# 前　　言

　われわれの 21 世紀 COE プログラム「東アジア世界の人文情報學研究教育據點」では "漢字文化の全き繼承と發展" をスローガンに掲げている。漢字による文化資源を全面的に繼承するためには、漢字文獻の各種データベースをコンテンツの面で量的に充實させるばかりでなく、更に新しい技術の開發による利便性の向上を目指すことが求められるのは當然であろう。しかも今後の漢字文獻データベースは、少なくとも東アジアの漢字使用圈における廣範な利用者を前提にして設計されねばならない。そのためには技術開發の面で、中國との交流を推進し、相互理解に努めること、ひいては日中間において存在し得る幾つかの問題について、その解決はともかく、現狀を正確に認識しておくことが緊急の課題であると思われる。そこで漢字文獻のデータベース化において中國で中心的な役割を果たしている中國國家圖書館とのあいだで、日中共同シンポジウム「漢字文獻資料庫的新技術」を開催することで、この課題に絲口を見いだそうと考えた。

　幸いにして國家圖書館の陳力副館長、及び嚴向東國際交流處長をはじめとする同館スタッフの熱心な協力があり、約半年に及ぶ準備の末、2005 年 1 月 22 日（土）、嚴しい寒さのなか、國家圖書館にほど近い湖北大厦を會場として開催の運びとなった。漢字文獻データベースの技術的な問題に關して、日本側 4 名、中國側 4 名の報告が行われ、それらを中心に率直かつ活潑な意見交換を行うことが出來たことは、極めて有意義であった。本册子にはこれらの報告の改訂稿を、當日の報告順に從って掲載する。ただ中國側報告者の内、中國國家圖書館科研處處長孫一鋼（SUN Yigang）氏の「數字文獻處理的標準規範研究」については、同氏からの申し出によりここに收錄することが出來なかったのは殘念である。

　漢字文獻データベースの技術について日中の研究者のあいだで直接に意見交換を行う機會はこれまでも決して多くなかった。今後、この種の試みが繼續して行われ、東アジア的規模で漢字文化繼承の具體的問題が一層活潑に討議されることを願ってやまない。

2005 年 7 月 17 日

高田時雄

# 目　次　Table of Contents

# Text-Searchable Image and Its Applications

Koichi Yasuoka

Documentation and Information Center for Chinese Studies,
Institute for Research in Humanities, Kyoto University

## 1　Introduction

Since 1996 proposal of the Council for Science, the university libraries in Japan have progressed "The Digital Library Project". Nowadays the union catalogue database of the university libraries (NACSIS-CAT) is almost completely equipped, and we can easily find any books and magazines in the libraries through the database on the Internet. But we are still far and away from the goal of "The Digital Library Project", which is the digitalization of all the books and the magazines in the libraries. The university libraries have only made displays of images of the rare books without their digital texts, their digital tables of contents, or their digital indices. The digital libraries in Japan now are not "libraries" but something like "museums", since they don't give us the way to "read" the books digitally.

In this paper the author represents the concept of text-searchable images and its applications. The author shows two formats, Portable Document Format and Scalable Vector Graphics, to actualize text-searchable images, and also shows a JavaScript-based program "`ttext-kanbun`" to produce text-searchable images in these formats. The author contributes this paper toward the true progress of the digital "libraries".

## 2　Text-Searchable Images

In this section we examine two formats, Portable Document Format (PDF) and Scalable Vector Graphics (SVG), to actualize text-searchable images.

### 2.1　PDF for Text-Searchable Images

The author has studied long time about text-searchable images using PDF [2]. And Adobe adopted some results of the study into PDF-1.4 [3] as "transparent text". Now we have two ways to actualize text-searchable images

using PDF. The one is to put a transparent text upon an image, and the other is to put an image upon a text written in white characters. The former way is only available with the browsers of PDF-1.4 and after, and the latter way PDF-1.2 and after. In this paper we use the latter way for backward compatibility.

PDF can represent both images and texts, but has some limitations on its format. PDF supports only two compression methods for color images, that are JPEG and ZIP. PDF supports several character-sets for CJK texts, Adobe-Japan1-6 [7] (including 14663 漢字 characters), Adobe-GB1-4 [1] (including 27629 汉字 characters), Adobe-CNS1-4 [5] (including 17625 漢字 characters), and Adobe-Korea1-2 [6] (including 4620 漢字 characters) under Japanese, mainland Chinese, Taiwanese, and Korean circumstances, respectively. We need "Japanese Language Pack" to read and search PDFs written in Adobe-Japan1-6 character-set, so as mainland Chinese, Taiwanese, and Korean. This means that these character-sets are incompatible with one another, and that PDFs for text-searchable images actually cannot get across the borderlines. In this paper we use JPEG for color images and Adobe-Japan1-6 character-set for texts to produce text-searchable images with PDF.

## 2.2 SVG for Text-Searchable Images

Tomohiko Morioka has studied about text-searchable images using SVG [4]. He actualized a text-searchable image to put an image upon a text. But in this paper we put a transparent text upon an image to actualize a text-searchable image using SVG.

SVG can include both images and texts, but the most contemporary viewer "Adobe SVG Viewer 3.0" has some limitations. SVG supports any kind of formats for color images, but the viewer suports only JPEG, PNG, and GIF. SVG supoorts any text-encodings but prefers UTF-8. In this paper we use JPEG for color images and UTF-8 for texts to produce text-searchable images with SVG.

# 3 Experiment and Result

The author wrote a JavaScript-based program "`ttext-kanbun`" to produce text-searchable images using PDF or SVG. "`ttext-kanbun`" runs on Internet Explorer 6 under Microsoft Windows XP.

We, members of COE21-project at Institute for Research in Humanities, Kyoto University, tried to make text-searchable images of 大唐西域記 (ex-橘寺-collection) with "`ttext-kanbun`" (Figure 1). We prepared 319 JPEG images for 大唐西域記, where each image has 2100×1950 pixels and total size of all images is 196807821 bytes, and its text written in UTF-8 consisting of 104725 characters (3138 different).
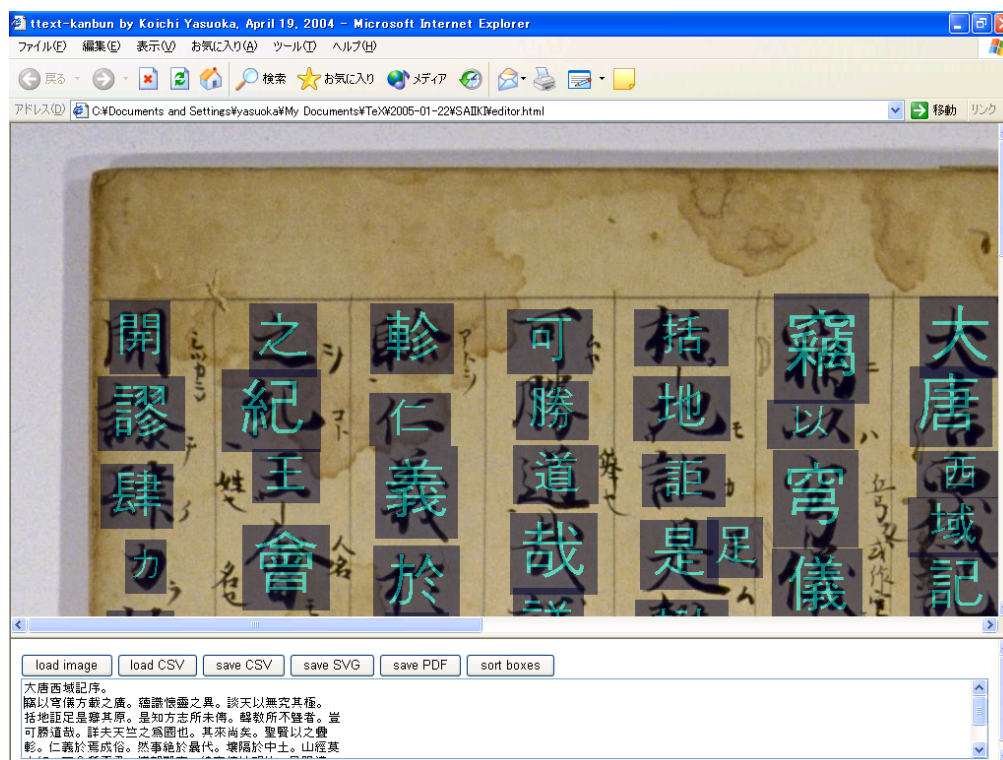


Figure 1: Snapshot of "`ttext-kanbun`"

First we produced text-searchable images using PDF (Figure 2). The total size of 319 PDF files was 202662390 bytes, 2.97% increasing from original JPEG images. We couldn't write 390 characters out of 104725 using PDF since they were not included in Adobe-Japan1-6. The 390 characters consisted of 51 different characters shown in Table 1. Then we combined the 319 PDF files into a multi-page PDF. The file-size of the combined PDF was 202440575 bytes, 2.86% increasing from original JPEG images.

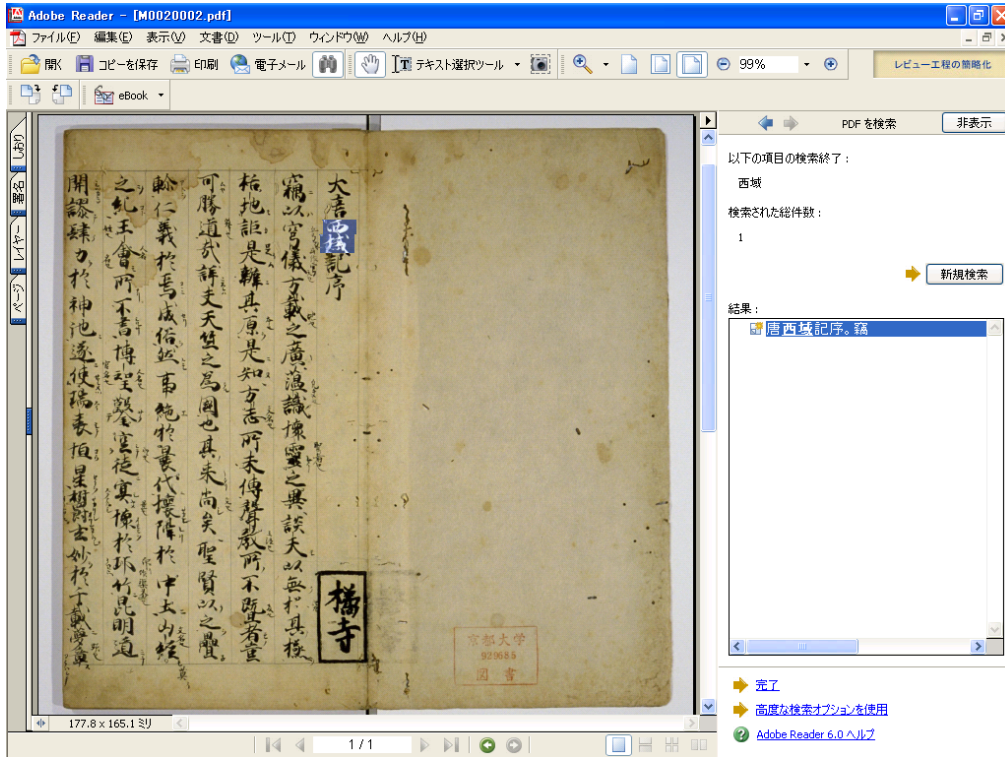Second we produced text-searchable images using SVG (Figure 3). The

Figure 2: Searching "西域" on "Adobe Reader 6.0"



Table 1: Characters not in Adobe-Japan1-6

Figure 3: Searching "西域" on "Adobe SVG Viewer 3.0"

| 猻 匲 旿 懰 捽 攱 毦 煩 絹 芉 茻 袠 卧 餕 |
| --- |

Table 2: Invisible characters on "Adobe SVG Viewer 3.0"

total size of SVG files and JPEG images was 203752662 bytes, 3.53% increasing from JPEG images only. All characters could be represented in SVG files, but 14 characters shown in Table 2 couldn't be displayed on "Adobe SVG Viewer 3.0", since the Viewer didn't support Unicode Plane-2 fonts.

