

Corpora as Digital Humanities Tools for Learning Foreign Languages

Iryna Dilay

irynadilayi@mail.montclair.edu

Language teaching practice is becoming more learner-centered

*As an assistant tool in language learning, **corpus** makes an easy and quick analysis of the greatest amounts of linguistic data possible, and the learners are provided with a new approach to learn a language independently.*

What is a corpus?

CORPUS is a collection of pieces of language text in electronic form, selected according to external criteria to represent, as far as possible, a language or language variety as a source of data for linguistic research (Sinclair 1996, 2004).

Corpora are **the main language data base.**



Pedagogic utility

Sinclair: “*One does not study all of botany by making artificial flowers*” (1991, 6)

1987 - the publication of Collins COBUILD English Language Dictionary, the first corpus-based dictionary for learners.

1988 – Johns’ paper on the use of corpus-derived and corpus-based materials in the language classroom

1990 – Sinclair’s Collins COBUILD English Grammar

Biber et al., 1998; Hunston, 2002; Kennedy, 1998; McEnery & Wilson, 2001; McEnery et al., 2005, in press; Meyer, 2002; Partington, 1996; Stubbs, 1996, 2001; Tognini-Bonelli, 2001; Tomlinson, 1998, etc.

Useful links

- <http://corpus.byu.edu> - Mark Davies' site: COCA, COHA, BNC, GloWbE, CORE, Now, etc
- <http://www.wordandphrase.info/> - COCA based
- http://corpus.byu.edu/coca/files/Teaching_Vocabulary_Through_DDL.pdf
- <http://www.lextutor.ca/> Compleat Lexical Tutor
- <http://www.linguee.com/> Linguee – a bilingual multidictionary supported by parallel translations.
- <http://projects.ael.uni-tuebingen.de/backbone/moodle/> Backbone - Pedagogic Corpora for Content & Language Integrated Learning
- <http://www.corpustool.com/> UAM Corpus Tool (to analyze your own corpus)
- <http://www.papyr.com/applets/concordancer/>

BENEFITS OF CORPORA FOR LANGUAGE TEACHING

- More accurate descriptions of language than textbooks/intuitions
- Exposure to contextualized, meaningful language in 'real' usages
- Examples of specific registers/genres of language
- Reference tool for independent/autonomous language investigation and learning
- Cited from: Jonathan Smart, Northern Arizona University



Pedagogical use

- **Corpus-based direct ESL/EFL teaching**
- **Corpus-based indirect ESL/EFL teaching**
- **Self-study**
- **Syllabus design**
- **Materials development**
- **Language testing**
- **Teaching LSP**
- **Teaching academic English**
- **L1-focused corpora – developmental (CHILDES)**
- **L2-focused corpora – learner corpora**



Corpus-based vs Corpus-driven approaches

Elena Tognini-Bonelli (2001) distinguishes between a corpus-based and corpus-driven approaches.

- **The corpus-based approach** (CBA) is a method that uses an underlying corpus as an inventory of language data. From this repository, appropriate material is extracted to support intuitive knowledge, to verify expectations, to allow linguistic phenomena to be quantified, and to find proof for existing theories or to retrieve illustrative samples. It is a method where the corpus is interrogated and data is used to confirm linguistic pre-set explanations and assumptions. It acts, therefore, as additional supporting material.
- **The corpus-driven approach** (CDA) is a methodology

Corpus-based (corpus-informed) learning

- **Deductive Approach**

Through a deductive approach, the definition and explanations of a word can be tested by running a **concordancer**. The teacher and the students are then exposed to a large number of contexts of the same word to confirm what is learned so as to refine the original generalizations. This, therefore, assists the student in exploring the language in great detail and thereby gaining further insights into its structure.

