

AUTOMATIC COMPILATION OF MODERN CHINESE CONCORDANCES

Syunsuke UEMURA*, Yasuo SUGAWARA*
Mantaro J. HASHIMOTO**, Akihiro FURUYA***

*Electrotechnical Laboratory, 1-1-4 Umezono, Sakura, Ibaraki 305, JAPAN
**Tokyo University of Foreign Studies, 4-51-21 Nishigahara, Kita, Tokyo 114, JAPAN
***Tokyo Metropolitan University, 1-1-1 Yakumo, Meguro, Tokyo 152, JAPAN

An automatic indexing experiment in Chinese is described. The first very large volume of modern Chinese concordances (two sets of one million-line KWIC index) has been compiled and materialized automatically with a modified kanji printer for Japanese.

INTRODUCTION

This paper describes an experiment to compile Chinese concordances automatically. A very large volume of KWIC indexes for modern Chinese (one million lines per set) has been compiled successfully with a kanji printer for Japanese. This paper discusses the purposes of the experiment, selection and input of the Chinese data, some statistics on Chinese characters (vs. kanji) and the concordance compilation process. Finally, examples from the computer-generated concordances are shown.

THE PURPOSES

The idea of machine-processing modern Chinese data originally came from Professor Yuen Ren Chao, Agassiz Professor Emeritus of Oriental Languages at the University of California at Berkeley, before one of the authors (Hashimoto) took over the directorship of the Princeton Chinese linguistics project. Chao served as the chief of the advisory committee to the project since its foundation. The idea, in short, was: so much has been said about the Chinese pai-hua-wen -- a written language of modern China -- yet nobody has ever clarified what it really was, i.e.; what the basic vocabulary was, what the major syntactic structure was, etc.: in other words the every detail of the reality of pai-hua-wen. Certain quantitative surveys were done before us, but even the most extensive one in those days was based on data consisting of no more than 100,000 characters. In addition, the selection was very poorly done -- most of the materials were primary school textbooks. We did not believe that school textbooks reflected the reality of the language, even in its written form. We chose one digit more than the previous one, namely 1,000,000 characters, though for various reasons, the actual data contained in our tape include several thousands more than one million [1, 2].

After completion of the computer input and editing of the million-character file at Princeton, researches towards statistical aspects of the data have been conducted [4]. As stated in [4], tables of character frequency can tell us various aspects of the Chinese, such as the basic character set, transient states of character strings and so on. This can be summarized as the first step of computer-processing modern Chinese data. However, in order to understand the reality of a language, besides statistics, concordances are the necessities which illustrate the contexts where and how those characters are used.

On the other hand, computer applications to Chinese have very limited background so far. No computer-generated concordances on Chinese have been reported yet. Thus the concordance generation project would not only be valuable to the understanding of Chinese pai-hua-wen, but also contribute to the development of the methodology to manipulate Chinese automatically. Consequently, a project to compile concordances of the Princeton million-character file was conducted at the Electrotechnical Laboratory during 1977-1979. This constitutes the second important stage of computer-processing modern Chinese.

THE CHINESE DATA

The Input of the Original Data

The first phase of the data input was done in Taiwan during 1969-1972 with a Chinese character keyboard, designed by Cheng Chin Kao -- a Chinese teletype Machine (manufactured by the Oki Denki Co., Ltd.). The code was converted into the Chinese standard telegraphic code in Waltham, Massachusetts at a computer company. The greatest difficulty, in addition to ordinary proofreading, consisted in the conversion of the so-called "combination characters" of the C.C.Kao system: any character not found in the Kao keyboard was punched so that part of it (normally the "radical") was represented by a character having the same radical in the keyboard, and another by a character having the same "signific". Necessary flags were of course attached to these "combination characters", yet the key punchers selected those constituent characters quite at random, sometimes

disregarding the position of a radical within a character, so that the results were often a hopeless mess.