# 4 Conclusion

In this paper the author has represented the concept of text-searchable images and its actualization using PDF or SVG. The author wrote a JavaScript-based program "`ttext-kanbun`" to produce such text-searchable images. As a result we have found that only 3% to 4% file-size increase is needed to add texts on JPEG images. The author now distributes "`ttext-kanbun`" at `http://coe21.zinbun.kyoto-u.ac.jp/~yasuoka/ftp/program/` and is pleased to help anyone to produce such text-searchable images.

# References

[1] Adobe-GB1-4 Character Collection for CID-Keyed Fonts, Technical Note #5079, Adobe Systems (November 2000).

[2] Koichi Yasuoka and Tokio Takata: Digital Rubbings — Their Past and Future, 2001 Pacific Neighborhood Consortium Proceedings (January 2001), ECAI Rubbings Work Session.

[3] Adobe Systems Incorporated: PDF Reference third edition — Adobe Portable Document Format Version 1.4, Addison-Wesley (December 2001).

[4] 守岡知彦: ポスト文字コード時代の文書処理技術に関する展望, 全国文献·情報センター人文社会科学学術セミナーシリーズ, No.12 (November 2002), pp.59-70.

[5] Adobe-CNS1-4 Character Collection for CID-Keyed Fonts, Technical Note #5080, Adobe Systems (May 2003).

[6] Adobe-Korea1-2 Character Collection for CID-Keyed Fonts, Technical Note #5093, Adobe Systems (May 2003).

[7] Adobe-Japan1-6 Character Collection for CID-Keyed Fonts, Technical Note #5078, Adobe Systems (June 2004).

# 古籍数字图书馆的中文信息处理技术综述

孙卫
中国国家图书馆

## 目 录

3500



1600

1959

504                                    469

[1]

1942    9

30                        39

---
[1]

11

1946

1924

.

80        1090

1965        1

.                                        40

570

1944        4

34                31

757

850

(  868  )                    1

6        1

                                    "

        "

        1899              1907

                                    "          "

    100

                        3500                          100

                                20        70

748

        20        80

                                0525              014/015

CCDOS                                    Windows  3.1

Windows 3.2

Windows 95　Windows

97　Windows 98　Windows ME

Windows　NT　4.0　　　　Windows　2000

Windows　XP　　　　　　　　Unicode

1

+

2

1

2

3

Unicode                                          -

2                       3

3   -3.5                      4                          5                          4

7                                      6              5    6

7          10    6

8        5                                    9          2

Unicode 4.0 / ISO 10646-2003          71000

4

GB13000      GB2312/GBK/Big5

ISO10646

2                          2003
3          2000
4                    2002                    1949
3    3    5              ——
5
6
7 http://140.111.1.40/start.htm
8
9
.

20902　　　　ISO10646　　　　　　　A 6582

27494　　　　　　　　　　[10]　　　　　　　　　B　4

3

（OCR）

99.997%　　　　　　　1

5　　　　　　　5

:

1　　　B/S

―――――――――――――
[10]

B/S

2

3

4

ISO10646

150

TrueType 1

TrueType 2

TrueType n

65536

1

1

ISO10646

150　　　　　　　　　ISO10646

ISO10646

| GB2312/GB18030/GBK | TrueType |
| --- | --- |

| BIG5 | TrueType |
| --- | --- |

| ISO10646 | TrueType 1 |
| --- | --- |
| | TrueType n |

2

2



GB2312/GB18030/GBK

BIG5

ISO10646 | TrueType 1

TrueType n

3

3

ISO10646

11



OCR

ISO10646

ISO10646

ISO10646

4

4　　　　　　B/S

ISO 10646

TrueType　　　　　　　　　　B/S

11　　　　　Windows 2000　Windows XP

95%

+

5%

5

6

5

6

B/S



7

7

0.03%　OCR　　　　　　　　1　　5

12

13　　　　　7　　1

7

1

7　　1

12
13

8

8

1

图 9

图 9                                                ISO10646

——                              ——

2



图 10

图 10

11

11

3



12

12                                        OCR

**XML**

XML

1

2

3

4                                                     XML

　　　　　　　XML

5

　　　　　B/S                                    JAVA

　　　　　　　　　ISO10646-2003

　　　　XML

　　　　　　　B                                    XML

B/S                                        JAVA

ISO10646-2003

ISO10646 2003                7

1                        95%                                3500

4

6

[14]

3500

_____

14

B/S

1999

1

BIG5

GB                    ISO10646

2

3

1

2

ISO 10646

3

1

13

2

14

1　http://idp.nlc.gov.cn/

2 http://202.96.31.42:9080/ros/index.htm

3　http://202.96.31.42:9080/wenxian/

4

1

"　"　　　　　　　　　　　　　　　　　　　　10+3　　　10

3　　　　　　　——"bbbbvvvppp.xxx"　"bbbb"

BookID　"vvv"　　　　　　　　　"ppp"　　　　　　　　"xxx"

3　*.jpg　　　　　　　　*.xml

*.pdf

52　　3　　　　　　　　　　3



2　　　　　　　　*.xml

```
01.   <?xml version="1.0" encoding="utf-16" standalone="yes"?>
02.   <page pageid="9999001040" bookid="" relate="9999001040">
03.     <staffidx>
04.     </staffidx>
05.     <workflow>71</workflow>
06.     <info>
07.       <format count="6">
08.         <modulus id="imagewidth">1103</modulus>
09.         <modulus id="imageheight">1856</modulus>
10.         <modulus id="landscape">0</modulus>
11.         <modulus id="l2r">0</modulus>
12.         <modulus id="type">0</modulus>
13.         <modulus id="case">000</modulus>
14.       </format>
15.       <setting count="23">
16.         <modulus id="1" />
17.         <modulus id="2" />
18.         <modulus id="3" />
19.         <modulus id="4" />
20.         <modulus id="5" />
21.         <modulus id="6" />
22.         <modulus id="7" />
23.         <modulus id="8" />
24.         <modulus id="9" />
```

```
25.          <modulus id="10" />
26.          <modulus id="11" />
27.          <modulus id="12" />
28.          <modulus id="13" />
29.          <modulus id="14" />
30.          <modulus id="15" />
31.          <modulus id="16" />
32.          <modulus id="17" />
33.          <modulus id="18" />
34.          <modulus id="19" />
35.          <modulus id="20" />
36.          <modulus id="21" />
37.          <modulus id="22" />
38.          <modulus id="23" />
39.        </setting>
40.      </info>
41.      <data>
42.        <text count="10">
43.          <line id="1">
44.            <position>1070,190,966,1766</position>
45.            <content>                                    </content>
46.          </line>
47.          <line id="2">
48.            <position>966,190,863,1766</position>
49.            <content>                                    </content>
50.          </line>
51.          <line id="3">
52.            <position>863,190,760,1766</position>
53.            <content>                              </content>
54.          </line>
55.          <line id="4">
56.            <position>760,190,656,1766</position>
57.            <content>                 </content>
58.          </line>
59.          <line id="5">
60.            <position>656,190,553,1766</position>
61.            <content>                                </content>
62.          </line>
63.          <line id="6">
64.            <position>553,190,447,1766</position>
65.            <content>                                     </content>
66.          </line>
67.          <line id="7">
68.            <position>447,190,343,1766</position>
69.            <content>                                </content>
70.          </line>
71.          <line id="8">
72.            <position>343,190,240,1766</position>
73.            <content>                                     </content>
74.          </line>
75.          <line id="9">
76.            <position>240,190,137,1766</position>
77.            <content>                           </content>
78.          </line>
79.          <line id="10">
80.            <position>137,190,33,1766</position>
81.            <content>                           </content>
82.          </line>
83.        </text>
84.        <annotate count="0" />
85.      </data>
86.  </page>
```

xml

03-05

07-14

15-39                                             pdf

41-85

        42-83                                  <text>

                                        <line>

                                              <position>

                                      x     y

                                  <content>

                              84        <annotate>

        <line>

              <content>

      xml

<content>                                              57

      ➢

            "    "    "    "          -

            "    "    "    "          -

➢

"   "    "   "        -

"       "   "   "          -

"     "

➢

"      "    "    "          -

➢

"      "    "     "          -

"      "    "    "          -

➢

"   "    "   "

"     "                    -

"     "                       -

"     "                       -

"     "                    -

"   "    "   "        +2

"     "                       -

"    "                  90

"      "   "    "

"     "

"     "

3

太僕館在俗呼驥馬頭府館院在宛平縣東張家灣巡檢司

楊村在土城西北北關巡檢司弘仁橋巡檢司城南州

里三十提舉司宣課司鹽場

檢校批驗所張家灣抽分竹木局

西岸黃船廠料磚廠河南

北關竹木局大通關鐵猫廠錦永館

分守州東曹師府通州左衛

東通州神武中衛通州右衛

院在府東察院橫定邊衛南州治激

場在舊城東關城

## 三河縣

縣治洪武甲明亭重修建正甲明亭在縣東旌善亭門西陰

陽學 醫學 養濟院 社倉 都察院俱無僧

會司官在縣 預備倉 軍儲倉 草廠治俱東北縣

陰東察院治在縣西察院公署在縣社稷壇今可改魏部道府

館一所治在西縣 戶部分司一所治在東

縣治南門街中明亭治在縣東旌善亭治在縣西陰陽學

## 武清縣

縣治在城中

4

5



古蹟

長城通州志云國初得唐寺丕墓誌石抆長
城周南具銘曰圪然孤墳長城之東可知長
城自北綿亙而南至唐時長城西南遺址尚
存也其曰蒙恬所築則非是攷昌平山水記
云湿餘得渡渡南有長城道跡遼史順州南
有齊長城天保時所築今潮縣武清二志俱
載境內有古長城疑昔聯為一耳今按縣東

潮陰志略　二十七

北一里有城址南抵武清北接通州俗稱長
隄或古城之古蹟興
晾鷹召方興紀要云在潮縣西南二十五里
高数丈周一頃元時遊獵多駐於此舊志云
在德仁務遼主遊獵駐輦之所今尚存土阜
高二丈餘又有故鷹召志稱在縣西四里亦
遼主遊獵之所今不詳其廢
呼鷹召元史至大元年七月築呼鷹召於潮

6

古蹟

長城通州志云國初得唐寺丕墓誌石於長
城周南其銘曰圮然孤墳長城之東可知長
城自北綿亘而南至唐時長城西南遺址尚
存也其曰蒙恬所築則非是玫昌平山永記
云溫餘得渡渡南有長城道蹟遼史順州南
有齊長城天保時所築今漷縣武清二志俱
載境內有古長城疑昔聯為一耳今按縣東

漷陰志略　二十七

北一里有城址南抵武清北接通州俗稱長
隄或古城之古蹟興

瞭鷹召方與紀要云玄狐漷縣西南二十五里
高數丈周一頃元時遊獵多駐於此舊志云
在德仁務遼主遊獵駐蹕輦之所今尚存土阜
高二丈餘又有故鷹召志稱在縣西四里亦
遼主遊獵之所今不詳其廎

呼鷹召元史至大元年七月築呼鷹台於漷

# From Text to Information -- Small Steps towards a Knowledgebase of Tang Civilization

Christian Wittern

## 1. Introduction

The Knowledgebase of Tang Civilization is aiming at providing a comprehensive interlinked research tool to aid in any research that has some relation to Tang Civilization. To build such a Knowledgebase, work has started with a focus on historical texts, e.g. the two dynastic "standard" histories *Jiu Tang Shu* 舊唐書 by Liu Xu 劉昫 and others (945), *Xin Tang Shu* 新唐書 by Ouyang Xiu 歐陽修 and others (1060) as well as Sima Guang's 司馬光 *Zizhi Tongjian* 資治通鑑, the latter provides a convenient chronicle of events of the Tang years for our purpose.

