

# On Corpus-based Autonomous Learning for College English Learners

Lihua Cai

School of International Studies, University of Science and Technology Liaoning, Anshan,  
114051, China

Cailihuao412@126.com

---

## Abstract

With the rapid development of interdisciplinary research in computer science and linguistics, corpora of different types are employed in the second language acquisition field. For Chinese English learners, native corpora like BNC will offer comprehensive guidance in lexical acquisition process; Chinese learner corpus is supposed to be used to predict and avoid effectively some typical errors and mistakes produced by college students.

## Keywords

Corpus; autonomous learning; college English.

---

## 1. Introduction

The official document issued by Education Department of China in 2007, namely “Curriculum teaching requirements for College English”, put forward that modern information technology, especially IT, should be fully employed to change traditional teaching model to promote students’ integrated application skills in English and cultivate students’ autonomous learning ability as well as comprehensive cultural quality. With development of computer technology, online information and corpus resources are expected to be employed efficiently for Chinese college students to learn English in and after class. This study aims to explore the practical way of computer-aided and corpus-based English learning for college students in their autonomous learning process.

## 2. Autonomous learning and corpus

Autonomous learning can be regarded as a sort of learning attitude as well as some kind of independent learning ability. Cohen (1990) pointed out that the success of language acquisition lies in learners’ factors and their ability to make good use of various chances to learn the language. Autonomous learning involves both external environment and internal environment. The former is the prerequisite and physical basis of the latter one which includes learners’ attitude and competence (Zhang Dianyu, 2005). Generally speaking, the theory of autonomous learning covers the following aspects (Shu Dingfang, 2004): a) learner’s attitude; b) learner’s ability; c) environment.

Corpus refers to a collection of naturally occurring samples of language which have been collected and collated for easy access by researchers and materials developers who want to know how words and other linguistic items are actually used. A corpus may vary from a sentence to a set of written texts or recordings. In language analysis corpuses usually consist of a relatively large, planned collection of texts or parts of texts, stored and accessed by computer. A corpus is designed to represent different types of language use, e.g. casual conversation, business letters, ESP texts (Jack R, 2005). Actually a number of different types of corpuses may be distinguished according to different purposes, structures or targets, namely specialized corpus, general corpus or reference corpus, comparable corpora, learner

corpus and so forth. In this study, learner corpus and specialized corpus on line may be employed to promote Chinese learners' English learning.

### 3. Corpus-based autonomous English learning model

Compared to traditional learning model, corpus-based autonomous English learning model may have some inherent advantages due to its features in nature. The details are as follows in the section below.

#### 3.1 Learner corpus and errors in the model

In recent years some high-quality learner corpora have been built in China to help Chinese English researchers shed new light on second language acquisition process of Chinese English learners. In addition, relevant findings in these studies may be of great value for English teaching and learning activities. Take the Chinese Learner English Corpus (shorten for CLEC) for example, it is the most important and representative learner corpus of written English produced by Chinese English learners in Mainland China. It was established under the supervision of Professor Gui Shichun at Guangdong University of Foreign Studies and Professor Yang Huizhong at Shanghai Jiao Tong University. The compositions in CLEC were produced by Chinese learners at different stages: middle school students (St2), college sophomores (Band IV students, St3), college juniors (Band VI student, St4), English-major juniors (St5) and English-major seniors (St6). CLEC is claimed to be reliable due to its scientific sampling process. First, the samples are all original writings without any correction; second, the sampling proportion is balanced among learners at different levels; third, the samples are from diverse sources so that the corpus covers learners' written output widely enough.

Therefore, it can be used to help Chinese learners to acquire English in autonomous learning style as well as a guide for college English teachers in the teaching process. CLEC has been tagged according to different types of mistakes and errors (see the table below), which is convenient for index engine to search for KWIC (key word in context) information closely related to typical errors produced by Chinese English learners. Both software package like AntConc or self compiled programs can fulfill the target. College English learners will benefit a lot from the study of high-frequency errors extracted from the corpus, hence the promotion of English proficiencies and accuracy in both written and spoken English.

