# On the Automation of Chinese History Literatures

Ching-Chun Hsieh *

Zy-Kaan Ding

Yuan-Hua Wang

Chuai-Jung Yeh

Chi-Chou Shu

Shu-Fen Tung

Shi Lin

Fung-Hua Wang

Research Fellow, Inst. of Information Science,
Academia Sinica *

Computing Center, Academia Sinica

## Abstract

This paper contains two parts. In the first part, a brief introduction to the full-text processing system of Chinese history literatures will be presented. In the second part, we like to share some experience of processing history literatures with readers. Comments and discussions on the character set, retrieval mechanisms, mark-up language and further applications will be presented.

# Introduction

In order to study the feasibility of using computer as a tool for humanities study, a long-term project namely "Automation of History Literatures" (史籍自動化) has been launched since July 1st, 1984. Under this project, a software package called "Chinese Text Processor" (abbr. CTP) has been developed. It is a full-text data-base based system with controlled vocabulary, free-term, and text content retrieval mechanisms. CTP has three versions. Some highlights and characteristics of these versions are listed in Table 1.

CTP has three subsystems, they are : (1) Creation Module, (2) Retrieval Module, and (3) Application Module. The creation module has two major functions. The first one is a data-entry functional module which turns the original text into machine readable form called source text files. Then, the source text files will be processed by a text structurization module, the second functional module, and will produce structurized text files. A functional block diagram of the creation module is shown in Figure 1. In Figure 1, readers can easily find out that the original text after been processed by the structurization module is disassembled into four related files. The text string file keeps the original text string. The illustration file keeps the figures, graphes, equations and notes of the original text. The text structure file is a tree structured logical representation of the text content, ie. chapters, sections, paragraphs etc. And, lastly, the format file stores all type-setting information of the original text, such as paging, character font/size etc.

For mark-up rules used by CTP, please refer to Appendix A. In Appendix B, a BNF representation of the structurization parser is presented. For more information about CTP, please refer to reference [1] , [2] , [3] , and [4] .

| items \ version | CTP 1.0 | CTP 2.0 | CTP 3.0 |
|---|---|---|---|
| 1. Hardware | Dragon 570(A 16-bit, 8086 based micro) | micro-VAX II with Dragon 570 as terminal | 3B2/3B5/3B15 of AT&T with Dragon 570 as terminal |
| 2. Operating System | CP/M | VMS 4.1 | UNIX 5.0/BINIX |
| 3. language used for development | BASIC (Compiler) | C | C |
| 4. Query form | Manual Driven, Chinese | Manual Driven, English | Manual Driven, English |
| 5. retrieval machemism | controlled vocaburary search through hashing or indexing files | free-term search by pattern matching through text string | ①content search ②free-term search by pattern maching through text-structure trees |
| 6. Basic text element handled | page | page | paragraph |
| 7. text structure representation | inverted files | inverted files | trees |
| 8. Proximity function | by volume, by page | by volume, by page | by volume, by page, by sub-tree |
| 9. text tested (No. of characters) | 食貨志(史記至隋)(200K) | 食貨志(史記至隋)(200K) | 食貨志(全)(600K) 三民主義講稿(200K) 先總統蔣公言論集--文化與教育(150K) |
| 10. No. of users | Single user | from 2 to 4 | from 8 to 48 |
| 11. Typical Response time in Seconds | 120-200 | 40 | 60-90 |
| 12. Date of Issue | March,1985 | March,1986 | Sept. 1986 |