The first challenge in the attempt to make texts available for such a purpose is to transfer them in a digital format that can be used as a fundament for all further work. This transformation will need to reflect the basic structure of the text as well as provide a means to relate to the higher level semantic entities contained therein. With such a text in place, the next step is to isolate, analyze and normalize the information atoms the text relates, which are in this case names of persons mentioned in the texts, titles of books, temporal and spatial identifications (names of places, dates etc) but also administrative titles and acts and finally the events that these text are reporting.

In order to make all this useable in the Knowledgebase, however there is one further step needed: The individual information items need to be related to each other explicitly in some way. While this connection is of course contained in the textual sources, sometimes explicit, but at least impicit, it requires considerable effort to encode these relations in machine-readable form.

## 2. The *Zizhi Tongjian*

### 2.1. The text and its tradition

The *Zizhi Tongjian* 'Comprehensive mirror to aid in government' (1086) by Sima Guang 司馬光(1019-1086)is easily the most influential single work of chronological (編年體) historical writing in China, in influence on later works second only to the *Shiji* 史記 'Record of the Historian' by Sima Qian 司馬遷. The work has been frequently annotated, excerpted, and expanded; there have been critical examinations and a considerable number of sequels. Figure 1 shows the text available in the widely used Baina 百納 edition. As can be seen, the text is running through, with only occasional spaces to indicate a change of topic.

*Figure1.*

Sima Guang and his compilers[1] made a critical examination of the sources available to them and recorded instances where they had to choose from conflicting accounts. This was published separately as *Zizhi Tongjian Kaoyi* 資治通鑑考異 in 30 juan by Sima Guang himself. Among the annotations that where made to the text of the 資治通鑑 those by Shi Zhao 史炤 (1100-1160)[2] and Hu Sanxing 胡三省 (1230-1287)[3] are the most frequently used. Over time, they have been folded into the text, together with Sima Guangs *Kaoyi* to make up the text, as it is now most conveniently accessed through the punctuated edition from Zhonghua Chuju 中華書局. This text, which has also been the basis for our digitization, is shown in Figure 2.

---

[1] Apart from Sima Guang himself the three main editors where Liu Ban 劉攽(1023-1089), who was responsible for the Warring States, Qin and Han periods , Liu Shu 劉恕(1032-1078), who took charge of the period from the Three Kingdoms to Sui and Fan Zuyu 范祖禹(1041-1098) who edited the records for the Tang and Five Dynasties period.

[2] The *Shiwen* 釋文 commentary.

[3] 廣註 *Guangzhu* and 音註 *Yinzhu* commentary in 30 juan

*Figure2.*

## 2.2. Structural features of the text

As can be seen from <u>Figure 2</u>, the structure of this text is much more complex than the one seen in <u>Figure 1</u>. Using this as a basis for digitization seems a much more complex and unnecessary cumbersome undertaking. It has however the advantage to provide additional information about the text and the entities the text talks about; for this reason in a long-term perspective, it seemed most convenient to start out from the modern edition.

The text as it is published is divided into the records of 16 dynasties or kingdoms. Of these, the Tang Records 唐紀 in 81 juan is by far the most voluminous. For the Tang Knowledgebase, work did concentrate exclusively on this part of the text. Within the Tang Records, the division is by Emperor, then by era period and finally by year. The editors choose to treat a calendar year as the basic unit, although obviously some changes of emperor or era do not fall at the beginning of the year. Within a year, the narrative is divided into individual paragraphs, which mostly follow those white spaces of the earlier editions, as can be seen by comparing the 'Zhonghua Shuju' edition with the 'Baina' edition. In cases where the editors did feel a further division was necessary, they did introduce their own paragraphs and, to distinguish these two types of paragraphs, the ones that appeared already in the early woodblock edition (and might have been introduced for all we know by Sima Guang himself) have numbers attached to them, which run through for

one whole year. In a first approximation, such a paragraph can be seen as the most basic narrative unit, reporting one 'event'.

Within a paragraph, the notes of various types are distinguished from the running text by smaller point size, but they are not further separated neither by source nor by category. Other clearly distinguishable features of the text are waivy and straight lines, which indicate titles of texts and other named entities respectively.

## 2.3. Translating structure to markup

As described above, the structural features of the text visible in the layout point to semantic features. For the electronic text to be useful, these features have to be named and transposed into machine readable form. For the purpose of this project, the *Guidelines for Electronic Text Encoding and Interchange*[4] have been used and adapted to the special purposes of this project[5]. In most cases the transposition from structure to markup has been straightforward and could be applied in a semi-automated fashion. What proved to be most difficult was the introduction of distinctions that had no equivalent in the printed page. For example, the named entities marked with a straight line lumped together quite distinct features:

- Personal names (`<rm>`)[6]
- Names of Places (`<dm>`)
- Era names (`<y>`)
- Names of dynastic periods or kingdoms (`<dyn>`)

Identifying and assigning the proper categories for these items proved to be by far the most time consuming task. The difficulty here is that even with supporting authority files, which we gradually created during this project, the assignment itself still has to be checked carefully. Another problem was the desire to not only identify the categories, but also link to the authority file, which might have different names for a person and additional information.

The text in smaller print, which lumped together all annotations that had been added to the original text either by Sima Guang or by a range of later commentators and editors, can be distinguished in the following way:

- text critical notes (`<app>`)

[4] Edited by Michael Sperberg-McQueen and Lou Burnard. First published in 1994, the most recent revision is of 2002.

[5] Adaptability is one of the major features of these Guidelines, which allow them to be both very broadly applicable and yet specific enough for varying local requirements. The adoption in this case consisted mostly of removing unnecessary elements and making the content models more stringent, so that it was easier to check the markup.

[6] In parentheses the elements used to indicate this feature. See the Appendix for a correspondence of these notations to the standard TEI notation.

- sound glosses (`<gloss>`)
- cross references (`<xref>`)
- critic of sources (this is indicated by a reference to *Zizhi Tongjian Kaoyi*))
- additional information (provided by Hu Sanxing or others)
- citations (`<cit>`)

It is by no means always obvious, who is responsible for a given note and when it was introduced, although in some cases, the name is mentioned explicitly.

In applying the markup, the need to economize on the amount of effort that is required to create texts with markup adds additional limitations; for this reason not all desirable features have been introduced.

A short excerpt from the markup of the master files, showing a part from the text in Figure 2 is given in the Appendix.

# 3. Beyond markup

The purpose of applying this type of markup is of course to make available the text content for further processing. Within the context of the Knowledgebase, this means deriving information from the text, abstracting from the individual expressions used there.

One area where this is specifically important is with regard to the persons mentioned in the text. Depending on the context and their standing in the world, many persons are referred to by different names in the course of their life, and this is naturally also reflected in our sources. Additionally, the chronological account and the several layers of notes to refer to persons differently in many cases.

In addition to simply marking textual features as being names of places, persons and so on, work has began to identify persons and relate their names; this is done separately from the text that constitutes the 'resource layer' in what we call the 'information layer'. This is of course a significant task that even for a text of comparatively limited scope, as the Tang Records of the 資治通鑑, where for about the first half of the Tang period[7], at the moment we have identified more than 6000 different persons mentioned in the text[8]. This identification is complicated by the fact that names or identifications such as Prince of Qin (秦王) are used for different persons at different times; frequently persons are only mentioned by

---

[7] To be exact, at present the period from the first year of Wude 武德 (618) to the end of Tianbao 天寶 (756) has been analyzed.

[8] This figure, which should of course not be taken to be the final toll, does also include persons from earlier periods that are mentioned in the text, as well as the names of later commentators and other persons, whose names appear in the sources. Only very preliminary work has been done to further analyze and categorize these data so far and we expect to devote more time and effort on this aspect in the years to come.

their given name, which may increase the number of identical strings that still have to be assigned to different records.

While the amount of effort to identify persons in such a way is prohibiting high, on the other side, we hope that the resulting data can be used as a training corpus to further annotate texts in this way semi-automatically, both texts from the Tang period, where the data and its content in the information layer can be straightforwardly applied, but also texts from other periods, where heuristics of the textual context derived from these data can provide hints to what has to be considered a personal name.

To give an example, here is a list of names used in our sources for the second Tang emperor Taizong, who is by far most frequently mentioned in the period considered here.

- 世民 (237[9])
- 唐太宗 (4)
- 太宗 (215)
- 太宗文武皇帝 (1)
- 太宗文武聖皇帝 (1)
- 文武 (1)
- 文武皇帝 (1)
- 文武大聖大廣孝皇帝 (1)
- 文武大聖皇帝 (1)
- 文武皇帝 (1)
- 文皇帝 (2)
- 秦公世民 (1)
- 秦王 (4)
- 秦王世民 (76)
- 趙公 (1)
- 趙公世民 (3)

There are 17 different appellations for Taizong and it should be clear that it would be quite impossible to find all places where he is mentioned in the source text by trying to perform a simple search for his name.

There are of course countless more aspects in which even a comparatively simple markup like the one used here can be used to further extract and analyze the text to gain new information. It is for example interesting to see at which time which name was used to refer to a person, or which nominal title was conferred

---

[9] The figures in parentheses are the number of occurrences in the text of the *Zizhi Tongjian* from 618 up to the end of the era Tianbao (756).

at which time. Since our text is a chronological record, it is quite easy to derive this information from the structure of the text.[10] The following Table 1 gives the names as they have been used, broken down by year; it is also mentioned whether the name is mentioned in the narrative main text, or somewhere in the commentary.[11]

Table1.

| Year | Occ. total | Occ. in text (note); total: 556 |
|---|---|---|
| 618 (武德–1) | 36 | • 秦公世民 1 ( 0 )<br>• 世民 19 ( 1 )<br>• 趙公 1 ( 0 )<br>• 趙公世民 3 ( 0 )<br>• 太宗 0 ( 2 )<br>• 秦王 2 ( 0 )<br>• 秦王世民 7 ( 0 ) |
| 619 (武德–2) | 23 | • 秦王世民 9 ( 0 )<br>• 太宗 0 ( 1 )<br>• 世民 11 ( 2 ) |
| 620 (武德–3) | 58 | • 秦王世民 18 ( 0 )<br>• 世民 27 ( 0 )<br>• 太宗 1 ( 13 ) |
| 621 (武德–4) | 85 | • 秦王世民 16 ( 0 )<br>• 太宗 1 ( 11 )<br>• 世民 57 ( 1 ) |
| 622 (武德–5) | 66 | • 秦王世民 14 ( 0 )<br>• 世民 35 ( 1 )<br>• 太宗 4 ( 16 ) |
| 623 (武德–6) | 6 | • 秦王世民 5 ( 0 ) |

[10] Since the text does give exact dates for many events, we do in fact plan to make this information accessible in machine readable form, but it has not yet been applied to the master copy of our text. It should also be noted that the narrative is at places not strictly chronological, but does follow the events where necessary, so the purely positional computation does have its limitations.

