Additionally, it must be underlined that, as far as parallel corpora are concerned, these are bidirectional and offer information about source texts as well as their translations to facilitate comparison between languages (Lee 2010). One such project that allows for the creation of lexicons and also monolingual corpora in 14 languages is The Preparatory Action for Linguistic Resources Organisation for Language Engineering (PAROLE). It offers standards and specifications for cross-linguistic analysis (Lee 2010). As far as strictly bidirectional parallel corpora are concerned, Lee mentions the English–Norwegian Parallel Corpus (ENPC) and the English–Swedish Parallel Corpus (ESPC).

An interesting example of corpora are those that include multimodal information, including speech transcripts connected with original audio or video recordings. Following Lee (2010), this allows for research into such aspects as prosody, gestures, and situated discourse to name only a few. The Scottish Corpus of Texts and Speech (SCOTS) is often quoted as an example model of this kind of corpora with its 4 million written and spoken texts (Lee 2010) as is SPOKES (<http://spokes.clarin-pl.eu/>) which currently contains 247,580 utterances (2,319,291 words) in transcriptions of spontaneous conversations (Pezik 2015).

Additionally, another useful solution for gathering necessary linguistic data is offered by the almighty power of the Internet. Thus, the Web can be treated as a corpus that allows one to find relevant data. This corpus, as Lee (2010) points out, is either dynamic or static including information connected with one particular moment of use or information that is constantly updated for new language sources. Examples of this application of the Internet include Web concordancers (e.g. WebCorp, Web-KWiC, KWiCFinder) to make research into concordance, the

Linguist's Search Engine which can be used to examine syntactic structures on the basis of parsed trees and the static web corpus ukWaC where two billion English words are lemmatized and tagged for parts of speech (Lee 2010).

2.2. Learner corpora

Since the major focus of the present work is on the relationship between language corpora, corpus analyses and their possible applications in language teaching and learning, it must be emphasized that these pedagogical implications resulted in the appearance of non-native speaker corpora (including written and spoken learner language). The corpus released in 2002 by Granger, Dangneaux and Meunier serves as an illustration of this pedagogical trend. In Tribble (1997) or Aston (2002) one can read about corpora created by students which centre around either genres or topics of particular interest to the group of students. Further, Braun (2005) developed a corpus of spoken English – ELISA – on the basis of a collection of interviews. Following Widdowson (1991, 2003), ELISA incorporates the principle of pedagogical mediation and the entire corpus is consistent, as far as pedagogical conceptualization is concerned, with respect to annotation, enrichment and search procedures. Thus, it promotes authentic data for learners since it uses both a great deal of decontextualized textual data as well as context-dependent interaction data (Widdowson 2003). It is worth noting that the European Minerva project SACODEYL (2005-08) (Braun 2010, Hoffstaedter and Kohn 2009, Pérez-Paredes and Alcaraz-Calero 2009, Pérez-Paredes 2010, Widmann, Kohn and Ziai 2010) also uses ELISA's pedagogical approach to a great extent including the design and corpus tools.

However, there are corpora dedicated to students who learn foreign languages. An example of such corpus is the Longman's Learner Corpus based on data gained from ESL students. Later, as Meyer (2002) points out, this corpus was used to write a dictionary which included suggestions for students' common

mistakes and strategies on how to counteract them. This information is also useful for teachers. Also, Lee (2010) references the International Corpus of Learner English (ICLE) created on the basis of students' argumentative essays illustrating different English language backgrounds.

Two further interesting examples of learner corpora are the CHILDES database and the Polytechnic of Wales (POW) Corpus (Lee 2010). These are resources that focus on data from children acquiring their native language. These resources are known as developmental corpora and they can assist in research into the way language forms are developed during the process of learning a first language (Lee 2010).

Obviously, this referential function as far as language is concerned is also fulfilled by traditional reference grammars that offer advice on how to form grammatical constructions in accordance with the rules of language (largely a prescriptive approach). An example of this is the corpus-based research of Quirk, Greenbaum, Leech, and Svartvik, which was published in 1972 (Meyer 2002). These scholars were pioneers in using corpora of written and spoken language to explain grammatical constructions.

3. Data-Driven Learning (DDL)

Data-Driven Learning (DDL) seems to be the best solution for the development of metalinguistic knowledge and learner autonomy since this method applies authentic language materials "to empower both teachers and students to develop competences in moving away from mere surface features of a text to selecting and understanding meanings and structures" (Corino and Onesti 2019: 1). One of the first advocates of this method was Johns (1991) who compared every student to Sherlock Holmes discovering the intricacies and mysteries of a language. Similarly, Sinclair (2004) praises corpus-based teaching for the use of authentic language materials. Moreover, Cobb and Boulton (2015) highlight that what is most valuable to the method is

