

Full Text processing
of
Chinese Language

- an experimental system for studying Chinese History Literatures

Ching-Chun Hsieh

Research Fellow,

Institute of Information Science Academia Sinica, Taipei, ROC

Introduction

Full-text databases and associated processing technology are becoming increasingly important for library automation, on-line information service systems, and office information systems. They also play a key-note in developing new computer applications where automation of the filing and search procedures will free workers from these time-consuming tasks.

Some major research and development topics related to full-text database systems now underway include the studies of document representation, access methods for text, text analysis (in natural language), and the new hardware for text retrieval. An excellent review of text access methods and retrieval hardware was given by Faloutsos [ref. 1]. And, an excellent position paper about full-text databases was given by Tenopir [ref. 2]. Since full-text of document is usually in natural language form, the studies mentioned above are unavoidably language dependent, ie they may not directly applicable to those full-text processing systems which are not in English. This fact becomes more obvious while Chinese full-text database systems are confronted, although at present, only quite limited and scattered effort has been devoted to the development of full-text database systems in Chinese language.

In this paper, an experimental Chinese text processor(CTP) for studying history literatures will be presented. The system provides both free-text and controlled vocabulary search for the text. The goal of this study is to find out the feasibility of using computers, especially small or microcomputers, as a tool for history studies. Our experiment shows that by applying existing technology, such a system is highly feasible.

The System

The project of CTP is a part of a long-term research plan, namely Computer Applications in Humanities(CAH) conducted by the Computing Center of Academia Sinica, Taipei, Taiwan, ROC. The project CTP was started in July, 1984. At present, the CTP is still an on-going project in its second phase. Therefore, as time goes, the system is unavoidably subjected to change. For the time being, we have two operational prototypes on a mini and a 16-bit micro, respectively. ✓

A basic conceptual configuration of the CTP is shown in Figure-1. It can be considered basically as a workstation for humanity studies. As comparing the CTP to an information retrieval system, the major differences are: the abstract text-structure supporting mechanism, more supporting tools, and a reinforced application section which supports documentation generation functions, i.e. user's note-card system [ref. 7, 8], word processor etc. The most important supporting tool, at present stage, is considered to be the auto-indexing program. These major differences will be presented in the subsequent sections.

The CTP is designed to facilitate the access of full-text documents, creating user's own private index systems to both the original text and user's own note-cards, and the generations of new documents during/after studies.