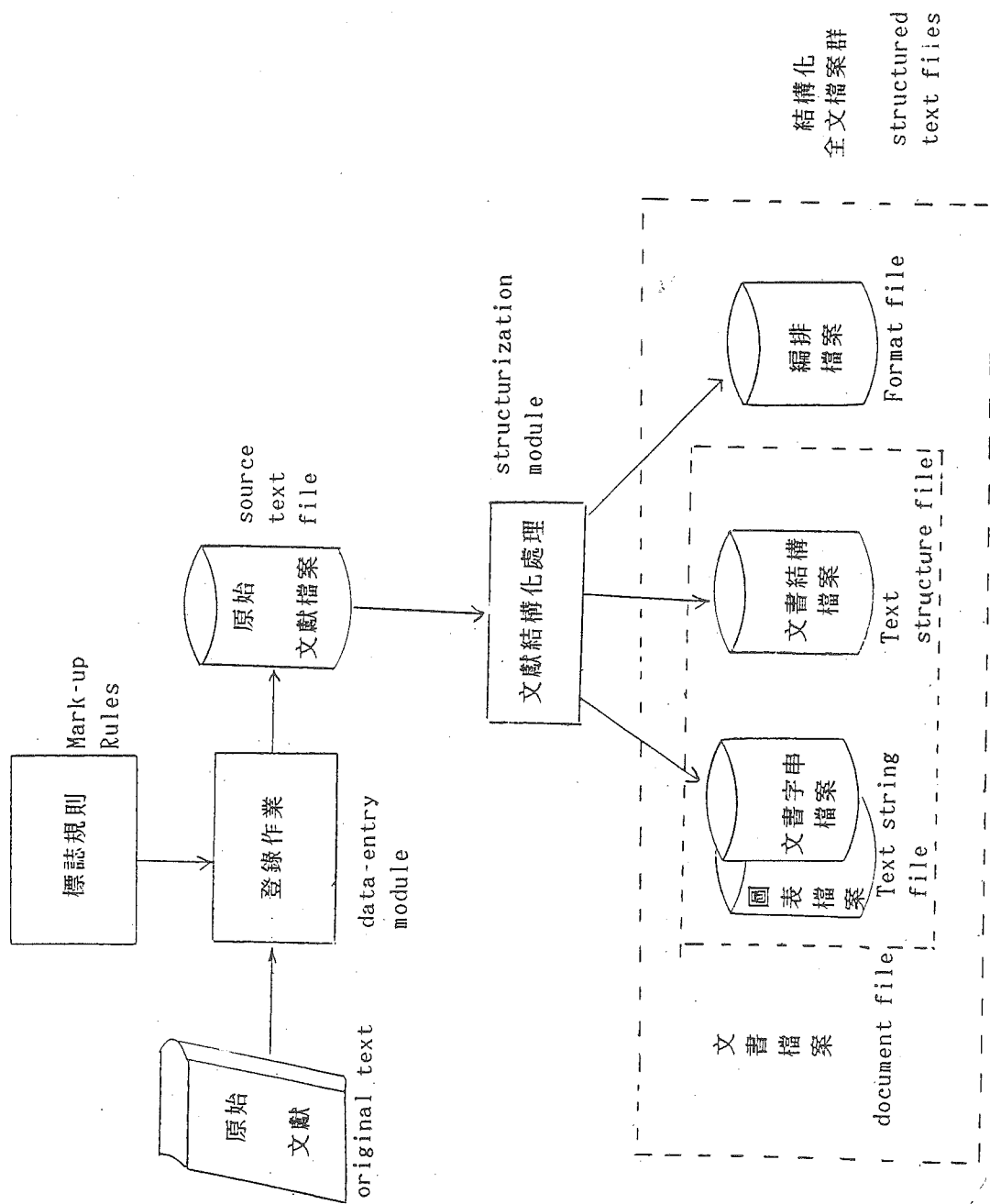Table 1. The characteristics of different versions of CTP

圖 1：全文處理系統之建檔作業部份方塊圖

Figure 1. A block diagram of the creation module

A functional block diagram of the retrieval module and the application module in shown in Figure 2. This part of the CTP can operate independently if structurized text files have already been produced by the creation module. In order to give you some ideas of the operations of the CTP, let me spend a few minutes in showing some slides of CTP for you.

