

*Sergey Zinin\**

**Keyword analysis of Chinese classics:  
“Thirteen Canons” and Zhuangzi**

**ABSTRACT:** This study applies keyword analysis to the Warring States Workshop (WSW) Ctexts corpus of classical Chinese texts. The WSW Ctexts corpus presents a good opportunity to extract a keyword character list of a limited number of linguistically and semantically related texts in a mid-size corpus. The initial goal of the author in this study is the production of the keyword lists, which could be further used for extraction of semantic information about texts. The contents of these lists depend on the type of scoring measure, “the keyness.” The lists could be combined into a synoptic table, arranged by a special set of thematic categories. This table could serve as a kind of a “semantic map” of the texts. Analysis of this synoptic table should help in estimating which scoring methods (log-likelihood (LL), chi-square (CHI), effect size, or pre-h-point (PHP) content characters) are more productive for content investigation of Chinese classics. The article demonstrates that, although some methods extract characters that describe topics, other methods could be better for selecting characters, characterizing stylistics. In addition, the synoptic table allows conducting some preliminary genre clustering of the texts.

**Content**

1. Introduction
  - 1.1. Keyword analysis
  - 1.2. Keywords and Chinese classics
  - 1.3. Research corpus
2. Keywords in Chinese classics
  - 2.1. Keyword character list compilation

---

\* Zinin Sergey, Warring States Project, University of Massachusetts, Amherst, USA; E-mail: [szinin@research.umass.edu](mailto:szinin@research.umass.edu)

2.2. Analysis of the top list of keyword characters and comparison with PHP content word lists

3. Key-keywords in Chinese classics

4. Comparison of outputs of keyword analysis methods

5. Discussion of results

6. Conclusion

Appendix

References

## 1. Introduction

### 1.1. Keyword analysis

Keyword analysis of text contents in corpus studies has been growing since the beginning of the 1990s. In this approach, keywords are defined as “word forms that occur in a text more frequently than expected by chance alone and are often closely connected to the overarching themes and genre of a text or set of texts” [Crvcek and Fidler, “Not all keywords are created equal,” 55].<sup>1</sup> They could be identified by “comparing the frequencies of words in a corpus with frequencies of those words in a (usually larger) reference corpus” [Baron et al., “Word Frequency,” 1]. It allows compilation of an ordered keyword list and discovering significant words. This list is “a useful tool for directing researchers to significant lexical differences between texts ... keywords can reveal a great deal about frequencies in texts which is unlikely to be matched by researcher intuition” [Baker, “The question is,” 8]. Keyword analysis allows studying such areas as thematic, stylistic, and linguistic analysis.<sup>2</sup>

Keyword analysis takes roots initially in theories of Firth and then Williams’s research of culturally significant words.<sup>3</sup> Word frequency lists

---

<sup>1</sup> “A standard technique in corpus linguistics and corpus-assisted discourse studies consists in automatically extracting keywords from corpora, to identify the content words that stand out in terms of frequency and keyness” [Gaspari and Venuti, “A golden keyword,” 131].

<sup>2</sup> “Keywords ought to identify words or phrases that may be specific to individual topics or questions asked, but they may also be revealing of authorial style, and if used by multiple authors in the corpus, could identify language features associated with a particular regional register” [Baker et al., “Triangulating Methodological Approaches,” 40], following earlier Scott’s “Key-words provide a useful way to characterize a text or a genre. Potential applications include: language teaching, forensic linguistics, stylistics, content analysis, text retrieval” [Scott, *Wordsmith Tools Version 3*, 70].

<sup>3</sup> Williams called these words “Keywords in two connected senses: they are significant, binding words in certain activities and their interpretation; they are significant, indicative words in certain forms of thought” [Williams, *A Vocabulary of*

have been used for content analysis for a long time, but using keyword lists is a comparatively new method, developed with advancement of personal computers at the beginning of the 1990s. The pioneers of this method often referred to the concept of “culture keywords” as a predecessor of their analysis [Scott, “PC Analysis of Key Words,” 233], and it is still treated as such by Stubbs [Stubbs, *Words and Phrases*, 145–194]. The term “keyword” had already took on several meanings by the time of development of keyword analysis [Bondi and Scott (eds.), *Keyness in Text*], but in the latter method, a vaguely defined “cultural significance” had been replaced by statistical significance. Therefore, in this context, “keywords are those whose frequency (or infrequency) in a text or corpus is statistically significant, when compared to the standards set by a reference corpus” [Bondi, “Perspectives on keywords and keyness,” 3].

In keyword analysis, the significance of a word as a potential keyword is measured by its “keyness score.” There are various keyness measures to implement keyword extraction and ranking. There are a few software packages in the market, which will allow extracting keywords automatically.<sup>4</sup> In this article, the author will consider three popular measures: log-likelihood (LL), chi-square (CHI), and recently introduced %DIFF relative frequency measure.<sup>5</sup>

Depending on the measure of keyness, keyword scores may be positive and negative.<sup>6</sup> Positive keywords could be roughly defined as “comparatively overused” words, comparing with word use in the reference corpus, and negative keywords will be “comparatively underused.” The specific measures of these properties will be described below.

There are more complex concepts built around keywords, most of which were developed by Mike Scott. One of them is the concept of key-keywords, which are keywords shared by a few texts of a corpus: a “key-

---

*Culture and Society*, 15]. As Richardson surmises, this concept is rooted in ideas of Valentin Voloshinov [Richardson, “Keywords Revisited”, 101]; Richardson’s interpretation has affected development of keyword analysis in the late 1990s. As Mike Scott notes in [Scott, *WordSmith Tools Version 5*, 165] “The term ‘key word,’ though it is in common use, is not defined in Linguistics,” as it has cultural theory roots.

<sup>4</sup> The most popular is Mike Scott’s WordSmith, which has served as the main force in popularizing the concept of keyword [Scott, *WordSmith Tools Version 3*].

<sup>5</sup> They are explained below.

<sup>6</sup> By Scott’s definition [Scott, *WordSmith Tools Version 3*, 71], “A word which is positively key occurs more often than would be expected by chance in comparison with the reference corpus. A word which is negatively key occurs less often than would be expected by chance in comparison with the reference corpus”.

keyword” is one which is “key” in more than one of a number of related texts. The more texts it is “key” in, the more “key key” it is. This will depend a lot on the topic homogeneity of the corpus being investigated.<sup>7</sup> Other related concepts are “associates” and “clumps.” An “associate” of a key-keyword X is “another keyword (Y) which co-occurs with X in a number of texts” [Scott, *Wordsmith Tools Version 5*, 158]. In addition, “clumps” is the name given to groups of keywords associated with a key-keyword [Scott, *Wordsmith Tools Version 5*, 161].<sup>8</sup>

The main reason for analyzing texts using keywords is the presumption that they express “aboutness,”<sup>9</sup> i.e., they allow understanding of text content, based on automatic extraction of frequent words. The key problem there is to select words that are relevant to the content of the text, and this is what the concept of “keyness” does. However, simply creating a list of these words is just the beginning. As Marina Bondi suggests [Bondi, “Perspectives on keywords,” 3], “identifying elements that are repeated to a statistically significant extent does not in itself constitute an analysis or an interpretation of the text or corpus. It does however point to elements that may be profitably studied and need to be explained.” Baker indicates one of the ways of further analysis—creating concordances

---

<sup>7</sup> “In a corpus of City news texts, items like bank, profit, companies are key key-words, while computer will not be, though computer might be a key word in a few City news stories about IBM or Microsoft share dealings” [Scott, *Wordsmith Tools Version 5*, 166].

<sup>8</sup> “The idea here is to refine associates by grouping together words which are found as key in the same sub-sets of text files. The example used to explain associates will help. Suppose the word wine is a key key-word in a set of texts, such as the weekend sections of newspaper articles. Some of these articles discuss different wines and their flavors, others concern cooking and refer to using wine in stews or sauces, and others discuss the prices of wine in a context of agriculture and diseases affecting vineyards. In this case, the associates of wine would be items like Chardonnay, Chile, sauce, fruit, infected, soil, etc.” [Scott, *Wordsmith Tools Version 5*, 161].

<sup>9</sup> This concept was popularized in the 1970s in the information sciences and content analysis by William Hutchins, and Scott related it to keywords, earliest in his *Wordsmith* manual of 1998 [Scott, *Wordsmith Tools Version 3*] and in the article of 2000 [Scott, “Picturing the Key Words,” 44]. Afterwards, it found popularity in corpus linguistics community. Hutchins discerned between document’s topic, summarization, and aboutness, where “concept of ‘aboutness’ ... associates the subject of a document not with some ‘summary’ of its total content but with the ‘presupposed knowledge’ of its text” [Hutchins, “The Concept of ‘Aboutness’”, 180].

based on keywords [Baker, “The question is,” 3], which is a future work for this author.<sup>10</sup>

It should be noted that keyword analysis is not a unique method of identifying text content by creating a list of significant words. Computational linguistics has been developing many other methods, e.g., topic extraction and topic signature (see [Kao and Poteet, *Natural Language Processing*]). Some of the methods implemented there also use LL and reference corpora, similar to keyword analysis.<sup>11</sup>

This author also investigated the content character list of this corpus, using “pre-h-point” (PHP) list method [Zinin, “Analysis of character-frequency lists”], and its results will be compared with the results of the keyword analysis.<sup>12</sup>

## 1.2. Keywords and Chinese classics

There are many studies on frequency of characters in classical Chinese texts, and a few of them analyze text content and genre attribution based on lexis.<sup>13</sup> However, although keyword analysis methods have been around for more than 20 years, this author did not find any research implementing keyword analysis to classical Chinese texts; therefore, this article may be one of the first studies.

Any frequency research on lexis of classical Chinese texts should choose the lexical unit of study. Up to the end of the 20<sup>th</sup> century, in Chinese and Western studies, this unit had been character. Increasingly, though, researchers have attempted to introduce words (monosyllabic and polysyllabic) as the lexical units. This approach has been working comparatively well for modern texts, but it has been hindered for classical texts by lack of modern authoritative and publicly available classic corpora with word segmentation. The few teams that implement word segmentation usually do not share their corpora resources. Calculating frequency

---

<sup>10</sup> Examples of such analysis could be found in articles of Karen Donnelly ([Donnelly, “Risk, chance, hope,”] and [Donnelly, “Dr. Condescending”]).

<sup>11</sup> However, the researchers in the keyword analysis area rarely compare these approaches with their methods.

<sup>12</sup> H-point divides the frequency list into two parts: PHP and post-h-point. For a frequency list, the function  $f(r)$  is introduced, where  $f(r)$  is the frequency for rank  $r$ . The “h-point can be defined as that point at which the straight line between two (usually) neighboring ranked frequencies intersects the  $y = x$  line,” i.e., where  $r = f(r)$ . The autosemantic words above h-point (PHP) tend to be “content words” [Popescu et al., *Aspects*, 24].