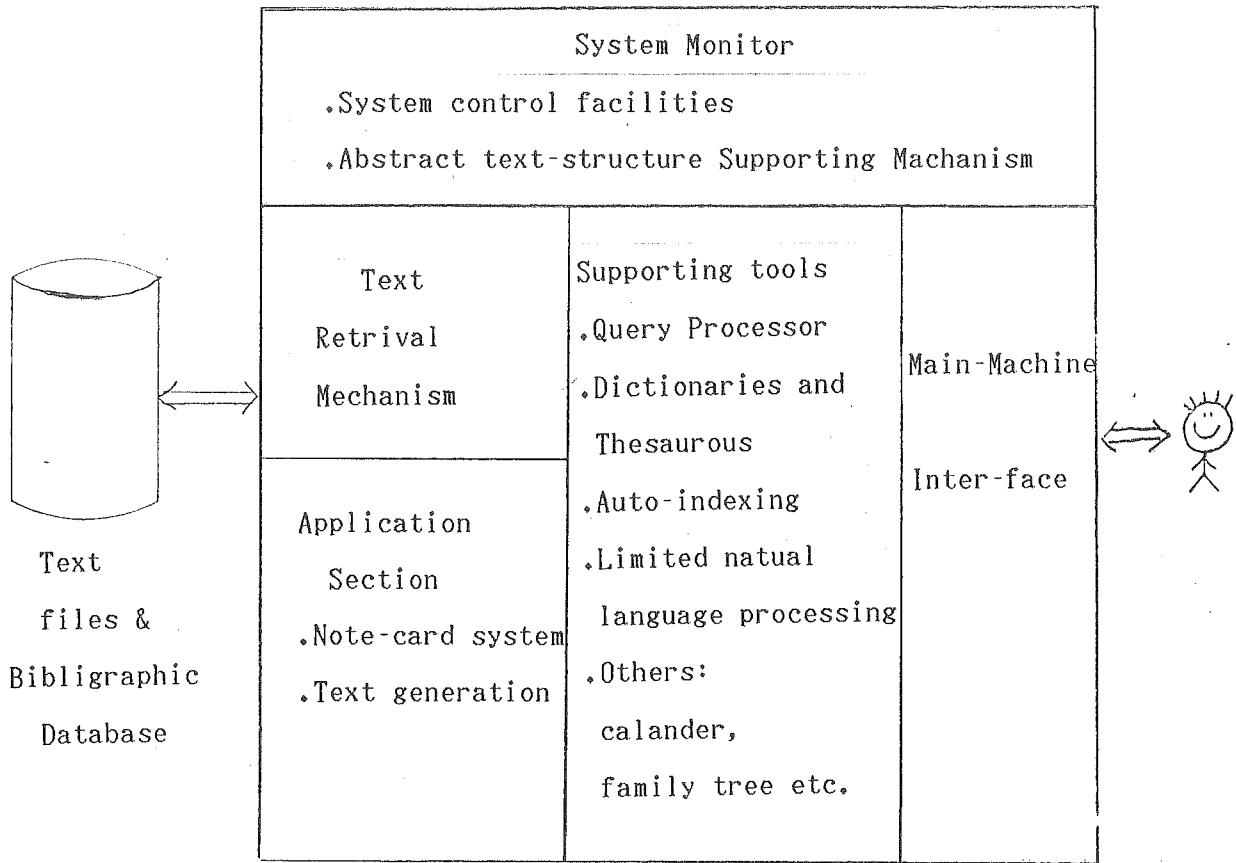


Figure-1 Base Configuration Concept of CTP

Text Representation

The object document in the full-text data-base under testing consists of all volumes of 食貨志 (for food and goods in meaning) from the history books of 24 Dynasties (廿四史) of ancient China. An example of 食貨志 is shown in Figure-2. It is of the form of an ordinary text book. In this paper, we will use full-text or text interchangeably to indicate such a type of document. A complete list of the 食貨志 and associated number of characters in each book is shown in Table-1. [ref. 3]

One of the reasons to choose 食貨志 as object document is that there is no tables and figures in them. Therefore, its logical structure is simple and it is very nature to use a tree-structured representation scheme in computer to describe the text organization of 食貨志. The tree-structured representation of 食貨志 is shown in Figure-3. In Figure-3A, a logical structure of a generalized text-tree is presented. It is implemented as a part of the system monitor. In Figure-3B, an example of the physical structure of 食貨志 is shown. Each terminating node of this graph represents a paragraph of the text, respectively. Of course, for a more detailed representation, sentences may be considered to be the terminating nodes. But for the moment, we choose paragraph as the smallest unit used in the internal representation of text. If any address to the sentence level is necessary, it will be carried out by a text-scan and analysis program within a chosen paragraph.

The advantage of separating the logical structure and physical structure of the text is obvious: the system is generalized to have the capability of handling more tree-structured texts.

宋史卷一百八十一

志第一百三十四

食貨下三

會子 鹽上

會子、交子之法，蓋有取於唐之飛錢。眞宗時，張詠鎮蜀，患蜀人鐵錢重，不便貿易，設質劑之法，一交二緡，以三年爲一界而換之。六十五年爲二十二界，謂之交子，富民十六戶主之。後富民質稍衰，不能償所負，爭訟不息。轉運使薛田、張若谷請置益州交子務，以權其出入，私造者禁之。仁宗從其議。界以百二十五萬六千三百四十緡爲額。

神宗熙寧初，立爲造罪賞如官印文書法。河東運鐵錢勞費，公私苦之。二年，乃詔置交子務于澶州。轉運司以其法行則鹽、礬不售，有害人中粗草，遂廢之。四年，復行於陝西，

志第一百三十四 食貨下三

四四〇三

史記卷三十

平準書第八

／止／

（一）國幣官表曰大司農，官有平準令。張大司農，官有平準令承者，以均天下郡國輸販，書則發之，賤則買之，貴賤相權輸，歸于京都，故名曰平準。（*）

漢初貧弱，漢興，撥秦之弊，丈夫從軍旅，老弱轉輸饑，作業劇而財匱，自天子不能具鈞駟，二而令民鑄錢，將相或乘牛車，齊民無歲蓄。三於是爲秦錢重難用，四更令民鑄錢，五一黃金一斤，六約法省禁。而不軌逐利之民，蓄積餘業以稽市物，物踊騰糶，七米至石萬錢，馬一匹則百物價高漲金。八

（二）國天子駕四馬，其色宜齊同。今言國家貧，天子不能具鈞色之駟馬。漢律作「醇駟」，醇與「純」同，純一色也。或作「辟」，非也。

（三）國如「澤」曰：齊婦無有黃鵠，故謂之齊民。若今言「平民」矣。晉灼曰：「中國被放之民也。」陸林曰：「無物可蓄藏也。」

平準書第八

一四一七

Figure-2: An example of the text of 食貨志

Name of the books 書 名	Number of characters 字 數
漢 書 (上、下)	3 6 4 1 1
史 記	1 4 3 4 3
舊 唐 書 (上、下)	2 4 8 0 7
唐 書 (一至五)	3 1 3 8 0
隋 書	1 2 0 5 3
魏 書	1 1 6 1 5
晉 書	1 1 0 8 5
舊 五 代 史	8 8 1 0
遼 史 (上、下)	4 9 5 6
金 史 (一至五)	5 2 3 2 0
宋 史 (上、下)	2 1 4 9 9 4
元 史 (一至五)	6 8 5 3 3
明 史 (一至六)	6 8 5 1 3
Total number of characters 共 計	5 5 9 8 2 0