The CTP version 3.0 is now ready to go to the public. In this paper, we like to present some worth-while considerations which we learned during the development of CTP. Your comments and discussions are hearty welcome.
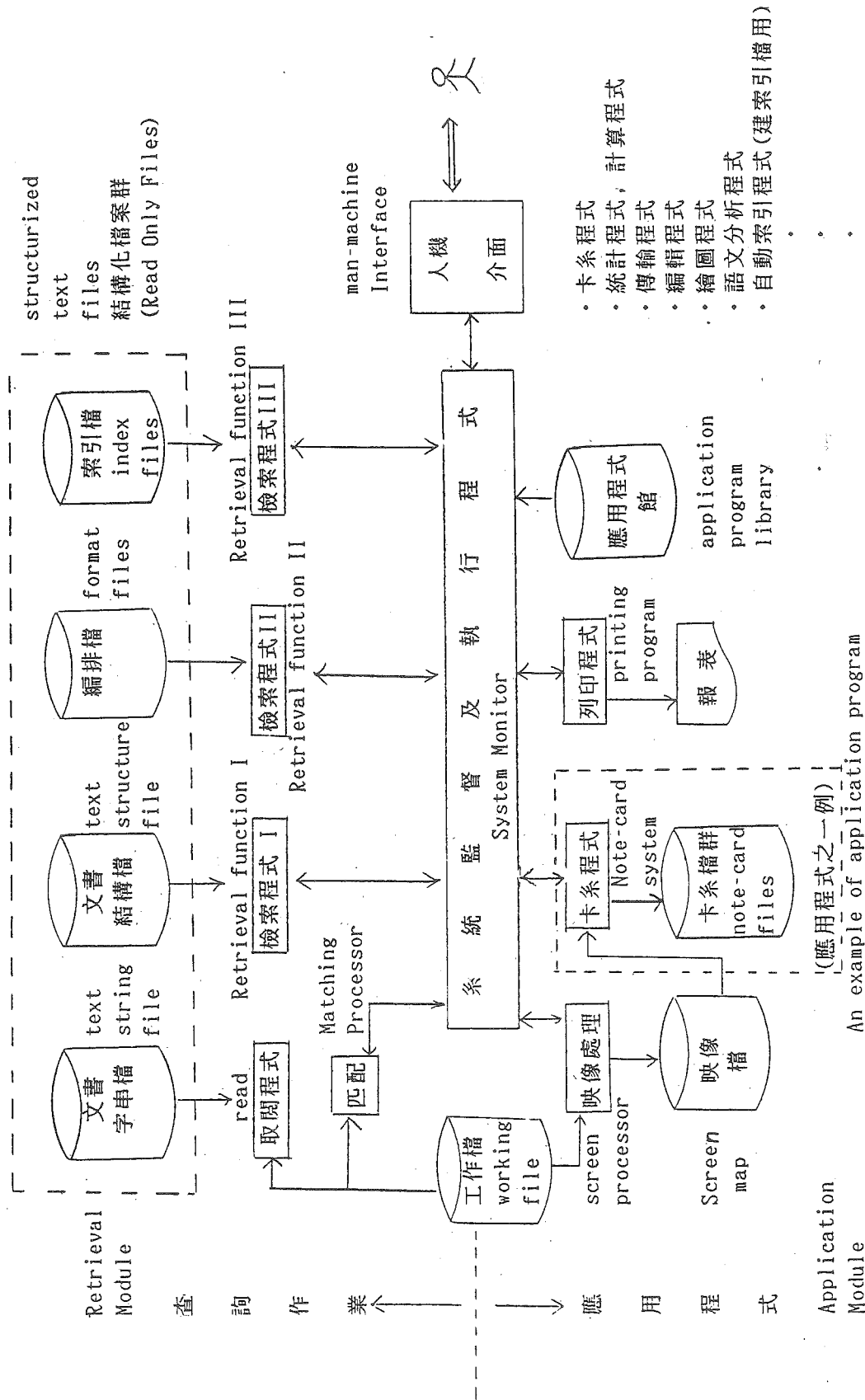
## Comments and Discussions on Processing History Literatures

While presenting history literatures in machine readable form, the information contained in the source text should not be distorted, added, and deleted by any means. But in real life, the above requirement can not be fulfilled with the present technology. Therefore, some efforts must be made to provide somewhat compromised solutions to this problem.

In this part of the paper, we will discuss the problems associated with character set, retrieval mechanisms, mark-up language, and furth applications in the subsequent sections.

### Character set

There are two character sets provided by vendor. One set contains about 13,000 characters which is identical to the set used by 通用漢字交換碼, a proposed standard set for commercial applications. The other set has more than 17,000 characters and it is a super set of the set mentioned above. In the first phase of our project while CTP 1.0 was developed,

Figure 2. A block diagram of the retrieval module and the application module

圖二 全文處理系統之查詢及應用部份結構方塊圖

the set of 13,000 characters was used. During the second phase of CTP 2.0, we shifted to the set of 17,000 characters. We have painful experience with these two character sets, because the coverage and completeness of these sets are obvious not capable enough for processing history literatures.

From Table 2 to Table 6, the reader can easily find that hundreds of characters are missing in these sets for processing history literatures. Therefore, we have to update the system whenever a character is to be added to the existing system. This is surely a time consuming and bothersome operation.

Another serious problem is the presence of variant forms of characters. The character patterns in ancient text can not always accommodate the present-day standard. So, in order not to distort the original text, we are asked to creat new characters whenever a mismatch of character pattern has been found. Even worse is in retrieval operation where both normal form and variant forms have to be used interchangeably. In other words, the capability of processing variant forms of characters is fatally needed for processing history literatures.

According to the facts we collected in Table 2 to Table 6, which is not limited to history literatures, we conclude that:

(1) The 通用漢字交換碼 is not capable enough for processing literatures.

(2) For processing literatures, CCCII [5] is highly recommended.