Key Word in Context (KWIC)

The screenshot shows the MonoConc Pro application window. The title bar reads "MonoConc Pro - [Concordance - [reason]]". The menu bar includes "File", "Concordance", "Frequency", "Display", "Sort", "Window", and "Info". The main text area displays a list of search results for the word "reason". The first result is a paragraph starting with "Labor will become cool partly because it has new leadership. The election of John Sweeney as AFL-CIO president has called forth a lot of new declarations of support for labor's goals from liberal activists and intellectuals. More important, AFL-CIO headquarters is buzzing with activity because the new crowd understands it has a small window &MD; perhaps a year &MD; to impress the country that it has some new ideas and new energy." The second result is a paragraph starting with "But the biggest reason labor will be important is that its basic...". The search results are displayed in a list format, with the word "reason" highlighted in blue in each entry. The status bar at the bottom of the window shows "318 matches", "Original text order", "Strings matching: reason", the file path "C:\Documents and Settings\bgreen35\My Documents\Corpora\NewYorkTimes\atwp1996\LW960102", "22,305,495 words", and the time "10:34 PM". The Windows taskbar is visible at the bottom, showing the Start button, several application icons, and the system tray.

<p>
Labor will become cool partly because it has new leadership. The election of John Sweeney as AFL-CIO president has called forth a lot of new declarations of support for labor's goals from liberal activists and intellectuals. More important, AFL-CIO headquarters is buzzing with activity because the new crowd understands it has a small window &MD; perhaps a year &MD; to impress the country that it has some new ideas and new energy.
<p>
But the biggest reason labor will be important is that its basic

... is growing evidence to suggest there is reason to worry. <p> Two months ago, the In ...
... ionhood." Hope? There is absolutely no reason for hope in Africa's most populous nati ...
... iters rarely bothered to note. The only reason we happen to know how early itrion had ...
... and new energy. <p> But the biggest reason labor will be important is that its bas ...
... for several years now, and never had a reason to regret the choice. Each year it seem ...
... lifornia efforts. <p> Of course, one reason California has been able to catch up in ...
... or teaching purposes, which is the main reason why this project was born, but today we ...
... d or dismissed. Tipler offers hope that reason may again serve as the handmaiden of fa ...
... r country during the Cold War, I see no reason morally, ethically, logically that you ...
... of the essential amino acids. For this reason, vegetarians must be sure to eat protei ...
... rector Sydney Pollack were an important reason why Paramount spent \$75 million making ...
... y uses to regulate itself." <p> One reason, Bagley said, is that since insurance c ...
... g on the situation, but say there is no reason the extended furlough will make things ...
... nd the Middle East, officials said. One reason is that Christopher travels on an Air F ...
... the much-debated site, if for no other reason than its selection has diverted attenti ...
... &UR; (End optional trim) <p> ``The reason that the ICC was created 100 years ago ...
... , the nation's top-rated station. ``The reason: I don't want our listeners thinking th ...
... petitioners, industry executives said. The reason is that Northrop Grumman and Westinghou ...
... Col. Richard Pernod, said there was no reason to inform anyone else in NATO or the ge ...
... nding by him. <p> ``We still have no reason to believe that Dan Dorfman has violate ...
... ock market reporting. But the immediate reason for the dismissal on Monday of the popu ...
... e network added that ``we still have no reason to believe he has violated any law or f ...
... anyone else. <p> ``We still have no reason to believe that Dan Dorfman has violate ...
... CNBC said Wednesday. ``... CNBC has no reason to follow Money magazine's approach in ...
... -Texas, instructs us that there is good reason to keep paying members of Congress whil ...
... y honchos later changed their minds. No reason they shouldn't have. Books portraying b ...
... . <p> ``People will pay for the same reason they pay for subscriptions to newspaper ...
... omics are the province of intellect and reason, not religion. <p> But it was his wr ...
... usly undermined the regime for a simple reason: Proud of one of the world's oldest civ ...
... d out, against our will and for no good reason, what right do they have to tell us whe ...

318 matches Original text order Strings matching: reason
C:\Documents and Settings\bgreen35\My Documents\Corpora\NewYorkTimes\atwp1996\LW960102 22,305,495 words 10:34 PM

Corpus-based learning

- **Establishing Word Frequency**

All the students are motivated to find out the frequency of some word in the unit, then the figures are gathered and the frequency list of each unit takes shape. The list helps the students have a sense of the distribution of the word in vocabulary and acts as an important guide in their vocabulary learning.

