

Corpus-based Study on Chinese College Students' Use of Formulaic Language in English Public Speaking

Yiming Zhang

Faculty of Arts, Anhui University of Finance and Economics, Bengbu, Anhui, China

Abstract

The teaching and learning of formulaic language in EFL setting has been systematically researched. However, up to date, little work has focused exclusively on EFL learners' use of formulaic language in English public speaking. This paper attempts to bridge this gap by studying the use of formulaic language of Chinese college students in English public speaking and further explore to what extent their use of formulaic language is different from that of native speakers. The data analyzed in the study are drawn from two corpora, i.e. the corpus of speeches of Chinese college students, and the corpus of speeches of native speakers selected from American Rhetoric. The results show that (1) Chinese college students use fewer polywords and institutionalized expressions than native speakers, both in number and in variety; (2) Chinese students tend to rely on sentence builders more than native speakers.

Keywords

Formulaic Language; English Public Speaking; Corpus; Chinese College Students; Contrastive Study.

1. Introduction

In recent years, researchers have shown an increased interest in the study of formulaic language (FL). FL investigation has become one of the most active fields of linguistic inquiry. It is developing at an unprecedented speed due to the considerable progress of corpus linguistic study. FL, also known as phraseologies, language bundles, or chunks, is ubiquitous in our life, ranging from "How are you?" to "Kill two birds with one stone." It can be vaguely defined as recurrent multiword patterns that are stored in one's memory and are used naturally without too much analysis. However, it is always a difficult task to clarify the definition of FL that can serve the purpose of linguistic study (Hunston 2002). Firstly, it is difficult to distinguish FL from non-formulaic one. Secondly, there is no clear boundary between different kinds of FL, which means that one chunk can often fall into different categories. Based on previous research (Biber 1999, Wray 2002) and the objective of this study, FL in this paper is defined as continuous or discontinuous word combinations expressing relatively complete meanings, composed of 2-6 words.

Over the past decades, numerous studies have been conducted on typical types of FL, efficient ways to learn it, its structures and functions, practical value, semantic and collocational differences across various groups, etc. For example, studies conducted by Biber et al. (2004), Schmitt (2004), Chen and Baker (2010) examined the most frequently used FL in spoken and written language. Boers (2005), Wang (2006), etc. concluded ways to the learn FL efficiently. Bible et al. (2004), Hyland (2008) categorized FL into a myriad of groups according to its functions. Ding (2005), Biber (2009), Ellis (2012), etc. examined the assumption that there should be a direct correlation between FL and foreign language learning. Research done by Deng (2013), Romer (2009), etc. compared the semantic and collocational differences between learners of various backgrounds. To date, however, little work has focused exclusively on the

use of FL in English public speaking. The research reported in this paper attempts to bridge this gap by studying how Chinese college students use FL to make English speeches, and how it differs from that of the native speakers.

This paper, which is exploratory and descriptive in nature, sets out to compare the way chunks are used by Chinese college students in English public speaking and by native speakers based on corpus investigation. It is composed of five sections. The first section discusses the research background, research objectives and significance of the thesis. In the second section, we offer a brief review of prior corpus-based and experiment-based research on the classification and acquisition of FL, as well as what earlier analysts have talked about its application, and semantic and collocational differences across various groups in particular. We then outline the procedure and methodology used in this study. After that we move on to show research results. Finally, a conclusion is drawn regarding the patterns of FL in public speeches delivered by Chinese college students as EFL learners and then we discuss the pedagogical implications of our findings.

2. Literature Review

Nattinger & DeCarrico (1992) first put forward a systematic way to classify FL according to its structures, as shown in Table 1. Biber et al. (1999) first defined FL as “word sequences that occur more than ten times per million words and are distributed in more than five texts”. Over the past decades, a large body of literature has investigated typical types of FL, efficient ways to learn it, its structures and functions, practical value, semantic and collocational differences across various groups, etc.

Table 1. Classification of FL

CLASSIFICATION	DEFINITION	EXAMPLE
Polywords	Fixed word combinations	by and large, a piece of cake, etc.
Institutionalized expressions	Sentence-level chunks which include proverbs, phatic language, etc.	no pains no gains, thank you, etc.
Phrasal constraints	A structure composed of fixed words, in which corresponding words and phrases can be filled as needed	two (days/weeks/months) ago, as (quickly) as possible, etc.
Sentence builders	Sentence-level chunks, which provide a structure for building the whole sentence and can be filled in with appropriate phrases or clauses as needed	I believe (that), it occurs to me (that), etc.