Controlled vocaburary search verses free-term search

It is a long argument between controlled vocaburary search and free term search since the full-text data-base became available to the public [6]. To our experience, we like to make the following comments:

A : The character set of 17,000 characters

B : The character set of 13,000 characters

Table 2. Missing Characters. The listed characters are used in 食貨志, but can not be found in the set of 13,000/17,000 characters

## Table 3. Variant forms.

The characters listed are used in 食貨志 and 史記。 Their corresponding variant forms are also shown.  But, those variant forms can not be found in the 17,000 set, and tense, they are also missing characters.

### Variant forms found from 史記

| No. | Normal | Variant |
|-----|--------|---------|
| 001 | 矼 | 矼 |
| 002 | 隸 | 隸 |
| 003 | 開 | 开 |
| 004 | 模 | 橅 |
| 005 | 起 | 起 |
| 006 | 弛 | 弛 |
| 007 | 崇 | 崇 |
| 008 | 廐 | 廐 |
| 009 | 污 | 污 |
| 010 | 弄 | 昪 |
| 011 | 煮 | 煮 |
| 012 | 廚 | 厨 |
| 013 | 亂 | 亂 |
| 014 | 氈 | 氈 |
| 015 | 敍 | 敍 |
| 016 | 騑 | 騑 |
| 017 | 鐏 | 鐏 |
| 018 | 楮 | 楮 |
| 019 | 廐 | 廐 |
| 020 | 宇 | 宇 |
| 021 | 朵 | 朶 |
| 022 | 杞 | 杞 |

### Variant forms found from 食貨志

| No. | Normal | Variant | No. | Normal | Variant |
|-----|--------|---------|-----|--------|---------|
| 001 | 鰥 | 鰥 | 023 | 頁 | 頁 |
| 002 | 棘 | 棘 | 024 | 遷 | 遷 |
| 003 | 渚 | 渚 | 025 | 衷 | 衷 |
| 004 | 箸 | 箸 | 026 | 侍 | 侍 |
| 005 | 筭 | 筭 | 027 | 毒 | 毒 |
| 006 | 盡 | 尽 | 028 | 持 | 持 |
| 007 | 号 | 号 | | | |
| 008 | 暄 | 暄 | | | |
| 009 | 恃 | 恃 | | | |
| 010 | 微 | 微 | | | |
| 011 | 都 | 都 | | | |
| 012 | 斥 | 斥 | | | |
| 013 | 羮 | 羮 | | | |
| 014 | 即 | 即 | | | |
| 015 | 時 | 時 | | | |
| 016 | 簒 | 簒 | | | |
| 017 | 者 | 者 | | | |
| 018 | 諸 | 諸 | | | |
| 019 | 報 | 報 | | | |
| 020 | 犀 | 犀 | | | |
| 021 | 赭 | 赭 | | | |
| 022 | 著 | 著 | | | |

Table 4. Missing Characters
    The following characters are used in 漢書/史記，but
    are not found in the set of 17,000 characters．
    (This table is not yet completed)

