# Design Principles and Functionality of Chinese Interlanguage Corpora: A Case Study of the HSK Dynamic Composition Corpus 2.0

**Baolin Zhang**

**Abstract**  Since the beginning of the twenty-first century, great progress has been made in terms of the construction of the Chinese Interlanguage Corpora and the essential role that these corpora has played in the study of Chinese second language teaching and research. However, there still exist some technical issues in terms of corpora design and functionality, such as the simplicity of the search function, difficulty searching for a certain interlingual phenomenon, and inconvenience caused by its not-so-user-friendly interface design. Furthermore, especially in current times, network safety has become an increasingly prominent issue and has resulted in a lack of operational corpora that satisfy the needs of the academic community. Under these circumstances, and in order to aim for Generation 2.0 which requires more delicacy and abundance, it has become necessary to adjust the design concepts and motivations behind the corpora. To ensure the uninterrupted operation and accessibility of the corpora, attentiveness and improvements in terms of system security are crucial. Meanwhile, optimizations and improvements of corpora features, especially the search function, are also essential for comprehensively meeting the needs of all users.

**Keywords**  Chinese interlanguage corpus · Design concepts · System security · Search function

B. Zhang (✉)
Research Institute of International Education of Chinese Language, Beijing Language and Culture University, Beijing, China
e-mail: zhangbl@blcu.edu.cn; baolin08@126.com

# 1   Introduction

## 1.1   *The Development of Corpus Construction and Its Applied Research*

Considered as the first Chinese interlanguage corpus in academia, the "Chinese Interlanguage Corpus System" was developed at the Beijing Language Institute in 1995. Zhang (2019: 86) notes that "despite there still being some problems with the corpus, such as the small scale of the corpus size, the limited breadth and depth in processing, and the lack of corpus retrieving speed (see Chen, 1996), the corpus system still holds a good reputation and practical value as a pioneering sharer in academia".

Several corpora have been built during the first decade of the twenty-first century, among which the most influential are the HSK Dynamic Composition Corpus (Beijing Language and Culture University), the Chinese Interlanguage Corpus for International Students (Jinan University College of Chinese Language and Culture, including written and spoken corpora), the Interlanguage Corpus for International Students (Sun Yat-sen University), and the Corpus of Chinese Interlanguage Error Analysis for Foreign Students (Nanjing Normal University).

During the second decade of the twenty-first century, more corpora were built as more Chinese teachers, experts, and scholars have devoted themselves to corpus construction. Some of these include the Chinese Interlanguage Corpus for Korean International Students (Ludong University), the Chinese Written Language Corpus for International Students (Beijing Chinese Language and Culture College), the Chinese Acquisition Corpus for Foreigners (Shanghai Jiaotong University), the Language Acquisition Corpus for Spoken Chinese (LAC/SC, The Chinese University of Hong Kong), the small-scale Foreign Student Oral Interlanguage Corpus (Suzhou University), the Corpus Based on Oral Telephone Examinations (Peking University), the Errors in Continuity of Chinese Characters Interlanguage Corpus(Sun Yat-sen University), the Global Chinese Interlanguage Corpus (led by BLCU and co-constructed by academia), the TOCFL Learner Corpus (Taiwan Normal University), and the Guangwai-Lancaster Chinese Learner Corpus (CLC, Guangdong University of Foreign Studies and Lancaster University, UK).

The construction and development of the Chinese Interlanguage Corpus have promoted the emergence of corpus-based research and achieved numerous important research results. Representative publications include Zhao et al. (2008), Zhang et al. (2008), Xiao et al. (2009), and Zhang et al. (2014). Taking the HSK Dynamic Composition Corpus as an example, 3858 research papers of various types involving this corpus were found in the China Knowledge Network (CNKI) database as of May 26, 2019 (Fig. 1).

These research papers mainly consisted of two categories: Master's thesis with a total of 2,929 papers and journal publications with a total of 732 papers (Fig. 2).