the substantial exposure to authentic language input in a controlled way. Furthermore, among the merits of DDL, Boulton (2016: 3) emphasizes the exploitation of the following elements/aspects: authenticity, autonomy, cognitive depth, consciousness raising, constructivism, context, critical thinking, discovery learning, individualization, induction, learning-to-learn, life-long learning, (meta)cognition, motivation, noticing, sensitization and transferability. However, it must be acknowledged that using DDL as an effective method requires time, practice, computer skills and most of all it must find favour with the students (especially those who do not feel comfortable with technological devices). Also, as Meunier (2011) points out, DDL necessitates considerable user investment in time and practice in order to use the data efficiently. As a result the role of a teacher changes from that of a sole authority possessing necessary knowledge to that of “a consultant, guide, coach and/or facilitator” (Suan Chong 2016). As far as students are concerned, whenever they attempt to solve language problems, they activate HOTS (higher order thinking skills), which will result in long-term knowledge retention and improved language skills (Corino and Onesti 2019: 2). Thus DDL, being a hands-on approach, provides opportunities for both teachers and students in indirect and direct applications of corpora in teaching and learning languages.

4. Discussion

As has been discussed above, there have been various types of corpora and different reasons for their creation. Without any doubt, language corpora are valuable language resources with multiple applications and the potential to fulfil different functions. However, the aim of this work is to see if corpus analysis (or working with corpora) can influence the teachers’ work and facilitate or enhance the process of learning. Thus, the assumption that is made for the sake of this article is that corpus analysis is a method of activating teachers and students. As follows,

the further discussion focuses on selected aspects connected with possible practical uses of corpus analysis in the teaching/learning process.

The first and foremost aspect of corpus analysis concerns the idea of the corpus as a source of knowledge about language itself. As a result, corpus analysis allows teachers and students to have access to various kinds of language data, depending on the corpus. Some of these corpora are open-source big-data resources, for example, for English the COCA – Corpus of Contemporary American English. If a given corpus is a current project, it is updated with actual uses of language, which makes it a more reliable and applicable resource.

4.1. Teachers

Without any doubt, the most obvious, and at the same time the most significant, function of a language corpus is that it provides knowledge about a language. As has been already mentioned, the purpose of the corpus dictates what kinds of texts are used to build it and, consequentially, what kind of language forms are to be expected.

The job of teachers constantly involves various kinds of interaction with their students. Beginning with lectures and classes through to meetings with their parents, this formal, and at the same time special, relationship always relies on cooperation. There are also physical representations of this cooperation in the form of tests, essays or exercises with a twofold role: on the one hand, they are proof of the students' level of knowledge and competences and, on the other hand, at the same time they provide evidence of mistakes and issues that have to be improved. Such evidence can be collected in a form of a corpus where only language data is gathered (without any personal detail). This collection can be further used to prepare additional teaching materials to revise the problematic issues. Additionally, the frequency and quantity of certain mistakes can prove the need for further reconsideration and revision of teaching syllabuses or

even software so they will be better suited to the real needs of the students.

Another issue connected with corpus analysis is inevitably related to the question of developing a teacher's competences and activating the process of teaching and learning. Some teachers meet the challenge of building their own corpus. In a practical sense this means first learning about the programmes and tools that can be helpful in creating such corpus (developing their computer skills, learning the new software necessary to build a corpus) and then collecting texts that provide language data for the corpus (developing research skills). However, teachers who do not want to build their own corpora can use resources which are already available and look for the necessary data in them (developing analytical skills). Yet, it must be also pointed out that the most demanding task for teachers is still to give focused directions to their classes and to guide their students through data discovery and interpretation since language corpora only provide language data without any analysis. Thus, the major responsibility of teachers (and later students) is to evaluate the information found.

As follows, creating such a corpus and later analysing it seems to be a way to activate teachers, because one of the main adversaries of every teacher is routine. To avoid routine, teachers attend various courses and trainings to raise their qualifications or to look for some alternative solutions for making their lessons or courses more interesting and inspiring to their students. This results in a situation where creating and analysing their own corpora is an additional instrument which allows teachers to break up the school routine and makes their job more attractive.

4.2. Students

Corpus analysis can be profitable for students as well. Introducing corpora as an alternative to dictionaries not only broadens learners' knowledge about possible language resources but also

offers a new, technology-oriented method of learning a language. Introducing learner corpora as educational projects is a worthwhile strategy since students are more motivated to work on language that comes from their own fields of interest. The benefit here is twofold: on the one hand, the student develops his or her language skills, and on the other, the student broadens his or her knowledge about a particular domain.

Furthermore, working with corpora and carrying out a corpus analysis is focused on two major tasks. The first is centred around the creation of a corpus by students. Such a corpus can include various kinds of texts, depending on its aim. To illustrate this idea, students could build a corpus of their own mistakes and another, referential corpus that represents the correct forms. Such corpora that function as reference resources will then include either their own texts with mistakes (genuine language productions) or texts which they collect from formal/standard resources. This is particularly useful for all kinds of revision and language drills that students can do on their own. An additional value from the perspective of a student is the fact that preparing and working with one's own corpus makes the whole process of learning highly personalized and autonomous and in consequence it allows for a significant amount of learning creativity and learning liberty.