| No. | Ch. | No. | Ch. | No. | Ch. | No. | Ch. | No. | Ch. | No. | Ch. | No. | Ch. | No. | Ch. | No. | Ch. | No. | Ch. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 001 | 丼 | 023 | �片 | 045 | 甍 | 067 | 紂 | 089 | 矔 | 141 | 傝 | 133 | 瀟 | 157 | 陮 | 177 | 檻 | 199 | 氺 |
| 002 | 庹 | 024 | 嵜 | 046 | 她 | 068 | 篝 | 090 | 蝍 | 112 | 逌 | 134 | 乏 | 156 | 窬 | 178 | 剋 | 200 | 壙 |
| 003 | 飍 | 025 | 倌 | 047 | 眹 | 069 | 毯 | 091 | 裛 | 113 | 竑 | 135 | 禮 | 157 | 关 | 179 | 獫 | 201 | 泩 |
| 004 | 鄙 | 026 | 燮 | 048 | 調 | 070 | 暢 | 092 | 路 | 114 | 蘇 | 136 | 敦 | 158 | 失 | 180 | 誠 | 202 | 駁 |
| 005 | 阢 | 027 | 罌 | 049 | 櫬 | 071 | 鸛 | 093 | 剌 | 115 | 斥 | 139 | 邀 | 159 | 挈 | 181 | 韁 | 203 | 驍 |
| 006 | 碧 | 028 | 陝 | 050 | 嫠 | 072 | 叟 | 094 | 雙 | 116 | 羨 | 138 | 賈 | 160 | 掋 | 182 | 擭 | 204 | 迋 |
| 007 | 緷 | 029 | 阢 | 051 | 浸 | 073 | 鯢 | 095 | 郗 | 117 | 錫 | 139 | 截 | 161 | 餓 | 183 | 搜 | 205 | 遍 |
| 008 | 闕 | 030 | 厭 | 052 | 經 | 074 | 斬 | 096 | 嶷 | 118 | 异 | 140 | 糒 | 162 | 耼 | 184 | 邃 | 206 | 振 |
| 009 | 芏 | 031 | 郝 | 053 | 潸 | 075 | 額 | 097 | 萆 | 119 | 即 | 141 | 桼 | 163 | 香 | 185 | 斅 | 207 | 岠 |
| 010 | 稿 | 032 | 棘 | 054 | 扑 | 076 | 塀 | 098 | 蠻 | 120 | 時 | 142 | 煥 | 164 | 剝 | 186 | 籥 | 208 | 埕 |
| 011 | 潰 | 033 | 睭 | 055 | 鱗 | 077 | 尽 | 099 | 雄 | 121 | 簒 | 143 | 乂 | 165 | 頦 | 187 | 廎 | 209 | 映 |
| 012 | 陀 | 034 | 霤 | 056 | 黀 | 078 | 鏠 | 100 | 姘 | 122 | 喜 | 144 | 翺 | 166 | 醁 | 188 | 棨 | 210 | 勝 |
| 013 | 偵 | 035 | 廬 | 057 | 壉 | 079 | 陵 | 101 | 袘 | 123 | 耇 | 145 | 撒 | 167 | 鼽 | 189 | 竈 | 211 | 崆 |
| 014 | 瘥 | 036 | 寧 | 058 | 湒 | 080 | 僆 | 102 | 鷇 | 124 | 諸 | 146 | 闌 | 168 | 廩 | 190 | 袤 | 212 | 侪 |
| 015 | 駬 | 037 | 茎 | 059 | 箐 | 081 | 顮 | 103 | 踰 | 125 | 報 | 143 | 燋 | 169 | 僑 | 191 | 爾 | 213 | 岠 |
| 016 | 褕 | 038 | 頷 | 060 | 徐 | 082 | 颿 | 104 | 憝 | 126 | 犀 | 148 | 迄 | 170 | 哀 | 192 | 罕 | 214 | 紖 |
| 017 | 騎 | 039 | 譌 | 061 | 鄁 | 083 | 谔 | 105 | 颿 | 127 | 芈 | 149 | 趲 | 171 | 郝 | 193 | 烰 | 215 | 獝 |
| 018 | 鉅 | 040 | 娷 | 062 | 銅 | 084 | 盠 | 106 | 雕 | 128 | 轎 | 150 | 湆 | 192 | 酾 | 194 | 敳 |  |  |
| 019 | 瑾 | 041 | 暤 | 063 | 嘽 | 085 | 崤 | 107 | 陪 | 129 | 肴 | 151 | 殹 | 173 | 嬴 | 195 | 胯 |  |  |
| 020 | 埋 | 042 | 掃 | 064 | 衙 | 086 | 嶧 | 108 | 恃 | 130 | 窒 | 152 | 昬 | 174 | 轇 | 196 | 鋎 |  |  |
| 021 | 輔 | 043 | 憋 | 065 | 邦 | 087 | 噎 | 109 | 微 | 131 | 劂 | 153 | 庖 | 175 | �爛 | 197 | 餞 |  |  |
| 022 | 敱 | 044 | 墟 | 066 | 傸 | 088 | 肴 | 110 | 都 | 132 | 軟 | 154 | 旭 | 176 | 蹁 | 198 | 屄 |  |  |

− 48 −

# Table 5. Missing character(normal form)
The following characters are used in Mandarin Daily News Dictionary but are not found in the set of 17,000 characters