More importantly, the rapid development of corpus construction and corpus-based applied research has propelled a shift for Chinese interlanguage and Chinese
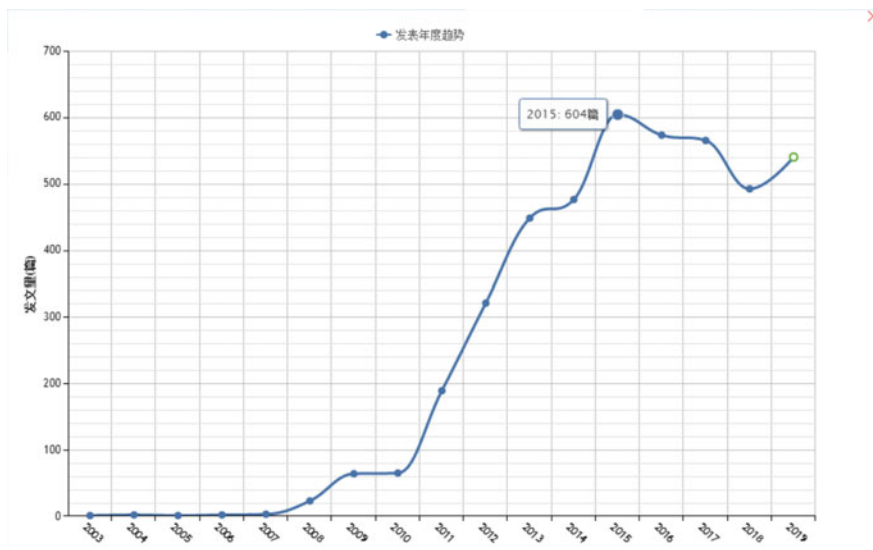
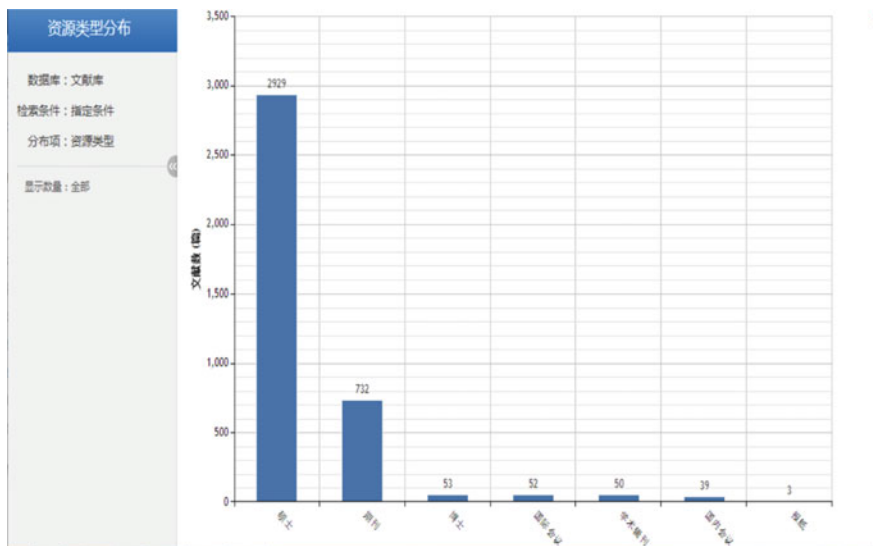**Fig. 1** Annual number of research papers based on HSK



**Fig. 2** Distribution of research resource types based on HSK

second language acquisition research. The field has seen a move away from small-scale, empirical, and speculative research toward large-scale, real data-based research that combines both quantitative and qualitative analyses. This shift has resulted in researchers being able to make conclusions with greater objectivity, universality, and stability.

## *1.2 Ontological Study of Chinese Interlanguage Corpus Construction*

The so-called ontology research refers to the research on the related theoretical issues of corpus construction. The reason for this naming is for the tendency of emphasizing the practice of corpus construction and despising theoretical discussion in the construction of corpus and emphasizing that theoretical research is an integral part of the construction of the corpus. Mainly it includes the following.

The overall design of a corpus occurs after the specific objectives of the corpus have been clarified and the necessity for its construction has been addressed. Addressing the feasibility of constructing the corpus involves (1) researching how to construct the corpus so that it meets the desired objectives; (2) clarifying its features; (3) determining its scale, materials, and structure; and (4) deciding annotation methods, principles of construction, and methods of application (Zhang & Xiliang, 2015). Representative research papers on overall corpus design include Chu and Xiaohe (1993) "The Basic Idea of Establishing a 'Chinese Interlanguage Corpus System'" and Cui and Zhang's (2011) "The Construction Plan of 'Global Chinese Learner Corpus'".

In addition to the overall design, labeling conventions are also an important component of corpus construction. The specification of corpus annotation schemes makes it easier to centralize what type of content is annotated, i.e. labeled, and how. The existence of a centralized annotation system is very important, especially in the case of corpus-based research paradigms. However, what to label can vary in corpus-building practices depending on the professionalism, knowledge, and hands-on experience of the project leaders, resulting in different functions and different use values of the corpus. Despite the importance of labeling conventions, not much research has been done in this area with the exception of several articles about the comprehensiveness of labeling (see Xiao et al., 2014; Zhang and Xiliang, 2018). Therefore, more extensive academic discussions or debates are needed to clarify the reasoning behind, and unify the understanding of, labeling conventions so as to form a labeling scheme that can be generally accepted by academia, thus promoting further development of corpus construction. Zhang et al. (2019)'s publication "Study on Standardization of Chinese Interlanguage Corpus Annotation" is one of the important related research projects in the field in recent years.