Moreover, students can benefit from the corpus analysis by using and examining prior existing corpora to find information and solutions to their particular language problems or to find applications of selected language forms. To illustrate the above issue one can refer to a case study where a student wants to consult a corpus (which then works as an outer standard language model) to learn and understand the differences in distribution and meaning between words of nearly the same meaning. This probably is a matter of intuition for native speakers but for learners of a foreign language, it may cause problems. The examples below focus on two English words *average* and *medium* in their adjectival and nominal functions and their Polish equivalent(s) since, as far as Polish is concerned, the form in an

adjectival function is the same for both *average* and *medium*. The following examples are from COCA (www.english-corpora.org/coca), NOW Corpus (www.english-corpora.org/nw) and PARARELA (<http://paralela.clarin-pl.eu>) and were retrieved between July and September 2020. Only two kinds of information from the corpora are being further scrutinized, namely the frequency (revealing the quantitative information) and the context (presenting qualitative information), since in the opinion of the present author these are the best and most accessible ways to show the differences between the two concepts in question.

4.3. COCA

The screenshot displays the COCA website interface. At the top, there are navigation tabs: SEARCH, FREQUENCY, and CONTEXT. Below the tabs, there are search options: ON CLICK: CONTEXT, TRANSLATE (?), GOOGLE, IMAGE, PRON/VIDEO, and BOOK (HELP). The search results show a frequency of 33635 for the word 'MEDIUM'. Below this, there is a table of search results with columns for year, source, and context. The table lists 20 results, each with a unique identifier, year, source, and a snippet of text containing the word 'medium'.

Year	Source	Context
1998	MAG CountryLiving	A B C teaspoon ground ginger 1/4 teaspoon salt 1 In heavy 9-inch skillet, heat oil over medium heat. With fork, pierce sausages all over several tim
2002	MAG CountryLiving	A B C , 1/4 cup water, and the spice bundle in a small saucepan set over medium heat. Cook, stirring occasionally, until the sugar melts and the liq
2015	MAG VegTimes	A B C # Preheat oven to 425F. Toss together tomatillos, oil, and oregano in medium bowl. Season with salt and pepper, if desired. Set aside. Sanda
2003	MAG AmerArtist	A B C be done with mineral pigments. His specific palette includes cadmium lemon, cadmium yellow medium , cadmium red light, permanent aliza
1994	FIC MassachRev	A B C the " mobster house." # Danny's father was a stocky man of medium height, with a rough voice deepened by decades of smoking. He spoke
1998	MAG CountryLiving	A B C overnight. 3 Prepare White Sauce: In heavy 2quart saucepan, melt butter over medium heat. Stir in flour, salt, and nutmeg; cook until bubbly
2007	ACAD TeachLibrar	A B C carried from the old world. Rigor and information fluency matter, no matter the medium -- so do excitement, engagement, and enthusiasm.
2007	SPOK NBC_Today	A B C meat rare. You say that's bad. MI-LEMPERT: Rare is bad, medium rare is bad. ROKER: Really? MI-LEMPERT: You don't want to
2013	MAG Prevention	A B C 30 to 40 minutes. Meanwhile, put 1 Tbsp oil in large skillet over medium heat. Add 2 tsp minced garlic and stir 1 minute, then add 1
2005	ACAD Education	A B C few caveats are in order. The presentation can be given to a small, medium , or very large group of parents. There are pros and cons for each
2009	ACAD CommCollegeR	A B C 14) and as small, less than 2,500 (n = 13); medium , 2,500 to 7,500 (n = 15); and large, more than
2006	MAG SouthernLiv	A B C Times have changed. Now moth orchids can be purchased growing in an orchid bark medium , sphagnum moss, a combination of these, or s
2012	MAG VegTimes	A B C in lime juice and zest. # 2 I Heat sunfloweroil in large skillet over medium heat. Cook eggplants 6 to 8 minutes in batches, or until golden bro
1996	ACAD AfricanArts	A B C with them. Most important, every Luba king is incarnated by a female spirit medium after death. Called Mwadi, such a medium inherits the c
2014	ACAD QuartRev/DistanceEd	A B C this article includes 3 narratives from students who were charged with using Twitter as a medium for sharing photographs and accompanyin
2002	MAG Sunset	A B C column into a sculptural element crafted from polished wood. Cabinets are constructed from inexpensive medium density fiberboard (MDF)
2017	TV Good Behavior	A B C HL - Feeling good? - No. - How you feeling? - Medium . Huh. Holy shit. I'm done. What? Why? All
2008	NEWS Denver	A B C // Salt and pepper to taste // Water as needed // Directions // In a medium bowl (or in a small food processor), whisk (or pulse)
2002	MAG PCWorld	A B C Pro/Wireless 5000, one of the first 802.11a network hubs. // The Message is the Medium Next-Generation Instant Messaging What is it? A wh
2017	MAG Nerdstz	A B C Whatever you're into, know that it's a gorgeous part of an engaging medium 's fantastic history.