The Selection of the Data

It was tried, at the selection of the data, to cover every conceivable category and style of writings in China since her modernization, the so-called May 5 Movement period, from ordinary novels to philosophical writings, from political speeches to newspaper articles, etc. etc. These categories and styles were classified and were assigned appropriate marks to show the genre. The partial list of these writings follow:

魯迅: 阿Q正傳	端木蕻良: 大江	陳独秀: 今日中國之政治問題
巴金: 家	張愛玲: 金鎖記	政務院: 中央人民政府政務院命令
茅盾: 子夜	趙樹理: 李家莊的變遷	薄儀: 我的前半生
老舍: 駱駝祥子	郭沫若: 金剛坡下	毛澤東: 實踐論
艾蕪: 石膏嫂子	梅蘭芳: 舞台生活四十年	林彪: 人民戰爭勝利萬歲
曹馮: 北京人	孫逸仙: 三民主義	劉少奇: 關於土地改革問題的報告
錢鍾書: 圍城	李濟: 殷墟建築遺存報告序	周恩來: 關於知識分子問題的報告
沈從文: 邊城	胡風: 文藝工作的發展及其努力方向	蔣介石: 行的道理

For a complete list of all these writings and of the genre marks, see [3]. All the proper nouns were so marked, as they may not correctly contribute to any statistical measurement of the written language except for these proper nouns themselves. These nouns were marked in the original texts by research assistants with enough command of the language to make correct judgment. Anything else, including punctuation marks of all sorts, in the texts were properly processed. Every sentence, including some vocative phrases, was numbered within the writing piece quite mechanically, though occasionally it was necessary for specialists to make certain judgment for segmenting sentences.

The Code System

The Chinese standard telegraphic code system includes some 9500 codes for Chinese characters. A code consists of a set of 4 digits, which represents one Chinese character. Among those 9500, 5231 have been used.

Statistics

Statistical analysis of this million-character file can be found in [4]. Some additional statistics are provided here. Fig. 1 shows the 10 most frequently used characters with their frequencies. These 10 characters occupy 17.1% of the total amount. Fig. 2 is a table of character frequencies vs. the number of character types. Fig. 3 shows the cumulative percentage of character occurrences as a function of the number of character types (in descending order of frequency). It indicates, for example, only 92 characters represent 47% of the entire data. There are 1170 characters each

of which are used more than 100 times and they occupy 92.8 % of the whole data.

Character	Frequency
的	46531
一	18077
是	17874
了	16390
不	16138
我	12827
在	11096
人	11057
有	10717
他	10332

Fig. 1. List of High Frequency Characters

Frequency	No. of Character Types
- 10001	10
10000 - 5001	13
5000 - 3001	32
3000 - 2001	37
2000 - 1001	106
1000 - 501	176
500 - 301	208
300 - 201	191
200 - 101	397
100 - 81	150
80 - 61	230
60 - 41	294
40 - 21	574
20 - 11	563
10 - 1	2250

Fig. 2. Frequency Distribution of Chinese Character Types

CHINESE CHARACTERS VS. KANJI

Chinese characters were imported into Japan sometime in the 5th century. Since then, they have been extensively used with a few additional characters created in Japan (this modified set of Chinese characters is called "kanji"), although hiragana and katakana (two sets of pure Japanese characters with their origin also in the forms of Chinese characters) were invented early in the 9th century.

"Chinese characters for daily use" established by the Ministry of Education for modern Japanese includes a 1850 kanji set, however several thousand more are still in use especially for proper nouns. The Japanese Industrial Standard (JIS) "Code of the Japanese Graphic Character Set for Information Exchange (C6226)" established in 1978 includes a 6349 kanji set, hiragana, katakana, Roman alphabet, Greek letters, Russian letters and other symbols. The kanji set is grouped into 2 levels, the first level a 2965 kanji set and the second level a 3384 kanji set. This means some 3000 kanji are considered to be enough for basic information exchange in Japanese. In this experiment, the kanji printer system T4100