Studies conducted by Schmitt (2004), Chen and Baker (2010), and Biber et al. (2014) mainly focused on identifying the typical types of FL used in spoken and written language. For example, Biber et al. (2004) have found that most FL used in spoken language is for building sentences, while most FL in academic papers is polywords or phrasal constraints. Subsequent studies have shown that the systematic differences of FL not only exist in different texts, but also in different language users. The research of Chen & Baker (2010) examined the assumption that academic writers use more polywords than non-academic ones, and EFL learners use more sentence builders than native speakers. Biber et al. (2014) concluded that authors of academic papers used more polywords and phrasal constraints than non-academic ones, and high-level EFL learners used more of them than low-level EFL learners. These studies indicate that structural differences in FL have become an important criterion for distinguishing different texts and authors.

Studies done by Wood (2002), Wray (2002), Boers (2005), Wang (2006) have shown several viable ways to facilitate the learning of FL. For example, Boers (2005) performed a computer-aided experiment to analyze the alliteration in FL and concluded that salient phonological patterning (e.g., alliteration) played a significant role in learning and memorizing FL.

There has also been an increasing amount of literature on the investigation of FL from a functional perspective. According to different functions of FL, Biber et al. (2004) categorized chunks into stance bundles, discourse organizers, referential expressions, and special conversational bundles. Research has shown that there are significant differences between spoken and written English in the use of these four types of FL. More stance bundles are used in speaking, while more referential expressions are used in academic writing. Hyland (2008) proposed a method to classify FL in academic papers from a functional perspective, dividing FL into research-oriented, text-oriented and participant-oriented.

Findings of Weinert (1995), Howarth (1998), Ding (2005), Biber (2009), Ellis (2012) analyzed the application and practical value of FL, noting that there should be a direct correlation between FL and foreign language learning. For instance, Ding (2005) carried out a host of investigations into the number of chunks used in students' English speaking and writing to explore the connection between a student's English level and the frequency of his using FL, finding that phraseology can help students to speak and write English in a more natural and accurate way. These studies have been largely concerned with the ways to boost FL learning and have also provided evidence for the positive effects of FL learning upon language acquisition.

More recently, studies have focused on the semantic and collocational differences of FL across various groups of people. Several studies have been conducted by Guan & Zheng (2005), Ädel & Erman (2012), Deng (2013) on the different use of FL between native language learners and non-native ones. For instance, Deng (2013) compared the preference of each group of people for certain types of chunks. Other studies have been largely concerned with the comparison of the use of FL between students and professional researchers in academic writing (Cortes, 2004; Romer, 2009). For example, Romer (2009) has found that some chunks are rarely used in academic writing, be they in the papers of students or professional researchers, but are very common in other forms of writing.

The abovementioned studies have well documented and have given important insights into FL structures and functions, as well as their correlation with the users' language levels or backgrounds in terms of academic writing, etc. One potential limitation, perhaps, is that they may have not adequately studied how chunks are used by EFL learners in speeches, i.e., English public speaking. Therefore, the present study, drawing on insights from previous studies, employs a corpus approach to study the use of FL of Chinese college students in English public speaking. We aim to find out the features of the use of FL of Chinese college students in English public speaking as compared with native speakers. We presume that the frequency and STTR of FL used by Chinese college students are both significantly different from those of native speakers.

3. Methodology

This paper sets out to investigate the semantic and collocational features of Chinese college students' use of FL in English public speaking based on corpus data, providing answers to the following three questions: (1) What are the semantic and collocational features of Chinese college students' use of FL in English public speaking? (2) What are the semantic and collocational differences of the use of FL between Chinese college students and native speakers in English public speaking?

3.1. Corpus

3.1.1. Research Corpus

The research corpus in this study is a self-built 130,800-word corpus which contains 300 speeches delivered by contestants in the two most influential English public speaking contests in China, (i.e. “FLTRP Cup” National English Speaking Contest and “21st Century Cup” National English Speaking Competition). It covers a wide range of topics including education, environment, mass media, personal development, society, etc. The 300 speeches can be subdivided into three groups: 100 speeches delivered by contestants who get first prize, 100 delivered by those who get second prize, and 100 delivered by those who get other prizes. They can well represent the speeches made by college students in different parts of China. The corpus has 130,800 tokens, and the Type-Token Ratio (TTR) is 5.22.