[11] Since this table is just ment to illustrate the idea, some years between the beginning of Zhenguan 貞觀 and the end of Tianbao have been skipped to save some space.

| Year | Occ. total | Occ. in text (note); total: 556 |
|---|---|---|
| | | • 世民 1 ( 0 ) |
| 624 (武德–7) | 46 | • 太宗 2 ( 7 )<br>• 秦王世民 5 ( 0 )<br>• 世民 32 ( 1 ) |
| 625 (武德–8) | 1 | • 秦王世民 1 ( 0 ) |
| 626 (武德–9) | 57 | • 秦王世民 1 ( 0 )<br>• 世民 46 ( 3 )<br>• 太宗 4 ( 3 ) |
| 627 (貞觀–1) | 7 | • 太宗 0 ( 5 )<br>• 秦王 1 ( 1 ) |
| 628 (貞觀–2) | 7 | • 太宗 1 ( 6 ) |
| 649 (天寶–8) | 2 | • 太宗 1 ( 0 )<br>• 文武大聖皇帝 1 ( 0 ) |
| 754 (天寶–13) | 2 | • 太宗 1 ( 0 )<br>• 文武大聖大廣孝皇帝 1 ( 0 ) |

# 4. Conclusions

It should be clear that this is just one example of how by abstracting from information contained in the text new layers of information can be constructed. While only a very simple example has been given, other rather easily gained information items are the co-occurrence of persons and places participating in an event (as described in one paragraph). Also, formal records of astronomical events, sacrifices and excursions of the Emperor can rather easily be harvested into machine readable form. At times, this might resemble approaches common from data mining methodology, there are nevertheless differences. Since the research target is information to be derived from narrative texts, all irregularities and inconsistencies of natural language and literary style have to be taken into account. Data mining and analytic tools are thus most usefully employed to give a rough dataset, which is then refined by a human researcher.

For this reason, the goal of constructing the *Knowledgebase of Tang Civilization* does not primary lie in constructing a comprehensive repository of information and resources pertaining to the Tang, but to provide the researcher with an

infrastructure that assists her in getting answers to her questions from the sources, and sometimes even getting answer for questions she did not ask.

# 5. References

C. M. Sperberg-McQueen Lou Burnard *Guidelines for Electronic Text Encoding and Interchange P4*, TEI Consortium Oxford, Providence, Charlottesville, Bergen, 2002.

Biaodian Zizhi Tongjian Xiaozu 楳點資治鑑小組*Zizhi Tongjian* 資治通鑑, Beijing 1956.

# 6. Appendix

## 6.1. Markup Example

The following is an example of how a part of the text in the sections shown in Figure 1 and 2 has been encoded in our files. The section startes with the header for the year, but then skips right directly to paragraph 17 on page 5880.

The following elements have been abbreviated: <dm> for <name type="place">, <rm> for <name type="person">, <y> for <date>, <dyn> for <name type="dynasty">. Instead of linebreaks, <lb> indicates places where breaks occur in our internal master documents.

```
<div><head n="武德-3">三年<note place="inline">(庚辰、六二
0)</note></head>
  (...)
  <div n="17"><p>夏，四月，丙申，上祠<dm>華山</dm>；壬寅，還<dm>長
安</dm>。<lb n="5880-175"/><note place="inline">華，戶化翻。從，還宣
翻。</note></p></div>
  <div n="18"><p>置<dm>益州道</dm>行臺，以<dm>益</dm>、<dm>利</dm>、
<dm><lb n="5880-200"/>會</dm>、<dm>鄘</dm>，<dm>涇</dm>、<dm>遂</dm>
六總管隸焉。<note place="inline"><dm>益州</dm>，<dyn key="ch174">隋
</dyn>之<dm>蜀郡</dm>。<dm>利州</dm>，<dyn key="ch174">隋</dyn><lb
n="5880-225"/>之<dm>義城郡</dm>，<dyn key="ch166">梁</dyn>之<dm>黎
州</dm>，<dyn key="ch110">晉</dyn>之<dm>晉壽</dm>，<name>蜀</name>
之<dm>漢壽</dm>，<dyn key="ch129">漢</dyn>之<dm>葭萌</dm>也。<lb
n="5880-250"/><dm>會州</dm>，<dyn key="ch174">隋</dyn>之<dm>涼川縣
</dm><dm>會寧鎮</dm>，<name>西魏</name>之<dm>會州</dm>也。<dm>鄘州
</dm>，<dyn key="ch174">隋</dyn>之<dm><lb n="5880-275"/>上郡</dm>，
<name>西魏</name>之<dm>敷州</dm>，<dyn key="ch169">後魏</dyn>之<dm>
北華州</dm><dm>中部</dm><dm>敷城郡</dm>也，<y to="中" value="太和">
太和中</y><lb n="5880-300"/>為<dm>東秦州</dm>。<dm>涇州</dm>，<dyn
key="ch174">隋</dyn>之<dm>安定郡</dm>。<dm>遂州</dm>，<dyn
key="ch174">隋</dyn>之<dm>遂寧郡</dm>，<dyn key="ch129">漢</dyn>之
```

&lt;dm&gt;&lt;lb n="5880-325"/&gt;廣漢縣&lt;/dm&gt;也 。是時&lt;dm&gt;益州&lt;/dm&gt;行臺所統，起&lt;name&gt;蜀&lt;/name&gt;，跨&lt;dm&gt;隴&lt;/dm&gt;而東北 。&lt;/note&gt;&lt;/p&gt;&lt;/div&gt;

   &lt;div n="19"&gt;&lt;p&gt;&lt;rmkey="r00280"&gt;劉武&lt;lb n="5880-350"/&gt;周&lt;/rm&gt;數攻&lt;dm&gt;浩州&lt;/dm&gt;，為&lt;rmkey="r01611"&gt;李仲文&lt;/rm&gt;所敗 。&lt;note place="inline"&gt;數，所角翻 。敗，補邁翻 。&lt;/note&gt;&lt;rmkey="r06179"&gt;&lt;lb n="5880-375"/&gt;宋金剛&lt;/rm&gt;軍中食盡；丁未，&lt;rmkey="r06179"&gt;金剛&lt;/rm&gt;北走，&lt;rmkey="r01602"&gt;秦王世民&lt;/rm&gt;追之 。&lt;/p&gt;&lt;/div&gt;

   &lt;div n="20"&gt;&lt;p&gt;&lt;rmkey="r03081"&gt;羅士&lt;lb n="5880-400"/&gt;信&lt;/rm&gt;圍&lt;dm&gt;慈澗&lt;/dm&gt;，&lt;note place="inline"&gt;&lt;cit&gt;&lt;title&gt;隋志&lt;/title&gt;：&lt;q&gt;&lt;dm&gt;河南郡&lt;/dm&gt;&lt;dm&gt;壽安縣&lt;/dm&gt;有&lt;m ref="c5880-401"&gt;慈澗&lt;/m&gt; 。&lt;/q&gt;&lt;/cit&gt;&lt;cit&gt;&lt;title&gt;水經註&lt;/title&gt;：&lt;q&gt;&lt;dm&gt;新安&lt;/dm&gt;有&lt;dm&gt;&lt;lb n="5880-425"/&gt;孝水&lt;/dm&gt;，&lt;dm&gt;孝水&lt;/dm&gt;東十里有水，世謂之&lt;m ref="c5880-401"&gt;慈澗&lt;/m&gt; 。&lt;/q&gt;&lt;/cit&gt;&lt;/note&gt;&lt;rmkey="r02590"&gt;王世充&lt;/rm&gt;使&lt;rmkey="r04210"&gt;太子玄應&lt;/rm&gt;&lt;lb n="5880-450"/&gt;&lt;app resp="【 章 】"&gt;&lt;lem&gt;救&lt;/lem&gt;&lt;rdg wit="十二行本；乙十一行；孔本"&gt;拒&lt;/rdg&gt;&lt;/app&gt;之，&lt;lb n="5880-475"/&gt;&lt;/note&gt;&lt;rmkey="r03081"&gt;士信&lt;/rm&gt;刺&lt;rmkey="r04210"&gt;玄應&lt;/rm&gt;墜馬，&lt;note place="inline"&gt;刺，七亦翻 。&lt;/note&gt;人救之，得免 。&lt;/p&gt;&lt;/div&gt;

   &lt;div n="21"&gt;&lt;p&gt;壬子，&lt;lb n="5880-500"/&gt;以&lt;dm&gt;顯州道&lt;/dm&gt;行臺&lt;rmkey="r05490"&gt;楊士林&lt;/rm&gt;為行臺尚書令 。&lt;note place="inline"&gt;去年正月，&lt;rmkey="r05490"&gt;楊士林&lt;/rm&gt;降 。&lt;lb n="5880-525"/&gt;&lt;/note&gt;&lt;pb n="5881-"/&gt;&lt;/p&gt;&lt;/div&gt;

   &lt;div n="22"&gt;&lt;p&gt;甲寅，加&lt;rmkey="r01602"&gt;秦王世民&lt;/rm&gt;&lt;dm&gt;益州道&lt;/dm&gt;行臺尚書令 。&lt;/p&gt;&lt;/div&gt;&lt;/div&gt;

100081

zniu@nlc.gov.cn

# The Research and Design of Multilingual Process System

Zhendong Niu
China Digital Library Corp. Ltd. Beijing, China 100081
zniu@nlc.gov.cn

Abstract: This paper discusses the multilingual processing system. It mainly focuses on the infrastructure of multilingual process system and the related key techniques including thesaurus, classification schemes etc.