<sup>13</sup> The author omits here a review of this literature, as it is not immediately relevant to the subject of keyword analysis. These works were reviewed in another article by the author [Zinin, “Vocabulary Richness”].

of characters, instead of words, remains the mainstream method in frequency research on Chinese classics. Most of the available articles in this area are character-frequency studies, not word frequency studies. This study belongs to the former category. It creates frequency lists using characters, not words. This author supports an opinion that, when it comes to the analysis of the contents and genre, using character as the unit of study should not considerably affect the results.

### 1.3. Research corpus

The corpus in this research is the Warring States Workshop (WSW) Ctexts open-source corpus, introduced in an earlier work of the author (see [Zinin, “Pre-Qin Digital Classics”])<sup>14</sup>. Despite the existence of a few available digital academic corpora, there is a lack of publicly available research open-source corpora in classical Chinese. Texts in digital academic corpora vary in textological decisions, which complicates comparisons of results. This is the reason why an open-source corpus, compiled on the basis of Creative Commons licensed digital texts, was used for this research.<sup>15</sup> It offers reproducibility of experiments, at the cost of a few textological flaws. The author believes that, because of the statistical nature of the research, potential text issues should not affect considerably its results.

This study applies keyword analysis methods to investigate the WSW corpus and analyzes the results. The WSW Ctexts corpus presents almost an ideal opportunity to analyze keyword lists of a limited number of linguistically and semantically close texts in a mid-size corpus. Each of the 14 texts, ranging in size from the Xiao Jing to the Zuo Zhuan, could be analyzed against the rest of the corpus<sup>16</sup>.

The WSW corpus is very well suited for keyword analysis research. Researchers who studied the relationship between the number and size of texts in reference corpora and keywords state that “the answer to the question ‘what is the ideal size of a reference corpus’ is five,” and that “a

---

<sup>14</sup> This author expresses gratitude to Bruce E. Brooks for continuing support of the project.

<sup>15</sup> This resource, created by the author, is based on Creative Commons digital versions of texts of Wikisource. It contains 12 texts of the Thirteen Classics (the Chun Qiu, the Gongyang Zhuan, the Guliang Zhuan, the Li Ji, the Lun Yu, the Mengzi, the Shi Jing, the Shu Jing, the Xiao Jing, the Yi Li, the Zhou Yi, and the Zhou Li). It does not include the Er Ya, but it includes the Zhuangzi (added as a balancing text). It also allows treating the Chun Qiu and the Zuo Zhuan as separate texts (instead of combined the Chun Qiu Zuo Zhuan), so there are 14 texts in total.

<sup>16</sup> The corpus includes also two small texts (the Guo Dian and the Mao Shi) which are not analyzed in this article.

reference corpus does not need to be more than five times larger than the study corpus” [Berber-Sardinha, “Comparing Corpora,” 12]. The WSW corpus complies with all of these criteria.<sup>17</sup> Looking at the role of genre variety, Goh, who analyzed the relationship between keywords and genre composition of reference corpora (RC), states, “genre difference of spoken and written RCs is not an important factor in keyword calculation” [Goh, “Choosing a reference corpus,” 249]. Therefore, even if WSW texts varied in genre, it should not affect the results. Considerable diachronic differences between texts could have affected the results [Goh, “Choosing a reference corpus,” 254], but it is not the case for the WSW corpus.

Some of the texts are comparatively large (e.g., the Zuo Zhuan and the Li Ji), but according to researchers who analyzed the effects of corpus and text sizes on keyword analysis, these could be considered within the norm.

The author applied previous statistical methods for content analysis of the corpus. In a previous study, the concept of the h-point was applied to identify most significant characters, describing text content. In that study [Zinin, “Analysis of character-frequency lists”], the significant content words were identified and then distributed over categories. These categories have been designed to make historical genre texts prominent.

This study is mostly of explorative character. The main goal of the author in this study is to identify keyword lists that could be used to extract semantic information about texts. The content of these lists depends on methods of scoring “keyness.” These lists will be combined into a synoptic table, arranged with a special set of thematic categories. This table could serve as a kind of a “semantic map” of the texts. Analysis of this synoptic table should help in identifying which methods are more successful. The article will not only identify LL and %DIFF keywords (and key-keywords) but also compare them with character lists from the h-point method. Better methods of keyness scoring could be applied to analyze other texts.

Many of the texts of the WSP Ctexts corpus are collections of smaller works, and, in reality, themselves are mini-corpora. But even smaller texts that have been assigned an “authorial figure” to them (like the Meng-zi or the Lun Yu), are not products of an “identifiable original author” [Boltz, “Why So Many Laozi-s?”, 10], as was well-known to the Chinese philological tradition. The recent discovery of manuscripts of Warring States and Han period instigated new analysis of textological structure of received canons. Many classics should be viewed as dynamic, flexible collections of smaller fragments. Two concepts should be noted

---

<sup>17</sup> Only one text, the Zuo Zhuan, is about five times smaller than the entire corpus.

as indicating productive research directions. The first is the idea of a “growth text” by Bruce Brooks [Brooks, “Before and After Matthew”, 1] and the second is the idea of “building blocks” by William Boltz [Boltz, “Why So Many Laozi-s?”, 12]. Brooks utilizes stylistic analysis to trace trajectory of text building, and Boltz stresses the role of “thematically based selection” [Boltz, “Why So Many Laozi-s?”, 12] in assembling texts from “building blocks”. Keyword analysis may provide useful information for any of these approaches.

This study is striving to be an open-source, reproducible, and verifiable effort. All source data should be available for download; in addition, the accompanying site contains all extended original results, which would not fit into the article’s limited space (see Appendices for links). The author is grateful to Bruce E. Brooks for continuing support of the project.

## 2. Keywords in Chinese classics

Keyword lists are created by assigning significance values (“keyness scores”) to words or characters and then selecting the most significant ones. Mostly, keyness scores are statistically significant values obtained by calculating LL and CHI scores. Some researchers (Gabrielatos and Marchi, “Keyness”) consider the statistical significance approach not the most effective measure and offer their own %DIFF measure, based on “effect size” approach. In this study, the author computed the values of all three measures (LL, CHI, and %DIFF) and compared the results.

### 2.1. Keyword character list compilation

LL keywords for the WSW corpus were calculated using formula<sup>18</sup>

$$-2 \ln \lambda = 2 \sum O_i$$

In case of a two-by-two matrix,

$$G2 = 2 * ((a * \ln(a/E1)) + (b * \ln(b/E2)))^{19}.$$

The list of keyword characters, produced by this method, is very large and needs filtering to obtain meaningful results. It is well known

<sup>18</sup> See [Rayson and Garside, “Comparing Corpora,” 3], and UCREL <http://ucrel.lancs.ac.uk/llwizard.html>.

<sup>19</sup> Where “a” is the frequency of a character in Corpus 1 (Reference corpus), and “b” is the frequency of a character in Corpus 2 (corpus under testing); “c” is the number of characters in Corpus 1 and “d” is the number of characters in Corpus 2. “a+b” will be the total number of a character in both corpora, and “c+d” is the number of all characters in both corpora. In these terms, expected values E1 (for Corpus 1) and E2 (for Corpus 2) will be  $E1 = c * (a+b) / (c+d)$  and  $E2 = d * (a+b) / (c+d)$ . See UCREL <http://ucrel.lancs.ac.uk/llwizard.html>.



that most rare words and characters go to the top of lists created using statistical significance. Most researchers recommend discarding of rare characters, often starting with minimum three counts.<sup>20</sup> In this study, owing to the size of the texts, only characters with frequency count four or higher were included in the resulting character lists. Another cutoff criterion is the p-value: it should be low enough. The value of  $p = 0.05$  is standard, which corresponds to an LL critical value of 3.84. In this study, the minimum of the LL values is 11.0<sup>21</sup>, so with one degree of freedom (d.f. = 1), p-value should be less than 0.001.<sup>22</sup> This is a more rigorous criterion than a p-value of 0.05. Even if these criteria are quite rigorous, for some large texts, the list of keyword characters could reach hundreds of entries. Finally, in keyword analysis, function words are usually removed from a keyword list [Archer, “Does frequency really matter?” 2–3]. Therefore, function (or “empty”) characters were removed from the resulting keyword character lists, with some exceptions.<sup>23</sup> The full list of LL keywords could be downloaded from the project’s Github site.<sup>24</sup>

CHI keywords were calculated using the contingency table and the formula from [Baron et al. Word Frequency, 44], see Table 1<sup>25</sup>.

Table 1. Contingency table for the CHI test

	Corpus 1	Corpus 2	Total
Frequency of feature	a	b	a+b
Frequency of feature not occurring	c	d	c+d
Total	a+c	b+d	N=a+b+c+d

<sup>20</sup> For example, Scott recommends three counts as default, but it depends on text’s size: “The default setting of 3 mentions as a minimum helps reduce spurious hits here. In the case of short texts, less than 600 words long, a minimum of 2 will automatically be used” [Scott, *WordSmith Tools Version 5*, 177].

<sup>21</sup> The critical value for 0.001 is 10.83 (DOI: UCREL <http://ucrel.lancs.ac.uk/llwizard.htm>).

<sup>22</sup> “Any word with a LL greater than or equal to 6.63 ( $p < 0.01$  for 1 d.f.) was considered key, and any word with a frequency less than 5 in either the Innsbruck Letter Corpus (before or after standardization) or the BNC sample was removed from the key word list,” [Baron et al, “Word frequency,” 55–56]. There was no intention in this study to look for very low p-values.

<sup>23</sup> Function words do not go into keyword lists, but could go if their usage is “function words can occur in a keyword list, if their usage is strikingly different from the norm established by the reference text” [Archer, “Does frequency really matter?” 3].

<sup>24</sup> See the Excel spreadsheet “chinese\_classics\_keywords\_log\_likelihood\_short”. See Appendix for the file structure’s explanation.

<sup>25</sup> Usage of terms a, b, c and d is different there than in case of LL calculations.

The CHI statistics ( $X^2$ ) will be calculated as follows:

$$X^2 = N(ad-bc)^2/(a+b)(c+d)(a+c)(b+d)$$

The criteria for frequency counts and p-value for this list were the same as for LL (11 and 4), function characters were also removed.<sup>26</sup>

A comparison of lists (see Table 2) shows that, at those criteria, there are more CHI keyword characters than LL keyword characters. However, the contents of the two lists are very close, with mostly the same characters but in a slightly different order. A more detailed analysis shows that, indeed, LL keyword list (with a few exclusions) is a subset of CHI keyword list. Therefore, CHI data will not be discussed in this article, as LL keyword list is representative enough well for statistical significance methods to be compared further with %DIFF keyword characters.

Finally, %DIFF value of keyness had been calculated by the formula (Gabrielatos and Marchi, Keyness: Appropriate metrics, 12):

$$((NF \text{ in SC} - NF \text{ in RC}) \times 100) / (NF \text{ in RC})^{27},$$

where NF = normalized frequency, SC = study corpus, and RC = reference corpus

As with LL and Chi lists, only characters with an absolute frequency of more than three are included.<sup>28</sup> As to the %DIFF value cutoff, it was also taken as 11, simply for conformity with LL and Chi numbers.<sup>29</sup> Table 2 contains comparative numbers for numbers of keywords retrieved by these methods.