3.1.2. Reference Corpus

The reference corpus consists of 55 speeches delivered by native speakers selected from a speech corpus, American Rhetoric (AR). The 55 speeches were carefully selected in order to ensure the representativeness and make the size of the reference corpus commensurate with the research one. The token of the corpus amounts to 130800, and the TTR is 5.29. The basic information of the two corpora is shown in Table 2.

Table 2. Basic information of the two corpora

CORPUS	TOKEN	TYPE	TTR
Observed corpus	130800	6824	5.22
Reference corpus	130872	6925	5.29

3.2. Research Procedure

The research procedure can be summarized as follows:

In the first step, we categorized FL into four categories, i.e. polywords, institutionalized expressions, phrasal constraints, sentence builders according to their structures and meanings based on the study of Nattinger & DeCarrico (1992), as shown in Section 2. According to their study, FL bears various functions. The abovementioned four categories of FL share some functions while differ in others. Specifically, polywords can help interlocutors to reduce the stress of immediate analysis and language decoding in communication and can well indicate the language proficiency of users. Not being able to proficiently use polywords is one of the major reasons why a lot of EFL learners sound unnatural and even confusing. Institutionalized expressions are the most effective FL to keep the flow of a conversation and to achieve phatic function. Phrasal constraints can serve as glue that hold a text together, thus greatly increase the coherence of speaking. Compared with the other three categories of FL, sentence builders can better help users to reduce language processing time, because a long sentence can be easily constructed through simple structures of sentence builders. This kind of FL best eases the stress of the brain when someone is required to respond and formulate sentences quickly.

Then we used the N-Grams model in AntConc to search for chunks composed of 2-6 words, and then excluded those which appeared fewer than 5 times. After that, we manually screened out the word combinations that were not qualified as FL, such as “think the”, “that we can”, “and I believe the”, “all thank you for your”, etc.

In the final step, we conducted a Chi-square test to explore the differences in the use of FL across polywords, institutionalized expressions, phrasal constraints and sentence builders.

4. Results and Discussion

The results of the contrastive analysis and discussion are as follows:

4.1. Chinese College Students' Use of FL

Table 3 shows the results obtained from the corpus analysis of the Chinese college students' use of FL in English public speaking.

Table 3. Features of the use of FL in the research corpus

CLASSIFICATION	TOKEN	TYPE	TTR	STANDARD FREQUENCY
Polywords	385	41	10.65	29
Institutionalized expressions	712	17	2.39	54
Phrasal constraints	993	76	7.65	76
Sentence builders	2158	98	4.54	165
Total	4248	232	5.46	324

As Table 3 shows, in terms of the frequency of the use of FL, the number of chunks in the research corpus reaches 4248 in total, and the type amounts to 232 in total. The uneven distribution of the four kinds of FL indicates that there is a preference for the use of certain kinds of FL among Chinese college students in English public speaking.

Specifically, it is observed that sentence builders, with a token of 2158, are the most frequently-used FL accounting for 50.80% of the total chunks, while polywords are used the least frequently, with a token of 385, accounting for only 9.06%. The underlying reasons for this difference can be twofold. First, both the structures and meanings of polywords are relatively fixed and are often used in more restricted contexts compared with other types of FL. This might make EFL learners feel less certain about the appropriateness of the use of polywords. Therefore, they choose to cut down the number of polywords to avoid mistakes. Furthermore, it can be noted that because the contestants need to finish writing those speeches in a relatively short period, they are more likely to select the types of chunks which can help them compose sentences quickly, such as sentence builders.

For TTR, institutionalized expressions have the lowest one of 2.39 while polywords have the highest one of 10.65; the TTR of phrasal constraints and sentence builders are 7.65 and 4.54 respectively. A closer look at the two special types of FL, i.e. sentence builders (the most frequently-used chunks) and institutionalized expressions (chunks of the lowest TTR) shows that "I think" is the most frequently-used chunk in sentence builders, occurring 1145 times and accounting for 53.06% of sentence builders; "Thank you" is the most frequently-used chunk in institutionalized expressions, occurring 612 times and accounting for 85.96% of institutionalized expressions.