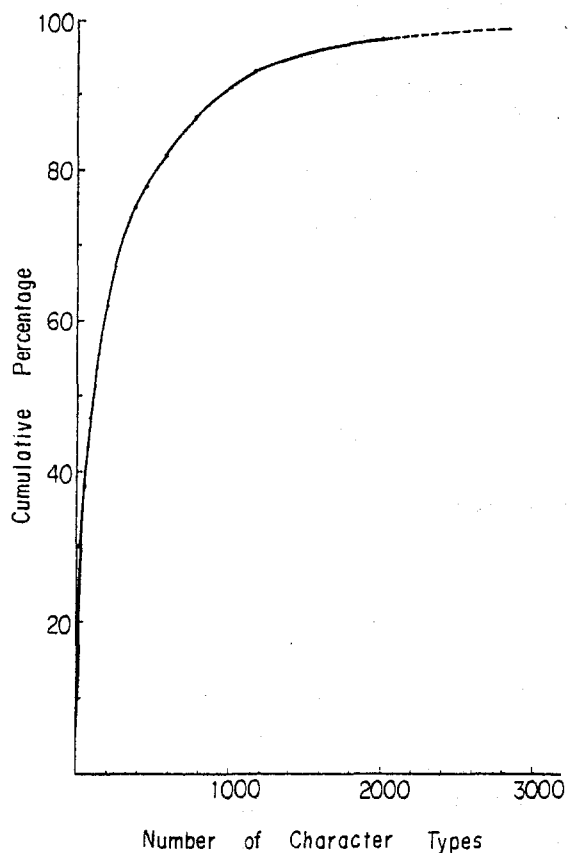


Fig. 3. Cumulative Percentage of Character Occurrences as a Function of the Number of Character Types

(Syowa Zyoho, Co., Ltd.) was used. A total of 8182 characters was available for this printer including 7360 kanji, hiragana, katakana, Roman alphabet, and other miscellaneous symbols. The system was developed 5 years before the establishment of JIS C6226.

As mentioned before, the million-character file included 5231 different Chinese characters. Among them, 295 were found to be unprintable (because they were not found in the T4100 system). The fonts of those 295 characters were designed and incorporated into the T4100 system. Later, when JIS C6226 was established, some of those 295 characters were found in the second level of the kanji set, namely 藝 (frequency 773), 門 (581), 黨 (563), 隨 (345), 證 (343), 缺 (189), 餘 (178), and 絲 (158). Fig. 4 shows the frequency of the remaining 287 characters. Their total frequency numbers 1100, which is 0.1% of the million-character file. This fact indicates that Chinese characters and kanji still overlap closely in modern Chinese and Japanese. (It should be noticed that the simplified Chinese characters are out of this scope since they did not exist at the so-called May 5 Movement period.)

THE CONCORDANCES

Besides the text itself, the Princeton million-character file contained information on the title, the author, the sentence numbers, and other miscellaneous editorial symbols (such as

Frequency	No. of Character Types
554	1 (狗)
228	1 (桌)
134	1 (碰)
128	1 (觸)
100 - 51	7
50 - 31	7
30 - 21	8
20 - 11	21
10 - 5	37
4	11
3	34
2	41
1	117

Fig. 4. Frequency Distribution of Chinese Characters which are not Found in the Kanji Set

marks to indicate proper nouns). Extensive work had to be done to interpret and reform editorial symbols. Fig. 5 shows the edited text sentences from the million-character file. After this editorial step and incorporation of Chinese character fonts to the T4100 kanji printing system, the concordance compilation process was started. Since we have had experience with the automatic compilation of one-million line concordances in Japanese [5], not many technical difficulties were encountered, except some malfunctions of our old kanji printer. Discussions on the salient features of those Chinese concordances follow.

Key Words

KWIC index style has been adopted as the form of Chinese concordances, since it is one of the most fundamental styles for computer-generated concordances. Because there is no clear segmentation of words in Chinese, and because one character represents a fairly sizable amount of information, each character was chosen as a "key word". Furthermore, no elimination of "non-key words" were made. Every character (including punctuation) was chosen as a key character. In this sense, the concordance may be named as "All characters in context" index. Consequently, one million character data required one million lines of index.