Table-1: All the volumes of 食貨志 and their number of characters

Access Method

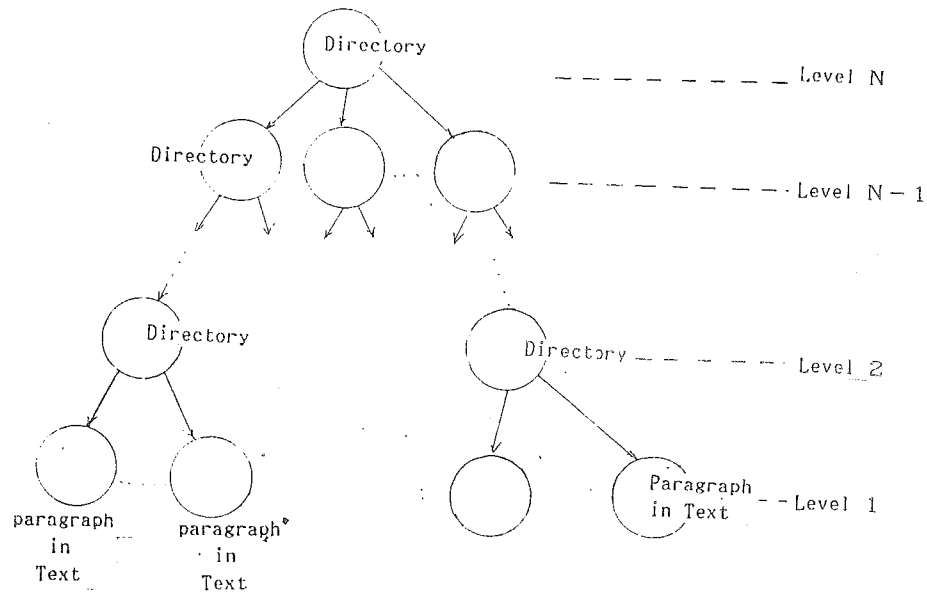
The system provides both the controlled vocabulary search and free-text search of the text database. At present, for controlled vocabulary search, only single key-word search is available. The keyword sets that the system provides include: the name of person(人名、字號等), the name of place, the name of official titles(職官名), the name of books appeared in the text, the expressions of time in terms of an era, a dynasty or an emperor(朝代, 時代, 帝號, 年號, 名號名等) and 10 classes of special terminologies such as 賦稅類(the class of tax terms), 鹽鐵類(the class of salt and metals terms) etc. A more detailed example can be found in Table-2. [ref. 3]

For each set of keyword, an inverted index file which maps each keyword onto the related paragraphs of the text is generated in the system generation phase. A program that automatically generates these inverted index files will be described in the next section. At present, only a manual-driven man-machine interface is available for users to retrieve text by these system keywords. No logical operators are available at this stage.

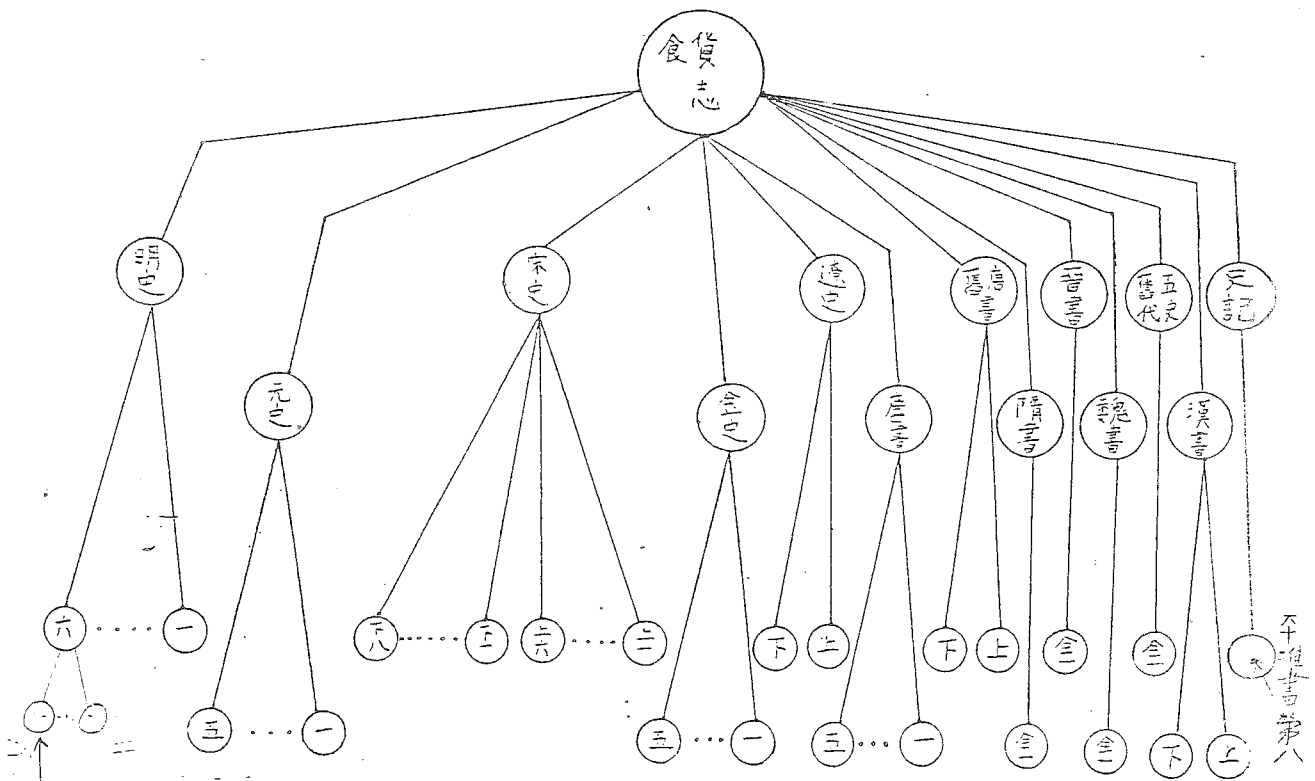
Besides these keywords provided by the system, one important theme of designing this system is to support user defined keyword searching facilities through the application section. In other words, this system supports two categories of keyword; one is in public domain to be used by all users and provided by the system, and another is merely created and used by a user in private manner.