Keywords: multilingual process, thesaurus, classification schemes, multilingual retrieval system

**1**

---

90

1. 1



Multilingual MenUSE

for EPOQUE

1.1

" 863"                                                                                                          "
   "

## 2

[1-4]　　2 1



**2.1**

2 1　,

Web Services,

# 3

## 3.1



**3.1**

[5-6]

## 3.1

## 3.2 　　　　　　　　　　　　　　　　　　　　　　　　　　　　　,

XML

**3.2**                                          **Tree View**

# 4

                                classification schemes

40                                  1998        "            "
                "                                "
                        "                                "        2000            "            "

| | | | | | |
|---|---|---|---|---|---|
| | | 中国图书馆分类法简表 (第四版) | | | |
| A | 马、列、毛泽东思想、邓小平理论 | P5 | 地质学 | TD | 矿业工程 |
| B | 哲学、宗教 | P7 | 海洋学 | TE | 石油、天然气工业 |
| C | 社会科学总论 | P9 | 自然地理学 | TF | 冶金工业 |
| D | 政治、法律 | Q | 生物科学； | TG | 金属学与金属工艺 |
| E | 军事 | Q1 | 普通生物学 | TH | 机械、仪表工业 |
| F | 经济 | Q2 | 细胞生物学 | TJ | 武器工业 |
| G | 文化、科学、教育、体育 | Q3 | 遗传学 | TK | 能源与动力工程 |
| H | 语言、文字 | Q4 | 生理学 | TL | 原子能技术 |
| I | 文学 | Q5 | 生物化学 | TM | 电工技术 |
| J | 艺术 | Q6 | 生物物理学 | TN | 无线电电子学、电信技术 |
| K | 历史、地理 | Q7 | 分子生物学 | TP | 自动化技术、计算机技术 |
| N | 自然科学总论 | Q81 | 生物工程学（生物技术） | TQ | 化学工业 |
| O | 数理科学和化学 | [Q89] | 环境生物学 | TS | 轻工业、手工业 |
| O1 | 数学 | Q91 | 古生物学 | TU | 建筑科学 |
| O3 | 力学 | Q93 | 微生物学 | TV | 水利工程 |
| O4 | 物理学 | Q94 | 植物学 | U | 交通运输 |
| O6 | 化学 | Q95 | 动物学 | V | 航空、航天 |
| O7 | 晶体学 | Q96 | 昆虫学 | V1 | 航空、航天技术的研究与探索 |
| P | 天文学、地球科学 | Q98 | 人类学 | V2 | 航空 |

**4.1**

"

"

"                             "

knowledge organization

thesaurus controlled
vocabularies classification code
subject headings

[7]

**5**

Berkeley

1    Jerome Yen. Multimodel and multilingual informedia – The iVIEW system. 2[nd] CCDL. Beijing. 2004.

2    Michael R. Lyn, Edward Yau, Sam Sze. A multilingual, multimodal digital video library system. Proceedings of the 2[nd] ACM/IEEE-CS joint conference on Digital libraries. P 145-153. ACM Press. 2002

3    Akira Maeda. Multilingual information processing for Digital libraries. Department of Computer Science, Ritsumeikan University 2001

4    Nuno Freire. Integration of Multilingual Classification Systems with the Dienst digital library system. http://www.ercim.org/ws-proceedings/DEL0S8/freire.pdf.

5    A Steven Pollitt. The key role of classification and indexing in view-based searching. IFLA'97 Copenhagen. Aug 31 – Sept. 2. 1997

6    C S Li A S Pollitt & M P Smith. Multilingual MenUSE – A Japanese front-end for searching English language database and vice versa. 14[th] BCS IRDG information retrieval colloquium, Lancaster 1992 Tony McEnery & Chris Paice (eds). Springer Verlag. Pp 14-37.

7    Zhendong Niu Mingkai Dong Jie Zhang Huaming Chen. A knowledge based solution for Digital libraries. 2[nd] CCDL. Beijing. 2004.

# Character Processing Based on Character Ontology

MORIOKA Tomohiko

## 1  Introduction

We use characters as a basis for data representation in computers, and as a tool for communication over computer networks. Computer programs are made of characters, we exchange mails that are sequences of characters, and many of the contents available over the Internet are realized in the form of characters.

Currently, in the field of information processing, characters are defined and shared using coded character sets. Character processing based on coded character sets, however, has two problems:

1. Coded character sets do not always contain a necessary character

2. Characters in coded character sets have fixed semantics

To resolve the problems, I proposed "Chaon" model which is a new model of character processing based on character ontology. In Chaon model, characters are defined, represented and processed according to its own character databases. Characters in Chaon model are independent from coded character sets for information interchange, and semantics of the characters stored in the database can be freely added or altered.

To realize the character processing based on Chaon model, I started "CHISE (Character Information Service Environment)" project with some other members. In CHISE project, we have developed some systems and databases to edit/process/print characters and texts. CHISE project is an open source project, so the results are freely distributed. We are realizing character processing environment based on Chaon model. In this paper, I explain an overview of the current state of character processing technology in CHISE project.

## 2  Chaon model

In Chaon model of character representation, a character is not a code point of a coded character set, but a set of the features it has. Characters are represented as character objects, and character objects are defined by character features. Character objects and character features are stored in character databases, and a character can be accessed using its feature as a key.

There are various information related with characters, so we can regard various things as character features, for example, shapes, phonetic values, semantic values, code points in various character codes. Figure 1 shows a sample image of a character representation in Chaon model which indicate a character "吉".



Figure 1: sample of character features to indicate "吉"

In CHISE model, each character is represented by a set of character features, so we can use set operations to compare characters. Figure 2 shows a sample of a Venn diagram of character objects. The diagram indicates that there are common semantic and phonetic values between characters "言", "云" and "謂" even if they don't have the same glyph, and characters "云" and "雲" have the common semantic and phonetic values in China.



Figure 2: Venn diagram of characters

As we have already explained, a coded character set (CCS) and a code point in the set can also be character features. Those features enable exchanging a character information with the applications that depend on coded character sets. If a character object has only a CCS feature, processing for the character object is the same with processing based on the coded-character model now we are ordinarily using. Namely we can regard the coded-character model as a subset of Chaon model.

# 3 Overview of CHISE system

Character processing based on Chaon model is to represent each character as a set of character features instead of a code point of a coded character set and process the character by various character features. It indicates that character processing system based on Chaon model is a kind of database system to operate character ontology. So the major targets of CHISE project are (1) character database systems, (2) character database contents and (3) CHISE based applications. CHISE Project is working on the targets and provides some results as free software.

As character database systems (and language bindings), following implementations are available:

**libchise** is a library to provide fundamental features to operate character database

**XEmacs CHISE** is a Chaon implementation based on XEmacs [9] (extensible text editing environment with Emacs Lisp interpreter) [Figure 3]

**Ruby/CHISE** is a Chaon implementation based on Ruby [8] (scripting language)

**Perl/CHISE** is a Chaon implementation based on Perl (scripting language)

For the Chaon implementations, currently two database contents are available as follows:

**CHISE basic character database** is a general character database attached to XEmacs CHISE

**CHISE-IDS** is a database about shapes of Ideographic characters

As CHISE based applications, we need at least editing system, printing system and character database utilities. As editing system, XEmacs CHISE is available. As printing system, the following programs are available:

**Ω/CHISE** is a multilingual type setting system based on Ω [5] (a multilingual TEX)

Figure 3: XEmacs CHISE

**chise2otf** is a converter to process CHISE texts with pLATEX[7] + OTF package [6]

**chise-tex.el** is like chise2otf, but it is implemented as a coding-system of XEmacs CHISE (written by Emacs Lisp)

As character database utilities, there are some emacs lisp programs, Ruby scripts attached to Ruby/CHISE and Perl scripts attached to Perl/CHISE.

In addition, the following system is available:

**Kage** is an automatic Ideographic glyph generating system [10]

Currently Kage is used in Ω/CHISE and chise2otf.

## 4 XEmacs CHISE

XEmacs CHISE (Figure 3), Ruby/CHISE and Perl/CHISE provide operations about character features based on Chaon model. In this paper, I describe a overview of them in the XEmacs CHISE case as an example.

### 4.1 Character representation

XEmacs CHISE represents a character as a character object. A character object is a first class object, just like a character in XEmacs, and completely a different type from the integer.

For programmers of Emacs Lisp applications, a character is represented by a set of character features. The character feature is a pair of a feature key and its value. The feature key must be a symbol and the value can be any Lisp object, that is, a character object itself can be a feature value of other characters, which makes it easy to represent relation network among characters. Namely a character is represented by an association-list of Lisp. In XEmacs CHISE, such kind of association-list which represents a character by its features is named "character-specification (char-spec)".

Each character also has an unique identifier called "character-id" although it is usually hidden from Emacs Lisp users.

## 4.2   Character object related functions

In order to define and handle a character and character features, XEmacs CHISE provides the following built-in functions.

**Function** define-char (*char-spec*)

> defines a character object that has a set of character feature *features*, and returns the object. *char-spec* should be an association-list.

> **[Example]**

```
(define-char
 '((name               . "CJK RADICAL MEAT")
   (general-category    symbol other)
   (bidi-category      . "ON")
   (mirrored           . nil)
   (ideographic-radical . 130) ; 肉 (radical)
   (ideographic-strokes . 0) ; (body strokes)
   (total-strokes       . 4) ; (total strokes)
   (<-ideographic-component-forms
    ((=ucs               . #x8089)     ; 肉
     ))
   (=ucs                 . #x2EBC)     ; 月
   (->subsumptive ; (included variants)
    ((=gt                . 37857)
     (=gt-pj-6           . #x3879)
     (=daikanwa          . 29237)
     )
    ((=ucs@unicode       . #x2EBC)
     )
    ((=big5-cdp          . #x8A73)
     )
    ((=big5-cdp          . #x8958)
     (=gt-k              . 00417)
     (=gt-pj-k1          . #x377D)
     ))
   ))
```

**Function** get-char-attribute (*character feature **&optional** default-value*)

> returns a value of a character feature specified by the key *feature* of a character object *character*.

> If the value of *feature* is not defined, *default-value* is returned. If *default-value* is omitted, nil is returned.

> **[Example]**

```
(get-char-attribute ?\u2EBC 'name)
→ "CJK RADICAL MEAT"
```

**Function** put-char-attribute (*character feature value*)

> adds or changes the value of a feature of a character. This function sets a Lisp object *value* to a value of the character feature specified by the key *feature* of a character object *character*.

> **[Example]**

```
(get-char-attribute ?あ 'foo)
→ nil
(put-char-attribute ?あ 'foo 1)
→ 1
(get-char attribute ?あ 'foo)
→ 1
```

**Function** remove-char-attribute (*character feature*)

removes character feature *feature* from character object *character*.

**Function** find-char (*char-spec*)

retrieves the character that has specified features *char-spec*.

XEmacs CHISE provides a map function for character features also. This function aims at finding characters with certain character features or processing characters using its character features.

**Function** map-char-attribute (*function feature*)

This function maps *function* over entries in *feature* (an association-list). *Function* is called with two arguments, a key and a value in the list repeatedly, until all the pairs in *feature* is used up.

**Function** char-attribute-alist (*character*)

returns the features of the character *character*. Every feature of a character is retrieved by this function.

**Function** char-attribute-list ()

returns the list of all existing character features except coded character sets.

XEmacs CHISE has on-memory database per each process besides the CHISE character database shared in the CHISE environment. The on-memory database works as a kind of cache memory for the external database. If a character feature is not found in the on-memory database, the feature value is read from the external database and the value is stored into the on-memory database. If a character feature is found in the on-memory database, XEmacs CHISE does not access the external database.

Modification functions for characters or their features of XEmacs CHISE, such as put-char-attribute, work for on-memory database. However it is volatile. So XEmacs CHISE has a function to save the character data.

**Function** save-char-attribute-table (*feature*)

saves each character's value of character feature *feature* into the CHISE character database.

XEmacs CHISE has functions to clear character data in the on-memory database to be able to reread from the CHISE database.

**Function** reset-char-attribute-table (*feature*)

clears character feature *feature* of every character in the on-memory database to be able to reread each value of the *feature* from the CHISE database.

**Function** reset-charset-mapping-table (*coded-charset*)

clears decoding-table of *coded-charset* in the on-memory database to be able to reread from the CHISE database.

XEmacs CHISE has a function to read every character's value of a specified feature from the CHISE database at a burst.

**Function** load-char-attribute-table (*feature*)

reads every character's value of a specified *feature* from the CHISE database at a burst.

In addition, there is a function to register a character feature.

**Function** mount-char-attribute-table (*feature*)

registers character feature name *feature* as a target to read from the CHISE character database.

By the way, Ruby/CHISE and Perl/CHISE don't have on-memory database, so written character data are written into the CHISE database directly. So there are no functions to sync between on-memory database and the CHISE database.