| No. | ch. | No. | ch. | No. | ch. | No. | ch. | No. | ch. | No. | ch. | No. | ch. | No. | ch. | No. | ch. | No. | ch. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 001 | 亠 | 023 | 筭 | 045 | 鈌 | 067 | 驫 | 089 | 搢 | 111 | 挀 | 133 | 曆 | 155 | 蛼 | 177 | 瞽 | 198 | 㥯 |
| 002 | 亡 | 024 | 蒽 | 046 | 鎰 | 068 | 騲 | 090 | 犾 | 112 | 矗 | 134 | 盉 | 156 | 鎈 | 178 | 乚 | 199 | 丨 |
| 003 | 亖 | 025 | 堷 | 047 | 矙 | 069 | 鴯 | 091 | 獾 | 113 | 姘 | 135 | 生 | 157 | 鎋 | 179 | 禩 | 200 | 川 |
| 004 | 爻 | 026 | 瘦 | 048 | 鈌 | 070 | 攉 | 092 | 荄 | 114 | 粜 | 136 | 三 | 158 | 鐙 | 180 | 羅 | 201 | 川 |
| 005 | 乚 | 027 | 菜 | 049 | 餶 | 071 | 凂 | 093 | 刜 | 115 | 嗒 | 137 | 嘞 | 159 | 鑝 | 181 | 裼 | 202 | X |
| 006 | 丁 | 028 | 蘩 | 050 | 骰 | 072 | 菶 | 094 | 昳 | 116 | 收 | 138 | 六 | 160 | 鍍 | 182 | 覷 | 203 | 8 |
| 007 | ㄒ | 029 | 蕪 | 051 | 鬐 | 073 | 殯 | 095 | 硬 | 117 | 氕 | 139 | 閩 | 161 | 鐵 | 183 | 覸 | 204 | 扌 |
| 008 | ㄠ | 030 | 蒜 | 052 | 醲 | 074 | 頭 | 096 | 駘 | 118 | 吽 | 140 | 轉 | 162 | 䢀 | 184 | 忍 | 205 | 髑 |
| 009 | ㄍ | 031 | 蜊 | 053 | 騳 | 075 | 算 | 097 | 氵 | 119 | 啾 | 141 | 鞓 | 163 | 圊 | 185 | 刬 | 206 | 覷 |
| 010 | ㄗ | 032 | 蛻 | 054 | 閒 | 076 | 鱘 | 098 | 達 | 120 | 爪 | 142 | 輔 | 164 | 忺 | 186 | 尗 | 207 | 宮 |
| 011 | 艸 | 033 | 旹 | 055 | 韇 | 077 | 鵷 | 099 | 三 | 121 | 嗯 | 143 | 靜 | 165 | 皿 | 187 | 蓬 | | |
| 012 | 冫 | 034 | 鉋 | 056 | 挍 | 078 | 殲 | 100 | 扡 | 122 | 扴 | 144 | 態 | 166 | 誣 | 188 | 彡 | | |
| 013 | 亞 | 035 | 合 | 057 | 懐 | 079 | 鬽 | 101 | 簙 | 123 | 廖 | 145 | 淺 | 167 | 蛇 | 189 | 辰 | | |
| 014 | 丅 | 036 | 鄩 | 058 | 臁 | 080 | 鬽 | 102 | 搶 | 124 | 爛 | 146 | 父 | 168 | 蜶 | 190 | 㠸 | | |
| 015 | 疒 | 037 | 壟 | 059 | 捨 | 081 | 鐦 | 103 | 搘 | 125 | 狹 | 147 | 騎 | 169 | 裛 | 191 | 瑨 | | |
| 016 | 癢 | 038 | 鎂 | 060 | 燭 | 082 | 鑆 | 104 | 轟 | 126 | 磅 | 148 | 觞 | 170 | 釱 | 191 | 世 | | |
| 017 | 嶝 | 039 | 瞽 | 061 | 嗞 | 083 | 鑝 | 105 | 耰 | 127 | 芘 | 149 | 餀 | 171 | 鉅 | 192 | 亖 | | |
| 018 | 睑 | 040 | 鎧 | 062 | 壅 | 084 | 薛 | 106 | 鰲 | 128 | 脛 | 150 | 餡 | 172 | 鐸 | 193 | 𠃌 | | |
| 019 | 眶 | 041 | 蜳 | 063 | 鎌 | 085 | 之 | 107 | 劉 | 129 | 徉 | 151 | 鯉 | 173 | 鉾 | 194 | 乛 | | |
| 020 | 鮑 | 042 | 蜚 | 064 | 乀 | 086 | 辟 | 108 | 滈 | 130 | 捏 | 152 | 藍 | 174 | 鐃 | 195 | 亖 | | |
| 021 | 攺 | 043 | 醹 | 065 | 陽 | 087 | 驫 | 109 | 皀 | 131 | 搭 | 153 | 藨 | 175 | 酺 | 196 | 搀 | | |
| 022 | 咽 | 044 | 醋 | 066 | 隋 | 088 | 搴 | 110 | 蒦 | 132 | 搭 | 154 | 薐 | 176 | 麿 | 197 | 摠 | | |

Table 6. Variant forms of characters found in the Mandarin daily News Dictionary and not found in the set of 17,000 characters. (to be continued)

Table 6. (to be continued)

Table 6.

(1) free-term search is more suitable for history literature processing than the controlled vocaburary search.

(2) In order to improve the system performance, we suggest to establish a controlled vocaburary index by the following :

· Set a counter for each free-term whenever being used for search.

· For a definite period of time, the above counts will be examined. If a term is found to be frequently used, then, it will be considered as a key word for search, and hence, its index shall be established.

In other word, we think a combination of both retrieval mechanisms will optimize the performance of the system. Besides, we suggest to establish the controlled vocaburary search on the foundation of free-term search and the frequency of usage of each terms.