term classes	Example																					
田制類	井 阡 度 世	田 陌 田 業	之 之 之	制 法 制	口 世 莊 荒	分 業 田	之 田	上 下 荒 營	田 田 田	墾 大 田												
水利類	纓 戍 屯 井	田 田 田		占 代 井 王	田 田 田		湯 籍 民 名	沐 田 田	邑 田	歸 公 私 公	受 上 中 下	田 田 田	步 晦 夫 屋									
賦稅類	什 初 履 十	一 稅 晦 一	晦 之 稅		賦 田 什 租	斂 租 五 稅	而 稅 一		三 什 田 口	十 一 租 賦	而 稅	一	算 民 租 田	錢 租 稅 租	什 賦 田 口	五 斂 租 賦						
雜稅類	課 租 兩 田	庸 稅 稅	調 法	秋 財 課 庸	稅 賦 役		調 租 庸	賦 調 法		兩 稅 稅 租	稅 錢 錢 賦		上 留 送 征	供 州 使 役	折 賦 正 貢	納 役 賦 賦	租 租 地 租	庸 、 調				
戶口類	六 漏 逃 雜	部 戶 戶 營	民 戶	賦 戶 男 婦	役 口 夫 人		奴 良 奴 男	婢 夫		老 小 癯 殘			樂 男 三 蔭	遷 女 長 附	隣 里 黨 復	長 長 長	力 里 三 屯	役 黨 長 民	戶 門 民 金 沙 遷			
錢財類	比 四 沈 比	輸 文 郎 輸 錢 錢		五 大 當 小	錢 錢 錢		五 貨 鼓 內	銖 泉 鑄 府		鑄 錯 朋 元	錢 刀 龜		白 鑄 赤 白	鹿 錢 側 金	幣 半 三 盜 復	兩 銖 鑄 小	錢 錢	赤 紫 赤 白	側 錯 側 金	之 錢 錢 三 品		
鹽鐵類	鹽 鹽 鹹 兩	法 麩 醃 池		顆 私 末 刮	鹽 鹽 鹽 鹹		煎 煎 竈 官	鹽 鍊 戶 鹽		鹽 糴 食 鹽	務 鹽 鹽 法		煎 青 青 白	鹹 白 鹽 鹽	鹽 食 大 甜	鹽 鹽 次 冷 鹽						
市糴類	六 除 輕 平	幹 貸 重 糴		均 常 贖 五	輸 平 禁 均	銅	入 甘 太 甘	粟 泉 倉 泉 倉		入 耐 均 平	殺 金 輸 準		耦 輓 常 拜	犁 犁 平 爵	買 贖 輕 顧	爵 禁 重 租						
漕運類	漕 漕 漕 含	運 米 吏 嘉 倉		河 柏 集 鹽	倉 倉 倉		轉 武 洛 河	運 牢 口 陽	倉 倉 倉	歌 上 綯 歲	支 填 舟 運	江 關	船 船		漕 水 車 運	輓 運 漕	太 北 陸 運	倉 丁	八 宿 轉 轉	遞 場 輸 運	水 飛 水 通	陸 輓 運 運
非財經	五 西 胡 鴻	帝 戎 部	三 皇	長 渾 鸞 六	城 和 坊			百 華 勇 寬	保 人 士 鄉	鮮 卑	八 十 軍 侍	丁 二 士 官	兵 丁 兵		偃 隆 玳 六	武 基 瑁 官 樓		夏 洛 突 吐	人 陽 厥 谷 渾	華 殿 東 修	人 最 宮 文	

Table-2: The 10 classes of special terminologies used as keyword set for text retrieval.



3A: A Logical Structure of the text tree



3B. The physical structure of 食貨志

Figure-3: Tree-Structured representation of text

The text-tree also provides a way of retrieval. It behaves like a content table of a book. Besides, it also provides proximity functions that will limit the search of text within a reasonably restricted scope.