Contexts

One of the deficiencies of the KWIC index style is that the context each line can show is limited to its line length. We could afford 55 characters for the context. Since one or two Chinese characters represent a word, this length

can accommodate more than 30 words of information in English.

Reverse Sorted Index

Two types of KWIC index have been produced. One is for the normal type, in which all lines are sorted in the ascending order of the Chinese standard telegraphic code of key characters (plus 7 succeeding characters). Fig. 6 shows an example page from this type of index. The other is the so called "reverse sorted" index. The major key for this type is the same as that of the normal type. The minor sort keys are, the characters immediately preceding the major key. Thus all lines for one key character are listed in the ascending order of the code for the character immediately preceding the key character and so on. Fig. 7 shows an example page from the reverse sorted concordance.

CONCLUDING REMARKS

The two sets of modern Chinese concordances can be reached at the National Inter-university Research Institute of Asian and African Languages and Cultures, Tokyo University of Foreign Studies. It should be noted that a concordance of one million lines amounts to over 25,000 pages (actually it counts for 27,341) or 50 volumes of a 5cm-width paper file. Before printing the whole index, engineers recommended linguists to use COM technique, but in vain. A microfiche version should have been produced for portability. Analysis of the concordances have just got off the ground. The resulting papers are expected to follow.

作家	作品	文番号	原文
艾蕪	石青	1	早上太陽仍像往天一樣，把晴美的陽光抹上滿峽的樹林，叫帶露的樹葉草葉都亮得耀人的眼睛。
		2	只是石青嫂子的心上却陰暗極了，陰暗得像夏季烏雲滿佈的天空一樣，隨時都會兩點似的落下淚來。
		3	看見屋裡踢倒的板凳，打爛的燈，再看見門前地裡一片亂踏的足跡。
		4	菠菜的葉子，踩來變成爛泥；番茄踩成一灘一灘的紅漿。
		5	那些紅漿很使石青嫂子疑心，怕是夜來扭扭的時候，他身上流出來的血。
		6	對河山腰上的汽車公路，一乘長途汽車馳過以後，便比平日還要靜寂，簡直靜寂得可怕，滿山禿露的亂石，在陽光下面更加顯得蒼老醜陋，彷彿一些生癩瘡的禿頭似的。
		7	人工整過的公路，隱藏在亂石裡面，一種原始的荒涼的氣氛，越發強烈地流露出來。
		8	有石青的時候，她從來沒有感到過她這間山峽中唯一的茅屋是孤獨的，寂寞的，可怕的。
		9	她只覺得面臨小河，背靠山嶺的一片斜坡，給予她無限的繁忙和勞碌。
		10	她終天頭上包著一張藍布帕子，不是拿鋤頭挖地，鐮刀割草，就是手腕上掛個籃子，採摘什麼東西。
		11	晚上星子都現在山峽的高空了，樹林茅屋全隱藏在輕霧裡面，小的孩子，坐在門前哭著媽媽的時候，她還在地裡摘著苦瓜豇豆或是茄子辣椒，準備明天一早挑到五里以外的鎮上去賣，好換點米回來。