Construction standards are another aspect of corpus construction that requires attention. Research in this area is a response to the lack of unified standards in the

construction of Chinese interlanguage corpora and the great arbitrariness in the practice of database construction. Construction standards are a summary of the experience in Chinese interlanguage corpus construction. These standards draw on various lessons, consolidate academic theories of corpus construction, designate the levels of corpus construction, and have important guiding significance for corpus construction (Zhang & Xiliang, 2015). Not much research has been carried out in this area yet, and Zhang and Xiliang's (2015) "On Construction Standards of Chinese Interlanguage Corpus" is a relatively comprehensive and systematic discussion of this issue.

A final noteworthy component of corpus construction is that of software systems. In the past, the construction of Chinese interlanguage corpora focused more on the size, composition, and labeling conventions and lacked attention with regard to the development of management and retrieval/search systems. In fact, the development of software systems, including management and retrieval systems, plays a very important role in enhancing the practical functions of a corpus and improving the level of corpus construction.

## 1.3 Existing Problems

The development and progress of the Chinese interlanguage corpus are undoubtedly huge, and it has been widely recognized by the academic circles. However, the problem is also obvious, which not only determines the construction level of the corpus, but also affects the application research based on the corpus. Mainly it is manifested in the following aspects:

(1) Annotated content is not comprehensive and cannot meet the needs of teaching and research in many aspects. For example, the "HSK Dynamic Composition Corpus" only has entries from intermediate and advanced learners and can only be used for static horizontal research, not for vertical research of the acquisition process. In addition, only character, word, sentence, and writing composition errors are annotated, making the corpus suitable for error analysis but not performance analysis.

(2) The search and retrieval function is too simple, resulting in the inability to search for some important language phenomena and limitations with regard to the functions of the corpus. For example, it is not possible to search for sentences using the "是……的" *Shi……de* construction, the "连……也/都……" *Lian……ye/dou……* construction, semi-fixed collocations, nor sentences where "离" *Li* is used as a separable word.

(3) The functional design is not user-friendly and is inconvenient. For example, you cannot automatically download the corpus results generated by a specific query; wrong recordings, mislabeling, or omissions found while browsing the corpus cannot be fixed or amended by users; thoughts, comments, and suggestions from users cannot be relayed to the creator.

(4) The network security is not up to industry standards and cannot be accessed outside specific school networks, which seriously affects the use of the corpus. This is largely attributed to the fact that the corpus was developed ahead of its time, and the programming language and technology used in its original construction have now become outdated. This has resulted in security loopholes in the system and the failure of the corpus to meet the open requirements. Consequently, the corpus could not continue to be accessed by domestic and foreign users, diminishing its application in Chinese language teaching and related research.

Faced with the four major issues outlined above, we have made various adjustments to the corpus to improve its functionality and applicability for Chinese language teaching and related research. We regularly patched the system which resolved some of the issues preventing the corpus from being available on the campus network and to those on campus. However, the patched system does not meet the needs of domestic and foreign academia. The entire HSK corpus was copied to the BCC corpus so that everyone could at least browse the HSK corpus. However, BCC is a native language corpus with a different retrieval method from the HSK corpus. It is still difficult to meet the needs of academia due to the inconvenience and inability to search in a mislabeled corpus.

## *1.4 Solutions*

Faced with Academia's urgent need for corpora despite their shrinking number, we decided to redevelop the HSK corpus' software system by using the current mainstream programming language. We did this in order to continue and better serve Chinese teachers, scholars, researchers, graduate students, and Chinese learners domestically and globally.

The task of developing the new system was approved and commenced on January 5, 2018. On February 11 of the same year, the new system was complete and deployed to the server. On March 28, the system was officially opened to the public following trial operations and debugging. The new version of the system is known as the "HSK Dynamic Composition Corpus (Version 2.0)", and can be accessed at hsk.blcu.edu.cn. The new system ranks high in safety performance, which has allowed it to remain open to the public. As a result, we have achieved the goal of rebuilding the corpus system and the HSK corpus can continue to serve Chinese language teaching and research around the world.

## 2  Design Principles

In view of the various problems in the construction of the Chinese interlanguage corpus, and in accordance with the purpose of building the corpus, we have formulated the basic principles for the redevelopment of the corpus system.