For the free-text search, users are always been asked to define the scope of the search through the text-tree. The free-terms can be any length of consecutive characters. After a successful free-term search, the user may ask the system to update his/her private index file by adding the newly found entry.

As a final remark of this section, we like to point out that the nodes of text-tree are excellent entry points for inverted index files. Linking the text-tree to the indexing system can provide all information about the organization and associated attributes of the text to the retrieval system. And thus, greatly enhanced the performance of the system. For instance, a field in a node descriptor may be used to describe the length of the text covered by that node. And from this information, we may calculate precisely the time needed for a free-term search in advance. As another example, a broader-term in an thesaurus may be assigned to a node of the text-tree automatically if its successor nodes are indexed by most of the terms under that broader-term.

An auto-indexing method for Chinese text

The importance of auto-indexing ability to an information retrieval system is obvious and has been discussed in many papers [ref. 1, 2]. For Chinese text database systems, the auto-indexing ability is even more important and can be considered as a necessity for overcoming the feasibility requirements of implementing such a system. In this research, a fairly simple auto-indexing method has been proposed and, to our surprise, the results are especially encouraging. The method is described as follows.

In Figure-4, a block diagram of the auto-indexing system is presented. The inputs to the auto-indexing system are the text file and the keyword file. The first part of the auto-indexing system is a matching program that scans each character in the text file and matches it against a tree of words [ref. 4, 5] in the keyword file. An example of a tree of keywords are shown in Figure-5. As you see that all of the keywords that share the same leading characters will be organized in a tree. Therefore in the keyword file, there is a forest of keyword trees.

By this arrangement, the time needed for matching can be reduced. Let n denotes the number of character in text and T denotes the average time required to match a character against a tree of keywords, then, the total time required to build the index file is nT . Of course, the result is application dependent. It depends upon the content of the text and the complexity of the trees of keywords. For the text of ancient document, we found that each tree is rather simple, most of them are trieval, i.e. only one or two words in a tree. Since hashing technique can be easily applied to located the related tree for a specific character, the value of nT is

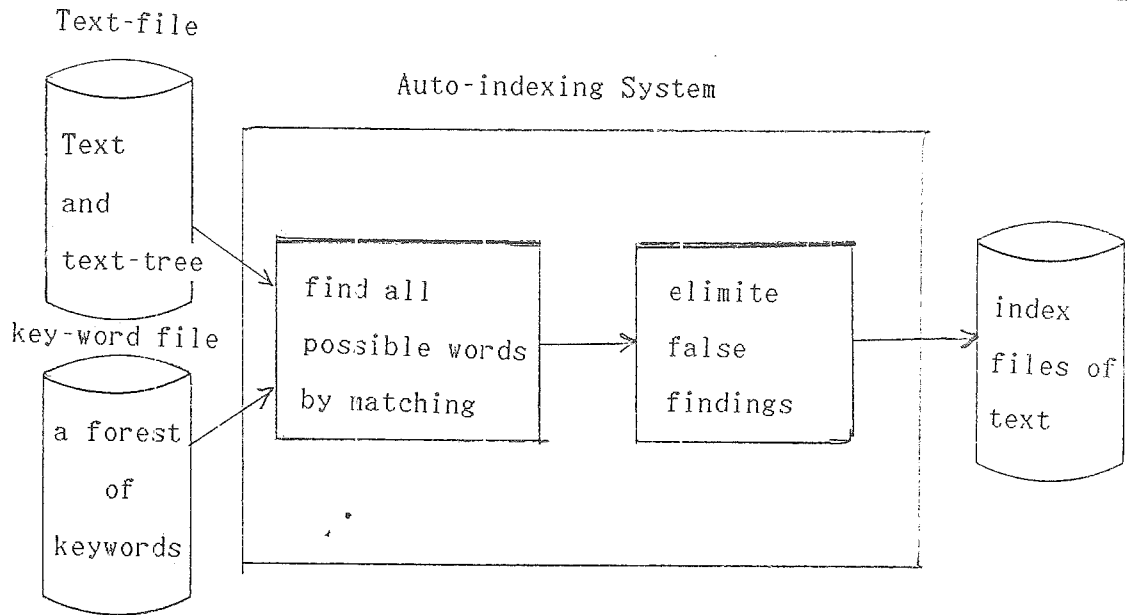


Figure 4 : A block diagram of Auto-indexing method

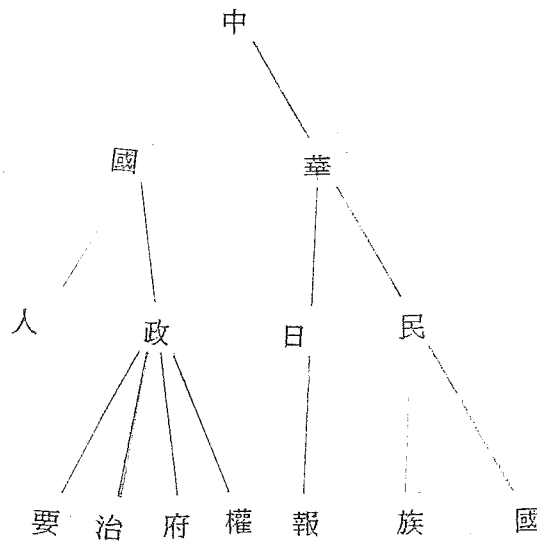


Figure 5 : An example of a tree of words, that share the same leading characters. The node with a circle indicate a termination of a word.

much less than the time required by a straight-forward matching method whose matching time required is $n.m.l.t$, where l represents the average length of keywords, m the total number of keywords in the forest and t the time required to match two characters.