Figure 1

<https://www.english-corpora.org/coca/>

Upon analysing the examples above, the students find that as far as *medium* is concerned, it is used in the corpus 33,635 times. They can observe that *medium* (meaning intermediate, in-between) as a modifier is used with such concepts as size (3, 10, 18, 11), heat (1,2, 6, 9), height (5), density (16), colours (4), or mood (17) - thus such concepts whose understanding is a matter of scale or gradeability. As far as the nominal function is concerned, *medium* is used to mean 'a means, a channel of transfer' (7, 12, 14, 15, 19, 20). Tracing the examples confirms the students' intuitions and gives them an insight into the definition of the specific content of the terms.

For *average*, COCA presents the following data:

The screenshot shows the COCA interface with the search term 'average' and a frequency of 115286. Below the search bar, there is a table of search results. The table has columns for rank, year, source, and context. The context column shows various sentences where the word 'average' is used.

Rank	Year	Source	Context
1	2012	BLOG addictinginfo.org	A B C Congressional Budget Office. Last year's increase was 4%. Compare that to the average 12% annual inflation rate during the previous 40 years. http://1.1
2	2012	BLOG addictinginfo.org	A B C America's working families. http://bit.ly/eSE4F # Under Obama, tax rates for average working families are the lowest they've been since 1950. http://bit
3	2012	BLOG addictinginfo.org	A B C cars. http://bit.ly/gCukV # He announced the development of a huge increase in average fuel economy standards from 27.5mpg to 35.5mpg starting in 2
4	2012	BLOG contracrastimes.com	A B C drastically. In the first seven months of 2012, California added jobs at an average rate of 23,000 a month. But in August and September, the average gain
5	2012	BLOG contracrastimes.com	A B C an average rate of 23,000 a month. But in August and September, the average gain was 6,800. # " We are seeing signs of a visible slowdown in
6	2012	BLOG crawfishboxes.com	A B C the two levels was a less than ideal 4.37. At Oklahoma he posted an average Game Score of 48.47 in 15 starts. At Corpus Christi he posted a Game
7	2012	BLOG crawfishboxes.com	A B C for a lefty is pretty good. Andy Pettitte, since 2002, has an average fastball of 88.8 miles per hour, according to FanGraphs so there's a chance
8	2012	BLOG crawfishboxes.com	A B C Houser # On the other hand, with one or two more MPH on his average fastball, Houser could be a very good power pitching prospect. With a hard
9	2012	BLOG crawfishboxes.com	A B C one. He has 20-25 HR potential, and he should hit for a decent average while walking at a respectable rate. Here's the problem, though: He
10	2012	BLOG crawfishboxes.com	A B C believe that he's not done growing. He has 20-25 HR potential, above average speed, and a strong arm with a quick release. He has quick hands
11	2012	BLOG cryptomundo.com	A B C become less and less convinced that these so called experts know much more than the average joe. Let's face it, so far they have tried fireworks, crying
12	2012	BLOG ...wencias.blogspot.com	A B C think outside that System? Can we believe they know much of anything about what average Americans are going through right now? # All of these showc
13	2012	BLOG ...tionalgeographic.com	A B C 'd expect. How big can a gray wolf get? In Yellowstone, the average weight of adult male wolves ranges between 100 and 120 pounds. The average weight
14	2012	BLOG ...tionalgeographic.com	A B C the average weight of adult male wolves ranges between 100 and 120 pounds. The average weight of adult female wolves ranges between 84 and 93 pou
15	2012	BLOG ...tionalgeographic.com	A B C Yellowstone has a variety of wildlife that the wolves can feast on year round in average temperatures that generally do not drop below -10F. Realistically i
16	2012	BLOG ...tionalgeographic.com	A B C 7, 10:22 am # I have a zoo pass and visit the zoo on average twice a week. Of all the animals at the zoo, the wolves are
17	2012	BLOG gantdaily.com	A B C new system, surpassing the 30-point mark for the third time and improved his assists average with six dimes. # The 'Black Mamba' also teamed up with k
18	2012	BLOG showbits.net	A B C films, and nearly a third of the total list. (Trivia: The average year of release of the 100 films is 1963; the average year of release of the 30 films i've seen is 1972. Both years
19	2012	BLOG showbits.net	A B C Trivia: The average year of release of the 100 films is 1963; the average year of release of the 30 films i've seen is 1972. Both years
20	2012	BLOG ...uality.wordpress.com	A B C issues are economic issues. # While LGBT persons tend to have more education on average than the general population, evidence suggests that they mal

Figure 2

<https://www.english-corpora.org/coca/>

As has been exemplified above, *average* appears in COCA 115,286 times. Taking its function as a modifier, among 20 examples above *average* (meaning estimated on given data, approximated, being representative of) modifies such nouns as inflation (1), family (2), economy (3), fastball (7), American (12), year (19) and concepts such as rate (4), gain (5), score (6), speed (10), or weight (13). In the nominal function, *average* is used in only one example (9). Other interesting uses of *average* are represented by phrases like *the average Joe* (11) and *on average* (20).