### 2.1  Aim

Our aim in the construction of Chinese interlanguage corpora has always been to serve the teaching and research of Chinese as a foreign language. At the Third International Symposium on Construction and Application of Chinese Interlanguage Corpus held in the Summer of 2014, our focus became to proactively and wholeheartedly serve Chinese language teaching and research all over the world. It is precisely under the guidance of this understanding and aims that the HSK corpus, whether it be version 1.0 which was completed and launched at the end of 2006, version 1.1 which was upgraded in August 2008, or version 2.0 which has been recently re-developed, was made available free of charge and without delay to users all over the world.

### 2.2  Principles

There are three core principles with regard to the re-designing of the HSK corpus: (1) ensuring reliable and secure operation; (2) ensuring that the functions meet user demands; and (3) ensuring a fast, simple, and user-friendly experience.

The reconstruction of the corpus software system was mostly due to network security issues. Therefore, the first requirement for the reconstruction of the corpus is that there should be no or minimum security risks. It must be ensured that the corpus can operate normally and continue to serve the academic community uninterruptedly. Specifically, first of all, the new corpus system must not have any high-risk and medium-risk vulnerabilities, and low-risk vulnerabilities should be kept to a minimum as much as possible so that it can successfully pass the security inspections implemented by relevant departments and units. Secondly, when there are high-risk and medium-risk vulnerabilities, it can respond quickly and solve the problem in time, so as to ensure that the corpus is normally opened and not closed. This is a new problem brought about by the rapid development of information technology in the Internet era, and corpus builders must pay close attention to this problem.

The second consideration that came into play when redesigning the HSK corpus was the need to ensure that the functions of the corpus were computationally powerful enough to meet user demands. The HSK corpus version 1.0 and version 1.1 are products of the 1.0 generation, which were built during the initial period of Chinese interlanguage corpora construction (cf. Zhang, 2019). These two versions embody

defining characteristics of that time in that they are simple and large scale, not fully functional, and have difficulty satisfying users' demands in many aspects. For example, in these earlier versions, it is possible to search for the usage of the separable word "合" *He*, but not for the usage of "离" *Li*. One can search for some sentence constructions with marker words such as the "把" *Ba* construction and the "比" *Bi* construction, but not for sentences with the "是……的" *Shi……De* construction or the "连……也/都……" *Lian……Ye/Dou……* construction, because these sentence constructions require two search terms.

These deficiencies in functionality may lead to incomplete research conclusions since the generated results, limited by the corpus' lack of functionality, do not comprehensively or accurately reflect second language speaker use. That is, in cases where the corpus does not have the ability to process the queried entry, as is the case with separable words like "Li", the phenomenon cannot be fully analyzed. Furthermore, the value of the corpus cannot be realized for relative corpus research under such conditions because certain sentence patterns cannot be retrieved. The new corpus system solves these problems and facilitates users to explore various language phenomena so as to better serve Chinese language teaching and research.

The third aspect of attention was to ensure the new design was a fast, simple, and user-friendly interface. The imperfect design and inconvenience of use with regard to version 1.0 and version 1.1 of the HSK corpus were also sources of problems. For example, the queried results could be downloaded automatically which resulted in negative user feedback. Users noted that the huge amount of queried results could only be downloaded manually page by page, resulting in sore wrists. Furthermore, users could not adjust the quantity of output results when browsing; they could not communicate with the administrator and give feedback in a timely manner when they encountered problems; they could not make corrections for errors they found in the corpus entry and annotation, so the errors continued to exist and cause problems for other users. The new system also solves these problems and is more user-friendly, allowing users to use the corpus more conveniently and to correct any errors they may find.

## 3 Functional Design

In order to solve the various existing problems in the construction of the corpus and make up for its shortcomings, the functions of the corpus should be improved in terms of corpus retrieval, presentation, data statistics, maintenance, message feedback, automatic download, etc., so that it can better serve Chinese teaching and research.

## *3.1 Search*

The basic way for users to use the corpus is corpus retrieval. From a user's point of view, the value of the corpus lies in the retrieval, presentation, and accessibility of the corpus. What they care about is whether the search function can retrieve the results that they need, and whether it can provide the convenience of collecting and retrieving data for their own teaching and research work.

The search function of a corpus should include the retrieval of specific characters, words, phrases, and sentences; the retrieval of annotated content; the retrieval of special sentence constructions, fixed and semi-fixed structures, compound sentences, and the usage of separated words such as "Li"; collocation searches; and the ability to retrieval data based on parts of speech.[1] The search parameters of a corpus should be constructed based on the characteristics of a language user's nationality, gender, age, composition topics or oral topics, and scores; the search function should help users gain access to the error corpus, the correct corpus, and all the entirety of the corpus.