---

<sup>26</sup> The full list of CHI keywords is contained in “chinese\_classics\_keywords\_CHI\_short” Excel spreadsheet and could be downloaded from the Github site. See Appendix for the file structure’s explanation.

<sup>27</sup> 0.0000001 was added to avoid division by zero.

<sup>28</sup> The full list of %DIFF keywords is contained in “chinese\_classics\_keywords\_diff\_short” Excel spreadsheet and could be downloaded from the Github site. See Appendix for the file structure’s explanation.

<sup>29</sup> The meaning of %DIFF is different from that of statistical significance, but, according to the author’s estimation, depending on text size, a similar LL value will be a bit higher than 10. There will be generally almost twice more %DIFF values than LL or Chi.

Table 2. Comparative list of keyword characters' numbers. V is the number of types in a text, N is the number of tokens, LL column contains the numbers of LL keyword characters, %DIFF is the number of %DIFF keyword characters, and CHI is the number of CHI keyword characters.

<b>Text<sup>30</sup></b>	<b>V</b>	<b>N</b>	<b>LL</b>	<b>%DIFF</b>	<b>CHI</b>
XJ	374	1800	33	73	40
ZY	1035	13354	165	299	190
LY	1365	15928	97	256	108
CQ	941	16791	146	261	159
SHU	1911	24539	262	523	293
SHI	2833	29622	462	838	537
MZ	1895	35372	168	480	198
GL	1595	40836	163	387	179
GY	1641	44232	190	438	219
ZL	2215	49462	360	684	395
YL	1551	53929	320	493	341
ZHZ	2968	65251	352	771	400
LJ	3050	99007	325	918	358
ZZ	3236	178564	437	1042	447

Charts on Figure 1 and Figure 2 display the relationships between the numbers of keywords and lengths of texts and vocabulary size, respectively. Owing to the differences in the type of measures, it is not possible to compare the numbers of keywords directly, but these charts display a good visual correlation between ratios of keywords and text characteristics.

---

<sup>30</sup> See the list of abbreviations in the Resources section.

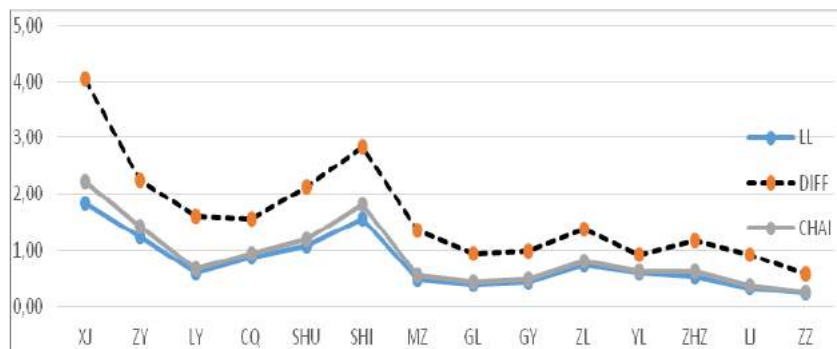


Figure 1. LL, %DIFF, and CHI keyword numbers/text length ratios, ordered by text length

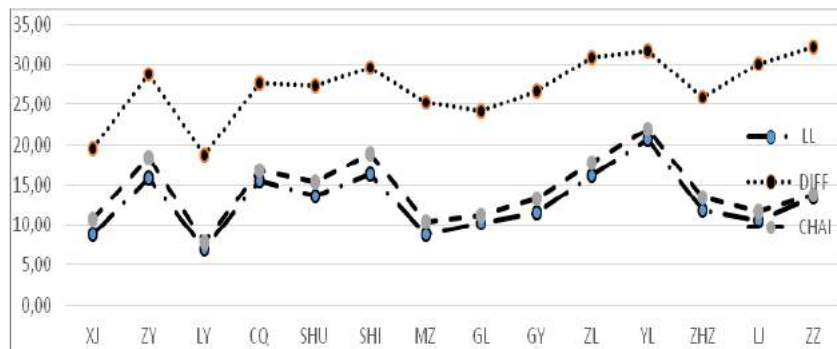


Figure 2. LL, %DIFF, and CHI keyword numbers/vocabulary ratios, ordered by text length

The ratios of all three scores correlate well. We can see that the ratio for keyword number/text length is generally falling with length increase, and this for keyword number/vocabulary probably depends on other text characteristics.

## 2.2. Analysis of the top list of keyword characters and comparison with PHP content word lists

In keyword analysis, it is often hard to handle all words/characters, even if strict criteria were imposed. For topic and genre analysis, most researchers recommend to select only the top part of keyword list, from

50 to 100 words/characters. The sizes of sample numbers are arbitrary and do not usually reflect text characteristics.<sup>31</sup>

This author earlier implemented another approach to the topic analysis, based on the concept of h-point, introduced by Popescu.<sup>32</sup> H-point should reflect texts' inner characteristics, and PHP lists of content word/characters should reflect texts' content. It may be beneficial to compare all three methods of selecting content characters based on their frequency: PHP content characters and top lists of LL and %DIFF scores. Although the h-point approach provides a "native" number of top characters, the LL and %DIFF methods lack such measure. The author decided to take 50 or 100 top characters from LL and %DIFF lists, depending on which number is closer to the corresponding number of PHP characters for this text. See Table 3 for the specific numbers.

Table 3. Matching LL and PHP keyword characters

#	1	2	3	4	5	7	8
Text	CQ	ZZ	GY	GL	SHI	SHU	ZY
Matching positive in all keywords	45	69	45	43	16	22	25
Total	49	125	68	63	37	37	37
% matching to all KW	91,84	55,2	66,18	68,25	43,24	59,46	67,57
Number of pre-h-point characters	60	182	94	91	66	67	51
Selected top KW	50	100	50	50	50	50	50
Matching positive in top keywords	36	48	35	35	9	14	14
Total	49	125	68	63	37	37	37
% matching to top KW	73,47	38,4	51,47	55,56	24,32	37,84	37,84
#	9	10	11	13	14	15	16
Text	ZL	YL	XJ	LJ	MZ	LY	ZHZ
Matching positive in all keywords	50	63	5	45	23	17	28
Total	74	81	7	91	43	29	56
% matching to all KW	67,57	77,78	71,43	49,45	53,49	58,62	50
Number of pre-h-point characters	96	114	20	136	78	50	101
Selected top KW	100	100	all	100	50	50	100
Matching positive in top keywords	44	54	5	30	16	16	23
Total	74	81	7	91	43	29	56
% matching to top KW	59,46	66,67	71,43	32,97	37,21	55,17	41,07

<sup>31</sup> "The vast majority of studies do not examine all keywords, but the top X (usually the top 100)" [Gabrielatos and Marchi, "Keyness," 3].

<sup>32</sup> See: Zinin, "Analysis of character frequency lists".

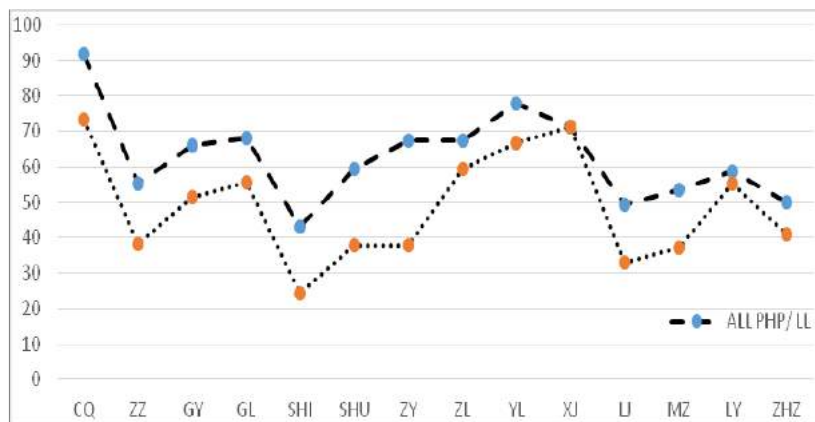


Figure 3. Ratio of PHP content characters, matching all and top LL keyword characters

Do PHP content characters match LL characters? Table 3 displays this information. For almost all texts, PHP content characters are on the top 100 most frequent list, as PHP characters are very frequent in text, and LL characters do not need to be very frequent. The study shows that, with the exception of the Shi Jing, more than 50% of the PHP characters belong to the LL lists. It means that the LL lists contain very frequent non-function characters, too. However, if the top 50 or 100 LL keyword characters are selected, and compared to the PHP characters, there is less significant overlap (from 30% to 50%) between the lists of PHP characters and LL characters. The top lists of LL keyword characters include not too many very frequent characters; otherwise, they would be overlapping more with the PHP lists.

Table 4. Matching %DIFF and PHP keyword characters

	1	2	3	4	5	7	8
	CQ	ZZ	GY	GL	SHI	SHU	ZY
DIFF							
matching	45	71	52	46	23	24	26
LL matching DIFF	91,84	56,8	76,47	73,02	62,16	64,86	70,27
REDUCED DIFF	12	1	5	4	1	3	8
	24,49	0,8	7,35	6,35	2,7	8,11	21,62
	9	10	11	13	14	15	16
	ZL	YL	XJ	LJ	MZ	LY	ZHZ
DIFF							
matching	37	64	6	52	25	20	32
LL matching DIFF	50	79,01	85,71	57,14	58,14	68,97	57,14
REDUCED DIFF	11	19	5	1	2	4	1
	14,86	23,46	71,43	1,1	4,65	13,79	1,79

Analysis of the overlapping of %DIFF list and PHP characters shows even a higher percentage of PHP characters on the complete %DIFF lists than on the complete LL lists. It should be caused by the size of the %DIFF lists. Analysis of the top part of the %DIFF lists shows that they do not practically overlap with the PHP characters, and the reason is the lower (than that for LL) frequencies of those top characters on the %DIFF lists.<sup>33</sup>

### 3. Key-keywords in Chinese classics

The notion of key-keyword is a derivation of the keyword concept. Key-keywords are (Scott, “PC Analysis of Key Words”, 237) “words which are key in a large number of texts of a given type.” There could be positive and negative key-keywords, as well. Table 5 presents the top part of the list of positive key-keywords for the WSW corpus, ordered by the number of texts to which they belong.<sup>34</sup>

Table 5. Most frequent key-keywords in the WSW Ctexts corpus

Character	Number of texts	Texts
天	7	shu, zhouyi, xiaojing, guodian, lijì, mengzi, Zhuangzi
民	7	shi mao, shu, zhouli, xiaojing, guodian, lijì, mengzi
道	7	shi mao, zhouyi, guodian, lijì, mengzi, lunyu, Zhuangzi
子	6	chunqiu, zuozhuan, gongyang, lijì, mengzi, lunyu
故	6	zuozhuan, shi mao, xiaojing, guodian, lijì, Zhuangzi
無	6	zuozhuan, shi, shu, zhouyi, mengzi, lunyu
萬	6	shi, shu, zhouli, guodian, mengzi, Zhuangzi
行	6	zhouyi, xiaojing, lijì, mengzi, lunyu, Zhuangzi
三	5	chunqiu, zhouyi, zhouli, yili, lijì
來	5	chunqiu, gongyang, guliang, shi, zhouyi
公	5	chunqiu, zuozhuan, gongyang, guliang, shi mao
古	5	shi mao, shu, lijì, mengzi, Zhuangzi
國	5	zuozhuan, gongyang, guliang, shi mao, zhouli

<sup>33</sup> For example, for the Chun Qiu, the h-point is 60; that is, all PHP characters have absolute frequencies more than 60. However, on the %DIFF list for the Chun Qiu, only 12 characters have a frequency equal or higher than 60. This situation is even worse for matching in other texts.