Thus, in studying only one corpus students can see the differences between the two terms in question, in their frequency as well as in the selection of words that they are used with. So, beginning with the frequency of terms and following on to their context, the students can learn how distinct these two words are and how they should be used.

4.4. NOW

The screenshot displays the NOW Corpus interface. At the top, there are navigation tabs: SEARCH, FREQUENCY, CONTEXT, and ACCOUNT. Below the navigation, there are search options: ON CLICK, TRANSLATE (?), GOOGLE, IMAGE, PRON/VIDEO, and BOOK (HELP). The search results are shown in a table with columns for RANK, DATE, SOURCE, and TEXT. The word 'MEDIUM' is highlighted in the text column.

RANK	DATE	SOURCE	TEXT
2	20-12-28 US	appleinsider.com	A B C a network connection at the time of software installation to function. Options for " Medium Security " and " None " are also available, with the former n
3	20-12-28 US	advances.sciencemag.org	A B C Commons Attribution license, which permits unrestricted use, distribution, and reproduction in any medium , provided the original work is properly cite
4	20-12-28 US	Chicago Tribune	A B C ended up pursuing a career in tech, she rediscovered her love for the artistic medium last December. # " I felt the urge to paint again, " she
5	20-12-28 US	Digital Photography Review	A B C getting some much needed relief with the double punch of the R6/R5 but Nikon's medium term prospects do not look good. I am not saying that Nikon
6	20-12-28 US	consequenceofsound.net	A B C his 1970 album Abraxas. Of course, Santana also has a " Smooth " medium blend named for his ubiquitous 1999 hit with Rob Thomas. # In case a
7	20-12-28 US	twincities.com	A B C suspects are described as black men, between 20 and 30 years old, with medium builds. One suspect is between 5-feet-9-inches and 6 feet tall. The oth
8	20-12-28 US	avclub.com	A B C hair, " and " good at math " highlighted what makes the niche-embracing medium of podcasting comedy so singular and special. Rodgers' and Boosters
9	20-12-28 US	avclub.com	A B C a great podcast, it's the best show going on internet culture in any medium , and this specific episode is better than the cherry on a whipped cream sun
10	20-12-28 US	Associated Press	A B C engagement through authentic conversation and high-quality content. They help their clients make Facebook a medium for a more personalized exper
11	20-12-28 US	Chicago Tribune	A B C was published in 2018. # The subject matter also has surfaced in a different medium , being dissected by Ziemia on a podcast released earlier in Decer
12	20-12-28 US	CIO	A B C years, as what was previously enterprise-class will become available and accessible for small and medium businesses. # For example, AI previously req
13	20-12-28 US	Forbes	A B C " People fainted in the presence of the Beatles because sightings of them in any medium were exceedingly rare. # Still, it brought up a question. How c
14	20-12-28 US	The Verge	A B C used to the idea. They also tend to work better on TV because the medium allows for longer-form storytelling that can navigate the intricacies of multi
15	20-12-28 US	nbcchicago.com	A B C night and Wednesday morning. # Most of the Chicago area is under a " medium " threat for hazardous travel, while some southern suburbs and parts
16	20-12-28 US	MarketWatch	A B C API platform market are as follows: # Minimum adoption of telecom API platform between medium and small operators. # The major driving factors of
17	20-12-28 US	The Guardian on MSN.com	A B C of people need to be vaccinated in order to achieve population immunity. In the medium term there will be pockets of the population in which the inf
18	20-12-28 US	PR Newswire	A B C anytime and anywhere. Be it a small office, a home office, a medium scale business or a large enterprise, the iR1643i adapts perfectly to any business e
19	20-12-28 US	advances.sciencemag.org	A B C Creative Commons **25:336:TOOLONG license, which permits use, distribution, and reproduction in any medium , so long as the resultant use is not fc
20	20-12-28 US	PR Newswire	A B C . please visit https: **25:124:TOOLONG # About My Business AdvocateMy Business Advocate is a medium through which small and medium-sized busi
21	20-12-28 US	avclub.com	A B C throughout its chronological filming process. The premise is simple: Six friends invite a medium into their Zoom chat to conduct a seance, which takes
22	20-12-28 US	Forbes	A B C # Going forward, we expect revenues to stay weak in the near to medium term, and if the comonarr is not able to control expenses, we believe

Figure 3

<https://www.english-corpora.org/now/>

The data above reveals that in the NOW Corpus *medium* appears 330,494 times (a number which considerably exceeds the use of *medium* in COCA). In the function of a modifier this word is used with such nouns as security (2), term (5, 22), blend (6), builds (7), business (12, 18), threat (15), or operator (16). As far as its nominal use is concerned, it is instantiated in examples 3, 4, 8, 9, 10, 11, 13, 14, 19, 20 and 21 of the above table. Thus, when compared to COCA, in the NOW Corpus (or at least in the first 20 examples) one can see that *medium* is used more as a noun than as an adjective. Additionally, the selection of nouns that *medium* modifies in NOW is different in quality from the ones that are modified in COCA, namely they are no longer nouns that require scalar modifiers.