# 5 Character features

## 5.1 Categorization

In the character processing based on Chaon model, it is important to analyze characters and their various properties and behaviors and represent them as character features. We can find sundry properties and behaviors of characters, and we can use infinite kind of character features. However common character database requires a guideline about character features. So we think that it is feasible to regard each character feature as an abstraction of an operation for characters.

In the point of view, character features can be categorized like following:

1. general character properties (such as descriptions of dictionaries)

2. mappings for character IDs

3. relations between characters

For example, radicals, strokes and phonetic values can be classed into the category 1, code points of UCS [2] can be classed into the category 2 and relations between character variants can be classed into the category 3.

Information of the category 2 is used for processing about character codes, such as code conversion. Processing about character codes consists of two kind of operations: encoding and decoding. To encode a character by a CCS is to get the CCS feature's value in the character. To decode a code-point of a CCS is to search a character whose value of the CCS feature is the code-point. Processing about character codes should be fast, so the CHISE character database has special indexes for decoding.[1]

For the processing about character variants, information of the category 3 is used.

## 5.2 Description for complex information

For development of a general purpose character database, we may find some cases that there are different kind of usages, purposes, applications, sources, interpretations, theories, etc. so it is hard to chose one feature value and we want to provide alternative values. In that cases, we may want to add metadata, such as sources of the values. To resolve the problem, we have to introduce structured feature value or structured feature name (key).

To represent structured feature values, a format named "character reference (char-ref)" is used in CHISE. It is a kind of property-list of S-expression (Lisp), property name indicates kind of metadata and property value indicates its data. As a special property name, :char is reserved to indicate a character which is added the metadata. Currently :sources is defined to indicate information source.

CHISE also has a format to represent structured feature names. In the structured feature names, "domain identifiers" and/or "metadata identifiers" are added to ordinary (base) feature names. The format is defined as following definitions:

<concrete feature name>
    := <base feature name> @ <domain identifier>
    |   <concrete feature name> / <domain identifier>

<metadata feature name>
    := <concrete feature name> * <metadata identifier>

For example, when total strokes is represented by character feature total-strokes and ucs is used as a domain identifier, concrete feature name is total-strokes@ucs. When source is represented by metadata identifier sources, total-strokes@ucs's source is represented by metadata feature name total-strokes@ucs*sources.

If there is a correspondence between different kind of features, such as radical and body-strokes, we can represent the correspondence by a domain identifier. For example, when radical is represented by ideographic-radical and body-strokes is represented by ideographic-strokes, two concrete feature names

    ideographic-radical@ucs
    ideographic-strokes@ucs

are corresponding.

---

[1]For the special treatment, we distinguish the category 1 and 2, but it seems that there are no essential differences.

## 5.3　Inheritance of character definition

If we construct a large scale character database including a lot of character variants, inheritance of character definition is good way to avoid to write a lot of common features. So CHISE introduces four special features to represent parent and child relations:

<−**subsumptive** defined character is a child of each character indicated by its value

<−**denotational** likewise

−>**subsumptive** each character indicated by its value is a child of the defined character

−>**denotational** likewise

# 6　Database contents

Character database is a fundamental part of the Chaon character representation model. Users can, of course, freely define or modify characters by adding new character features, but a rich and accurate database would be a great place to start, and it will also attract new users. We have thus developed a standard character database for Chaon implementations. Currently two database distributions are available:

1. CHISE basic character database

2. Database about structure information of Ideographs (CHISE-IDS)

The former is a basic character database attached in XEmacs CHISE which is realized by a collection of define-chars while the later is a database for Ideographs to represent information about shapes which is represented by "Ideographic Description Sequence (IDS)" format defined in ISO/IEC 10646-1:2000 [2].

Structure information of Ideographs is information about combinations of components of Ideographs. A lot of Ideographs can be represented by a combination of components, so the information is a useful. It is not only representation of abstract shapes, but also related with semantic values and/or phonetic values. So we planned to develop a database about structure information of Ideographs for every Ideograph which consists of combination of components. In the 2001 fiscal year, we realized CHISE-IDS database with supporting of "Exploratory Software Project (未踏ソフトウェア創造事業)" run by IPA (Information-technology Promotion Agency, Japan) [11]. Currently it supports CJKV Unified Ideographs and Extension A of ISO/IEC 10646-1:2000 [2] and Extension B of ISO/IEC 10646-2 [3]. We are also working for representative glyph image of JIS X 0208:1990 and Daikanwa Dictionary.

Before we developed CHISE-IDS database, there are some databases including structure information of Ideographs: CDP database [1] by Academia Sinica, database about gaiji (private used character) used in CBETA and "Konjaku Mojikyo" [4]. These databases use original formats so it is not easy to convert to IDS format. Konjaku Mojikyo is a proprietary software so their data are not opened for the public. In the view of licence, CDP database and CBETA database are available with free software licenced under the term of GPL while Konjaku Mojikyo is not. So we converted CDP database and CBETA database to IDS and integrate them with CHISE database.

Currently CHISE basic character database (is a part of XEmacs CHISE) and CHISE-IDS package are distributed separately. However CHISE-IDS package provides an installer to integrate CHISE-IDS database files with CHISE basic character database. CHISE-IDS package also have some utility programs to use structure information of Ideographs in XEmacs CHISE:

**ids.el** IDS parser

**ids-read.el** utility to read CHISE-IDS database files into XEmacs CHISE

**ids-dump.el** utility to dump structure information of Ideographs stored in XEmacs CHISE (represented by character feature ideographic-structure) into CHISE-IDS database format

**ids-util.el** utility to convert structure information of Ideographs into other representative glyph images corresponding with specified domains

**ids-find.el** utility to search Ideographs by components

Currently, ids-find.el has two commands to search Ideographs by components named: ids-find-chars-including-components (= ideographic-structure-search-chars) and ids-find-chars-covered-by-components.

If you type

M-x ideographic-structure-search-chars [CR] components [CR]

Ideographs that have every component are displayed. (Figure 4)



Figure 4: result of M-x ideographic-structure-search-chars [CR] 木金 [CR]

If you type

M-x ids-find-chars-covered-by-components [CR] components [CR]

Ideographs that consists of one or more usage of every component are displayed. (Figure 5).



Figure 5: result of M-x ids-find-chars-covered-by-components [CR] 木金 [CR]

# 7    Conclusion

I have described a brand new character processing model Chaon and overview of CHISE project and character processing in CHISE systems such as XEmacs CHISE. Chaon character representation model is powerful and

radical enough to solve the problems that the present coded character model has, and the implementation of XEmacs CHISE or other CHISE systems have shown that the model is feasible enough to build a application on.

Chaon model gives users freedom to create, define and exchange characters of their need, as it is easy to change character databases or modify character features dynamically. XEmacs CHISE provides a good framework to experiment the character representation. With the CHISE basic character database, XEmacs CHISE can handle various characters including characters defined in Unicode. Even if a character is not defined in Unicode, users can add it into CHISE database to define it by its features. Users can handle each character based on their point of view or policy. For example, XEmacs CHISE provides some unification rules or mapping policies about Unicode. With the CHISE-IDS database, users can search Ideographs easily. This method is also available for non-Unicode characters.

Currently, CHISE project provides basic elements to process Chaon based text: text editor (XEmacs CHISE), scripting languages (Ruby/CHISE and Perl/CHISE) and type setting system ($\Omega$/CHISE). You can try KNOP-PIX/CHISE which is a DVD bootable GNU/Linux system including XEmacs CHISE and $\Omega$/CHISE. Its image is available at `http://kanji.zinbun.kyoto-u.ac.jp/projects/chise/dist/KNOPPIX/`.

CHISE project is an open source project, so its results are distributed as free software. Information about CHISE project is available at:

- `http://cvs.m17n.org/chise/`

- `http://kanji.zinbun.kyoto-u.ac.jp/projects/chise/`

These WWW pages, various programs and data are managed by CVS (a kind of revision control system), so users can get the latest snapshot. There are mailing-lists about CHISE project: for English and Japanese. If you are interested in CHISE project, please join to the lists.

# References

[1] 漢字庫. `http://www.sinica.edu.tw/~cdp/zip/hanzi/hanzicd.zip`.

[2] International Organization for Standardization (ISO). *Information technology — Universal Multiple-Octet Coded Character Set (UCS) – Part 1: Architecture and Basic Multilingual Plane (BMP)*, March 2000. ISO/IEC 10646-1:2000.

[3] International Organization for Standardization (ISO). *Information technology — Universal Multiple-Octet Coded Character Set (UCS) – Part 2: Supplementary Planes*, November 2001. ISO/IEC 10646-2:2001.

[4] 今昔文字鏡. http://www.mojikyo.com/.

[5] The Omega typesetting and document processing system. http://omega.cse.unsw.edu.au:8080/.

[6] Open Type font 用 VF. http://psitau.at.infoseek.co.jp/otf.html.

[7] Ascii 日本語 TeX(pTeX). http://www.ascii.co.jp/pb/ptex/.

[8] The object-oriented scripting language Ruby. http://www.ruby-lang.org/.

[9] XEmacs. http://www.xemacs.org/.

[10] 上地 宏一. 漢字フォント自動生成サーバ "影 KAGE" の構築 — 文字コードの枠組みを越える次世代漢字処理の提案 —. **漢字文献情報処理研究**, 3:143–147, 2002.

[11] 守岡 知彦 and クリスティアン・ウィッテルン. 文字データベースに基づく文字オブジェクト技術の構築. In **情報処理振興事業協会 平成 *13* 年度 成果報告集**. 情報処理振興事業協会, 2002. http://www.ipa.go.jp/NBP/13nendo/reports/explorat/charadb/charadb.pdf.

# 基于字符集的中文信息处理

翟喜奎

（国家图书馆 数字图书馆管理处）

**文摘**：根据在古籍文献处理时所遇到的中文信息处理问题提出了基于字符集的中文信息处理。基于字符集的中文信息处理应包括汉字排序，汉字可有多种排序标准(拼音、部首、笔画、四角号码等等)；规范检索，用简体、繁体、异体字都能统一检索；汉字输入方法；支持大字符集的显示等等。探讨基于字符集的中文信息处理方法。

**关键词：** 中文信息处理　汉字字符集 汉字排序 汉字属性标准

## 一、中文字符集的发展

### 1．GB 2312-80《信息交换用汉字编码字符集－基本集》

中国于1980 年3 月颁布了第一个汉字编码字符集标准，即GB 2312-80《信息交换用汉字编码字符集－基本集》。该标准符合ISO 2022编码体系结构。1981 年5 月1 日开始正式实施，它奠定了中国中文信息处理技术的发展基础。

随着GB 2312-80的颁布，中国颁布了相应的《15×16、24×24、32×32汉字点阵字模集及字模数据集》。所谓点阵字形，就是以点的形式来表现的字符或汉字的形态。15×16点阵字形，可以表示《信息交换用汉字编码字符集－基本集》中的绝大部分汉字。由于15×16的点阵字形只能表示横向笔画和竖向笔画都不超过八笔的汉字。如果一个汉字的横向笔画或者竖向笔画超过了八笔（如"量、酬"等字），在点阵字模就容纳不下。在《基本集》中，这样的汉字共有138个，只好压缩笔画做变通处理。15×16点阵字形适于屏幕显示，作校对之用。

24×24点阵字形，可以全部表示《基本集》中的6763个汉字的笔形结构，用不着压缩笔画，而且能够保持横细竖粗的宋体风格，适用于针式打印、喷墨打印，是一种很有使用价值的点阵字形。