The frequency of *know* across varieties

SECTION	ALL	US	CA	GB	IE	AU	NZ	IN	LK	PK	BD	SG	MY	PH	HK	ZA	NG	GH	KE	TZ	JM
FREQ	307 762 2	753 024 497	214 497	622 175	149 328	243 370	121 207	137 269	620 90	756 92	515 23	724 54	700 83	787 59	496 56	669 18	850 18	555 46	598 12	450 22	641 79
PERMIL	1,63 3.06	1,94 6.76	1,59 1.63	1,60 5.14	1,47 8.07	1,64 2.08	1,48 9.20	1,42 3.50	1,33 2.89	1,47 3.55	1,30 4.86	1,68 5.97	1,68 2.92	1,82 1.08	1,22 7.58	1,47 5.12	1,99 3.57	1,43 2.77	1,45 6.65	1,28 0.49	1,62 1.92

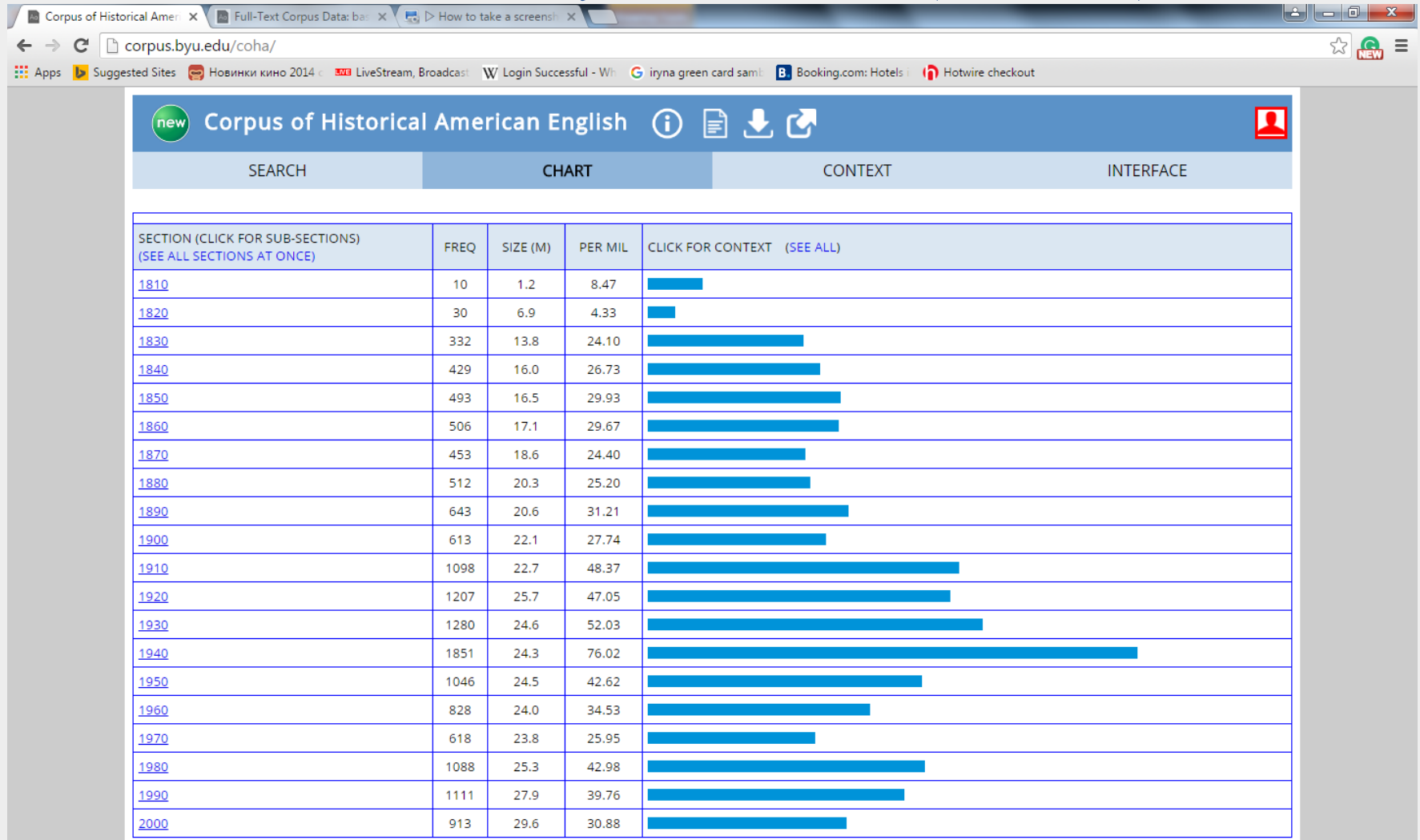
Table 1. Distribution of *know* across GloWbE varieties