As for the recall and precision functions of the retrieval method are concerned, an example is given in Table-3. It can be easily proved that according to our indexing method, the index file made by human professionals is a sub-set of the index file found by the machine. According to Table-3, the percentage of the false indexes are very low. It is usually less than 5%.

The second part of the auto-indexing system is a text analysis program that will eliminate some possible false indexes. The text analysis program consists of a rule-based analyzer [ref. 6] which will eliminate some false findings by examining the syntax information obtained during matching. An example is given in Figure-6. Since the number of false indexes are considerably small, this correcting process may be omitted in some applications.

```
#---- 以下為 auto-index 所建立的資料 ---- #
#
#      621 Mar 11 09:04 his2azzbk.ind      #
#      156 Mar 11 09:04 his2azzbk.iov     #
#     2622 Mar 11 09:04 his2azznm.ind     #
#      416 Mar 11 09:04 his2azznm.iov    #
#     5589 Mar 11 09:04 his2azznu.ind     #
#      104 Mar 11 09:04 his2azznu.iov    #
#     2070 Mar 11 09:04 his2azzof.ind     #
#       0 Mar 11 09:04 his2azzof.iov     #
#     6279 Mar 11 09:04 his2azzpl.ind     #
#      728 Mar 11 09:04 his2azzpl.iov    #
#     1587 Mar 11 09:04 his2azzyr.ind     #
#      572 Mar 11 09:04 his2azzyr.iov    #
```

Totle memory space : 20744 bytes

```
#--- 以下為 partial auto-index 所建資料 ---#
#
#      621 Mar 11 09:54 his2mzzbk.ind    #
#      156 Mar 11 09:54 his2mzzbk.iov    #
#     2622 Mar 11 09:54 his2mzznm.ind    #
#      416 Mar 11 09:54 his2mzznm.iov    #
#     5589 Mar 11 09:54 his2mzznu.ind    #
#       52 Mar 11 09:54 his2mzznu.iov    #
#     2070 Mar 11 09:54 his2mzzof.ind    #
#       0 Mar 11 09:54 his2mzzof.iov    #
#     6279 Mar 11 09:54 his2mzzpl.ind    #
#      208 Mar 11 09:54 his2mzzpl.iov    #
#     1587 Mar 11 09:54 his2mzzyr.ind    #
#      624 Mar 11 09:54 his2mzzyr.iov    #
```

Totle memory space : 20197 bytes

The difference : 547 bytes (about 2.6%)

Table-3 : A comparison of Auto-indexing and human corrected partial auto-indexing of all index files in 遼史食貨志.