The search function should be simple, convenient, and easy to use.

### 3.1.1 General Search of Strings

This is the basic retrieval function of the corpus which allows one to search for specific characters, words, phrases, and sentences in the corpus. Generally speaking, this function is available in any corpus. As far as the HSK corpus is concerned, search parameters can be set based on factors such as the candidate's nationality, composition topic, certificate level, test time, and test score.

It should be noted that there are two "composition scores" in the retrieval conditions, which can indicate the selection range of the two scores before and after. For example, the first score is set at 60 and the latter one is set at 80, which means the retrieved corpus results originate from compositions from a score range of 60 to 80.

Below are examples of entries for specific characters, words, phrases, and sentences.

Take "帮" *bāng* as an example for word query (Fig. 3).

The word query takes "帮助" *bāngzhù* "help" as an example (Fig. 4).

Phrase query taking "帮助别人" *bāngzhù biérén* "helping others" as an example (Fig. 5).

Sentence query taking "我们应该帮助别人" *wǒmen yīnggāi bāngzhù biéren* "we should help others" as an example (Fig. 6).

The usage query of the separable word "离" takes "帮忙" and "见面" as examples. A space must be added between the two components of the separable word (e,g, 帮[space]忙, Fig. 7; 见[space]面, Fig. 8) in order to generate relevant corpus results.

---

[1] The role of part-of-speech retrieval is of great significance to the construction of the corpus and the use of the corpus by users. However, it is a pity that the HSK corpus does not realize this function. The "Global Chinese Interlanguage Corpus" fulfills this function.

**Fig. 3** Results produced searching for the character "帮" *bāng*



**Fig. 4** Results produced searching for the character "帮助" *bāngzhù* "help"



**Fig. 5** Results produced searching for the phrase "帮助别人" *bāngzhù bié ren* "help others"



**Fig. 6** Results produced searching for the sentence "我们应该帮助别人" *wǒmen yīnggāi bāngzhù biéren* "we should help others"

### 3.1.2 Sentence and Text Search

The HSK corpus offers an exhaustive collection of composition errors made by foreigners during the composition section of the Advanced Chinese Proficiency Test

**Fig. 7** Results produced searching for the separable words "帮忙"



**Fig. 8** Results produced searching for the separable words "见面"



**Fig. 9** Results produced searching for the "把" Ba sentence

with a focus on five areas: characters, words, sentences, text, and punctuation marks. Among them, errors in characters, words, and punctuations can be retrieved in either character strings or word and vocabulary lists. Errors in sentences and text can be queried using the sentence and text retrieval function.

Errors sentence retrieval takes "把" Ba sentence as an example (Fig. 9).

See the figure below for the error text (Fig. 10).

The above two search methods are available in the 1.0 and 1.1 versions of the corpus. These methods can solve retrieval problems such as when searching in an error corpus.

**Fig. 10** Results produced searching for the error text

### 3.1.3 Advanced Search

In version 2.0 of the HSK corpus, two advanced search features have been added: (1) search parameters for specific conditions, and (2) the ability to search for word collocations. The additions of these features have further enhanced the functionality of the HSK corpus as more results can be retrieved in a single search. In addition, it is also possible to search for the separable word "离" in the new 2.0 version.

The method of generating results based on specific search parameters in the HSK 2.0 corpus is suitable for retrieving specific sentence patterns, semi-fixed structures, and complex sentences with two marker words. The reason for the relatively powerful retrieval capabilities of this type of search is the use of regular expressions. Regular expressions[2] are quite common and general methods for corpus retrieval, yet they are relatively unfamiliar to linguistic professionals with a liberal arts background. The HSK corpus is easy to use and suitable for students of non-STEM majors due to the changes made based on the theoretical and practical background of liberal arts students. These changes include a liberal arts transformation of regular expressions, the simplification of mathematical equations into frame structures, and ensuring successful searching through filling mark words in corresponding positions.

For example, the search of sentences containing the "是……的" construction (Fig. 11) and the "连" construction (Fig. 12).

Fixed structure retrieval with "爱……不……" (Fig. 13) and "一……就……" as examples (Fig. 14).

Take "或者……或者……" *or…or…* as an example for complex sentence retrieval (Fig. 15).