The frequency of *think* across varieties

SECTION	ALL	US	CA	GB	IE	AU	NZ	IN	LK	PK	BD	SG	MY	PH	HK	ZA	NG	GH	KE	TZ	JM
FREQ	282 426 6	745 325	190 665	684 665	131 165	247 302	116 530	103 226	443 59	507 01	360 26	680 18	552 12	541 65	428 54	537 44	539 50	372 56	395 06	285 54	410 43
PERMIL	1,49 8.63	1,92 6.85	1,41 4.79	1,76 6.35	1,29 8.29	1,66 8.61	1,43 1.74	1,07 0.47	952 .25	987 .03	912 .38	1,58 2.75	1,32 5.82	1,25 2.42	1,05 9.42	1,18 4.71	1,26 5.06	960 .99	962 .12	812 .12	1,03 7.23

Table 1. Distribution of *think* across GloWbE varieties

Democracy over time (COHA)



Democracy in COCA

Corpus of Contemporary American English

SEARCH CHART CONTEXT OVERVIEW

SECTION (CLICK FOR SUB-SECTIONS) (SEE ALL SECTIONS AT ONCE)	FREQ	SIZE (M)	PER MIL	CLICK FOR CONTEXT (SEE ALL)
SPOKEN	6326	109.4	57.83	
FICTION	489	104.9	4.66	
MAGAZINE	4034	110.1	36.64	
NEWSPAPER	4895	106.0	46.20	
ACADEMIC	12924	103.4	124.96	
1990-1994	8386	104.0	80.64	
1995-1999	5734	103.4	55.43	
2000-2004	5613	102.9	54.53	
2005-2009	4391	102.0	43.03	
2010-2015	4544	121.6	37.38	

New coinages: *Brexit* in Now Corpus



Collocations: adj + *rain* in COCA

		FREQ		ALL	%	MI	
1	<u>HEAVY</u>	1388		45597	3.04	5.72	
2	<u>COLD</u>	464		66639	0.70	3.59	
3	<u>TROPICAL</u>	287		8466	3.39	5.87	
4	<u>POURING</u>	260		5519	4.71	6.35	
5	<u>TORRENTIAL</u>	238		416	57.21	9.95	
6	<u>FREEZING</u>	205		5009	4.09	6.14	
7	<u>DRIVING</u>	185		36030	0.51	3.15	
8	<u>STEADY</u>	182		15670	1.16	4.33	
9	<u>WET</u>	140		20525	0.68	3.56	
10	<u>PURPLE</u>	100		10476	0.95	4.04	
11	<u>GENTLE</u>	77		10032	0.77	3.73	
12	<u>TEMPERATE</u>	67		1329	5.04	6.44	
13	<u>HEAVIER</u>	64		4331	1.48	4.67	
14	<u>ICY</u>	55		4373	1.26	4.44	

From ‘three Ps’ to ‘three Is’

Presentation – practice – production =corpus-based

Illustration – interaction – induction (Carter and McCarthy (1995) =corpus-driven:

- a) looking at real data;
- b) discussing and sharing opinions and observations;
- c) making one’s own rules.