<sup>34</sup> It includes key-keywords found in at least seven texts (including the Mao Shi and the Guodian, which are otherwise not considered in this study). The full list of positive and negative keywords (appearing at least from three texts) could be found at the online reference Excel spreadsheet “chinese\_classics\_keywords\_log\_likelihood\_short” (See Appendix for Resources).

圍	5	chunqiu, zuozhuan, gongyang, guliang, zhouli
地	5	zhouyi, zhouli, guodian, liji, Zhuangzi
姜	5	chunqiu, zuozhuan, gongyang, guliang, shi_mao
孔	5	shi, liji, mengzi, lunyu, Zhuangzi
孰	5	gongyang, guodian, mengzi, lunyu, Zhuangzi
己	5	chunqiu, gongyang, guliang, lunyu, Zhuangzi
思	5	shi, shi_mao, xiaojing, mengzi, lunyu
怨	5	zuozhuan, shi_mao, shu, mengzi, lunyu
政	5	zuozhuan, shi_mao, shu, zhouli, lunyu
敗	5	chunqiu, zuozhuan, gongyang, guliang, guodian
方	5	shi, shu, zhouli, liji, Zhuangzi
明	5	shu, zhouyi, xiaojing, liji, Zhuangzi
曰	5	zuozhuan, shu, zhouyi, mengzi, lunyu
正	5	chunqiu, gongyang, guliang, zhouyi, yili
歸	5	chunqiu, zuozhuan, gongyang, guliang, shi
父	5	chunqiu, gongyang, guliang, xiaojing, liji
王	5	chunqiu, zuozhuan, shi_mao, shu, mengzi
百	5	shi, shu, zhouli, liji, mengzi
義	5	shi_mao, xiaojing, liji, mengzi, Zhuangzi
詩	5	zuozhuan, shi_mao, xiaojing, liji, mengzi
變	5	guliang, shi_mao, zhouyi, liji, Zhuangzi
足	5	guodian, liji, mengzi, lunyu, Zhuangzi
身	5	xiaojing, guodian, liji, mengzi, Zhuangzi
邦	5	shi, shu, zhouli, guodian, lunyu
雨	5	chunqiu, gongyang, guliang, shi, zhouyi

If many characters are “overused” in several texts together, relative to the reference corpus, there could be some semantic affinity between these texts. Depending on the share of these texts in corpus, positive key-keyword, participating in many texts, may be “underused” in other texts, i.e., be a negative key-keyword for those texts. If some texts have affinities based on instances of positive and negative key-keywords, it could be that they fall together in semantic sense.

For example (see Table 6), characters *tian*, *min*, and *dao* are positive key-keywords in philosophical and fiction texts and negative in historical texts, and characters *zi* and *gong* are positive key-keywords in historical texts and negative key-keywords in philosophical and fiction texts.



Table 6. Text distribution of several key-keywords

Character	Positive	Negative
天	shu, zhouyi, xiaojing, liji, mengzi, Zhuangzi	chunqiu, zuozhuan, gongyang, guliang, zhouli, yili
民	shu, zhouli, xiaojing, liji, mengzi	zuozhuan, gongyang, guliang, Zhuangzi
道	zhouyi, liji, mengzi, lunyu, Zhuangzi	zuozhuan, gongyang, shi, zhouli, yili
子	chunqiu, zuozhuan, gongyang, liji, mengzi, lunyu	shi, shu, zhouyi, zhouli, yili, Zhuangzi
公	chunqiu, zuozhuan, gongyang, guliang	shi, shu, zhouyi, zhouli, yili, liji, mengzi, lunyu, Zhuangzi

It is possible to expand beyond just a list of key-keywords, and Scott introduces [Scott, “PC Analysis of Key Words”, 238] “associates”, as “words found to be key in the same texts as a given key key word”<sup>35</sup>. This notion is instrumental for the concept of (Scott, “PC Analysis of Key Words”, 241) “clump of associates”, which “is a set of associates formed by co-occurrence in the same texts which gave rise to associates.” Clumps could be used for defining a semantic area for groups of texts.

There is another way to find out if key-keywords could be used as a feature to group texts together. A matrix of texts and key-keywords was prepared<sup>36</sup>, where only the presence/absence of a character was entered and then clustered, using the free online tool ClustVis.<sup>37</sup> The resulting graphic representation of text/key-keyword clusters (so-called HeatMap) is presented in Figure 4.

<sup>35</sup> As Scott comments [Scott, PC, 239], “This notion of Associate is strikingly close to the early 1950s and 60s discussions of collocation”. Therefore, this paper will not address neither associates (nor clumps) nor collocations.

<sup>36</sup> See the list of these key-keyword under “Top key-keyword character for clustering experiment” in Appendix.

<sup>37</sup> <http://biit.cs.ut.ee/clustvis/>

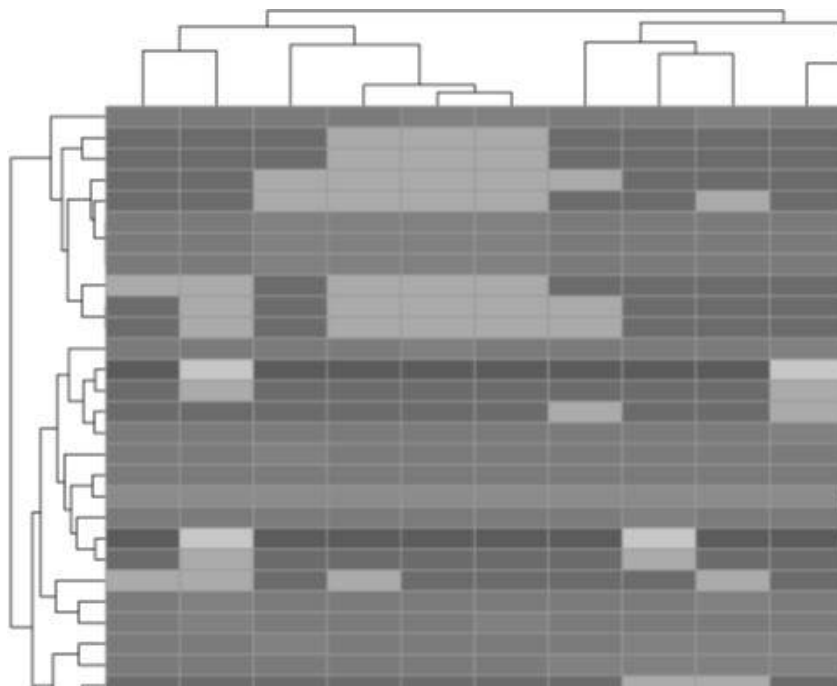


Figure 4. HeatMap of key-keyword clusters

There are two aspects of the ClustVis's HeatMap. On one side, it clusters texts according to their term (key-keywords) vectors; on the other side, it clusters terms according to their text vectors. Term clustering will not be analyzed in this article. Text clustering produces, most prominently, a historical group, and then, it groups together the Shi Jing and the Shu Jing (and the Zhou Li) and philosophical texts (the Lun Yu, the Meng-zi, the Xiao Jing, and the Zhuang-zi). Ritualistic texts are scattered, though. This opens another avenue for using key-keywords for genre and semantic analysis.

#### 4. Comparison of outputs of keyword analysis methods

The main benefit of applying keyword analysis, as well as h-point methods to frequency lists of Chinese classics, is in compiling concise lists of characters that could be used for topic and genre analysis of corpus texts, which could be called a semantic map of the texts. Having three instances of lists allows us to estimate the effectiveness of the applied methods. In addition, it may happen that some methods could be better suited for presenting specific aspects of texts. The breakdown of lists by texts and by categories is presented in Appendices.

Characters will be broken down by a set of categories, which were earlier utilized by the author for PHP content words. The category sets were not derived through topic discovery or modeling; they have been designed by the author and defined by the author's analysis of the PHP content characters of the Thirteen Classics.<sup>38</sup> The author decided to select categories that expose semantic dimensions of the classics (CQ, ZZ, GY, and GL, as well as SHU, LJ, YL, and ZL), as historical texts and, especially, chronicles. These categories have been informed mostly by the Chun Qiu PHP content characters' list. The main categories are Numeric, Calendrical, Social, and Politico-Geographical. The remaining characters were distributed across Part-Of-Speech (POS) categories: Nouns, Verbs, Adjectives/Adverbs, and finally Miscellanea. Practically all PHP characters of the Chun Qiu could be assigned to non-POS categories, with the exclusion of some verbs.<sup>39</sup> The order of text also reflects this idea — on the top, there are properly historical texts, and other texts follow in their internal id number order. The most important element is the lists of keyword and content characters. They could be reorganized into any other set of categories.

In addition, the texts below are grouped according to genre categories, presumed by the author.<sup>40</sup> One of the possible usages of character-

---

<sup>38</sup> This approach is similar to what Karen Donnelly implements [Donnelly, "Dr. Condescending," 91] "Initial analysis was carried out using Wordsmith Tools to elicit the top 100 lexical keywords from each corpus, which were then grouped thematically in order to allow comparison across the 3 corpora and guide selection for further study using collocations and concordance lines."

<sup>39</sup> Many Chinese characters (words) could play the role of various POS — noun, verb, adjective, etc. This phenomenon is well known as "POS variation or categorial ambiguity (i.e., ambiguity in grammatical category)" [Tsou and Kwong, "Some basic and Salient Linguistic Features," 7]. There are various solutions for this situation. Character (or word) could be assigned to all possible categories, or only one, dominating, category could be assigned to it. Tsou and Kwong offer an interesting log-ratio measure for the latter solution [Tsou and Kwong, "Some basic and Salient Linguistic Features," 8], but it requires a POS-marked corpus. In this study, only one category was assigned uniformly for all texts. It could be not optimal, as it may depend on text, which syntax category the character belongs. The author considers it a temporary provisional solution. The characters could be re-shuffled in other studies, if necessary.

<sup>40</sup> These texts were classified in their own categories, starting from the Qi Lüe catalogue (except the Zhuangzi). Most of them are collections of various texts, sometimes from different genres. As Lisa Raphals defines, "The Chinese Classics are a group of texts of divination, history, philosophy, poetry, ritual and lexicography that have, to a significant extent, defined the orthodox Ruhist (Confucian) tradition of China. Since the Song dynasty (960–1279), they have consisted of the

frequency lists is attempting to find similarities between lists and relating them to the genre categories. A pure numeric approach (described above), consisting in clustering vectors, created based on the number of characters in topic/lexical categories, is not conclusive. It only showed that “historical” texts could be grouped together, but the category system was pre-built to favor historical terms. Some preliminary suggestions will be made below, but this should be the subject of a special investigation.