32×32点阵字形比24×24点阵字形更能体现宋体风格，能完整地表现汉字的笔锋，使撇笔和捺笔自然婉转，舒畅流利，字体质量较高。

1992年中国颁布了矢量汉字的字模集及数据集：

◆　GB/T 13844-1992图形信息交换用矢量汉字 单线宋体字模集及数据集
◆　GB/T 13845-1992图形信息交换用矢量汉字 宋体字模集及数据集
◆　GB/T 13846-1992图形信息交换用矢量汉字 仿宋体字模集及数据集
◆　GB/T 13847-1992图形信息交换用矢量汉字 楷体字模集及数据集
◆　GB/T 13848-1992图形信息交换用矢量汉字 黑体字模集及数据集

1992年、1993年中国颁布了高精度点阵汉字标准：

◆　GB/T 14242-1993 信息交换用汉字64×64点阵黑体字模集及数据集
◆　GB/T 14243-1993 信息交换用汉字64×64点阵楷体字模集及数据集
◆　GB/T 14244-1993 信息交换用汉字64×64点阵仿宋体字模集及数据集

- GB/T 14245-1993 信息交换用汉字64×64点阵宋体字模集及数据集
- GB/T 14717-1993 信息交换用汉字128×128点阵宋体字模集及数据集
- GB/T 14718-1993 信息交换用汉字128×128点阵黑体字模集及数据集
- GB/T 13443-1992 信息交换用汉字128×128点阵楷体字模集及数据集
- GB/T 13444-1992 信息交换用汉字128×128点阵仿宋体字模集及数据
- GB/T 14719-1993 信息交换用汉字256×256点阵宋体字模集及数据集
- GB/T 14720-1993 信息交换用汉字256×256点阵黑体字模集及数据集
- GB/T 13445-1992 信息交换用汉字256×256点阵楷体字模集及数据集
- GB/T 13446-1992 信息交换用汉字256×256点阵仿宋体字模集及数据集

## 2．编码字符集的繁体字和简体字对应编码

1984 年"全国计算机与信息处理标准化技术委员会"提出编码字符集的繁体字和简体字对应编码的原则，并做出了制定六个信息交换用汉字编码字符集的计划。这六个集分别命名为基本集（GB2312-80）、第一辅助集(辅一)、第二辅助集(辅二)、第三辅助集(辅三)、第四辅助集(辅四)、第五辅助集(辅五)。其中，基本集、辅二集、辅四集是简体字集，辅一集、辅三集、辅五集分别是基本集、辅二集、辅四集的繁体字映射集，且简/繁字在两个字符集中同码(个别简/繁关系为一对多的汉字除外)。这六个集均采用双七位编码方式，但为了避开ASCII表中的控制码，每个七位只选取了94 个编码位置。所以每张代码表分94 个区和94 个位。其中前15 区作为拼音文字及符号区或保留未用，16 区到94 区为汉字区。第一辅助集(辅一)、第二辅助集(辅二)、第三辅助集(辅三)、第四辅助集(辅四)、第五辅助集(辅五)分别于1987年至1991年颁布。

- GB 12345-90《信息交换用汉字编码字符集——第一辅助集》
- GB 7589-87《信息交换用汉字编码字符集——第二辅助集》
- GB13131-1991《信息交换用汉字编码字符集——第三辅助集》
- GB 7590-87《信息交换用汉字编码字符集——第四辅助集》
- GB13132-1991《信息交换用汉字编码字符集——第五辅助集》

## 3．GB13000《信息技术通用多八位编码字符集》

1990年中国颁布了GB13000《信息技术通用多八位编码字符集》。

随着国际间的交流与合作的扩大，信息处理应用对字符集提出了多文种、大字量、多用途的要求。1993年国际标准化组织发布了ISO/IEC 10646-1《信息技术通用多八位编码字符集第一部分体系结构与基本多文种平面》。我国等同采用此标准制定了GB 13000.1-1993。该标准采用了全新的多文种编码体系，收录了中、日、韩20902个汉字，是编码体系未来发展方向。

## 4．GBK 编码字符集

1995 年12 月完成GBK 规范。GBK 编码是GB2312-80 国标码的扩充（其中GB 表示国标，K 表示扩展）。但是GBK编码本身不是国家标准。该编码规范完全兼容GB2312-80。

### 5．GB18030-2000《信息交换用汉字编码字符集基本集的扩充》

2000年3月中国颁布了国家标准GB18030-2000《信息交换用汉字编码字符集基本集的扩充》是我国继GB2312-1980和GB13000-1993之后最重要的汉字编码标准，是未来我国计算机系统必须遵循的基础性标准之一，该标准是国家强制性标准。在中国大部分计算机系统仍然采用GB 2312编码。GB 18030与GB 2312一脉相承，较好地解决了旧系统向新系统的转换问题，并且改造成本较小。从我国信息技术和信息产业发展的角度出发，考虑到解决我国用户的需要及解决现有系统的兼容性和对多种操作系统的支持，采用GB 18030是我国目前较好的选择，而GB 13000.1更适用于未来国际间的信息交换。考虑到GB 18030和GB 13000的兼容问题，标准起草组编制了GB 18030与GB 13000.1的代码映射表，使得两个编码体系可以自由转换。同时，还开发了GB 18030基本点阵字型库。

### 6．台湾字符集

《中文资讯交换码 CCCII》。CCCII 是 Chinese Character Code for Information Interchange 的缩写，是经台湾中研院中美会及国科会等单位支持，于 1979 年 12 月集合由台湾图书馆学者、文字学家及电脑专家组成"国字整理小组"提出的汉字编码。经过一些修改，被美国采纳为美国国家标准 ANSI Z39.64-1989，同时也被称为东亚字符编码（EACC）用于图书馆书籍目录方面。

《通用汉字标准交换码 CNS 11643》。1983 年 10 月，台湾科学委员会、教育部国语推行委员会、中央标准局、行政院主计处电子资料处理中心共同制定了《通用汉字标准交换码》（Chinese Ideographic Standard Code for Information Interchange，简称 CISCII 码），经试用修订，1986 年 8 月 4 日由台湾中央标准局公布为法定标准，标准编号为 CNS 11643。这一标准于 1992 年 5 月 21 日重新修订公布，更名为《中文标准交换码》（Chinese Standard Interchange Code）。1995 年 1 月 4 日，台湾中央标准局又公布了 CNS 11643-1《中文标准交换码使用方法》。

上述两个汉字字符集，CNS 11643 为通用的标准交换码，适用面较广。CCCII 使用面相对较窄，相当于行业规范。

BIG-5 码是 1984 年台湾资讯工业策进会根据《通用汉字标准交换码》制订的编码方案。

### 7．日本字符集标准

1978 年，日本政府公布了日本工业标准 JIS C 6226-1978《信息交换用汉字字符集》。该标准于 1983 年进行了修订，新增加了 4 个汉字，并将该标准编号改为 JIS X 0208-1983。

1990 年日本发布了第二个日本汉字编码字符集标准 JIS X 0212-1990，作为日本汉字交换码辅助集。

1993 年日本发布了第三个日本汉字编码字符集标准 JIS X 0221-1993，该标准是有 20,902 个汉字的编码标准。接着，有了 1996 年的《表外汉字字体表试案》，日本政府在公布该方案时，并且声明，它"是法令、公用文书、报纸、杂志、广播电视等一般社会生活中，使用表外汉字字体的依据"，"这个字体表将明治以来传统的印刷文字字体（并非《康熙字典》字体的本身，而是以《康熙字典》为依据作成的明治以来的铅字字体即《康熙字典体》）置于印刷标准字体的地位。"

目前，日本最新汉字编码字符集标准是 JIS X 0221-1:2001。

### 8．韩国字符集标准

1987 年韩国制定了韩国标准编码字符集 KS C 5601-1987，共有 8224 个字符。1991 年韩国制定了编码字符集的辅助集 KS C 5657-1991，增收汉字 2856 个。

### 9．ISO/IEC 10646 与 UNICODE

国际标准化组织（ISO）于 1984 年开始研究、制定《信息技术 通用多八位编码字符集（UCS）》国际标准，即 ISO/IEC 10646。1993 年 5 月，该标准的第一部分：体系结构与基本多文种平面（即 ISO/IEC 10646.1）正式发布。前后经历了九年的时间。ISO/IEC 10646-1 的第二版更加全面，即 ISO/IEC 10646-1：2000。与 ISO/IEC 10646-1：2000 等同的工业标准是 Unicode3.0，Unicode 是 Universal Code 的简称，即统一编码。除了作为 ISO/IEC10646 编码的一种称谓外，Unicode 同时还是由 HP，IBM，APPLE，MICROSOFT 等一些国际知名企业组成的一个联盟的名称。该联盟的主要宗旨就是要通过市场手段推进多文种的统一编码，因此称为 Unicode。它的广泛使用将会使得软件开发费用大幅度降低，开发更为快捷，可共享资源更为丰富，使用者的投入也将大幅度降低，便于推广。软件可以实现一个版本的世界范围内通用，不再需要多个版本、多种语言的产品了。目前兼容和支持该标准的已有许多大型厂商，如微软、苹果、SUN、甲骨文等国际性软件公司。

目前，ISO/IEC 10646 国际标准的最新版本是 2003 年修订的 ISO/IEC 10646:2003 等同的工业标准是 Unicode 4.0。

## 二、中文字符集的基本内容

### 1．GB 2312-80《信息交换用汉字字符集—基本集》

GB2312-80《信息交换用汉字字符集—基本集》 收录汉字信息交换用的基本图形字符，采用一字一码的原则，实现简化汉字6763 个，总计7445 个图形字符。具体包括：一般字符202 个，序号60 个，数字22 个，英文字母大小写共52 个，日文平假名169 个，希腊字母大小写共48 个，俄文字母大小写共66 个，汉语注音符号37 个，一二级汉字共6763 个。其中：一级常用汉字3 755 个，按照汉语音排序；二级非常用汉字3008 个，按照偏旁部首排序。该标准的制定和应用为规范、推动中文信息化进程起了很大作用。
- 双字节编码
- 范围：A1A1～FEFE
- A1-A9：符号区，包含682 个符号
- B0-F7：汉字区，包含6763 个汉字。

### 2．《汉字国标扩展规范GBK》

《汉字国标扩展规范GBK 》，在MS Windows 9x/Me/NT/2000、IBM OS/2 的系统中广泛应用。是GB2312 国标码的扩充。它是国家技术监督局1995 年为中文Windows 95所制定的新的汉字内码规范（其中GB 表示国标，K 表示扩展）。该规范在字汇一级上支持ISO10646 和GB13000 中的全部中日韩（CJK）汉字，并与国家标准GB2312-80 信息处理交换码相兼容。
- 双字节编码，GB2312-80 的扩充，在码位上和GB2312-80 兼容
- 范围：8140～FEFE（剔除xx7F）共23940 个码位
- 包含21003 个汉字，包含了ISO/IEC 10646-1 中的全部中日韩汉字。

### 3．GB 12345-90《信息交换用汉字编码字符集—第一辅助集》

国家标准GB1234-90《信息交换用汉字编码字符集—第一辅助集》于1990 年发布，是与基本集对应的繁体字字符集，共收图形字符7583 个，其中前15 区除收集了GB 2312-80 中前15 区内收的全部字符外，又增收了35 个竖排标点符号和汉语拼音符号共717个字符。从