It is important to note that this search method is still based on form. This means that results will be generated as long as there are set search terms in the corpus,

---

[2] Regular expression is a kind of logical formula for string manipulation. It uses some pre-defined specific characters and the combination of these specific characters to form a "rule string". This "rule string" is used to express the pair of characters and is a kind of filtering logic for strings.

**Fig. 11**  Results produced searching for the "是……的" construction



**Fig. 12**  Results produced searching for the "连……也……" construction



**Fig. 13**  Results produced searching for the "爱……不……" construction



**Fig. 14**  Results produced searching for the "一……就……" construction



**Fig. 15**  Search results generated for complex sentence retrieval

**Fig. 16** Collocation retrieval example: "汉语" *Hànyǔ* Chinese left collocation situation

yet these results may not actually be the entire range of linguistic data housed in the corpus. For example, sentences "爱情不是长久不衰的", "对朋友的爱不是对父母的爱也不是对爱人的爱" and the semi-fixed structure "爱……不……" have nothing to do with each other, even though there are "爱" and "不" in the sentence.

The ability to search for word collocations allows one to search for the co-occurring words preceding or following a word and their frequency. In this way, one can find out what words are collocated to the left or right of a certain word, count the corresponding collocation frequency, and sort them in descending order of frequency. This is a significant data retrieval method because it provides usage of words by demonstrating both the frequency and information obtained before and after a word. The generated results are equivalent to the "Word Collocation Dictionary", which can serve as an important reference for Chinese language teaching.

Taking "汉语" *hanyu* "Chinese" as an example, the most frequent collocate on the left side is "学习" *xuexi* "study/learn", with a frequency of 585 (Fig. 16); the second most frequent word is "学" *xue* "study", with a frequency of 523. These are the two most frequent collocations. It can be seen that "学习汉语" *xuexihanyu* "learn Chinese" and "学汉语" *xuehanyu* "study Chinese" are the two collocations that learners use the most and have the best mastery. From the perspective of Chinese language teaching, these collocations are also the most important words to be taught to learners and should be the focus of teaching. The frequency of "对" *dui* "right" on the left is 48, while "觉得"*juede* "think/feel" only appears 9 times. The most frequent collocate on the right side is the auxiliary "的" *de*, with a frequency of 491(Fig. 17); a comma is used more often at the end of sentences on the right with a frequency of 344; the frequency of "有" *you* "have" after "汉语" *hanyu* "Chinese" is 28; and the frequency of "越来越" *yuelaiyue* "more and more" after "汉语" *hanyu* "Chinese" is only 4.

Details can be seen as follows.

**Fig. 17** Collocation Search example: "汉语" Hànyǔ Chinese right collocation situation

## 3.2  Corpus Presentation

In order to provide as much convenience as possible for users' teaching and research, background information regarding the corpus language data and its authors should be presented. In addition, this information should also accompany the search results generated by users when they use the corpus. Furthermore, the retrieved labeled results should conform to the original corpus in terms of form (e.g. composition, audio, video, etc.). In order to adapt to different browsing habits of users, it should be possible for users to adjust the number of results displayed on each page.

The background information of the language data refers to the author's nationality, gender, test time, composition topics, subjective oral and composition test scores, objective listening, reading, and comprehensive expression test scores, their total test results, and obtained certificates. This background information plays an important role in studying and judging learners' acquisition of Chinese. For examples of background information, see the green fonts in Figs. 13 and 14.

The following is a data entry of an original composition (Fig. 18) followed by the annotated full text in the corpus (Fig. 19).
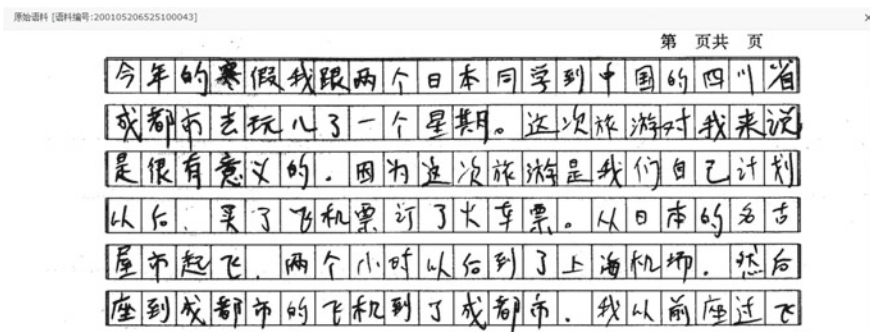


**Fig. 18**  An original composition in the corpus

Fig. 19 Annotated full text in the corpus



Fig. 20 Users can set the number of results displayed per page

Users can set the number of results displayed per page according to their own reading habits (Fig. 20).