Fig. 5. An Example from the Edited Text

，我國生產資料所有制的社會主義革命，已經基本完成。	這	是一個巨大的變革。但是，僅有經濟制度上的這個變革；	劉少奇在北	5 6
國是一個社會主義的大國，但又是一個經濟落後的窮國，	這	是一個很大的矛盾。要使我國富強起來，需要幾十年艱苦	毛澤東關於	6 0 2
的詩就是一種不負責的東西了，不負責的東西是好的嗎？	這	是一個很重要的問題，所以，第一種主張，就側重在這種	聞一多詩與	1 8
和人民都在實際上懂得這個問題。為什麼現在又有人覺得	這	是一個新問題呢？這是因為過去國內外的敵鬥爭很尖銳	毛澤東關於	1 3 6
如何把地下發掘的資料，與傳下來的紀錄資料連綴起來。	這	是一個普遍的考古問題。也是在中國區域，考古家所面臨	李濟 再談	1 9 8
倒是蹲在家裡坐冷板凳好些。那第三科科長沒有管理他，	這	是一個沉悶而少話說的青年人，油黑的面孔上生著幾顆面	沙丁 代理	2 7
頭上坐著，沐浴著藍色的霧，漸漸地感到了老年的沉重。	這	是一個沒有月色的初夜。沒有遊人。衰草裡也沒有蟋蟀的	何其芳遲暮	7
受人注意而後竟受到擱頭，這便是一個解脫，不革命。	這	是一個真實，我們應該有勇氣來承認這真實，承認這失敗	茅盾 從牯	1 1 3
們要進行大規模的建設，但是我國還是一個很窮的國家，	這	是一個矛盾。全面地持久地厲行節約，就是解決這個矛盾	毛澤東關於	5 8 4
中一般。三面都是山，像半個環兒擁著；人如在井底了。	這	是一個秋季的薄陰的天氣。微微的雲在我們頂上流著；岩	朱自清溫州	4 0
底變局。中國人本來城裡人，到此時忽然成為鄉下人了。	這	是一個空前底變局。這是中國人所遇到底，一個空前底挫	馮友蘭辨城	6 8
關於正確處理人民內部矛盾的問題，	這	是一個總題目。為了敘述的方便，分為十二個小題目。在	毛澤東關於	1
表表示出來的。我們仍以下面這個三段論定理為例QQ。	這	是一個複合命題，它有三個原始命題Q。這三個原始命題	金岳霖論推	5 5 2
，這自從工業革命之後在西洋所發生的那一套生活方式，	這	是一個豐裕的經濟。我並不覺得自己配談西洋文化，我缺	費孝通中國	1 7 1
沒看清那是誰，已經把那人抓過來摔在地上。他斷定了	這	是一個賊。「多多頭。打死我也不怨你，只求你不要說出	茅盾 春蚤	3 2 5
來，坐在炕沿：「柴火燒沒了。」女人說，睜老趙一眼。	這	是一個跟他吃盡千辛萬苦，也不抱怨的好心眼的小個子女	周立波暴風	4 8 0
們開始看到了現實主義的創作要求底多方面的發展現象。	這	是一個重要的轉換期，因為，在社會學的意義上說，思想	胡風 文藝	2 7
放在枕邊的皮包拿過來，從中取出一個摺子，他叫她瞧。	這	是一個銀行的摺子，淨存一欄上，七千五百二十三元。外	張天翼報復	3 4 5
從前採用的黨內鬥爭方法叫做「殘酷鬥爭，無情打擊」。	這	是一個錯誤的方法。我們在批評「左」傾教條主義的時候	毛澤東關於	9 9
髻上那朵「蛋花」，跟鵝毛一塊掉在「蛋簾」的迎兒上。	這	是一個隆重的儀式。千百年相傳的儀式。那好比是誓師典	茅盾 春蚤	2 4 0
不巧撞著了一個巡路的小長毛，當時沒法，只好殺了他，	這	是一個「結」。然而從老通寶懂事以來，他們家替這小長	茅盾 春蚤	3 3