As far as *average* is concerned, NOW offers the following data:

The figure displays two screenshots of the NOW Corpus (News on the Web) interface. The top screenshot shows the search results for the word 'AVERAGE', with a frequency of 1704126. The bottom screenshot shows a list of 36 examples of 'AVERAGE' usage in various contexts, including sports, economics, and science.

Example	Source	Context	Frequency
17	19-10-21 US	Montgomery Advertiser	A B C
18	13-09-30 IN	Cricknet Country	A B C
19	16-08-10 IN	Times of India	A B C
20	14-07-27 US	Sleep Review	A B C
21	20-12-12 GB	thetimesonline.com	A B C
22	16-04-13 NZ	sportal.co.nz	A B C
23	17-06-28 PH	ABS-CBN Sports	A B C
24	18-02-26 CA	Yahoo News	A B C
25	13-04-11 CA	driving.ca	A B C
26	19-08-22 SG	tnp.sg	A B C
27	10-02-04 GB	NHS Choices	A B C
28	14-02-04 ZA	Mail & Guardian Online	A B C
29	20-01-30 CA	kitco.com	A B C
30	17-02-03 MY	malaysiandigest.com	A B C
31	18-10-25 ZA	MyBroadband	A B C
32	19-12-19 US	PLOS	A B C
33	17-11-14 ZA	eProp.co.za	A B C
34	12-01-25 US	Agricultural Research	A B C
35	13-08-14 CA	Globe and Mail	A B C
36	19-01-15 NG	Guardia	A B C

Figure 4

<https://www.english-corpora.org/now/>

According to the data in the above tables, the frequency of the word *average* in NOW is 1,704,126 times. As far as the application of *average* as an adjective is concerned, it is used with such nouns as age (26, 27), American (29), estimate (30), projection (31), escalations (33), temperature (34) or gain (35). Thus, if compared to COCA, there are two similar examples (American, gain), and the rest of examples differ. The nominal uses of the word *average* are represented in examples: 19 and 21. Yet, what draws the attention is the verbal use of *average* as provided in examples: 17, 18, 20, 22, 23, 24, 25, 28, 32 and 36, the use that has no representation in the examples from COCA. So, differences between *medium* and *average* as presented in NOW in terms of their semantic quality do not seem so obvious as in COCA. However, on the whole (and as the above analysis shows) comparing data from different corpora adds additional information for students looking to find solutions to language intricacies.

4.5. Paralela

It is highly probable that the examples described above do not provoke any questions for native speakers who, without any problems, master the qualitative differences between *medium* and *average*. Yet, these qualitative differences are the most difficult for non-native speakers who frequently look for equivalent terms in their mother tongues. Such a situation is exemplified below where *average* and *medium* have the same equivalent in Polish- *średni* (in its basic form).