### 5. Discussion of results

This study will not analyze topics of specific texts in detail because of the lack of space. However, it will discuss briefly keyword characters (and content characters) extracted by all three methods and how these sets of characters differ in relation to established categories, for each of the texts. The following section briefly discusses relevant parts of the synoptic table<sup>41</sup>.

#### Historical texts

##### 1. Chun Qiu

Table 7. Comparative keyword chart for the Chun Qiu

Method	Nu- meric	Calen- drical	Social	Politico- Geo- graphical	Nouns	Verbs	Adjec- tive/ Adverb	Misc
PHP	十二 三七 四八 六五 九	月年夏 春冬秋	公人侯 子師孫 王夫伯 叔曹	齊晉宋鄭 衛楚陳邾		有伐來 葬盟奔 莒歸蔡 殺杞帥 侵	大正	
LL	十七 二八	月冬春 秋夏年 癸辛亥 酉巳壬 午庚丑 乙丁	公侯伯 孫卒杞 叔曹	齊宋衛晉 鄭邾莒楚 蔡薛陳		伐帥師 來葬盟 奔侵		

following thirteen texts” (Raphals. “Chinese Classics,” 135). In the current article, the lexicographic text (the Er Ya) was omitted, the divinatory text (the Yi Jing) was classified as philosophical, and the Shu Jing and the Shi Jing were classified as fiction, instead of singling out just the Shi Jing as “poetry.”

<sup>41</sup> The complete table could be restored from these parts, or found in GitHub resources, see Appendix.

%DIFF		酉癸冬 巳卯亥 丑月壬 秋春乙 辛 丙戊午 丁夏庚	杞曹帥 郎	鄆薛邾滕 邲鄩宋洮 衛邾媯	麟蠡李 薑霄頓 款零灌 彊翦纒 羯	莒蔑句		
-------	--	---	----------	---------------------	-------------------------------	-----	--	--

In the Numeric section, only the PHP and LL methods bring numeric characters on top lists (see Table 7), and there are more characters in the PHP section than in the other measures.<sup>42</sup> In the Calendrical section, all methods bring up season names, but LL and %DIFF bring many cyclical signs on the top, unlike PHP, which could be explained by the lower frequency of cyclical signs. Therefore, LL and %DIFF methods seem to be considering dating by cyclical signs more prominent than by numbers. In the Social category, PHP and LL are close.<sup>43</sup> %DIFF only includes a few family terms on the top of its list. In the Politico-geography section, all methods bring in names of states or kingdoms, with some variations. The major difference between methods could be observed in Nouns and Verbs categories. The PHP and LL conceptual framework is consisting mostly of Social and Politico-geographical terms.<sup>44</sup> %DIFF output is practically missing any topical verbs.

Interpreted from the point of view of PHP content characters, the Chun Qiu could be perceived as a text, where there are social actors, related to some territories and states, who conduct some political actions, referred to seasons, month, and years. A similar interpretation could be given on the basis of LL. The type of interpretation, allowed by PHP and LL output, places the Chun Qiu into a genre category of historical chronicle. However, %DIFF only tells that the text is about some dates and states, but misses on actors and their actions.

<sup>42</sup> %DIFF method does bring numbers to the keyword character list, but they do not get on the top list, e.g., 十 gets into position 77 on CQ %DIFF, below cut-off of 50.

<sup>43</sup> LL missing such important categories as 人 and 王.

<sup>44</sup> However, it includes many high-frequency verbs. A few verbs overlap between PHP and LL (e.g., 伐, 來, and 葬).

## 2. Zuo Zhuan

Table 8. Comparative keyword chart for the Zuo Zhuan

Meth- od	Nu- meric	Calen- drical	Social	Politico- Geo- graph- ical	Nouns	Verbs	Adjec- tive/ Adverb	Misc
PHP	二三十 一五四 六	月年日 秋春冬 夏	子公人 君侯王 師夫伯 氏叔孫 歸臣民 季父司 軍孟仲 士帥尹 蔡	晉楚齊 國鄭衛 宋陳吳 周城秦 趙魯邑	禮事文 書天德 朝罪門 馬亂女 謀	曰為有可 伐盟謂死 殺行出成 告言知入 敢來聞欲 立見命奔 敗用亡取 生執辭求 政信令主 懼獻問宣 還獲乘棄 服	大寡今 武小難	是吾 未下 上中 右
LL			子叔氏 師公君 諸軍尹 侯季武 臣伯孫 貳族嬖	晉楚鄭 吳趙秦 韓邑國 衛魯荀 虢遠	產罪禍 書文圖 鮑亂范 卻賂駟	曰請懼盟 宣寵討敗 謀告棄叛 奔歸許亡 召洩殺怒 臧遂伐逐 囚死門賦 襄獲逃昭 免敢真逞 偪	寡樂敞 穆難平 魏俘絳 疆戎崔	吾
%DIF F			尉蟠媯 姚甥嬖 褚頡姑 潘	穎郎遠 濫號郟 鄂郟邯 邲巴	范帑孟 裔胙旆 青躒魋 衷蚡莜 祚紇諺 駸龙筆 檣猓賂 產儕儵 檣綫殿 鮑駮廚 繡董	洩紆悛逞 緝謗寵犒 偪瞞墮闕 狃瘡懼迂 蕪寔庇縛 蒐門囚傲 輶袂劫掠 句	俘絳暉 徹扶訴 齶鞏壘 麗樂瑰 鋼尪綿 楯	余

The Zuo Zhuan contains numeric and calendrical characters only in the PHP list, and even these numbers are lower than those in the Chun Qiu (see Table 8). All methods deliver more characters in Verbs and Social area

than in the Chun Qiu. %DIFF considerably expands its number of nouns. PHP and LL bring up nouns, too. In addition, there are characters in the Adjectives/Adverbs section. It is clear that nouns, produced by the PHP and LL methods, are important historical narrative terms,<sup>45</sup> which characterizes the Zuo Zhuan as a historical text. Although there is some overlapping with the %DIFF output, its nouns are less significant. The same could be said about Verbs. LL and PHP brought up characters, describing important historical actions, whereas %DIFF does not. All this allows, based on keyword and content characters, qualifying the Zuo Zhuan as a historical narrative, but may not be specifically a chronicle.

### 3. Gongyang Zhuan

Table 9. Comparative keyword chart for the Gongyang Zhuan

Method	Numeric	Calendrical	Social	Politico-Geographical	Nouns	Verbs	Adjective/Adverb	Misc
PHP	二三十 一五四 七	月年日 秋春冬 夏	子公人 君侯王 師夫伯 氏叔孫 歸季父 帥曹婁	晉楚齊 國鄭衛 宋陳邾	書天	曰為有可 伐盟殺出 言入來立 取執葬莒 譏弑稱	大正然	未是 吾
LL	十	月春秋 冬年夏 朔	公婁侯 曹伯孫 卒師季 帥歸蔡 杞	齊宋邾 衛晉鄭 陳桓	書災	譏言稱貶 弑伐葬記 諱莒來滅 盟侵殺	隱	曷奈
%DIFF	什頃		婁袁侄 鱗郎盱 繆郟	閩濤盾 系鄰甌	鷓葵潔 蝓翬鞍 孛書 獲	譏貶錄記 篡托諱剽 逡巧運弑 褒稱 諛	昧煬咍 斐隱憊 曠	曷 遯奈

The Gongyang Zhuan is closer to the Chun Qiu than to the Zuo Zhuan, even though in the LL list it has very reduced set of Numeric terms (see Table 9). Both PHP and LL produce rich sets of Social and Political terms, and both lists are practically missing Nouns. In Nouns, the characters that are available are similar to what the Zuo Zhuan list produces. The verbs also have a rich set of terms. Therefore, on the basis of these lists, the text could be characterized as a historical text with considerable time indications.

<sup>45</sup> There are overlapping frequent terms, such as *luan*, *zui* and *wen*.

#### 4. Guliang Zhuan

Table 10. Comparative keyword chart for the Guliang Zhuan

Meth- od	Nu- meric	Calen- drical	Social	Politico- Geo- graphical	Nouns	Verbs	Adjec- tive/ Adverb	Misc
PHP	二三十 一五四 七	月年日 秋春冬 夏	子公人 君侯王 師夫伯 叔孫歸 父帥蔡 曹	晉楚齊 國鄭衛 宋陳邾	事天	曰為有 可伐盟 殺出言 入來辭 葬莒弒 志	大正內	是未
LL	十	月春 秋冬夏 年日朔 辛癸	公侯卒 曹孫帥 師蔡姬 杞伯	齊宋衛 鄭邾晉 陳楚狄 滕		正伐葬 來弒言 盟莒志 侵殺奔 戰敗稱 繒		
%DIF F		酉朔	伉謚郟 曹侄	郟鄆郟 鄆桷郟 鄆薛滕 邾	鴉伶雩 頰驟累 贈翬螽 纘鬚彊	繒巧嫌 借挈搜 逾斥弒 羅唁崩 錄壅斲 嫁葬髡	謹卑 媯	

The Guliang Zhuan is similar to the Gongyang Zhuan; their Calendrical, Social, and Politico-geographical sections are very close (see Table 10). It also should be considered a historical narrative.

#### Ritualistic texts

#### 5. Li Ji

Table 11. Comparative keyword chart for the Li Ji

Meth- od	Numer- ic	Calen- drical	Social	Politico- Geo- graphical	Nouns	Verbs	Adjec- tive/ Adverb	Misc
PHP	二三十 一五四	月年日 夏	子公人 君侯王 夫臣民 士母婦 孔	楚國	禮事文 天德朝 門世賓 喪東樂 方道廟 義衣宗 位	曰為有可 謂死行出 成言知入 敢立見命 用生執主 問食祭教 學哭拜居 服	大小正 內貴外 然明反 長尊	是未 上中



LL		日	親君婦 父夫孔 妻舅士 祖 櫛	鄉	禮喪故 廟樂杖 衣冠天 氣麻尸 孝宗義 音棺肉 節情聲 齒事身 本功器 佩鬼味 堂墓臭 體繭粥 總禘經 踴衾 謁	哭祭服衰 學斂食練 製修哀教 養容問祔 順浴飲發 行附居沐 除誌殯吊 寢帶	親敬尊 貴幼誠 玄大素 賤明祖 疏	
%DIFF			婿紳婢	蕢剡	爾禪絳 胎總緹 粥綏維 蠟笏洞 旒尸杖 縞紛帨 觸庠翠 衿綬腎 駟閣尋 縱竽蔥 鷹闌臭 衾菜痛 麻犢雁 蚤棺醜 洽禘禘 簞紉幃 疫穉墀 麤箭苦 咳經謁	誌製禮詘 僕練鋪漱 減祔涂沅 溜措紐聖 綴伸殯掃 煎吊浴哭	臚暖恒 蕤糜褻 菲腥峻 暗揄歡 況	毋

The Li Ji has content character lists that are very different from the lists of historical texts (see Table 11). The PHP method extracts some Numeric and Calendrical terms, but any politico-geographical terms are practically missing. The Li Ji has a Social section, even though not so big as it is in historical texts. The main difference between those texts is that all methods bring in for the Li Ji many Nouns, Verbs, and Adjectives/Adverbs. There is a difference between how both the PHP and LL methods work and how %DIFF method works. The first two methods extract, on the top, meaningful terms, describing the conceptual frame of the text. The %DIFF

method also extracts important terms for the Li Ji, but these terms are not “conceptual.” They rather characterize the text stylistically.