To sum up, the results above suggest that the most important two categories of FL used by Chinese college student in speeches are sentence builders and institutionalized expressions. The students tend to make speeches using these chunks.

4.2. Comparison between Chinese College Students and Native Speakers

To conduct Chi-square test, we calculated standard frequency on a basis of 1,000,000 words. Table 4 shows the differences of the frequency (FREQ) and standard frequency (SF) between Chinese college students and native speakers in English public speaking when they use FL. Table 5 shows the differences of the TTR and Standard Type-Token Ratio (STTR) between the two groups.

Table 4. Comparison of FREQ and SF between the two groups

CLASSIFICATION	CHINESE STUDENTS		NATIVE SPEAKERS		χ^2	p
	FREQ	SF	FREQ	SF		
Polywords	385	2900	668	5100	607.43	0.00*
Institutionalized expressions	712	5400	236	1800	1806.53	0.00*
Phrasal constraints	993	7600	1362	10400	439.51	0.00*
Sentence builders	2158	16500	2057	15700	20.20	0.00*
Total	4248	32400	4323	33000	5.69	0.02*

Table 5. Comparison of TTR and STTR between the two groups

CLASSIFICATION	CHINESE STUDENTS		NATIVE SPEAKERS		χ^2	p
	TTR	STTR	TTR	STTR		
Polywords	10.65	58.87	10.19	56.34	0.08	0.77
Institutionalized expressions	2.39	14.01	13.59	79.67	48.63	0.00*
Phrasal constraints	7.65	47.64	8.89	55.38	0.50	0.48
Sentence builders	4.54	26.16	3.53	20.33	0.80	0.37
Total	5.46	38.23	4.64	32.51	0.37	0.55

As Table 4 shows, in terms of the total frequency of the chunks, there is a significant difference ($\chi^2=5.69$, $p=0.02\leq 0.05$) between Chinese college students and native speakers. Chinese students cannot use as many chunks as native speakers, both in number and in variety. This indicates that in general, Chinese college students' ability to deploy FL is weaker than that of native speakers. It can be seen that compared with native speakers, Chinese college students use fewer polywords and phrasal constraints but more institutionalized expressions and sentence builders.

The SF of polywords across the two groups is significantly different ($\chi^2=607.43$, $p=0.00\leq 0.01$). This type of FL is used by native speakers 2200 more times every 1,000,000 words than by Chinese college students. In English public speaking, among the four types of FL, polywords can best help speakers deliver speeches in a more natural and accurate way. Therefore, English speeches delivered by Chinese students are, in all likelihood, less natural and idiomatic than those made by native speakers. Upon further study, we have found more evidence to prove this assumption. We counted the frequency of a few polywords with certain functions and then calculated the standard deviation (σ). The results of the polywords introducing opinions are shown in Table 6. It can be observed from Table 6 that the standard deviation of Chinese students (10.80) is higher than and almost twice as much as that of native speakers (5.60). This

indicates that native speakers use polyword in a wider range and are capable of replacing one polyword with many others that express the same meaning to balance the frequency of polywords being used. It is obvious that Chinese college students rely on “in my opinion” and “as far as I’m concerned” too much in their speaking, making their language less natural and idiomatic, while native speaker deploy many other chunks to express the same meaning, such as “as for me” (6 times), “from where I stand” (6 times), “from what I can see” (5 times). Furthermore, we have also found that Chinese students make more mistakes in polywords compared with native speakers. For example, “in the campus” is used by Chinese college students up to 11 times in total while the correct expression should be “on the campus”.

Table 6. Comparison of opinion-introducing polywords between the two groups

POLYWORD \ GROUP	CHINESE STUDENTS	NATIVE SPEAKERS
In my opinion	31	22
As far as I’m concerned	17	13
From my perspective	12	10
From my point of view/ In my view	6	9
To my knowledge/way of thinking/mind	4	7
As for me	1	6
From where I stand	0	6
From what I can see	0	5
TOTAL	71	79
STANDARD DEVIATION (σ)	10.80	5.60

The SF of institutionalized expressions is also significantly different ($\chi^2=1806.53$, $p=0.00\leq 0.01$). Chinese students use 3600 more institutionalized expressions every 1,000,000 words than native speakers. It can be contributed to one of the two reasons below: (1) Chinese college students are better at using institutionalized expressions. (2) Chinese college students depend on a few institutionalized expressions, leading to a lack of variety. Based on the results shown in 4.1, we believe that there should be more plausibility in the latter reason. Chinese college students tend to rely on certain institutionalized expressions like “Thank you” too much, while native speakers are able to employ more varieties of FL to realize phatic functions.