處」起來的。韓愈談他自己做古文，「惟陳言之務去。」	這	是一句最緊要的教訓。語言跟著思想情感走，你不肯用俗	朱光潛咬文	8 6
的片子就是今天晚上……您不是早就說過？……」阿寶摸不清	這	是一回甚麼事，粉臉太太驟然添上了一臉怒色，圓胖的鼻	王統照小紅	2 0 9
旁邊還有一株牽牛藤在晚風里微微擺動它的柔軟的腰肢。	這	是一幅靜的，美麗的，幻想的圖畫。我不覺痴痴地望著它	巴金 窗下	1 5 0
初見你的時候，那就老得多了。」這是一件很巧的事。「	這	是一張尺多寬的小小的橫幅，馬孟容君畫的。上方的左角	朱自清溫州	1
主義，似乎行政事務做得多了，所以就放鬆了思想工作。	這	是一所三樓三底兩夾廂的上海式樓房。鳳二爺住樓上的客	梅蘭芳舞台	4 3
麼理智都要躲避的，他未經考慮一下就和你訂婚了。可是	這	是一方面的理由，但不是全面的理由，甚至不是主要的理	周揚 整頓	1 7 7
未經考慮一下就和你訂婚了。可是這是一時的，你記著，	這	是一時的，你記著，這是一時的，過後，熱情一冷了下去	張天翼報復	2 3 9
抗侵略和反抗壓迫的人民戰爭。人民戰爭必將蓬勃發展，	這	是一時的，過後，熱情一冷了下去，他會覺得他自己滑稽	張天翼報復	2 3 9
大珍，帶上門出去。沿著荷塘，是一條曲折的小煤屑路。	這	是一條客觀規律，既不以美帝國主義者的意志為轉移，也	林彪 人民	2 3 6
地築壘的臂，要將他的妻孥奴子推到不可知的安全地方。	這	是一條幽僻的路；白天也少人走，夜晚更加寂寞。荷塘四	朱自清荷塘	5
當什麼都好了，就可以不費氣力享受現成的幸福生活了，	這	是一條悠長而異樣的路。宿羊山雖遙遙在望，但它的倩影	蕭乾 宿羊	2 9
他並不覺得這有什麼說不過去；有時候揣他，他還覺得	這	是一種不實際的想法。我國少數民族有三千多萬人，雖然	毛澤東關於	3 7 7
主要的創作力量移用到文藝行政工作方面去的危險傾向，	這	是一種優越，那些拉破車的根本就用不上電石燈。現在，	老舍 駱駝	4 5 6
翼戴政府，擁護政府，天下清平，門第亦同享安泰之樂，	這	是一種喧賓奪主，殺雞求卵的辦法。沒有作家，沒有作品	周揚 整頓	1 8 2
以，第一種主張，就側重在這種宣傳的效果方面，我想，	這	是一種大氣度。他們在政治上常抱一種領先的姿態。西漢	錢穆 中國	3 0 0
有了空前廣泛的發展，許多民族獨立國家已經建立起來。	這	是一種對於詩的價值論者。好些人念一篇詩時是不理會他	聞一多詩與	1 9
從朝鮮歸來的人，會知道你正生活在幸福中。請你意識到	這	是一種巨大的世界力量。一切被壓迫民族的獨立自主的要	劉少奇在北	3 3
的作家；假如你為小資產階級訴苦，便幾乎罪同反革命。	這	是一種幸福吧口巴，因為只有你意識到這一點，你才能更	魏巍 誰是	1 1 6
領導知識分子進行文化建設，我們就天然地不會犯錯誤。	這	是一種很不合理的事。現在的小資產階級沒有痛苦麼？他	茅盾 從牯	1 3 4
樣的死一次。這已超過了忍受的問題，在他們，十分覺得	這	是一種很危險的想法。而在有些地方，我們有一些同志正	周恩來關於	3 5 2
6 6 3 8	這	是一種恥辱。這裡一直不會早下去的，但是，他們等不得	端木良大江	8 9
	這	(這是一這是)		