Data-driven learning – a learner is a **researcher! (Johns (1991), Leech (1997), Flowerdew (1993), etc.)**

Corpus-driven (data-driven) learning

- **Inductive Approach**

The study of some key words can be conducted through this approach. The teacher can ask the students to observe and discuss the concordance before explaining the word. Working with corpora does not necessarily require students sitting at computer terminals. The students are motivated in their learning of English and at the same time, learn some important skills that would make their learning less intimidating and keep them in regular contact with English.

- **Establishing Colligation, Collocation and Semantic Prosody**

The computerized corpus, with sufficient, authentic and typical data, enables users to conduct comprehensive and overall research into collocational patterns. Once patterns of non-native deviance have been discovered, students can be explicitly made aware of these patterns, and they will eventually be able to modify their linguistic behavior into a more native-like direction.



Example: to bear _____resemblance
CONCORDANCE:

- All of this **bears no** resemblance to the Europe of today. ...
- Queuing up at the polls at election time **bears a strong** resemblance to a gathering of the herd at a...
- Most **bear little** resemblance to their wild ancestors.
-which is under fire for a tea kettle ad that **bears an uncanny** resemblance to Adolf Hitler.
- Their stories **bear a striking** resemblance.

Learning vocabulary

In isolation:

- through naming objects;
- through translation;
- through definitions;
- through associations: synonyms, antonyms, word families, etc.

In context:

“You shall know a word by the company it keeps” (Firth).

Lexical Approach by Lewis (1993) – learning a word in lexical chunks (collocations).



Concordance

“Learning a word from either a short definition or a single sentence context tends to produce inert lexical knowledge that does not facilitate the word's comprehension in a novel context, while learning a word from several contexts, with or without a definition, tends to produce rich, transferable knowledge”
(Mezynski, 1983).

What is the missing word?

The prime minister didn't offer any (...) alternatives.

What is the word now?

The prime minister didn't offer any (...) alternatives.

Plans can not be expected to continue losing money year after year and remain (...) .

Mr. Jeffries could be a (...) candidate for mayor as early as 2017.

They are totally not economically (...) in rural areas.

Only in rare cases is the embryo (...), and the condition can be lethal for the mother.



Stevens' Experiment (1991)

The task was to have students recall a known word to fill a gap in a text, which was either a single gapped sentence or a set of gapped concordance lines. Learners retrieved words from memory more successfully when cued by the concordance lines in spite of their chopped-off nature.

Learning Synonyms

The use and differentiation of synonyms is usually a problem for many learners of English. The best way to help them to deal with this problem is exposing them to ample examples of real English, which can be obtained easily via a large English corpus. In this aspect, **concordances** can play a unique role here.

Establishing connotations:

- *Childlike, youthful, childish, young, immature, juvenile*
- *Chubby, fat, plump, overweight*
- *Confident, secure, proud, egotistical*

Semantic prosody of synonyms

‘Whichever way you look at it, it is a fact that “the meaning of a word can often be illuminated by the other words which it tends to co-occur with” (Wierzbicka 1987, 21), so that e.g.

comparing the adverbs which the verbs *rebuke*, *reprimand* and *reprove* tend to co-occur with, will yield important clues as to the semantic differences between them:

rebuking tends to be done *sharply*, whereas *reprimanding* tends to be done severely; only *reproving* can be done *gently* but cannot be done *sharply*, *severely*.



Compleat Lexical Tutor

Go to

<http://www.lextutor.ca/>

Learning grammar

Questioning the NORM!

- Progressives with stative verbs
- Prepositions
- Verb complementation (Infinitive or Gerund?)
- Reflexive pronouns
- Countable vs uncountable nouns
- Subject–predicate agreement
- Clauses
- Articles
- ...



Error analysis

- **Misuse of Words**

The purpose of error analysis is to seek the causes of the errors and the law of learning a language, but finding out the representative problems or errors requires the amount of the learner's output to be as large as possible. Hence, the analysis of the learners' interlanguage would be much more efficient if the outputs are transferred into an electronic corpus accumulatively from learners and then explored via a computer program.

Leaner Corpora

Learner corpora are sources of error analysis for different L2 learners, especially useful for educators.

Learner Corpus Association

(<http://www.learnercorpusassociation.org/>).