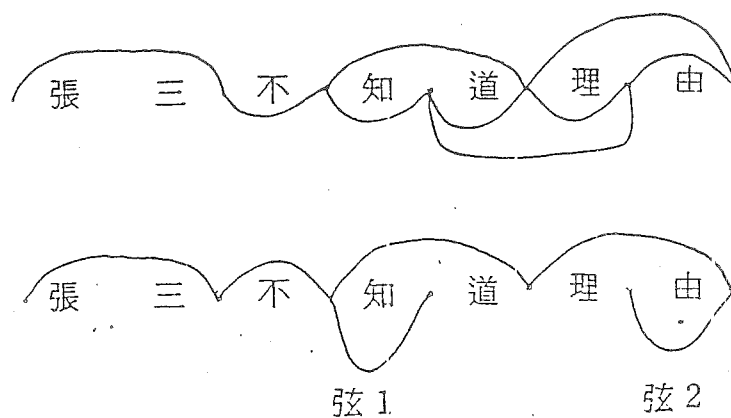


Figure 6 An example of eliminating false indexes by text analysis

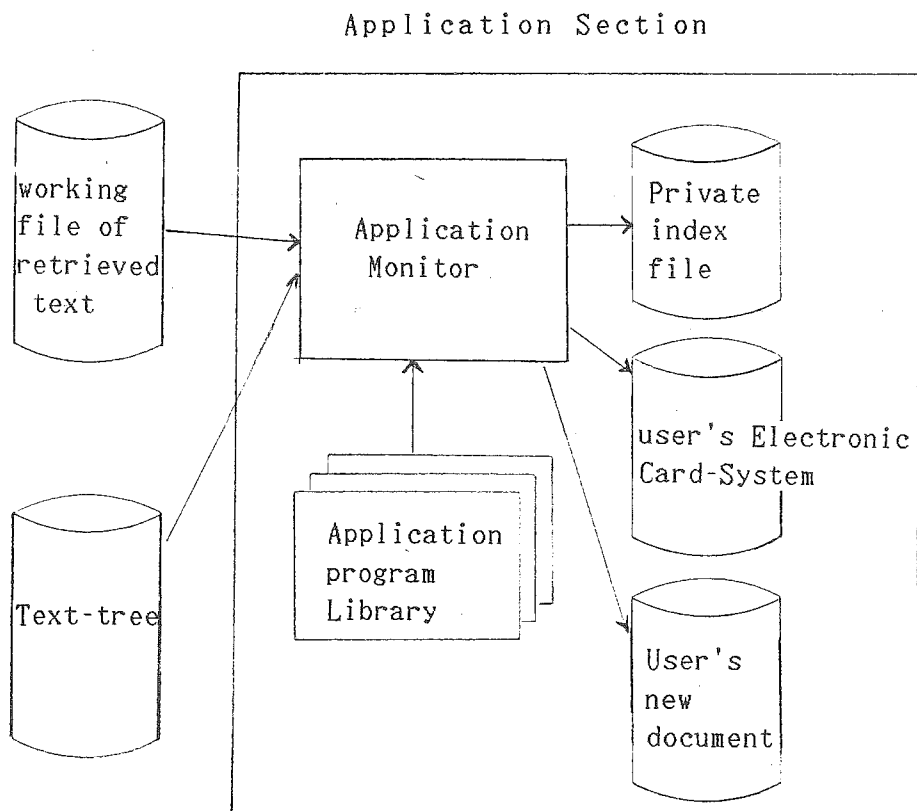


Figure 7 : A block diagram of Application Section and its operational environment.

Application section

In our system, an application section is designed for interfacing the retrieved text to user's applications. At present, the application section provides two major functions; to establish user's own private indexing system and to generate a card-system after user studied the text. A block diagram that shows the application section and its related system parts is shown in Figure-7.

The application section will be activated by user's request during retrieval. It provide users a manual-driven interface. Through the manual, the user can establish his/her own private index file, and generates a deck of cards on which a memo or a digest of the original text can be written. After a deck of cards has been established, the user can use them to create new document, such as essays, study reports, etc.

Conclusion

For illustration, two examples of search are given in Figure-8 and Figure-9, respectively.

By the experience of project CTP, we like to state the following comments as the conclusion of this paper :

1. There is a high demand of automating tool books usually referenced by humanist. Most of these tools are knowledge-based in nature. And thus, the automation involves the application of knowledge-based systems and AI technology. This is a field that open to us for future development.
2. Natural language processing capability will directly enhance the system performance and improve the friendliness of man-machine interface. Even with very limited natural language processing techniques, such as morphological analysis, do help a lot.
3. The representation of text is still worth doing research.
4. Application section needs further enhancement. Word processor, text editor spread-sheet, note-card system, statistical package, etc, can greatly facilitate the power of text-processing.

Acknowledgement

The author would like to express his sincerely thanks to the Council for Culture Planning and Development, Executive yan (行政院文化建設委員會) for supporting this project.

Figure-8 An example of search by "董仲舒"

代號	內容	書名	頁數	段數
1	董仲舒	漢書	1137	2
2	仲舒	漢書	1137	2

請鍵入始印代號：2 後按換行鍵
 請鍵入終印代號：2 後按換行鍵
 請鍵入列印格式：1(0:直印,1:橫印)後按換行鍵
 請鍵入列印元件：0(0:顯示幕,1:印表標)後按換行鍵