## 3.3 Statistical Analysis

All the language data within the HSK corpus has been statistically analyzed. The data produced by the statistical analyses provide an overview of the corpus, including the total number of characters, words, composition topics, and text (Fig. 21); all kinds of error data related to characters, words, sentences, text, and punctuation (Fig. 22); characters and word information sorted by factors such as year, country, and obtained HSK certificate (Fig. 23). The data is especially useful for studying learners' acquisition of Chinese and can serve as an important source of reference for Chinese language teaching.

Examples of statistical charts are listed as follows.

**Fig. 21** An overview of the corpus



**Fig. 22** All error data types related to characters, words, and sentences



**Fig. 23** Characters information sorted by year

## *3.4 Others Functions*

In addition to the above functions, the corpus also provides functions such as crowd-sourcing maintenance, message feedback, personal workspace, automatic down-loading, and adding related resources to further enhance the ability to serve teaching and research.

(1)   Crowdsourcing maintenance

Large-scale interlanguage corpora are usually annotated manually and rely on hundreds of annotators. Inconsistent annotation and labeling, and even errors and omissions, are inevitable. Although quality monitoring is carried out during the process, it still cannot solve the problem completely. According to the concept of crowdsourcing, allowing users to correct errors and omissions in the input version and marked version of a corpus is an effective way to improve the quality of that corpus.

Specific instructions: double-click on the corpus entry to open the dialog box for modification and editing. Edit the entry then click submit and wait for it to update. It will then go through backend screening and once finished will be published and replace the original entry. The screening process is carried out by the corpus administrator who decides whether the modification made by the user should replace the original entry. This is an especially important step in the process which helps avoid inadvertent and incorrect changes made by users. In this way, crowdsourcing can effectively improve the quality of a corpus' transcriptions and annotations so that it can better serve the majority of users.

(2)   Feedback

Users will inevitably encounter various questions that need to be answered in order to better use the corpus. In addition, they may have comments and suggestions that are important for the construction and improvement of the corpus. Thus, a medium for communication and feedback is necessary for effective communication between corpus constructors and users. The approach taken by the HSK corpus was to add a "feedback message" function to facilitate communication, exchange, and discussion between users and the corpus builder (Fig. 24).

From a practical point of view, this function serves a good communicative purpose.

(3)   Personal workspace

Setting up a "personal workspace" in a corpus is a good idea and can have many practical functions. For example, users can maintain their own information, input clerks can extract corpus entries for input and transcription, annotators can label entries, and users can perform corpus-assisted analyses, research, and even write papers. In short, it can be a working platform for builders to build a corpus and for users to conduct related research. At present, the functions of the personal workspace in the HSK corpus are still inadequate and should be enriched.

**Fig. 24** Sample feedback message

(4)  Automatic language data download

In response to the inconvenience reported by users of the corpus in the past, version 2.0 of the HSK corpus has been equipped with an automatic download feature. The generated results can be downloaded automatically (limited to 500). This fast and convenient feature forgoes the labor of manual copying and downloading each entry.

(5)  Increase and accumulation of the relevant resources

The inclusion of practical resources closely related to Chinese language teaching in the corpus can provide users with great convenience in teaching and research. For example, in our study, we compared the characters used in the corpus with the characters used in the "Chinese Proficiency Vocabulary and Chinese Character Level Syllabus". We found that there are 2905 Chinese characters in the syllabus, and 3904 different Chinese characters in the HSK corpus. Learners mastered 999 more Chinese characters than stated in the syllabus. Among the 3905 Chinese characters, there are 2778 characters in the syllabus accounting for 71.16% and 1126 characters outside the syllabus accounting for 28.84%.

The comparison with the "List of Frequently Used Modern Chinese Characters" (1988) for native speakers shows that there are 3500 frequently used characters in the word list, which are divided into 2500 frequently used characters and 1000 sub-frequently used characters. Comparing the 3904 characters in the corpus with the "Frequently Used Characters List", there are a total of 3153 characters in the table, with 2452 frequently used characters accounting for 98.08% of 2500 frequently used characters and 701 sub-frequently used characters accounting for 70.1% of the 1,000 sub-frequently used characters.

Based on these studies and findings, we have compiled "A Comparison List of 2500 Frequently used Characters with HSK by Phonetic Errors", "A Comparison List of 2500 Frequently used Characters with HSK by Total Frequency ", and "A Comparison List of 2500 Frequently used Characters with HSK by Error Frequency", which have been put into the statistical information as an important reference for teaching Chinese characters (Fig. 25).