ICLE (The International Corpus of Learner English),
Louvain, Belgium (written)

ITAcorp (The International Teaching Assistants Corpus),
Pennsylvania State University, USA (spoken)

ICCI (The International Corpus of Crosslinguistic
Interlanguage), Tokyo University of Foreign Studies,
Japan (written)

CLEC (The Chinese Learner English Corpus), China
(written)

Importance of Learner Corpora

The interface between L1 and L2 materials.

Frequency information from native-speaker corpora alone is not sufficient to inform curriculum and materials design. Rather, 'it is important to strike a **balance between frequency, difficulty and pedagogical relevance**. That is exactly where learner corpus research comes into play to help weigh the importance of each of these' (Meunier 2002: 123).

Shortcomings and caveats of using corpora

- Learners are overwhelmed by the extremely large number of examples generated by their searches.
- Giving undue prominence to what is simply frequent at the expense of rarer but more effective or salient expressions (Cook 1998:61)
- Corpus data are traces of texts rather than discourse and must be recontextualized in language teaching (Widdowson 1990,2000)
- Corpus samples may have to be adapted when used with low levels or young learners, and when corpus examples contravene or offend sociocultural norms and customs.
- The need for learners and teachers to acquire new skills and assume new roles in the studying process.
- The need to make corpora more user-friendly.

Conclusion

In a learner-centered classroom, students are seen as being able to assume a more active and participatory role than in traditional approaches.

Corpora can be very useful tools for language learners, since they allow these learners to quickly and easily see how native speakers use the language in a wide variety of naturally occurring texts. Corpus-based research and teaching has the potential to empower non-native teachers and researchers, since native speaker introspection is no longer considered the one infallible source of insights into language structure and use.

Corpus use is not meant to replace existing teaching

Bibliography

- Aijmer, K. (Ed.). (2009). *Corpora and language teaching*. Amsterdam, The Netherlands: John Benjamins.
- Boulton, A., Carter-Thomas, S., & Rowley-Jolivet, E. (Eds.). (2012). *Corpus-informed research and learning in ESP: Issues and applications*. Amsterdam, The Netherlands: John Benjamins.
- Campoy-Cubillo, M. C., Fortuño, B. B., & Gea-Valor, M. L. (Eds.). (2010). *Corpus-based approaches to English language teaching*. London, UK: Continuum.
- Flowerdew, J. (2009). Corpora in language teaching. In M. H. Long & C. J. Doughty (Eds.), *The handbook of language teaching* (pp. 327-350).. Oxford: Wiley-Blackwell.
- Hunston, S. (2002). *Corpora in applied linguistics*. Cambridge, England: Cambridge University Press.
- Jones, M., & Durrant, P. (2010). What can a corpus tell us about vocabulary teaching materials? In A. O’Keeffe & M. McCarthy (Eds.), *The Routledge handbook of corpus linguistics*. London, UK: Routledge.
- McCarthy, M. J. (2004). Using corpora in language teaching. *CALPER Digests*. Retrieved from <http://calper.la.psu.edu/publications.php>.
- McEnery, T. (2011). What corpora can offer in language teaching and learning. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning* (Vol. 2, pp. 364-380), New York, NY: Routledge.
- Meunier, F. (2002). The pedagogical value of native and learner corpora in EFL grammar teaching. In S. Granger, J. Hung and S. Petch-Tyson (eds), *Computer learner Corpora, Second Language Acquisition and Foreign Language Teaching*, pp. 119-142. Philadelphia: John Benjamins.
- Reppen, R. (2010). Using corpora in the language classroom. In B. Tomlinson (Ed.), *Materials development in language teaching* (pp. 35-50). Cambridge, UK: Cambridge University Press.
- Sinclair, J. M. (2004). (Ed.). *How to use corpora in language teaching*. Amsterdam, Netherlands: John Benjamins.
- Stevens. V. (1995). Concordancing with language learners: Why? When? What? *CAELL Journal*, 6(2), 2-10.
- Sun, Y., & Wang, L. (2003). Concordancers in the EFL classroom: Cognitive approaches and collocation difficulty. *Computer Assisted Language Learning*, 16, 83-94.
-