1史記 2漢書 3晉書 4魏書 5隋書 6舊唐書 7唐書 8舊五代 9 0
 【中文】【半形】 字根：_

人名查詢

請鍵入查詢書名： 後按換行鍵
 請鍵入始詢頁數： 後按換行鍵
 請鍵入終詢頁數： 後按換行鍵
 請鍵入查詢人名：董仲舒 後按換行鍵

請鍵入列印元件：0(0:顯示幕,1:印表標,2:不印,0:放棄)後按換行鍵
 找到了!!! 2 筆 08:52:28 08:53:55

1史記 2漢書 3晉書 4魏書 5隋書 6舊唐書 7唐書 8舊五代 9 0
 【中文】【半形】 字根：_

仲舒 漢書 頁1137 段 2

是後，外事四夷，內興功利，役費並興，而民去本，董仲舒說上曰：「春秋它穀不豐，至於麥禾不成則書之，以此見聖人於五穀最重麥與禾也。今關中俗不好種麥，是歲失春秋之所重，而真生民之具也，願陛下幸詔大司農，使關中民益種宿麥，令毋後時。」(一)又言：「古者稅民不過什一，其求易共；(二)使民不過三日，其力易足。民財內足以養老盡孝，外足以事上共稅，下足以畜妻子極養，故民說從上。(三)至桑則不然，用爾秋之法，收帝王之制，除井田，民得賣買，富者田連仞佰，貧者亡立錙之地，又鄙川澤之利，管山林之饒，(四)荒淫越制，踰後以相高，國有田君之尊，里有公侯之富，小民安得不困？又加月為更卒，已復為正，一歲地或一歲力役，三十倍於古；(五)田租口賦，鹽鐵之利，二十倍於古。(六)或耕桑之田，見稅什五。(七)故貧民常衣牛馬之衣，而食大豪之食，重以倉粟之更，刑戮妄加。(八)民愁亡聊，亡適山林，轉為盜賊，結衣半道，斷獄歲以千萬數。漢興，循而未改，古非田法難奪奪行，宜少近古。(九)限民名田，以澹不足，(一〇)塞并兼之路。鹽鐵皆歸於民，去奴婢，尊桑之祿，(一一)薄賦斂，省徭役，以寬民力，然後可善治也。」仲舒死後，功費愈甚，天下虛耗，人復相食。(一二)

(一) 師古曰：「宿麥，謂其直饒冬。」
 (〇) 師古曰：「印上一頁(2)印下一頁(9)停止？」

1史記 2漢書 3晉書 4魏書 5隋書 6舊唐書 7唐書 8舊五代 9 0
 【中文】【半形】 字根：_

官名查詢	
請鍵入查詢官名:	後按換行鍵
請鍵入始詢頁數:	後按換行鍵
請鍵入終詢頁數:	後按換行鍵
請鍵入查詢官名: 尚書	後按換行鍵
請鍵入列印元件: 0(0:顯示幕,1:印表機,2:不印,9:放棄)後按換行鍵 找到了!!! 33 筆 11:55:59 11:57:31	
1史記 2漢書 3晉書 4魏書 5隋書 6新唐書 7唐書 8舊五代 9	0
[中文] [半形]	字根:

代號	內容	書名	頁數	段數
17	尚書度支郎中	魏書	2859	2
18	尚書僕射	魏書	2860	2
19	尚書令	魏書	2863	4
20	尚書令	隋書	680	2
21	尚書左丞	隋書	676	4
22	尚書左丞	新唐書	2122	4
23	尚書左丞	唐書	1344	2
24	尚書省	舊唐書	2086	1
25	尚書省	新唐書	2089	2
26	尚書省	舊唐書	2118	1
27	尚書省	唐書	1369	3
28	尚書省	唐書	1372	1
29	尚書省	唐書	1400	5
30	尚書省	唐書	1402	1
31	尚書右僕射	舊唐書	2118	4
32	尚書比部	唐書	1347	3

總數(Y/N)

1史記 2漢書 3晉書 4魏書 5隋書 6新唐書 7唐書 8舊五代 9 0

[中文] [半形] 字根:

代號	內容	書名	頁數	段數
1	尚書	晉書	782	1
2	尚書	晉書	780	1
3	尚書	晉書	792	3
4	尚書	晉書	793	2
5	尚書	晉書	793	2
6	尚書	晉書	793	2
7	尚書	晉書	798	1
8	尚書	晉書	798	1
9	尚書	魏書	2852	2
10	尚書	魏書	2859	2
11	尚書	魏書	2862	2
12	尚書	魏書	2865	2
13	尚書	隋書	675	2
14	尚書	隋書	682	2
15	尚書	新唐書	2093	1
16	尚書	唐書	1402	4

總數(Y/N)

1史記 2漢書 3晉書 4魏書 5隋書 6新唐書 7唐書 8舊五代 9 0

[中文] [半形] 字根:

代號	內容	書名	頁數	段數
33	尚書省都事	唐書	1404	1

總數(Y/N)

1史記 2漢書 3晉書 4魏書 5隋書 6新唐書 7唐書 8舊五代 9 0

[中文] [半形] 字根:

請鍵入始印代號: 20 後按換行鍵
請鍵入終印代號: 20 後按換行鍵
請鍵入列印格式: 1(0:直印,1:橫印)後按換行鍵
請鍵入列印元件: 0(0:顯示幕,1:印表機)後按換行鍵

Figure-9 An example of search by "尚書"

References

1. Christor Faloutsos, "Access methods for text", Computing Surveys, ACM volum 17, NO.1, March 1985.
2. Carol Tenopir, "Full-text databases", Annual Review of Information Science and Technodgy, vol. 19, PP215-246, ASIS, 1984.
3. 毛漢光：史籍自動化第一年總報告，中央研究院，台北市，中華民國，1985.
4. C. C. Hsieh, "Experiments on the search and recognition of word from Chinese text", The first Asia-Pacific Conference on Library Science, Taipei, ROC, March. 1983.
5. 何文雄，中文斷詞的研究，國立台灣工業技術學院 工程技術研究所電子所 碩士論文，July, 1983.
6. 陳正佳，一套中文語法分析系統的研究與設計，國立台灣大學 資訊工程學研究所 碩士論文。June, 1985.
7. 林瑞華，卡系與資料整理，中央卡系統推廣中心，鳳山，高雄，台灣，中華民國，1978.
8. 梅棹忠夫原著，余阿勳、劉焜輝譯，知識誕生的奧秘，晨鐘出版公司，台北市，中華民國，1976.