**Fig. 25** Statistics example graph

## 4   Conclusions

Through the review of the construction of the Chinese interlanguage corpora and the discussion of existing problems, we put forward some new concepts and functions of corpora construction in order to improve the level of corpus construction and better serve teaching and research.

1. The objective and fundamental purpose of building corpora are to serve Chinese language teaching and scientific research all around the world. The premise of ensuring this function is to make sure that the corpora are always open to the public. This requires the corpora systems to be secure without any high-or medium-risk vulnerabilities. This is a new situation and a new problem brought about by the development of new information technology, which must be paid great attention to by corpora builders.

2. Improvements in corpora software systems can enhance the functionality of corpora and can better meet users' needs. The improvement and enrichment of search and retrieval methods to enable users to query some words, phrases, and sentences that were previously unavailable are one such example of functionality enhancement. The rich and practical statistical information has important reference value for teaching and research. A user-friendly interface and the design of certain humanized functions, such as the autonomous setting of the number of corpora presentations and automatic downloading, can provide users with convenience and improve their user experience. The function which allows users to modify the transcriptions and annotations by crowdsourcing allows for the continuous improvement in the quality of corpora annotations.

3. Users have the most say on what kind of functions a corpus should have. Their questions, comments, and suggestions in the process of using a corpus are of great significance to corpus construction and should be understood in time and given feedback as soon as possible. Therefore, it is particularly important for corpus builders to communicate with users and to maintain a smooth and effective medium for communication as provided by the "feedback message" function.

4. In the past, corpora were designed and built in a simple and extensive way, which was in the initial stage of corpora construction, or Corpora Generation 1.0. The development of the HSK Corpus 2.0 has made us realize the important

role of software systems. A good software system can make the corpus powerful, easy to use, and possess "fine and rich" characteristics. It also promotes corpora construction into Generation 2.0. The development and transition from the simple Generation 1.0 to the refined Generation 2.0 reflect the developing progress of Chinese interlanguage corpora construction. This is also an inevitable result of the technological progress of the times.

There are some important characteristic differences between Generation 1.0 to Generation 2.0. as noted below:

Corpora labeling: individual layer labeling → comprehensive labeling
Labeling mode: error labeling → error labeling + basic labeling
Search method: simple search → advanced search
Construction concept subcontracting → crowdsourcing
Research Paradigm: Error Analysis → Comprehensive Investigation of Interlanguage
Data view: individual data → big data

It can be said that with HSK corpus version 2.0, the construction of Chinese interlanguage corpora has entered Generation 2.0 from Generation 1.0, with 2018 being regarded as the first year of Generation 2.0.

# References

Chen, X. (1996). *Introduction to* "*Chinese interlanguage corpus system*"[C]. *Selected Papers of the Fifth International Chinese Language Teaching Conference*, Beijing: Peking University Press, pp. 459–467.

Chu, C., & Xiaohe, C. (1993). The basic idea of establishing the "Chinese interlanguage corpus system." *World Chinese Teaching, 3*, 199–205.

Cui, X., & Baolin, Z. (2011). "Global Chinese language learners corpus" construction plan. *Language Application, 2*, 100–108.

Xiao, X. et al. (2009). *Research on the difficulty and classification of Chinese sentence learning for foreign students.* Beijing: Higher Education Press.

Xiao, X., & Wenhua, Z. (2014). The comprehensiveness and classification of Chinese interlanguage corpus labeling. *World Chinese Teaching, 3*, 368–377.

Zhang, B. (2019). From 1.0 to 2.0_construction and development of chinese interlanguage corpus. *International Chinese Language Teaching Research* (4), 84–95.

Zhang, B. et al. (2019). *A Study on the Standardization of Chinese Interlanguage Corpus Annotation.* Beijing: Peking University Press.

Zhang, B., & Xiliang, C. (2015). On the construction standards of Chinese interlanguage corpus. *Language Application, 2*, 125–134.

Zhang, B., & Xiliang, C. (2018). New thoughts on the research of chinese interlanguage corpus annotation standards—also on the design of the "global Chinese interlanguage corpus, annotation standards[C]. *Selected Papers of the Third International Symposium on the Construction and Application of Chinese Interlanguage Corpus,* Beijing: World Book Publishing Company.

Zhang, B. et al. (2014). *A corpus-based study of Chinese sentence acquisition for foreigners*, Beijing: China Book Publishing House.

Zhang, B. et al. (2008). *Thematic research on Chinese vocabulary based on interlanguage corpus.* Beijing: Peking University Press.

Zhao, J. et al. (2008). *A study of Chinese syntax based on interlanguage corpus.* Beijing: Peking University Press.