Show / hide columns		Show	20	entries	First	Previous	1	2	3	4	5	Next	Last
Lp	English	Polish											
1	If laws regulating consumer transactions were harmonised throughout the EU , small and medium businesses (SMEs) and consumers in all Member States would benefit . <i>Parlament Europejski</i>	Na harmonizacji przepisów prawa regulujących transakcje konsumenckie w całej UE skorzystałyby wszystkie małe i średnie przedsiębiorstwa (MŚP) oraz konsumenci . <i>European Parliament</i>											
2	Madam President , Commissioner , just one topic should be the focus of attention when the summit is held between the European Union and Japan at the end of this month : the disaster that has devastated the people of Japan as a result of earthquakes , the tsunami and continuing radioactive contamination , and the concrete role that can be played by the EU , its Member States and individual citizens in helping to deal with the resulting problems in the short , medium and long term . <i>Parlament Europejski</i>	Pani Przewodnicząca , Panie Komisarzy ! <i>European Parliament</i>											
3	Madam President , Commissioner , just one topic should be the focus of attention when the summit is held between the European Union and Japan at the end of this month : the disaster that has devastated the people of Japan as a result of earthquakes , the tsunami and continuing radioactive contamination , and the concrete role that can be played by the EU , its Member States and individual citizens in helping to deal with the resulting problems in the short , medium and long term . <i>Parlament Europejski</i>	Jeden temat powinien znaleźć się w centrum uwagi , kiedy z końcem miesiąca będzie odbywał się szczyt pomiędzy Unią Europejską a Japonią : katastrofa , która zdruzgotowała Japończyków w wyniku trzęsień ziemi , tsunami i trwającego skażenia radioaktywnego oraz konkretna rola , jaką może odegrać UE , jej państwa członkowskie i poszczególni obywatele , pomagając uporać się z wynikłymi problemami w krótkim , średnim i dłuższym okresie . <i>European Parliament</i>											
4	These factors will have a significant impact on energy costs in the medium term , and we will need to assess their repercussions on Europe 's current environmental strategy . <i>Parlament Europejski</i>	Te czynniki w średnim okresie znacząco wpłyną na koszty energii i będziemy musieli ocenić ich reperkusje dla obecnej strategii środowiskowej Europy . <i>European Parliament</i>											
Lp	English	Polish											
1	That is why , given the average level of catches over the past three years , the reference tonnage has been increased from 55 000 tonnes to 63 000 tonnes . <i>Parlament Europejski</i>	Oto dlaczego , z uwagi na średni poziom połowów w ciągu ostatnich trzech lat , tonaż referencyjny uległ podwyższeniu z 55 000 do 63 000 ton . <i>European Parliament</i>											
2	I also welcome the proposal included in the report that aims to guarantee payment of the full monthly wage during maternity leave , which is 100 % of the last monthly salary or the average monthly salary . <i>Parlament Europejski</i>	Z zadowoleniem przyjmuję również zawartą w sprawozdaniu propozycję służącą zagwarantowaniu wypłacania pełnego wynagrodzenia miesięcznego w okresie urlopu macierzyńskiego , wynoszącego 100 % ostatniego wynagrodzenia miesięcznego lub średniej miesięcznych wynagrodzeń . <i>European Parliament</i>											
3	Farming comes way below the average and that has got to be addressed in any future direction in which we reform the common agricultural policy , but I welcome the document . <i>Parlament Europejski</i>	Działalność rolnicza plasuje się znacznie poniżej średniej i rozwiązanie tego problemu należałoby uwzględnić niezależnie od kierunku , jaki miałyby przybrać przyszła reforma wspólnej polityki rolnej , ale z zadowoleniem przyjmuję ten dokument . <i>European Parliament</i>											
4	An average Japanese girl , for example , can expect to live to the age of 83 . <i>Parlament Europejski</i>	Na przykład przewidywana długość życia dla przeciętnej japońskiej dziewczynki wynosi 83 lata , podczas gdy w Lesotho wyniosłaby 42 lata . <i>European Parliament</i>											
5	It is because the World Meteorological Organisation 's figures are clear - they are average figures . <i>Parlament Europejski</i>	Ponieważ liczby podawane przez Światową Organizację Meteorologiczną są jednoznaczne - są to średnie . <i>European Parliament</i>											
6	At the moment , the EU youth unemployment rate has reached 20 % , while the average EU school dropout rate is 16 % , and in some countries , like Portugal , it has reached 40 % . <i>Parlament Europejski</i>	Stopa bezrobocia wśród młodzieży sięgnęła w tej chwili w UE 20 % , natomiast średnio 16 % osób w UE przedwcześnie kończy naukę , a w niektórych krajach , jak na przykład w Portugalii , odsetek ten wyniósł 40 % . <i>European Parliament</i>											

Figure 5

<http://paralela.clarin-pl.eu/#search/pl/>

In cases similar to the one mentioned above, an option to solve the problem of differences between apparently semantic terms is offered by corpus analysis of the original language. Furthermore, in PARALELA a student can read the different ways in which the words in question function across languages.

On the whole, the above examples of sentences from selected corpora (COCA, NOW, PARALELA) offer a wide selection of illustrations for 'language-in-use' situations for the words *average* and *medium*. However, if a student looks for real-life language applications, a reference to a corpus seems justified. Native speakers intuitively know how to use language (especially fixed expressions) in a given context. Moreover context, as a language phenomenon, has not been researched through grammar books, coursebooks or handbooks for practising 'language-in-use' situations. In other words, language learners have to learn the contextual environment for particular expressions by heart, so a reference source to check if the learners' intuition prompts a correct solution is a useful tool.

5. Conclusions

As has been discussed above, corpus analysis is a useful tool to be applied in teaching and learning foreign languages. Furthermore, selected aspects, theories and examples of corpora prove that they are valuable language resources that, on the one hand, register language forms and, on the other hand, function as reference resources available via open access to a broad public.

Yet, the main question of this article concerns the issues of how corpus analysis can influence the process of teaching and learning foreign languages. The suggestion presented above is that corpus analysis is definitely a method of activating teachers and students to enhance the educational process of teaching and learning foreign languages both inside and outside of the classroom. Furthermore, an additional advantage of using corpus analysis is the fact that students are given freedom to work on materials that they themselves identify with, as well as to

pursue their interests in selected fields which allows for a great amount of autonomy in learning.