## 6. Yi Li

Table 12. Comparative keyword chart for the Yi Li

Method	Numeric	Calendrical	Social	Political-Geographical	Nouns	Verbs	Adjective/Adverb	Misc
PHP	二三一		子公人 君夫爵 賓婦賓		禮門東北 西面馬南 階屍位解 席阼酒堂 首幣弓醢 奠	曰為有出 告入立見 命取執主 獻拜升降 受祭射興 洗揖祝送 授答荅進 退酌稽佐 篚辭俎	大外長	下上 右左 反從
LL			司婦卒 爵賓		西北東階 解屍南興 席阼位耦 筵醢房堂 脯門豆幣 酒弓矢首 拾鼎衆面 巾束羞戶 肺楹菹奠 脊銅饌醴 匕	拜主升坐 降洗俎受 揖執射設 荅祝答篚 祭初酌介 送授取盥 佐立稽加 酬進辭徹 出反遂退 獻縮適臘 薦侑擯釋 躋踴復筭	實袒	左右

%DIF F			稊爵		解枋樞枳 斫鋼匕陔 塾觚笄階 胃饌飪拾 匜弣面鏃 冪髀肺腸 奠肫筵阼 脊槃壁楣 甗楹耦艇 茵菹屍房 脯樹西巾 韭羸酪箱 醢醢齋席 糶軫縶筭 縲膊	苔洗啐嘑 侑篚擗筭 俎縮複摺 升答揖拜 坐扱擯醮 盥酬臘酌 撲降 概醕	孺纒謾 奕績蕘 胖妥	毋第
-----------	--	--	----	--	--	--	------------------	----

The Yi Li's lists are consisting mostly of Noun and Verb characters, practically missing characters in Numeric, Calendrical, and Politico-Geographical categories (see Table 12). They could be characterized similarly to what was said about the Li Ji.

### 7. Zhou Li

Table 13. Comparative keyword chart for the Zhou Li

Method	Nu- meric	Ca- lendri- cal	Social	Politico- Geo- graph- ical	Nouns	Verbs	Adjec- tive/ Adverb	Mis- c
PHP	二 三 一 五 八 九 百	日 歲	子 人 王 夫 氏 民 司 士 師 史 賓 帥 客 官 徒	國 府 邦	禮 事 馬 寸 掌 方 物 車 喪 法 器 田 刑	曰 為 有 謂 行 入 命 用 政 令 食 祭 祀 治 禁 辨 鼓 受 教 服 啐	大 小 正 內 長 共 胥	下 上 中
LL	二 四 八 六 十 五 三 百 九	歲	士 徒 史 客 職 國 野 官 賓 軍 群 兵 司 仆 吏	邦 市 地 方 田 府	掌 祀 法 寸 事 車 尺 物 鼓 刑 旅 器 牲 財 喪 圭 節 旗 材 獄 金 馬 膳 輪 獸 弓 筋 數 角 罰 弊 鐘 畿	令 禁 治 詔 祭 屬 訟 舞 分 役 贊 政 參 任 縣 兇 守 裸 教 建 待 巡 奏 領 馭 敘 蹕 園 積 戒 辨	共 胥 均 大 隸	

<b>%DIF FF</b>		昨	仆史嬪 徒吏	府域邦 幽	掌帑弩畿 軹鉦壇琮 弊鏡駕褐 唇蜃瘍幅 版縛朝筋 虞筍弛部 塵轂澮鑊 牲柝諭甕 柎毳寸助 韶案醫法 尺敷珍鱗 爻泉擊輪 橈鐘甬鬱 纘鎮枚	馭頌蹕療 蕪裸覘迤 斫圍窆辨 搏揉詔禁 煮敘染訟 舛兇判舂 齋芟擾	侔液均 胥瑞倨 植隸	
--------------------	--	---	-----------	----------	---	---	------------------	--

The Zhou Li's lists are similar to those of the Yi Li and the Li Ji, but they contain many Numeric characters, unlike the latter two texts (see Table 13). However, it practically does not contain calendrical terms, so the origin of the numbers should be different. Similarly to the Yi Li, it contains many Social terms, which is not surprising for a ritual text.

### Philosophical texts

#### 8. Lun Yu

Table 14. Comparative keyword chart for the Lun Yu

Method	Numeric	Calendrical	Social	Politico-Geographical	Nouns	Verbs	Adjective/Adverb	Misc
PHP	三		子公 人君 夫孔		禮事仁 道路	曰為有可 謂行言知 聞見問學	好	是 吾 未
LL			子孔 友君	邦	仁路貢 道張顏 色樊善 勇恥 忠蔽淵 鯉	曰問學知 冉言聞回 見信政求 謂行無怨 沽敏	直佞遠 怡巍偃 騫遲枉	吾
%DIF F			與友	顛	鯉樊貢 蔽仁路 顏畔張 牆淵勇 恥禱色 躬	沽冉諒切 譬學敏誨 回欺問惑 恭	偃怡巍 騫希枉 便佞倦 遲愚泰 篤貧狂 驕儉直	斯

In the Lun Yu lists, calendrical and numeric characters disappear in all lists (see Table 14). There are comparatively less keyword characters for Social actors than in many Ritualistic texts. Politico-geographical area is also non-existent. The main keywords there are defined by catch-all sections of Nouns and Verbs, similarly to the Yi Li. However, the LL and %DIFF methods also produce many characters in a rich Adjective/Adverb area. The character Nouns lists of the PHP and LL methods capture well the main conceptual frame of the text, whereas the characters produced by the %DIFF method rather characterize it stylistically. All lists contain just one, but very important term — *ren*. Finally, verbs are very basic and reflect the discursive side of the text.

### 9. Mengzi

Table 15. Comparative keyword chart for the Mengzi

Method	Numeric	Calendrical	Social	Politico-Geographical	Nouns	Verbs	Adjective/Adverb	Misc
PHP	一百		子公人 君王夫 民父士 孟孔舜	齊國	事天仁 道樂心 義	曰為有 可謂行 言知聞 欲見問 食	天今然	是吾 未下
LL	百		孟堯子 民王奚 妻舜禹 孔霸	鄒	仁智善 義天心 耕飢道 性足里 惡獸	養云問 仕欲放 充無愛 居章餽	賢湯慕 誠汗賸 悅	
%DIFF			氓舜叟 孟堯嫂 霸	鄒邠庾	智界眸 奔謳庫 鬻飢烹 驩吠沼 供耕簞 鷓屑濱	餽搜孳 曆慕充 覺	訑煖賸 悅傑紕 惻潔壑 褐餓巍 沛餒 汗	沓

The Mengzi is a philosophical text, like the Lun Yu (see Table 15). It also has no keywords in the Numeric, Calendrical, and Politico-geographical sections. However, it has a wider list of Social terms.<sup>46</sup> Characters in the Nouns section are also close to the Lun Yu's lists, which means that they discuss the same subject area. The Verbs list is shorter and reflects a discourse structure.

<sup>46</sup> Including practically all of the Lun Yu's terms, their lists are close.

It should be noted that the absence of characters in the Numeric, Calendrical, and Politico-Geographical sections does not define the text's genre automatically. It only means that such text does not belong to historical literature. Its specific genre should be analyzed by content characters in the Nouns and Verbs sections; the Nouns play the most important role.

### 10. Xiao Jing

Table 16. Comparative keyword chart for the Xiao Jing

Method	Numeric	Calendrical	Social	Politico-Geographical	Nouns	Verbs	Adjective/Adverb	Misc
PHP			子人民父		事天孝			
LL			親民父母		孝事故詩 天身德思 義	云爭愛 教順移 治敢生 行	敬嚴悌 滿聖明 先	
%DIFF F			親父母民 兄家		孝蓋忠詩 身法聖思 事德宗刑 廟天地善 義	云移爭 愛順教 哀治養 敢故守 生失道	悌嚴滿 悅敬昔 安貴明 終先	莫

XJ is a typical philosophical text. There are no characters in the Calendrical, Numeric, and Geographical sections, and a few characters in the Social and Nouns sections (see Table 16). The LL and %DIFF methods add characters to the Verbs section, and %DIFF has a large Adjective/Adverbs section.

### 11. Zhou Yi

Table 17. Comparative keyword chart for the Zhou Yi

Method	Numeric	Calendrical	Social	Politico-Geographical	Nouns	Verbs	Adjective/Adverb	Misc
PHP	二三五 四六九		子人君		天象位道 咎彖利	曰有可行 用貞志孚 咎	大小正吉 凶明亨終	未下 上中
LL	九六				象利志頤 雷井夬蠱 巽兌坎	貞咎孚曰 悔無終厲 謙艮行蹇 順剝征涉 履震應噬 蒙遯渙	吉亨凶剛 柔元險光 壯需當攸 吝	

%DIFF					夬巽象拇 棟頤頻兌 蠱坎雷 牝利	艮嗑嗃啞 渙孚貞謙 噬遜咎剝 蹇悔窺婚 姤	吝亨媯剛 凶睽需吉 羸拯汔攸 妄漸惕壯 柔牀	
-------	--	--	--	--	---------------------------	-----------------------------------	------------------------------------	--

The Zhou Yi lists characterize it as not a historical text; however, the PHP method extracts several numbers, because of their frequent usage in the text (see Table 17). However, the LL and %DIFF methods ignore numbers. The text content is best represented by PHP nouns.

## 12. Zhuangzi

Table 18. Comparative keyword chart for the Zhuangzi

Method	Nu- meric	Ca- lendri- cal	Social	Politico- Geo- graphical	Nouns	Verbs	Adjec- tive/ Adverb	Misc
PHP	三一	日	子人 君王 夫民	國	事天道物神 方德心形聖 足名世樂義 身仁坤	曰為有可 謂死行出 成言知聞 欲見用生 問治	大今同 明然	是吾 未下 上中
LL	萬		尼堯 奚莊 墨倪	地	形物天道性 足精衆意心 俗身名神劍 機情盜處故 惡果德日海 顏世累仁枝 死耳聃聃崖 竅聖	知生遊化 謂始忘辯 見回問聞 惑似靜應 理通論解 況語觀存 治譽殉齧 避離調	邪然真 虛冥滿 老巧古 愚窮全 悲恬瞽 默	吾嘗 孰
%DIFF			倪	崑崙	崖陌筵扁影 毫軀竅涯捆 纏蛙形蟻湖 係狀機圓跂 樞聃屠齋机 脩杯尻精骸 蝟暮籠斛竿 符勿脛鑑垢 疵毀 鼈	攫憐闕避 齧悟娛擢 捐遊撓剖 掎掇臥波 攫響啍鈞 捶諛支調 彫	真恬瞽 僻徨徬 腐漫謬 默邪淡 豪怪俄 盲澹全 蓬怵道 悲俱伎 蹶	勉嚮 徧

The Zhuangzi is a regular non-historical text. Its Noun section contains many terms that qualify it as a philosophical text, especially in the PHP list (see Table 18).