Significant difference is observed in phrasal constraints as well ($\chi^2=439.51$, $p=0.00\leq 0.01$). Native speakers use 2800 more phrasal constraints every 1,000,000 words than Chinese college students. As is discussed in 3.2, phrasal constraints are a good indicator of the coherence of speaking. Therefore, it can be speculated that the speeches delivered by Chinese college students are less coherent than those by native speakers. We have conducted further study to observe the difference between Chinese college students with higher scores and those with lower ones in English public speaking contests and concluded that the former group generally use more phrasal constraints, both in token and in type, which indicates that phrasal constraints are useful in public speaking if speakers want to receive more recognition from the audience.

In terms of sentence builders, there is also significant cross-group difference ($\chi^2=20.20$, $p=0.00\leq 0.01$). This indicates that, comparatively speaking, Chinese students are not so adept at

using English to compose speeches within a short period, thus need to depend more on sentence builders, which can serve as great ready-made frameworks to help them save time and energy. We have found, upon further study, that the frequencies of “I think” and “I believe” are significantly different between the two groups. “I think” occurs 1145 times and accounts for 53.06% of sentence builders in the observed corpus, while it only occurs 430 times and accounts for 20.99% of sentence builders in the reference corpus. “I believe” has a frequency of 194 and accounts for 8.90% of sentence builders in the observed corpus while its frequency is merely 103 and accounts for 5.03% in the reference corpus. This finding suggests that compared with native speakers, Chinese students are likely to rely more on these two chunks in speeches.

As Table 5 shows, in terms of the total STTR of the chunks, there aren't significant difference ($\chi^2=0.37$, $p=0.55>0.05$) between Chinese college students and native speakers. This indicates that the overall richness of chunks is similar between the two groups. It can be seen that the only significant cross-group difference lies in the STTR of institutionalized expressions ($\chi^2=48.63$, $p=0.00\leq 0.01$). The STTR of institutionalized expressions of Chinese college students (14.01) is tremendously lower than that of native speakers (79.67), suggesting that Chinese college students tend to repeat some institutionalized expressions while native speakers are able to use a wider range of them with less repetition. This points to the feature of Chinese students' use of institutionalized expressions: “Large in quantity, but small in richness.”, as revealed in Wang (2009).

5. Conclusion

This paper has discussed: (1) the semantic and collocational features of Chinese college students' use of FL in English public speaking; (2) the semantic and collocational differences of the use of FL in English public speaking between Chinese college students and native speakers. It has been found that: (1) Chinese college students prefer using sentence builders and institutionalized expressions in English public speaking while they can't use polywords so proficiently as other categories of FL; (2) In general, Chinese students' use of FL is significantly different from that of native speakers in number and variety, but similar to that of native speakers in richness. Specifically, they cannot use as many polywords and phrasal constraints as native speakers, which makes their speeches less natural. Besides, Chinese college students are likely to depend on a few polywords and institutionalized expressions too much, leading to a lack of variety. Last but not the least, they have more preference for sentence builders, which can save their energy and time when they are required to deliver a speech with little preparation time. (3) Chinese college students tend to rely more on several chunks than native speakers, such as “I think”, “I believe” “thank you” “in my opinion”, etc.

The findings of this research have significant implications for the effective ways to improve Chinese college students' language fluency and coherence in English public speaking. Firstly, they are supposed to find alternative polywords and institutionalized expressions that express the same meaning or realize the same function in order to increase the frequency and variety of these chunks. Secondly, they are expected to acquire more phrasal constraints to increase the coherence of their speeches. Thirdly, they should avoid relying on sentence builders too much, while use more devices such as parallelism, asyndeton and tricolon to build sentences.