Fig. 6. A Page Example from the Chinese Concordance (Normal Style)

動的唯心論哲學家，他們也有一套唯物論與唯心論的統一之愛」，自從人類分化成為階級以後，就沒有過這種統一成唯心論統一於唯物論，而不是相反。由於這個對立統一事物的矛盾法則，即對立統一的問題，我們可以總起來說幾句。事物矛盾的法則，即對立統一的問題，得到了如下幾點認識。第一，要認真貫徹對立統一鬥爭的統一，於重大現實毛主席所提出的關於鬥爭與統一由各個有關部門直接負責，而由中央宣傳部負責進行統一一樣罵也是空的。現在是「陽春白雪」和「下里巴人」統一死。所以便懶著不動了。但是終於還是得走，想得到萬一是一水底月，鑽中天，但在自己的心裡能不能否認總含有萬一仰一仰，像畫家設完一層色那麼退後看看。然後，又逐一三個女兒的悲慘結局，你的懷疑慢慢會變成惆悵。在園丁臉細眉，白淨皮色，太太模樣的女人，年紀不過二十六七聲音亂叫了。但是那位作調人的警察却冷笑，扳著陳老七朱三太的上門討利息，他記起還有兩注存款，橋頭陳老七略顯渾濁，有出山泉水的意思。若溯流而上，則三丈五丈南京大學教授和我談話的地方，即離開左面的斷崖數十丈可以跨過似的。然而實景中並沒有石條，只是相距若干丈指著他們誇耀似地對我們說：「這兩位高僧，是我們方丈抱著羊竿，哀喚著「天妃」的慈聲。我的心舟在起落萬丈城基礎，同時却也有它的社會基礎。所以，你如把「張三生皺著高鼻上的皮，搖頭繼續說：「撒爛污。十個阿孺三鄉下你那些田，早早脫手得好。自從改了民國，接二連三的聲音向黨民問道。「我剛剛看見抬傷兵進城，接二連三玩兒去，誰忍耐站在先先生書桌前，晃著身子，背早上上繞上領導前進者，總是少數特殊分子，遺落在後追隨不上。假使沒有了頭顱，却還能做服役和戰爭的機械，世上是共產的事實，不是言論。歐洲之所以罵乎我們中國之上方面的研究，成為革命黨正確地決定其政治上和軍事上。社會主義已經成為世界體系。占人類總數三分之一以上得到的，有關這兩個基本問題的答案及解脫，證實了以上這時還不改變策略，它將坐困於為數至少在二十萬人以上，聯保主任跑來報告，說是索橋邊已經扣留下二十個以上彈會怎樣地鑽進他底身體，他也感到一種超乎壓力以上是人民大眾呢？最廣大的人民，占全人口百分之九十以上的革命的功利主義者，我們是以占全人口百分之九十以上的戀愛已經失去了莊嚴的時候，隨時隨地可以有接吻以上和他們在江北的地都管計在內，他們三家都有一千塊以上外。單純的過程只有一對矛盾，複雜的過程則有一對以上在小屯的建築遺址中，地基部份的夯土，有夯至十層以上那人頭後面突出一個硬硬的小凸餅，顯然是一個中年以上

的理論。我們前面已經批評到的黑格爾，就是突出地站在唯心的愛。過去的一切統治階級喜歡提倡這個東西，許多所謂聖的原理意義特別深遠而豐富，完全可以適用於說明唯心論與法則，是唯物辯證法的最根本的法則。列寧說：「就本來法則，是自然和社會的根本法則，因而也是思維的根本法則，是唯物辯證法原理，必須承警唯心論與唯物論是可以相互醜的基本原理，却有偉大的啟示性，完全可以應用來討論唯物監督。中央宣傳部應該經常檢查各部門各地區執行中央問題，是提高和普及統一的問題。不統一，任何專門家的機會，假使碰見水呢？於是從炎熱裡向前走。重傷的隨走的希望。因此，我自從十四歲以後便不願從速訂婚。我的父的軀開，把另一面的魚們攆齊。又往後仰身端詳了一番，回的樸實言語中，佻說中的古怪老頭和他的女兒重新復活過來的光景。我不會抽煙，我也沒有買什麼東西，我只求她換一的肩膀道：「我勸你少找點麻煩罷。到那該，中什麼用。你的二百元和張寡婦的一百五十元，總共十來塊錢的利息，都的深潭皆清澈見底。深潭中為白日所映照，河底小小寬石子的地方，我的確看到有一根不很大的石條伸出在空中，照相的兩個斷崖，我們所登的便是左面的斷崖。我想：這地方叫的師兄，年紀都快八十歲了，是從城裡某公館裡回來的。」思潮中震蕩時，母親。縱使你在萬里外，寫到「母親」兩的爸爸巴」說成「張三他爸爸巴」，或者你不說「他還沒來的血，有九個不清潔。毒來些，醜醜來些。要不得，要不得的打仗，倒嘗有一年間過，把地面上糟塌得不成樣子，中間的，不曉得有多少，」黨民激動地用顫抖的聲音說，「真可的多難的書？啊，飛。不是那在樹枝上矮矮的跳著的麻雀兒的，依然混同一色，那纔是社會整體之真骨幹。結果詩人，情形就何等地醒日呵口可，這時再不必用什麼制帽勳章來的，不是政治哲學，完全是物質文明。因為他們近來的物質戰略戰術方針的重要方法之一，是一切共產黨人都應當注的社會主義國家，已經組成了以蘇聯為首的社會主義陣營。的說法。整理田野考古發掘所得的資料，一個最迫切的問題義和團民的包圍之中，而其結果就是這個王朝的統治被粉災民了。才一走進屋裡，還來不及把醃肉掛向籬壁的竹釘的恐怖。電燈光閃閃地刺痛他底眼睛。「這燈光。」他煩躁的人民，是工人，農民，兵士和城市小資產階級。所以我們最廣大羣衆的目前利益和將來利益的統一為出發點的，所行為的。去年是不是我和你當了你未婚夫的面互相擁抱，的好地，條通和黃土包子還不管在內。街裡的福來德燒鍋，的矛盾。各對矛盾之間，又互相成為矛盾。這樣地組成客觀的；好些地方，把夯過了的土毀了再夯大力。這一營造方式的女人，挽著髮髻彭吉。她的心往下一沉。她知道那是金花