References

- Aston, Guy (2002). "The learner as corpus designer". In: Bernhard Kettemann and Georg Marko (eds.). *Teaching and Learning by Doing Corpus Analysis*. Amsterdam: Rodopi, 9–25.
- Boulton, Alex (2016). "Integrating corpus tools and techniques in ESP courses". *ASP*: 69: 113–137. DOI: 10.4000/asp.4826.
- Braun, Sabine (2005). "From pedagogically relevant corpora to authentic language learning contents". *ReCALL* Vol 17/1: 47–64.
- Braun, Sabine (2010). "Getting past 'Groundhog Day': Spoken multimedia corpora for student-centred corpus exploitation". In: Tony Harris, Maria Moreno Jaén (eds.). *Corpus Linguistics in Language Teaching*. Frankfurt am Main: Peter Lang.
- Cobb, Tom, Alex Boulton (2015). "Classroom applications of corpus analysis". In: Douglas Biber, Randi Reppen (eds.). *The Cambridge Handbook of English Corpus Linguistics*. Cambridge: Cambridge University Press, 478–497.
- Corino, Elisa, Christina Onesti (2019). "Data-Driven Learning: A scaffolding methodology for CLIL and LSP teaching and learning". *Frontiers in Education* 4/7. DOI: 10.3389/educ.2019.00007.
- Hoffstaedter, Petra, Kurt Kohn (2009). "Real language and relevant language learning activities: Insights from the SACODEYL project". In: Anton Kirchhofer, Jutta Schwarzkopf (eds.). *The Workings of the Anglosphere. Contributions to the Study of British and US-American Cultures*. Trier: WVT, 291–303.
- Johns, Tim (1991). "From printout to handout: grammar and vocabulary teaching in the context of data-driven learning". In: Tim Johns, Philip King (eds.). *Classroom Concordancing. English Language Research Journal* 4: 27–45. DOI: 10.1017/CBO9781139524605.014.
- Lee, David, Y. W. (2010). "What corpora are available?" In: O'Keeffe, Anne, Michael McCarthy (eds.). *The Routledge Handbook of Corpus Linguistics* (1st ed.). London; New York, NY: Routledge, 107–121.

- McEnery, Tony, Andrew Wilson (2011). *Corpus Linguistics: An Introduction* (2. ed., repr). Edinburgh: Edinburgh Univ. Press.
- Meunier, Fanny (2011). "Corpus linguistics and second/foreign language learning: exploring multiple paths". *Revista Brasileira de Linguística Aplicada* 11: 459–477. DOI: 10.1590/S1984-63982011000200008
- Meyer, Charles, F. (2002). *English Corpus Linguistics: An Introduction*. Available at <https://DOI.org/10.1017/CBO9780511606311>.
- O’Keeffe, Anne, Michael McCarthy (2010). "What are corpora and how have they evolved?" In: Anne O’Keeffe, Michael McCarthy (eds.). *Routledge Handbook of Corpus Linguistics* (1st ed.). London; New York, NY: Routledge, 3–13.
- Pérez-Paredes, Pascual, Jose, M. Alcaraz-Calero (2009). "Developing annotation solutions for online Data Driven Learning". *ReCALL*, 21/1: 55-75.
- Pérez-Paredes, Pascual (2010). "Corpus linguistics and language education in perspective: Appropriation and the possibilities scenario". In: Tony Harris, Mariá Moreno Jaén (eds.). *Corpus Linguistics in Language Teaching*. Frankfurt: Peter Lang, 53–73.
- Renouf, Antoinette (2007). "Corpus development 25 years on: from super-corpus to cyber-corpus". In: Roberta Facchinetti (ed.). *Corpus Linguistics 25 Years on*. Amsterdam: Rodopi, 27–49.
- Sinclair, John (ed.) (2004). *How to Use Corpora in Language Teaching*. Amsterdam; Philadelphia, PA: John Benjamins.
- Suan Chong, Chia (2016). "Ten innovations that have changed English language teaching". Available at <https://www.britishcouncil.org/voices-magazine/ten-innovations-have-changed-english-language-teaching>.
- Tribble, Chris (1997). "Improvising corpora for ELT quick-and-dirty ways of developing corpora for language teaching". In: Barbara Lewandowska-Tomaszczyk, Patrick Melia (eds.). *PALC’97 Practical Applications in Language Corpora*. Lodz: Lodz University Press, 106–17.
- Widdowson, Henry, G. (1991). "The description and prescription of language". In: James Alatis (ed.). *Linguistics and Language Pedagogy: The State of the Art*. Washington, DC: Georgetown University, 11-24.
- Widdowson, Henry, G. (2003). *Defining Issues in English Language Teaching*. Oxford: Oxford University Press.

Widmann, Johannes, Kurt Kohn, Ramon Ziai (2010). "The SACODEYL search tool: exploiting corpora for language learning purposes". In: Ana Frankenberg-Garcia, Lynne Flowerdew, Guy Aston (eds.). *New Trends in Teaching and Language Corpora. Proceedings of the TaLC 2008*. London: Continuum, 321-327.

Language resources

<https://www.english-corpora.org/coca/>

<https://www.english-corpora.org/now/>

<http://paralela.clarin-pl.eu>

Joanna Redzimska
ORCID iD: 0000-0001-7837-9397
Uniwersytet Gdański
Instytut Anglistyki i Amerykanistyki
ul. Wita Stwosza 51
80-308 Gdańsk
Poland
joanna.redzimska@ug.edu.pl