It has to be acknowledged that this study is limited in 2 aspects. First of all, the speeches in the research corpus are delivered only by those taking part in speech contests, students who are, relatively speaking, of higher English level, and thus may not be able to represent the whole group of Chinese college students very well. What's more, we haven't found feasible ways to evaluate the language fluency of speeches quantitatively, so it is hard to explore the correlation between the use of FL and the language fluency of speeches. Further work needs to be done to

establish whether Chinese college students' language can be more fluent and idiomatic with more chunks being used in speeches based on corpus that can better represent the whole group.

References

- [1] Adel, A. & Erman, B. (2012). Recurrent word combinations in academic writing by native and non-native speakers of English: A lexical bundles approach. *English for Specific Purposes*, (2), 81–92.
- [2] Biber, D. (2009). A corpus-driven approach to formulaic language in English. *International Journal of Corpus Linguistics*, 14(3), 275–311.
- [3] Biber, D., Conrad, S. & Cortes, V. (2004). If you look at...: Lexical bundles in university teaching and text-books. *Applied Linguistics*, (3), 371–405.
- [4] Biber, D., Gray, B. & Staples, S. (2014). Predicting patterns of grammatical complexity across language exam task types and proficiency levels. *Applied Linguistics*, (59), 1–31.
- [5] Biber, D., Johansson, S., Leech, G., Conrad, S. & Finegan, E. (1999). *Longman Grammar of Spoken and Written English*. Harlow: Pearson.
- [6] Boer, F. (2005). Finding ways to make phrase-learning feasible: The mnemonic effect of alliteration. *System*, 33, 225–238.
- [7] Chen, Y. & Baker, P. (2010). Lexical bundles in L1 and L2 academic writing. *Language Learning and Technology*, (2), 30–49.
- [8] Cortes, V. (2004). Lexical bundles in published and student disciplinary writing: Examples from history and biology. *English for Specific Purposes*, (4), 397–423.
- [9] Deng, Y. C. (2013). Usage patterns of formulaic sequences in Chinese EFL learners' oral production. *Foreign Languages and Their Teaching*. 7(3), 60–65.
- [10] Ding, Y. R. & Qi, Y. (2005). Use of Formulaic Language as a Predictor of L2 Oral and Written Performance. *Journal of PLA University of Foreign Languages*, 28(3), 49–53.
- [11] Ellis, N. (2012). Formulaic Language and Second Language Acquisition: Zipf and the Phrasal Teddy Bear. *Annual Review of Applied Linguistics*, (32), 17–44.
- [12] Guan, B. & Zheng, S. (2005). Recurrent word combinations in Chinese college students' English writing. *Modern Foreign Languages (Quarterly)*, 28(3), 288–296.
- [13] Hyland, K. (2004). *Genre and Second Language Writing*. Lansing: The University of Michigan Press.
- [14] Hyland, K. (2008). Academic clusters: Text patterning in published and postgraduate writing. *International Journal of Applied Linguistics*, (1), 41–62.
- [15] Nattinger, J. & DeCarrico, J. (1992). *Lexical phrases and language teaching*. New York: Oxford University Press.
- [16] Romer, U. (2009). English in academia: Does nativeness matter, *Anglistik: International Journal of English Studies*, (2), 89–100.
- [17] Schmitt, N. (2004). Knowledge and acquisition of formulaic sequences: A longitudinal study, In N. Schmitt (Ed.), *Formulaic Sequences: Acquisition, Processing and Use* (pp.157-160). Amsterdam: John Benjamins Publishing Co.
- [18] Wang, L. F. & Qian, J. (2009). A Corpus-based Study on Chunk Patterns of Chinese EFL Public Speakers. *Foreign Language Research*, 115(6), 115–120.
- [19] Wang, L. F. & Zhang, D. F. (2006). An overview of the developments of studies on L2 prefabricated chunk acquisition abroad. *Foreign Languages and Their Teaching*, 17(5), 17–21.
- [20] Weinert, R. (1995). The role of formulaic language in second language acquisition: A review. *Applied Linguistics*, 16(2), 180–205.
- [21] Wood, D. (2002). Formulaic Language in Acquisition and Production: Implications for Teaching. *TESL Canada Journal*, 20(1), 1–15.
- [22] Wray, A. (2002). *Formulaic Language and the Lexicon*. Cambridge: Cambridge University Press.