賀麟 論唯 1 5 3
毛沢東在延 3 6 4
賀麟 論唯 1 4 7
毛沢東矛盾 1
毛沢東矛盾 4 8 7
賀麟 論唯 1 4 5
賀麟 論唯 1 2 5
周恩來關於 3 6 5
毛沢東在延 2 7 9
端木良大江 5 3
郭沫若黑貓 1 1
老舍 黑白 2 0
師陀 果園 1 1 0
巴金 生與 8
茅盾 林家 8 2 4
茅盾 林家 2 0 2
沈從文迎城 5 6
豐子愷廬山 2 6
豐子愷廬山 2 1
郁達夫遲桂 4 0 6
謝婉瑩寄給 1 6
董同龢國語 7
吳組湘官官 1 6 4
張愛玲金鎖 5 7 4
巴金 家 2 1
徐志摩想飛 7
錢穆 中國 2 2 7
魯迅 春末 5 5
孫逸仙三民 4 2 5
毛沢東矛盾 3 5 1
劉少奇在北 3 0
李濟 再談 1 9 7
翦伯贊養和 1 2 2
沙丁 代理 2 0 0
巴金 家 7 8
毛沢東在延 1 3 8
毛沢東在延 2 7 5
張天翼報復 2 3 5
周立波波風 1 4 6
毛沢東矛盾 3 5 9
李濟 毀墟 5 7
張愛玲秧歌 3 3 2

4 1 0 4 的 (一的一上的)

- 1 6 3 0 5 -

Fig. 7. A Page Example from the Chinese Concordance (Reverse Sorted Style)

REFERENCES

1. Kierman, F.A. and Barber, E.: "Computers and Chinese linguistics", Unicorn, No. 3 (1968)
2. Boltz, W.G., Barber, E. and Kierman, F.A.: "Progress report on Pai-hua-wen computer count and analysis", Unicorn, No. 7, pp. 94-138 (1971)
3. Hashimoto, M.J., et al.: "A grammatical analysis of the Princeton million-character computer file", Bulletin of the Chinese Language Society of Japan, No.222, pp. 1-16,36 (1975)
4. Hashimoto, M.J.: "Computer count of modern Chinese morphemes", Computational Analysis of Asian and African Languages, No. 7, pp. 29-41 (1977)
5. Uemura, S.: "Automatic Compilation and Retrieval of Modern Japanese Concordances", Journal of Information Processing, Vol. 1, No. 4, pp. 172-179 (1979)