Table 1 Text sample in CLEC

<p>&lt;ST 6&gt; &lt;SEX ?&gt; &lt;Y ?&gt; &lt;SCH GFL&gt; &lt;AGE ?&gt; &lt;WAY 1&gt; &lt;DIC 2&gt; &lt;TYP  2&gt; Euthanasia should be Legalized in China      Euthanasia, or mercy killing,  means helping to hasten the death of a person who is badly suffering. According to  Joseph Fletcher, it includes people getting incurable disease [fm1,-] and people in a  helpless condition, such as trapping in a blazing fire. [sn8,s] In China, suicide is  legal, which means, people are legal to kill themselves in a helpless condition, so  what we consider [fm1,-] is only whether it is legal to end the life of a [np7,1-]  incurable patient. I am in favor of the legalization of euthanasia, though some  others against it.[sn8,s] Those who argue [fm1,-] that euthanasia is  inhumane. It is a false argument. Death, most of the time, is the end of long suffering  period. With the advanced medical technique [fm1,-] and equipment, human life  can be extended. On one hand, it is indeed a good thing to provide people health  when the diseases are curable; on the other hand, it is rather a bad thing to extend  people's suffering when the diseases are incurable.</p>
---

#### 3.2 BNC and Online corpora in the model

Large quantities of native speakers' authentic texts in English can be collected and processed to build up a corpus for reference. The British National Corpus (BNC) is a 100 million-word collection of

samples of written and spoken English from a wide range of sources, designed to represent a wide cross-section of British English from the later part of the 20th century. The written part which accounts for 90 percent, includes extracts from regional and national newspapers, specialist periodicals and journals for all ages and interests, academic books and popular fictions, published and unpublished letters and memoranda, school and university essays, among many other kinds of texts. In addition, BNC is encoded according to the Guidelines of the Text Encoding Initiative (TEI) to represent both the output from CLAWS (automatic part-of-speech tagger) and a variety of other structural properties of texts (e.g. headings, paragraphs, lists etc.). Full classification, contextual and bibliographic information is also included with each text in the form of a TEI-conformant header.

BNC can offer authentic guide for Chinese learners in the process of autonomous learning, especially for the collocation identification and distinction of synonyms in terms of semantic prosody and idiomatic level (see the figure below). Furthermore, other native corpora are available for free online such as COCA, COHA, BASE and so forth.

The screenshot shows the AntConc 3.2.1w software interface. At the top, there's a menu bar with File, Global Settings, Tool Preferences, and About. Below the menu is a toolbar with Corpus Files, Concordance, Concordance Plot, File View, Clusters, Collocates, Word List, and Keyword List buttons. The main window displays a KWIC concordance search results table. The table has columns for Hit (number), KWIC (the context of the word 'cause'), and File (the source file for each hit). The search term 'cause' is entered in the search bar at the bottom left. Other search options include Words (checked), Case (unchecked), and Regex (unchecked). The search results show 148 hits across 50 files. Buttons for Start, Stop, Sort, and Save Window are also visible at the bottom.

Hit	KWIC	File
1	ndered that the word Maidenhead did not appear to cause the inhabitants of that town any problems. Profes	B1002010.txt
2	using the Sex Pistols in a film. It was to be the cause of his downfall. His first plan was to have the f	B1002010.txt
3	d; he knew the pain which a broken marriage could cause. But as the relationship deepened, so a separatio	B1002010.txt
4	nce founding Island, Blackwell had championed the cause of Jamaican music in its different guises of ska,	B1002010.txt
5	Richard's impetuosity, now become more and more a cause of annoyance. As Virgin's accountant, Jack Claydo	B1002010.txt
6	ry had been put to some use over the years in the cause of promoting his own artists, not least the Sex P	B1002010.txt
7	, and his faithfulness and loyalty to the African cause were profound. His love of Africa was paralleled	B1002010.txt
8	as another manifestation of his commitment to the cause of African development. He attended the early mee	B1002010.txt
9	as been released on bail fails without reasonable cause to surrender to custody. Section 6(1) provides th	B1002010.txt
10	ebate on the measure will be taken to promote the cause of hanging. In the Commons vote in June last year	B1002010.txt
11	ployees, warned that bringing in the troops would cause long delays and unnecessary suffering. Army perso	B1002010.txt
12	lief would have been tragically wrong. The direct cause of the Clapham Junction accident was undoubtedly	B1002010.txt
13	lective liability which lies on British Rail. The cause of Mr Hemingway's uncharacteristic errors failing	B1002010.txt
14	lay room on Sunday November 27 undoubtedly direct cause of accident. Excessive overtime blunted his worki	B1002010.txt
15	he had provided valuable assistance to the allied cause as Prime Minister after Italy's surrender in 1943	B1002010.txt
16	the effective moisture content of soils and thus cause a decrease in biological productivity. He makes n	B1002010.txt
17	asis( bilharzia), a debilitating disease that can cause liver fibrosis and bladder cancer, which is trans	B1002010.txt
18	lower Yangtze if this project proceeds. Here, the cause will be a very different one to that associated w	B1002010.txt
19	The resulting loss in crop productivity is grave cause for immediate concern, even if it can be partiall	B1002010.txt
20	o Arnold( 1987), shifting cultivation is the main cause of deforestation, causing 70, 50 and 35 per cent	B1002010.txt
21	fallow period of only 7 years, which is likely to cause long-term nutrient depletion so that maintenance	B1002010.txt

