# A Retrieval System for Interlanguage Analysis

**Lung-Hao LEE [a], Li-Ping CHANG [b],**
**Bo-Shun LIAO[a], Wan-Ling CHENG[b] & Yuen-Hsien TSENG[a*]**
[a]*Information Technology Center, National Taiwan Normal University, Taiwan*
[b]*Mandarin Training Center, National Taiwan Normal University, Taiwan*
*samtseng@ntnu.edu.tw

**Abstract:** In this paper, we describe the development of a retrieval system that is designed for analyzing the interlanguage. We adopt the annotated TOCFL learner corpus as the target to explore the language acquisition for leaners of learning Chinese as a foreign language. An illustrative scenario is presented to demonstrate the functionalities of implemented prototype system. This system can be deemed as a computer-assisted tool for contrastive interlanguage analysis research.

**Keywords:** second language acquisition, learner corpora, Mandarin Chinese

## 1. Introduction

Learner corpora: the Longman Learners' Corpus, the International Corpus of Learner English (ICLE) (Granger, 2003), and the Cambridge Learner Corpus (CLC) (Nicholls, 2003), to name but a few, are important collection of foreign language learners' linguistic production for research of second language acquisition and foreign language teaching (Granger, 2002). To make learner corpora to be more useful, they must be annotated using defined error types for automatic or manual analysis (Díaz-Negrillo & Fernández-Domínguez, 2006).

From the viewpoint of engineering, annotated learner corpora can be employed to develop specific Natural Language Processing (NLP) systems for educational applications. For instances, computer-assisted essay writing (Milton, 1998), spell error checking (Yu et al., 2014; Tseng et al., 2015), and grammatical error detection/correction (Chodorow and Leacock, 2000; Izumi et al., 2003; Lee et al., 2013, Ng et al., 2014; Yu et al., 2014; Lee et al., 2015). From linguistic perspectives, interlanguage is the type of linguistic system used by the second-/foreign- language learners who are in the process of learning a target language. Contrastive Interlanguage Analysis (CIA) is the main methodology that combines research areas of corpus linguistics and second language acquisition (Granger, 2015). Comparing learner corpus with native speaker's usages, researchers can identify learners' incorrectly linguistic usages or overgeneralized situations (Ishikawa, 2009). In addition, linguistic features of different L1s learners can also be obtained from CIA researches (Chang, 2014).

In this work, we develop and implement a retrieval system to help researchers to analyze the interlanguage. Our system is flexible to meet information needs in terms of various searching conditions. Besides, the search results can be downloaded easily if needed.

## 2. The Retrieval System of Annotated Leaner Corpus

The learner corpus is mainly originated from the computer-based writing Test of Chinese as a Foreign Language (TOCFL). The writing test is designed according to the six proficiency levels of the Common European Framework of Reference (CEFR). Test takers have to complete two different tasks for each level. For example, for the A2 (Waystage level) candidates, they will be asked to write a note and describe a story after looking at four pictures. All candidates are asked to complete the writings on line. Each text is then scored on a 0-5 point scale. Score 5 means high-quality writings, score 3 is the threshold for passing the test, and so forth. There are 4,567 essays have been collected in the TOCFL learner corpus.

The native Chinese speakers are then trained and asked to label the grammatical error types of learners' writings using the tagging editor (Lee et al., 2014). For the purpose of studies in Chinese learners' interlanguage, hierarchical error tags are designed. One is target modification taxonomy, which includes mis-ordering (permutation), redundancy (addition), omission (deletion), and mis-selection (substitution). The other is linguistic category classification that consists of linguistic types, for example, noun, verb, preposition, specific construction, and so on. So far, 2837 essays with the score above 3 have been annotated. In total, there are 33, 497 error instances. The top 3 error tags are Sv (mis-Selection of verbs), Sn (mis-Selection of nouns), and Madv (Missing of adverbs). Their frequencies are 3838, 2252, 1714, respectively.

The searching functions of our retrieval system can be divided into two main parts: (1) *Basic search*: users can select the main types of error tags, i.e., modification types, and the linguistic categories. The levels of learner's language proficiency in CEFR and the scores of the learners' written essays can be chosen by ticking all that apply using checkbox. Searchers can also choose learner's mother-tongue language and types of writing styles. Besides, we also provide the concordance function to show the character contexts surrounding the search target in the search results. (2) *Advanced search*: when the search targets are determined, searchers can further filter the search results by including/excluding the characters occurring in the left-hand/right-hand sides. Moreover, the search results can be downloaded easily in plain text format for further research.


## 3. An Illustrative Scenario for Interlanguage Analysis

We present a scenario to illustrate the effectiveness of our developed retrieval system for interlanguage analysis. Take the '讓' (*rang4* 'to make') sentence for example, we can choose the main error type S and the sub-type *rang*. Figure 1 shows the searching results. We found that learners usually confuse '讓' (*rang4* 'to make') with '把' (*ba3* 'disposal marker'), '對' (*dui4* 'to someone'), and '給' (*gei3* 'to give'). If there is no error tag annotated in the corpus, even we search the keyword '讓' (*rang4* 'to make') and investigate one by one sentence to find the erroneous usages, but we cannot find the misused sentences with ba3 or dui4. With the help of this retrieval system, we can shorten the time efficiently. Moreover, we can limit the searching results into the specific word only, such as '把' (*ba3* 'disposal marker'), which will benefit to do deep observation and analysis. In addition to filtering function, we can also select the specific learners' attributes such as the learners' mother tongue or their proficiency. Take advantage of these functions, the analysis of interlanguages could be more easily and quickly done.