## The needs of standard mark-up language

A very good acticle that explains the concept and needs of a standard mark-up language can be found in the Appendix A of ISO standard DIS 8879 [ 7 ] . The mark-up language serves as a bridge that allows data sharing among many application areas such as full-text database, printing industry, word processing, etc. In our application, mark-up language is a necessity for text element recognition.

During the development of CTP, it is sorry to find out that no such a standard mark-up language for Chinese text is available. Therefore, we have to work out a replacement, because we do not have enough resource to develop a standard generalized mark-up language parallel to that of DIS 8879. This replacement, aiming for text element recognition, is the BNF representation of 食貨志 text shown in Appendix B. So far, this abstract BNF model has been successfully used to represent 食貨志, 孫中山先生三民主義講稿, and 蔣中正先生言論集. It seems powerful for certain "regular"

structured literatures. However, to design a generalized mark-up language for Chinese text is surely a necessity for furthes development in the Chinese text processing field.

Further Applications

In a traditional information retrieval environment, where information is organized as formatted record type, the most important function of such a system is retrieval. This situation is no longer true when full-text can be retrieved from a database. Enumerous applications can be developed upon retrieved text, to name a few, such as statistical analysis, computing, graphics, word processing/text generation, etc. Therefore, the development of full-text database and full-text processing technology provide an excellent opportunity to integrate information retrieval systems, text/word processing systems, and computation/analytical systems all together. Besides, many new applications may be created, such as a work-station for humanities study. Therefore, we conclude that full-text processing is a very promising field for us to looking for further development.

# References

1. 毛漢光：〈史籍自動化第一年總報告〉，中央研究院 歷史語文研究所。
   Institute of History & Language, Academia Sinica, Taipei, July 1985

2. 謝清俊等：〈中文全文處理系統的設計與製作〉，中央研究院 計算中心
   Computing Center, Academia Sinica, Taipei, Sept. 1986

3. 謝娟娟：〈中文全文處理系統使用手冊〉，中央研究院 計算中心
   Computing Center, Academia Sinica, Taipei, Sept. 1986

4. C. C. HSIEH, "Full-Text Processing of Chinese Language— —an
   experimental system for studying Chinese History Literatures" the SEAL
   sub-committee on Library Technology, AAS. 38th Aunual Conference,
   Chicago, U.S.A. March, 1986

5. 〈Chinese Character Code for Information Interchange〉 vol.1, vol.2
   Chinese Character Analysis Group, 行政院文化建設委員會，Taipei, 1980,
   1983.

6. Carol Tenopir, "Full-Text databases", 〈Annual Review of Information
   Science and Technology〉, vol.19, pp215-246, ASIS, Nov. 1984

7. ISO Draft International Standard DIS 8879, 〈Information Processing—
   Text and Office Systems—Standard Generalized Markup Language(SGML)〉,
   ISO/TC 97, International Organization for Standardization, Geneva,
   Switzerland, Oct. 1985.

# Appendix A

# Mark-up rules used by CTP

# 全文標誌規則

## 壹．通則

一、凡需利用本中心發展的「中文全文處理系統」(CTPS)來處理的文獻，都必須嚴格依照本規則的規定做資料登錄的工作。若否，CTPS將產生錯誤而無法正常工作。

二、依照本規則登錄的文獻中不能有外國文字、字母、及阿拉伯數字。本規則亦不能用之於圖形、表格、公式等。故凡有上列資料之文獻，不宜直接以本規則處理。

三、文獻登錄時以頁為單位，亦即每頁建一檔案，此檔案之名稱為：「xy．src」之格式，其中x為原始文獻之代碼，y為以數字表示之頁碼，而src表示原始檔案之類別，文獻之代碼由系統工程師設定，不可任意取用。

四、資料登錄時，應盡量依照原稿之版面排列形式。凡有空頁、空行、空白之處均須比照登錄。除在細則中有規定者外，不得任意改變版面之形式。

五、所有文字及符號之鍵入，除細則中另有規定者外，均取全形。

六、關於此標註規則之實務詳如細則。凡對此規則之使用有任何疑問時，請立刻與系統工程師連絡，不可以隨自己意思選擇處理之方式。

## 貳．細則

一、檔案結構

檔案之第一行應先以阿拉伯數字打入頁次(即通則三中之y)，然後按換行鍵，自第二行輸入正文。

二、正文

(1)若為空白頁，則其正文以 "a" 表示之。

(2)正文前至少有二個全形空白(bb)。

三、書、卷

　　⑴書卷起處加 "d"，迄處加 "e"。

　　⑵書卷之標題起處加 "f"，迄處加 "g"。

　　⑶標題群之間以一個或一個以上之全形空白隔開。

　　⑷每個標題之起處加 "l"，迄處不加特別標誌。

四、字形大小

　　小字起處加 "★"，迄處加 "■"。

五、註釋及校勘

　　⑴註釋起處加 "h"，迄處加 "i"。

　　⑵校勘起處加 "j"，迄處加 "k"。

　　⑶註釋和校勘同時出現處，校勘出處加 "※"。

# Appendix B
# BNF representation of the
# structurization module

# 食貨志的制式 Ｂ Ｎ Ｆ 表示法

(1) meta symbols of the BNF notation:

   〈　〉 ::= |　{ }*　{ }n　[ ]　" "

(2) constructs of the BNF notation:

| | |
|---|---|
| 〈…〉〈…〉 | means sequence expression concatentation |
| ::= | means "is defined as" |
| \| | means "or" |
| {…}* | means "repeatition option times" |
| {…}n | means "repeatition n times" n great than 0. |
| […] | means "optional elements" |
| "…" | means "terminal symbols" |

(3) 邏輯結構 (Logical Structure):

〈卷　　　〉 ::= 〈標頭〉 { [〈標題〉] 〈大段〉 }* [〈校勘記〉]

〈標　　頭〉 ::= "d" 〈卷頭〉 { [〈附卷頭〉] }* "e" "¥n"

〈卷　　頭〉 ::= 〈卷別〉 "¥n"

〈卷　　別〉 ::= {〈詞〉}*

〈附 卷 頭〉 ::= {〈詞〉}* "¥n"

〈標　　題〉 ::= "f" "¥n" { {〈全形空白〉}* 〈標題群〉 "¥n" }* "g" "¥n"

〈標 題 群〉 ::= 〈標題詞〉 [ { {〈全行空白〉}* 〈標題詞〉 }* ]

〈標 題 詞〉 ::= 〈詞〉 ["★" { 〈子標題〉 }* "▓"]

〈子 標 題〉 ::= {〈詞〉}*

〈大　　段〉 ::= 〈段〉 [〈註釋群〉] | 〈前敘正文〉 {〈引言內容〉}* 〈後敘正文〉
　　　　　　"l" {〈大段〉}* | 〈統計條列句子〉

〈段　　　〉 ::= {〈全形空白〉}2 〈正文直到換行後接兩個全形空白以前為止〉

〈正文直到換行後接兩個全形空白以前為止〉 ::= { 〈句子〉 }*

〈前敘正文〉 ::= [ {〈句子〉}* ] { 〈中文字〉 }* ":" [ "〔" 〈勘碼〉 "〕" ] "¥n"

〈引言內容〉 ::= [ {〈全形空白〉}0 { 〈中文字〉 }* "¥n" ]
　　　　　　{〈全形空白〉}4 { 〈中文字〉 }* "¥n"
　　　　　　{ {〈全形空白〉}2 { 〈中文字〉 }* "¥n" }*

〈後敘正文〉 ::= {〈全形空白〉}0 {〈句子〉}*

〈統計條列句子〉 ::= {"¥n"}* { {〈中文字〉}* "¥n" }*

〈註釋群〉 ::= "h" "¥n" {〈註釋〉}* "i" "¥n"

〈註　釋〉 ::= {〈全形空白〉}1 "〔" 〈註碼〉 "〕" 〈註釋正文〉

〈註釋正文〉 ::= {〈中文字〉}* "¥n" 〔 { {〈全形空白〉}5 {〈中文字〉}* "¥n" }* 〕

〈註　碼〉 ::= "○" | "一" | "二" | "三" | "四" | "五" | "六" |
　　　　　　　"七" | "八" | "九"

〈校勘記〉 ::= "j" "¥n" "校勘記" "¥n" {〈校勘〉}* "k" "¥n"

〈校　勘〉 ::= "〔" 〈勘碼〉 "〕" 〈校勘正文〉

〈校勘正文〉 ::= {〈中文字〉}* "¥n" 〔 { {〈全形空白〉}4 {〈中文字〉}* "¥n" }* 〕

〈勘　碼〉 ::= "○" | "一" | "二" | "三" | "四" | "五" | "六" |
　　　　　　　"七" | "八" | "九"

〈詞　　〉 ::= {〈中文字〉}*

〈句　子〉 ::= {〈中文字〉}* 〈句的標號〉

〈全形空白〉 ::= " "

〈中文字〉 ::= 〈宏碁公司出品之字形產生器DPC-24/60 之字集〉 |
　　　　　　　〈新造字〉 | 〈教育部發佈的標點符號〉