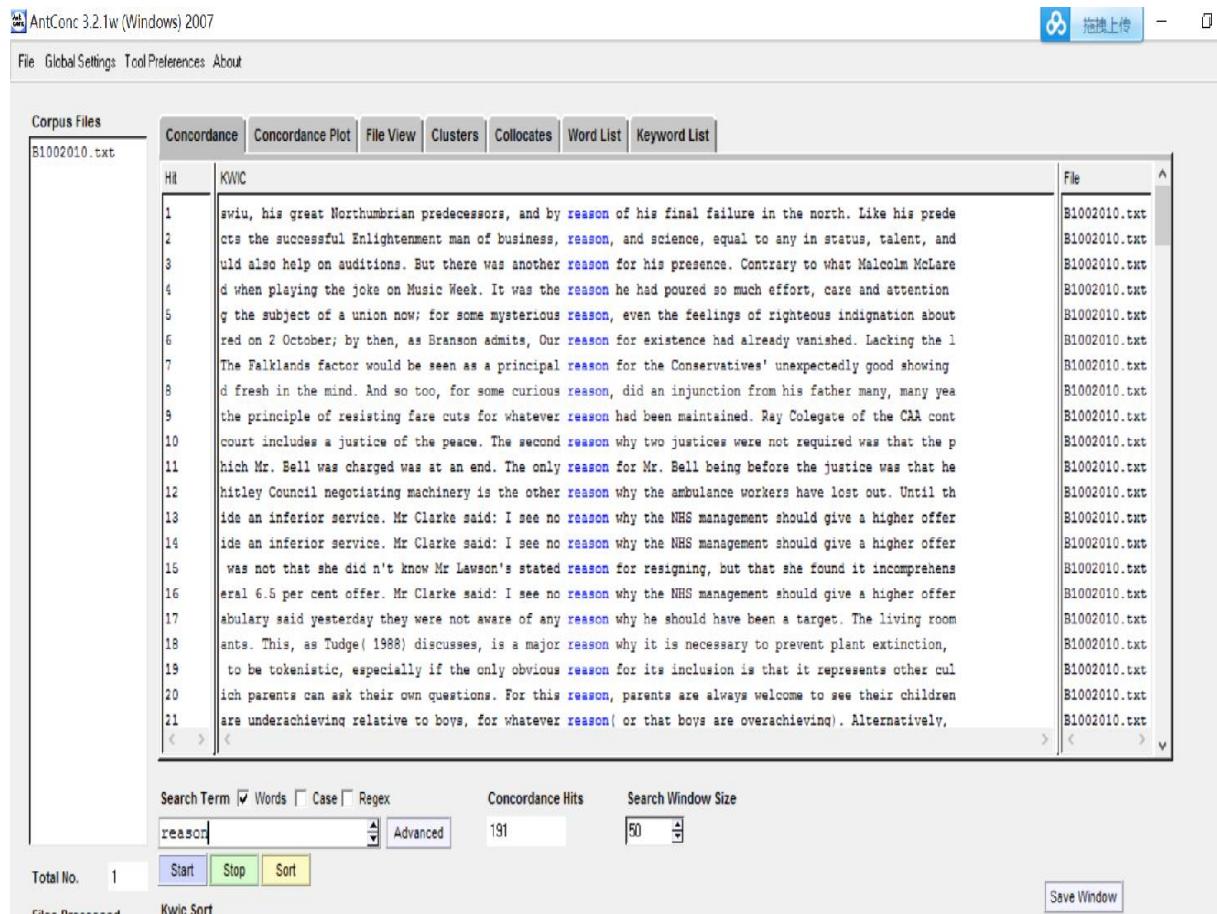


Figure 1. Collocation comparison of “cause” and “reason”

#### 4. Conclusion

To sum up, corpora of different types can be regarded as a valuable resource of language data in both native language and interlanguage. The former one will offer guidance and reference for language learners and the latter may help predict and avoid typical errors and mistakes that Chinese English learners tend to produce in the process of second language acquisition. Therefore, corpus-based model is supposed to be efficiently employed in autonomous learning process.

#### References

- [1] Cohen A. Language Learning: Insight for Learners, Teachers, and Researchers[M]. New York: Newbury House, 1990.
- [2] Zhang Dianyu. Strategy and Autonomous Learning of English Learning[J]. Foreign Language Teaching. 2005(1).
- [3] Shu Dingfang. Foreign Language Teaching and Reform: Problems and Strategies[M]. Shanghai: Shanghai Foreign Language Education Press. 2004.
- [4] Jack R. Longman Dictionary of Language Teaching and Applied Linguistics[M]. Beijing: Foreign Language Teaching and Research Press. 2005.