## Fiction

### 13. Shi Jing

Table 19. Comparative keyword chart for the Shi Jing

Meth- od	Nu- meric	Ca- lendri- cal	Social	Politico- Geo- graphical	Nouns	Verbs	Adjec- tive/ Adverb	Misc
PHP	四百	月日	子公人 君王歸 民孔	國	天德女心 山方思樂 南方山	曰為有可 行言來命 維	大憂	是中 予載
LL			牡		心胡思福 方桑隰女	載止采俾 憂式瞻 飛徂寐逝 烝陟懷	昊赫碩 肅烈淑 皇醉	予載
%DIFF				狃淇	獵駢駿鷲 幡劬零荏 黽坳鱗卉 騏遑監 芑靈繆鎬	躡維適觀 緝究儼咽 啄泮	駟粲芾 涓離悠 奕苑慘 昊子汎 嚶緜	兮噦 屈猗

The PHP method identifies a few Numerics and Calendricals in the Shi Jing (see Table 19). It has many Social terms in the PHP list, but almost nothing in this category in the LL and %DIFF lists. It is placed here in the “fiction” genre category, but this qualification needs to be elaborated later on the basis of keyword and content character analysis.

### 14. Shu Jing

Table 20. Comparative keyword chart for the Shu Jing

Meth- od	Nu- meric	Ca- lendri- cal	Social	Politico- Geo- graphical	Nouns	Verbs	Adjec- tive/ Adverb	Misc
PHP	三一 五四 百		公人王 民帝殷	周邦	事文天德	曰有言敢 命用作刑	大今小明	下上 惟子 汝厥
LL			帝 王 民 禹	殷邦	德命天艱 刑典底	克訓宅越 罰保休允 格念逸祗 恭嗣威率 迪敷猷	丕永明咸 誕欽庶	予惟 汝厥



%DIFF				岳洛岱	忱眷救底 彰豬洪艱 棐疇	迪遜敷猷 罔導冲佑 肇格祗俞 迓誥胤賚 宅訓逸	丕亮欽惇 愆訖俊誕 灼彥 爽	惟矧 朕厥 汝茲
-------	--	--	--	-----	--------------------	-------------------------------------	-------------------------	----------------

In the PHP list of the Shu Jing, there are no characters in Calendrical and in Politico-geographical, but there are some Numeric characters (see Table 20). There are characters in the Social section in the PHP and LL lists, but the Noun section is rather small in the PHP and LL sections, and Verbs are abundant. This could be interpreted as a text about social actors.

## 6. Conclusion

This paper has investigated the character-frequency lists of the WSW Ctexts corpus to identify significant keyword characters. Three methods were utilized to calculate keyness score: LL significance test, CHI test, and %DIFF. Comparison of the results of the first two methods showed that LL test is very close to the CHI test, and the latter was omitted in this study. The results of the %DIFF test were different, especially for the top parts of the lists.

Beside regular keyword character lists, this study created lists of positive and negative key-keywords for the corpus. Analysis of those lists, as well as clustering texts across the top key-keywords, has shown that key-keyword groups could be used to identify genre groups of texts. It may be even more beneficial to analyze clusters of key-keyword associates and “clumps,” which could be the subject of a future study.

The study further concentrated on the top parts of lists. Depending on the keyness score method, keyword filtering parameters, and text size, there could still be hundreds of keywords on lists for some texts. Although all of them could be relevant in understanding texts’ content, researchers in keyword analysis usually concentrate on the top parts of these lists.<sup>47</sup> There is no “native” measure to identify how many characters to select for the top list, but usually a number from 50 to 100 is recommended.

Previously, the author has conducted a study on selecting most significant characters on character-frequency lists, employing “PHP” method to identify content characters. It allowed, depending on text size, to identify from 30 to 100 content characters for each text and can be considered a native measure.

It is logical to compare results from the PHP content character list and keyword analysis methods. Taking as guide the number of PHP content

<sup>47</sup> The full keyword lists are still available, for this study, for analysis on the accompanying GitHub website, see Appendix.

characters for a text, the number of characters from the top of the keyword lists was defined, provisionally, either 50 or 100 characters. Therefore, for each text, three content character sets were created, and for each set, characters were broken down into a set of categories, which were earlier used in PHP content study, losing the frequency order. This latter set of categories was created, based on specific characteristics of the Chun Qiu; therefore, it allows tracing how well other texts fit the Chun Qiu topic model – the model of a historical chronicle. Such set of categories is biased toward the genre of the chronicle, but in any case, it allows getting some basic genre specifics of texts.

Production of a synoptic table with keyword character lists for corpus texts has been one of main goals of this study. This table allows comparing various methods of extracting content characters. It could be called a semantic map for the texts. Owing to the scope limitation of this article, the comparison of the specific lists was not conducted thoroughly, and it is a matter of further research. The main point in the current comparison was to identify how well texts fitted the category model of the Chun Qiu and to try to group texts according to characteristics of their content lists.

Comparison of sets of extracted significant content characters demonstrates that the PHP and LL keyword character lists are most beneficial for understanding the corpus text content, and seemingly, the PHP list gives a slightly better picture of the text content. The %DIFF method, though, could be useful for analyzing the stylistic characteristics of texts.

Suggested set of categories, designed for the identification of historical texts (chronicles), demonstrated to be effective for content lists developed using the PHP and LL methods. The %DIFF method was not that effective in selecting characters, characteristics for chronicles. The developed synoptic table of content characters will be used in further studies of individual text characteristics.

## Appendix

### Abbreviations for texts in the WSW Ctexts Corpus

Chun Qiu	CQ / chunqiu
Chun Qiu Zuo Zhuan	CQZZ / chunqiu zuozhuan
Gongyang Zhuan	GY / gongyang
Guliang Zhuan	GL / guliang
Li Ji	LJ / liji
Lun Yu	LY / lunyu
Mengzi	MZ / mengzi
Shi Jing	SHI / shi
Shu Jing	SHU / shu
Xiao Jing	XJ / xiaojing

Yi Li	YI / yili
Zhou Yi	ZY / zhouyi
Zhuangzi	ZHZ / zhuangzi
Zhou Li	ZL / zhouli
Zuo Zhuan	ZZ / zuozhuan

### **Top key-keyword characters for clustering experiment**

天子無行民道萬三來圍孔己方明曰正歸父百雨故公古國地姜孰  
思怨政敗王詩變足身邦義

### **Internet resources**

The corpus website:

<http://www.umass.edu/ctexts/index.php>

Username: ctexts, password: umass

### **GITHUB resources**

Github DOI: [https://github.com/wsw-ctexts/vocabulary\\_richness/](https://github.com/wsw-ctexts/vocabulary_richness/)

#### **chinese\_classics\_keywords\_log\_likelihood\_short.xlsx**

This file contains lists of LL keyword characters for the corpus. Positive and negative keywords for each text are contained in a dedicated sheet (e.g., for the Chun Qiu in the sheet “LL\_CQ”). It has three areas: area of positive keywords, area of negative keywords, and comparison with PHP content characters (the list of PHP synsets is contained in “PHP\_SYNSETS” sheet). Summary data is contained in the “KK\_LL\_STATS” sheet. The “KW\_Charts” sheet contains comparative charts of LL, CHI and %DIFF keywords. The “KW\_Summary” sheet contains a detailed comparison of PHP and LL characters (with a chart). The sheet “LL\_ALL” contains the synoptic list of all LL keyword characters of the corpus. The sheet “LL\_ALL\_KKW” contains the list of all key-keywords. The sheet “LL\_KKW\_SUMMARY” contains the list of key-keywords with lists of their texts attached. The sheet “LL\_KKW\_POS\_SUMMARY” contains a similar list of only positive key-keywords. The sheet “LL\_KKW\_NEG\_SUMMARY” contains a similar list of negative key-keywords.

#### **chinese\_classics\_keywords\_CHI\_short.xlsx**

The file contains data on CHI keyword characters of the corpus. The file has the same structure, as chinese\_classics\_keywords\_log\_likelihood\_short.xlsx, but only has “KW\_CHI\_CHARTS” and “CHI\_SUMMARY” sheets as summary data.

### **chinese\_classics\_keywords\_diff\_short.xlsx**

The file contains data on CHI keyword characters of the corpus. The file has the same structure, as `chinese_classics_keywords_log_likelihood_short.xlsx`, but is missing comparison with PHP content characters.

### **chinese\_classics\_keykeywords\_matrix.xlsx**

The file contains data on top LL key-keywords, in the form, used for clustering algorithms.

## **Literature**

Archer, Dawn. "Does Frequency Really Matter?" In *What's in a Word-List? Investigating Word Frequency and Keyword Extraction*, ed. by Dawn Archer, n.p.: Ashgate Pub. Limited, 2009, 1–16.

Archer, Dawn (ed.) *What's in a Word-List? Investigating Word Frequency and Keyword Extraction*. Digital Research in the Arts and Humanities. n.p.: Ashgate Pub. Limited, 2009.

Baker, Paul. "The question is, how cruel is it?" Keywords, Foxhunting and the House of Commons. In *Word Frequency and Keyword Extraction AHRC ICT Methods Network Expert Seminar on Linguistics 8 September 2006*, Lancaster: UCREL, 2006, 1–8.

Baker, Paul, Jesse Egbert, Tony McEnery, Amanda Potts and Bethany Gray. "Triangulating Methodological Approaches in Corpus Linguistic Research." In *Corpus Linguistics 2015, Abstract Book*, ed. by Formato, Federica and Andrew Hardie. Lancaster: UCREL, 2015, 39–41.

Baron, Alistair, Paul Rayson, and Dawn Archer. "Word Frequency and Key Word Statistics in Corpus Linguistics." *Anglistik: International Journal of English Studies* 20, no. 1 (2009): 41–67.

Berber-Sardinha, Tony. "Comparing Corpora with WordSmith Tools: How Large Must the Reference Corpus Be?" In *Proceedings of the Workshop on Comparing Corpora. Association for Computational Linguistics*, 9 (2000): 7–13.

Boltz, William G. "Why So Many Laozi-s?" In *Studies in Chinese Manuscripts: From the Warring States Period to the 20th Century*, edited by Imre Galambos, Budapest: Institute of East Asian Studies, ELTE, 1–32.

Bondi, Marina. "Perspectives on Keywords and Keyness: An Introduction." In *Keyness in Texts*, ed. by Bondi, Marina, and Mike Scott. Studies in Corpus Linguistics, 41. Amsterdam; Philadelphia: John Benjamins Pub. Co, 2010, 1–18.

Bondi, Marina, and Mike Scott (eds.) *Keyness in Texts*. Studies in Corpus Linguistics, 41. Amsterdam ; Philadelphia: John Benjamins Pub. Co, 2010.

Brooks, E. Bruce. "Before and After Matthew." In *Studies in the Didache*, ed. by Jefford and Draper (forthcoming).

Cheng, Winnie. "Corpora: Chinese-Language." In *The Encyclopedia of Applied Linguistics* ed. by Carol A. Chapelle, Oxford: Blackwell Publishing Ltd, 2012.