Figure 1. Searching results of rang4 sentence in the annotated TOCFL corpus

# 4. Conclusions and Future Work

This article describes our retrieval system that can be applied to analyze error types in annotated learner corpora. An illustrative scenario of this prototype system is presented for lnterlanguage analysis. We will further collect researchers' feedbacks and discuss with them to enhance its functions.

# Acknowledgements

# References

Chang, L.-P. (2014). Salient linguistic features of Chinese learners with different L1s: a corpus-based study. *International Journal of Computational Linguistics and Chinese Language Processing*, 19(2), 53-72.

Chodorow, M., & Leacock, C. (2000). An unsupervised method for detecting grammatical errors. *Proceedings of NAACL'00* (pp. 140-147). Seattle, Washington: ACL Anthology.

Díaz-Negrillo, A., & Fernández-Domínguez, J. (2006). Error tagging systems for learner corpora. *RESLA*, 19, 83-102.

Granger, S. (2002). A bird's eye view of learner corpus research. In Granger, S., Huang, J. & Petch-Tyson, S. (eds.) *Computer Learner Corpora, Second Language Acquisition, and Foreign Language Teaching*. Amsterdam & Philadelphia: Benjamins, 3-33.

Granger, S. (2003). The International Corpus of Learner English: a new resource for foreign language learning and teaching and second language acquisition research. *TESOL Quarterly*, 37(3), 538-546.

Granger, S. (2015). Contrastive interlanguage analysis: a reappraisal. *International Journal of Learner Corpus Research*, 1(1), 7-24.

Ishikawa, S. (2009). Phraseology overused underused by Japanese learners of English: a contrastive interlanguage analysis. *Phraseology, Corpus Linguistics and Lexicography*, 87-100

Izumi, E., Uchimoto, K., Saiga, T., Supnithi, T., & Isahara, H. (2003). Automatic error detection in the Japanese learner's English spoken data. *Proceedings of ACL'03* (pp. 145-148), Sapporo, Japan: ACL Anthology.

Lee, L.-H., Chang, L.-P., Lee, K.-C., Tseng, Y.-H., & Chen, H.-H. (2013). Linguistic rules based Chinese error detection for second language learning. *Proceedings of ICCE'13* (pp. 27-29), Bali, Indonesia: Asia-Pacific Society for Computers in Education.

Lee, L.-H., Lee, K.-C., Chang, L.-P., Yu, L.-C., Tseng, Y.-H., & Chen, H.-H. (2014). A tagging editor for learner corpus annotation and error analysis. *Proceedings of ICCE'14* (pp. 806-808), Nara, Japan: Asia-Pacific Society for Computers in Education.

Lee, L.-H., Yu, L.-C., & Chang, L.-P. (2015). Overview of the NLP-TEA 2015 shared task for Chinese grammatical error diagnosis. *Proceedings of the 2nd Workshop on Natural Language Processing Techniques for Educational Applications* (pp. 1-6), Beijing, China: ACL Anthology.

Milton, J. (1998). WORDPILOT: enabling learners to navigate lexical universes. *Proceedings of the International Symposium on Computer Learner Corpora, Second Language Acquisition and Foreign Language Teaching*, Hong Kong, China.

Ng, H. T., Wu, S. M., Briscoe, T., Hadiwinoto, C., Susanto, R. H., & Bryant, C. (2014). The CoNLL-2014 shared task on grammatical error correction. *Proceedings of CoNLL'14* (pp. 1-14), Biltmore, Maryland: ACL Anthology.

Nicholls, D. (2003). The Cambridge Learner Corpus – error coding and analysis for lexicography and ELT. *Proceedings of CL'03* (pp. 572-581), Lancaster, UK.

Tseng, Y.-H., Lee, L.-H., Chang, L.-P., & Chen, H.-H. (2015). Introduction to SIGHAN 205 bake-off for Chinese spelling check. *Proceedings of SIGHAN'15* (pp. 32-37), Beijing, China: ACL Anthology.

Yu, L.-C., Lee, L.-H., Tseng, Y.-H., & Chen, H.-H. (2014). Overview of SIGHAN 2014 bake-off for Chinese spelling check. *Proceedings of CLP'14* (pp. 126-132), Wuhan, China: ACL Anthology.

Yu, L.-C., Lee, L.-H., & Chang, L.-P. (2014). Overview of grammatical error diagnosis for learning Chinese as a foreign language. *Proceedings of the 1st Workshop on Natural Language Processing Techniques for Educational Applications* (pp. 42-47), Nara, Japan: Asia-Pacific Society for Computers in Education.