Cvrcek, Václav and Masako Fidler. "Not all keywords are created equal: How can we measure keyness?" In *Corpus Linguistics 2013 : Abstract Book*, ed. by Hardie, Andrew, and Robbie Love. Lancaster: UCREL, 2013, 55–57.

- Donnelly, Karen. "Risk, chance, hope — the lexis of possible outcomes and infertility." In *Corpus Linguistics 2013: Abstract Book*, ed. by Hardie, Andrew, and Robbie Love. Lancaster: UCREL, 2013, 67-69
- Donnelly, Karen. "'Dr Condensing' and 'Nurse Flaky': The Representation of Medical Practitioners in an Infertility Corpus." In *Corpus Linguistics 2015, Abstract Book*, ed. by Formato, Federica and Andrew Hardie. Lancaster: UCREL, 2015, 91-93.
- Gabrielatos, Costas. "Selecting Query Terms to Build a Specialized Corpus from a Restricted-Access Database." *ICAME Journal* 31 (2007): 5-44.
- Gabrielatos, Costas and Anna Marchi. "Keyness: Appropriate metrics and practical issues." In *CADS International Conference 2012. Corpus-assisted Discourse Studies: More than the sum of Discourse Analysis and computing?*, University of Bologna, 2012
- Gaspari, Federico and Marco Venuti. "A golden keyword can open any corpus: theoretical and methodological issues in keyword extraction." In *Corpus Linguistics 2015, Abstract Book*, ed. by Formato, Federica and Andrew Hardie. Lancaster: UCREL, 2015, 131-133.
- Goh, Gwang-Yoon. "Choosing a Reference Corpus for Keyword Calculation." *Linguistic Research* 28: 1 (2011): 239-256.
- Hutchins, William John. "The Concept of 'Aboutness' in Subject Indexing." In *Aslib Proceedings*, 30, n.p.: MCB UP Ltd, 1978, 172-181.
- Kao, Anne and Poteet, Steve R. (eds.) *Natural Language Processing and Text Mining*. London: Springer-Verlag, 2007.
- Kilgarriff, Adam. "Language Is Never, Ever, Ever, Random." *Corpus Linguistics and Linguistic Theory* 1, no. 2 (2005): 263-276.
- Lin, Chin-Yew, and Eduard Hovy. "The Automated Acquisition of Topic Signatures for Text Summarization." In *Proceedings of the 18th Conference on Computational Linguistics*, Association for Computational Linguistics, vol. 1, 2000, 495-501.
- Nenkova, Ani, and Kathleen McKeown. "A Survey of Text Summarization Techniques." In *Mining Text Data*, ed. by Charu C. Aggarwal and Cheng Xiang Zhai, Boston: Springer US, 2012, 43-76.
- Popescu, Ioan-Ioviț, and Gabriel Altmann. *Word Frequency Studies*. Berlin; New York: Mouton de Gruyter, 2009.
- Popescu, Ioan-Iovitz, J. Mačutek, and Gabriel Altmann. *Aspects of Word Frequencies*. Lüdenscheid, 2009.
- Rayson, Paul, and Roger Garside. "Comparing Corpora Using Frequency Profiling." In *Proceedings of the Workshop on Comparing Corpora*, Association for Computational Linguistics, 2000, 1-6.
- Rayson, Paul. "From Key Words to Key Semantic Domains." *International Journal of Corpus Linguistics* 13, no. 4 (2008): 519-49.
- Richardson, Kay. "Keywords Revisited: The Present as History." *Social Semiotics* 5, no. 1 (1995): 101-17.
- Scott, Mike. "PC Analysis of Key Words—and Key Key Words." *System* 25, no. 2 (1997): 233-245.

- Scott, Mike. *Wordsmith Tools Version 3*. Oxford: Oxford University Press, 1999.
- Scott, Mike. "Picturing the Key Words of a Very Large Corpus and their Lexical Upshots — or Getting at the Guardian's View of the World" in *Teaching and Learning by Doing Corpus Analysis*, ed. by B. Kettemann and G. Marko, Amsterdam: Rodopi, 2002, 43–50.
- Scott, Mike. *WordSmith Tools Version 5*. Liverpool: Lexical Analysis Software, 2010.
- Stubbs, Michael. *Words and Phrases: Corpus Studies of Lexical Semantics*. Oxford, Malden, MA: Blackwell Publishing, 2001.
- Tsou, Benjamin K., and Olivia Oi Yee Kwong. "Some Basic and Salient Linguistic Features Across Chinese Speech Communities from a Corpus Linguistics Perspective," In *The Oxford Handbook of Chinese Linguistics*, ed. by William S-Y. Wang and Chaofen Sun. Oxford: Oxford University Press, 2015, 1–23.
- Williams, Raymond. *Keywords: A Vocabulary of Culture and Society*. Oxford University Press, 2014.
- Zinin, Sergey. "Pre-Qin Digital Classics: Study of Text Length Variations". In Scholarly Reports of the Department of China of the Institute of Oriental Studies, Russian Academy of Sciences, issue 15, The 44<sup>th</sup> Conference "Society and State in China", vol. XLIV, pt.2, Moscow, 270–311, 2014.
- Zinin, Sergey. "Vocabulary richness of early Chinese texts: macroanalysis of the Thirteen classics and the Zhuangzi." In Scholarly Reports of the Department of China of the Institute of Oriental Studies, Russian Academy of Sciences, issue 20, The 46<sup>th</sup> Conference "Society and State in China", vol. XLVI, pt. 1, Moscow, 197–253, 2016
- Zinin, Sergey. "Analysis of character-frequency lists of Chinese classics and its application to content analysis and genre attribution." In Scholarly Reports of the Department of China of the Institute of Oriental Studies, Russian Academy of Sciences, issue 22, The 47<sup>th</sup> Conference "Society and State in China", vol. XLVII, pt. 1, Moscow, 254–299, 2017

## CONTENTS

### History and ethnography

- Diakova O.V.** Archaeological sources on the role of Polzevskaya culture in formation of the Tungus-Manchu community ..... 6
- Kukeev D.G.** Some features of the Oirat's resistance at the final stage of the annihilation of Zunghar khanate: "makhachin" ..... 13
- Golovachev V.Ts.** Illness and death of Paul Ibis, the hero and author of the "Excursion to Formosa": trans-temporal diagnosis ..... 23
- Krivokhizh S.V.** On emergence of the "sovereignty" concept in China ..... 36
- Nesterova E.I.** Chinese Embassy visit to Russia in 1870 and prospects of Russian-Chinese relations development (based on the documents of RI Foreign Policy Archive) ..... 40
- Borokh O.N.** *Xueyi* journal and dissemination of scientific knowledge in China in the first half of the twentieth century ..... 50
- Pozhilov I.E.** One's fundraising, the other's robbery (on "expropriation of exploiters" in the CCP's agrarian revolution) ..... 62
- Galenovich Yu.M.** The nation of China in the war against Japanese aggression in 1930s — 1940s ..... 72
- Zabrovskaia L.V.** Social security system in the People's Republic of China till the beginning of economic reforms (the 1950s–1980s) ..... 84
- Dmitriev S.V.** Three southern cities: sketches of history ..... 94
- Gracheva Yu.A.** The modern village life of the people Bouyei in Guizhou (based on the field research) ..... 105

### Source studies and historiography

- Starodubtseva N.S.** Review of the newspaper *Renmin ribao* publications about archaeological discoveries in China in 2017 ..... 122
- Blyumkhen S.I.** Hermeneutics of names and images in ancient Chinese texts: Bo Yi-kao, Bo Yi and Shu Qi ..... 133
- Popova G.S.** "The Preface to Writings" (*Shu Xu*) in "The Historical Records" (*Shi Ji*) by Sima Qian ..... 162

<b>Popova G.S.</b> <i>Shu Jing</i> (“The Classic of Writings”) at the turn of the II–I cent. BC (based on the materials of “The Historical Records” by Sima Qian) .....	183
<b>Popova G.S.</b> <i>Shu Jing</i> (“The Classic of Writings”) and <i>Yi Zhou Shu</i> (“The Lost Book of Zhou”): points of intersection .....	215
<b>Zinin S.V.</b> Keyword analysis of Chinese classics: “Thirteen Canons” and Zhuangzi .....	240
<b>Tišin V.V.</b> Some remarks on historical ethnonymy and toponymics of inner Asia (in the context of problem on <i>Xi</i> 霽 and <i>Bái-xí</i> 白霽) .....	279
<b>Vinogradova T.I.</b> What is the “Utopian monism” of Su Xun (based on the sinological card-file of academician V.M. Alekseev) .....	288
<b>Molodyakov V.E.</b> Jacques Bainville on the Japanese continental policy in China: the Manchurian incident and its consequences .....	296

### **Economics, politics and law**

<b>Leksyutina Ya.V.</b> Contemporary China’s contribution to UN peacekeeping .....	305
<b>Portyakov V.Ya.</b> A new round of the PRC’ economy opening-up .....	312
<b>Gruzinov I.I.</b> Possibilities and problems of modern Chinese ideology export .....	315
<b>Semenov A.A.</b> The role of small democratic parties in the PRC political system .....	321
<b>Sukhadolskaya L.L.</b> Brand with Chinese specifics as a new force of China economic development .....	330
<b>Dikarev A.D.</b> The legal discourse and political consequences of the Hague tribunal award on Philippines vs. China arbitration case .....	340
<b>Altantsetseg Noosgoi.</b> The Chinese-Mongolian relations of comprehensive strategic cooperation: desires and reality .....	357
<b>Kaygorodova N.A.</b> Some aspects of control over urbanization processes in the PRC on the verge of XX and XXI centuries .....	368



<i>Chubarov I.G.</i> Seeking for suitable urbanization model in China (Guangdong province case study) .....	376
<i>Magdalinskaya Yu.V.</i> Legal aspects of independent directors' institution formation and its establishment in the corporate practice in China .....	383
<i>Baldanova R.A.</i> On legal regulation of the illegal migration to the PRC .....	398
<i>Sazonov S.L.</i> Chinese aviation complex .....	403
<i>Ivanov S.A.</i> Prospects of cooperation between Heilongjiang province and Primorsky Krai .....	416
<i>Chen Xiao, Sazonov S.L.</i> On the priorities of the Russian-Chinese transport integration development .....	425
<i>Wu Zi, Sazonov S.L.</i> The main sources of financing of the Eurasian transport corridor .....	437
<i>Yu Tao, Sazonov S.L.</i> Russia and China should work together to develop "Ice Silk Road" in the Arctic region .....	450
<i>Gordienko D.V.</i> Changing the level of economic security of China in the implementation of the strategy of Economic Belt of the Silk Road .....	458
<i>Mukhamadiyeva Ya.I., Muratshina X.G.</i> Climate talks in China–EU negotiations .....	477
<i>Zhilkiyaev S.N.</i> From historian to Red inquisitor. Essay of political biography Wang Qishan .....	482

#### In memoriam

<i>Alpatov V.M.</i> In memory of S.Ye. Yakhontov .....	492
CONTENTS .....	511
目录 .....	514