16 区至91 区共收6866 个繁体汉字。目的在于规范必须使用繁体字的各种场合，以及古籍整理等。一级汉字数和二级汉字数都与GB2312-80 相同，另有103 个繁体字是属于简/繁为一对多的字（比GB2312 多103 个字，其它厂商的字库大多不包括这些字）。对于简/繁一对多的情况，则选一个最通用的繁体字码置于与基本集中该字相对应的码位，其余的则按拼音序编码于88和89 区。实际上可以将GB1234-90理解为GB2312-80的繁体字对应版。如果在同一套计算机系统中要支持简、繁体字共存，只能采取两个分别代表GB2312-80 和GB12345-90 的代码页。

### 4．GB 7589-87《信息交换用汉字编码字符集—第二辅助集》

GB 7589-87《信息交换用汉字编码字符集—第二辅助集》是作为基本集的补充而编制的，收入通用规范的简体汉字，收字7237 个，以201个部首为序排列，部首次序按笔画数排列，同部首字按部首以外的笔画数排列，同笔画数的字以笔形顺序(横、直、撇、点、折)为序。不收繁体字和被淘汰异体字。

### 5．GB 13131-1991《信息交换用汉字编码字符集—第三辅助集》

辅三集是辅二集对应的繁体字字符集，简繁体有一一对应关系，收字与辅二集相同。

### 6．GB 7590-87《信息交换用汉字编码字符集—第四辅助集》

辅四集是作为基本集的补充而编制的，均收通用规范的简体汉字，收字7039 个，都以201个部首为序排列，部首次序按笔画数排列，同部首字按部首以外的笔画数排列，同笔画数的字以笔形顺序(横、直、撇、点、折)为序。不收繁体字和被淘汰异体字，第二辅助集和第四辅助集共约有4200 多个字是经过类推简化得到的，提高了整个字符集的规范性，但降低了字符集的实用性。

### 7．GB 13132-1991《信息交换用汉字编码字符集—第五辅助集》

辅五集是辅四集对应的繁体字字符集，简繁体有一一对应关系，收字与辅四集相同。

### 8．GB 13000《信息技术通用多八位编码字符集》（ISO/IEC10646）

GB13000是ISO/IEC10646的等同标准。国际标准化组织为了将世界各民族的文字进行统一编码，制定了UCS 标准。根据这一标准，中、日、韩三国共同制定了《CJK 统一汉字编码字符集》，其国际标准号为：ISO/IEC10646，我国国家标准号为：GB13000。该汉字编码字符集就是通常人们所说的大字符集，它编入了20902 个汉字，收集了大陆一、二级字库中的简体字，台湾《通用汉字标准交换码》中的繁体字，58 个香港特别用字和92 个延边地区朝鲜族"吏读"字，甚至涵盖了日文与韩文中的通用汉字，满足了方方面面的需要。

### 9．GB 18030-2000《信息技术信息交换用汉字编码字符集基本集的扩充》

GB18030标准是中国政府于2000年3月颁布的最新中文汉字编码标准，是我国继GB2312-1980 和GB13000-1993 之后最重要的汉字编码标准，是未来我国计算机系统必须遵循的基础性标准之一。GB18030 收录了27484 个汉字。双字节部分收录内容主要包括GB13000.1 全部CJK 汉字20902 个、有关标点符号、表意文字描述符13 个、增补的汉字和部首/构件80 个、双字节编码的欧元符号等。四字节部分收录了上述双字节字符之外的，包括CJK 统一汉字扩充A 在内的GB 13000.1 中的全部字符。GB18030 编码空间约为160 万码位，目前已编码的字符约2.6 万。随着我国汉字整理和编码研究工作的不断深入，以及国际标准ISO/IEC 10646 的不断发展，GB18030 所收录的字符将在新版本中增加。

### 10. 统一码（Unicode）

统一码(Unicode)与ISO 10646 国际编码标准互相兼容。统一码是由一个名为Unicode学术学会的机构制订的字符编码系统。该系统是为了将世界上几十种紊乱的字符编码整合在一起，以期减少各电脑商开发国外市场时遇到的问题，美国各大电脑厂商组成了策进会，以推广一个世界通行的编码体制，以支持世界主要语文的书面文本的交换、处理及显示。Unicode学术学会的成员大部分为计算机软硬件的供货商。

在1991 年，国际标准化组织与Unicode 学术学会决定共同制订一套适用于多种语文文本的通用编码标准。自此以后，该两个组织便一直紧密合作，同步发展ISO 10646 国际编码标准及统一码。国际标准化组织提供ISO 10646 国际编码标准内的字符及编码资料，Unicode 学术学会则对这些字符及编码资料提出应用的方法以及语义资料作补充。ISO 10646 国际编码标准与统一码所包含的字符及使用的编码是相同的。统一码可被视为是ISO 10646 国际编码标准的实践版。因此，支持统一码的产品，亦支持ISO 10646 国际编码标准。由Unicode 学术学会制订的统一码3.0 版本，于2000 年2月正式推出。这个版本收纳了49,194 个来自世界各地不同语文的字符，其中包括27,484 个东亚的表意文字(汉字，汉字是经过CJK 整合的，即将中日韩文中相近的汉字用单一的编码，称为统一汉字Unihan，共2 万多个，但并不包含一些罕见的字，如康熙字典中的一些古字)。Unicode 编码有多种实现，常见的有UTF8, UTF16, UCS-2, UCS-4 等，统一码3.0 版本是与ISO/IEC 10646-1:2000 对应的版本。统一码3.1 版本于2001 年3 月推出。这个版本的主要特点是增加了44，946 个新字符，其中42，711 个为表意文字。连同统一码3.0 版本原有的字符，统一码3.1 版本共收录了94，140 个字符，其中表意文字总数超过70，000 个。统一码于2002 年推出的3.2 版本。虽然这个版本包括了1，016 个新字符，但其包含的表意文字则与统一码3.1 版本相同。统一码2003年推出最新4.0版本，统一码4.0版本与ISO 10646:2003 国际编码标准的现行版本完全对应，目前，在网络、Windows 系统和很多大型软件中得到应用。

## 三、基于字符集的中文信息处理

### 1. 问题的提出

数字图书馆，是面向未来互联网发展的信息管理模式。以数字资源的制作、存储、管理、传输和服务为主要特征的数字图书馆技术，是 21 世纪国际科技文化竞争的焦点之一。数字图书馆涵盖多个分布式、超大规模、可互操作的异构多媒体资源库群，面向社会公众提供全方位的知识服务。可以说，数字图书馆将实现对人类知识的普遍存取，并最终消除人们在信息获取方面的不平等。它既是知识网络，又是知识中心，同时也是一套完整的知识定位系统。

国家数字图书馆工程将在国家图书馆二期工程内建设国家数字图书馆国家中心，并通过应用系统开发实现数字资源采集、加工、处理、存储、归档、组织、发布和利用全过程。

地方志是我国所特有的一种文献形式，其中有 6，300 余种、120，000 余册。建国前的旧方志是国家图书馆独具特色的馆藏之一，所存文献数量与品质极高。旧方志数字化是一项重要的人文学术研究基础工程，利用计算机及网络技术进行深入的整理、开发，在当今数字化时代势在必行，它将大幅度地提高大众学习、认识中国古代地方文化的效率，即可以将学者的时间和精力从艰苦而繁琐的爬梳、翻检工作中解放出来，又可以向普通读者打开发掘、发现地方志宝藏的大门。地方志文献的数字化是全部中文文献数字化事业的一个复杂特例，是数字化图书馆事业的一个重要部分。

2004 年，国家图书馆将 50 万筒子页地方志文献进行文本化数字处理和版式还原。目的

是要实现地方志文献全文检索。共处理汉字约 2.3 亿，遇到 UNICODE4.0 之外的字符集集外字大约 4500 字。因此，在处理古籍文献时，所使用的汉字数量是很大的。除此之外，规范检索问题即用简体、繁体、异体字都能统一检索；检索到的结果进行汉字排序问题；大字符集汉字输入问题；大字符集汉字显示问题等等都要进行处理。这些都是中文信息处理中所遇到的实际问题。

基于字符集的中文信息处理应包括汉字排序，汉字可有多种排序标准(拼音、部首、笔画、四角号码等等)；规范检索，用简体、繁体、异体字都能统一检索；汉字输入方法；支持大字符集的显示等等。

## 2．进行汉字属性标准研究

为了解决基于字符集的中文信息处理即汉字排序、规范检索、汉字输入、汉字显示等问题，就要对汉字属性标准进行研究。目前，基于GB13000.1《信息技术通用多八位编码字符集》即(ISO/IEC10646.1-1993)、UNICODE1.0的汉字属性标准研究已经完成，解决的汉字数量是基本集20902个汉字。但是，基于UNICODE4.0（ISO/IEC10646：2003）的汉字属性标准研究，当前还是空白，要解决的汉字数量是扩充A集6582个汉字、扩充B集42711个汉字。要加速该方面的研究，满足数字图书馆资源建设以及实际应用工作的需求。

汉字属性标准研究的基本内容是汉字字型标准化、汉字标准发音、字型特征（包括汉字总笔画数量、汉字起笔至末笔笔形值、部首笔画数量、部首序号、部首外起笔至末笔笔形值、异体字数量、异体字字型等）、各种编码（包括四角号码、输入编码、其他汉字字符集编码等）以及构词和使用频度等。

## 3．汉字排序

汉字的排序方法是依据我国多年延续下来的传统规范和若干限定条件对汉字进行排序处理。即要使汉字象拉丁字符一样成为有序的集合，在计算机内能够进行比较、计算。我国目前使用的汉字排序方法主要有四种：部首法、汉语拼音法、笔画法和四角号码法。这些方法的规则如下：

部首法是以部首归并汉字的一种排检方法。它是先把汉字按其所属的部首归并集中。部首按笔画数多少排列先后顺序， 笔画数目相同的部首，依起笔笔形(横、竖、撇、点、折)排列先后顺序。同属一个部首的字，其先后顺序仍然是先按部首之外的笔画数排列， 部首之外的笔画数目相同的，再依起笔笔形顺序排列。

汉语拼音法是按照汉字发音和声调来归并排列汉字的一种方法。它的一般形式是：先按汉字的发音和声调来归并汉字，按字母的序列排序。音、调相同依笔画数多少排列。笔画数相同，再依起笔笔形(横、竖、撇、点、折)排列先后顺序。

笔画法是按照笔画数目及起笔笔形来归并排列汉字的一种方法。它的一般形式是：先按笔画数多少来归并汉字，笔画数相同，再依起笔笔形(横、竖、撇、点、折)排列先后顺序。

四角号码法是一种以数码来代表汉字四角的笔形并据此来排列汉字先后次序的方法。先按四角号码数多少来归并汉字。四角号码相同，依字中"横"笔的多少排列。"横"笔相同，依整体字的笔数排列。整体字的笔数相同，再依起笔笔形(横、竖、撇、点、折)排列先后顺序。

用计算机处理汉字排序问题的规则，见下表：

|  | 因素1 | 因素2 | 因素3 | 因素4 | 因素5 |
|---|---|---|---|---|---|
| 部首法 | 部首序号 | 部首外汉字笔数 | 部首外汉字起笔至末笔笔形值 | 内码 |  |
| 汉语拼音法 | 汉语拼音 | 声调 | 总笔画数 | 汉字起笔至末笔笔形值 | 内码 |
| 笔画法 | 总笔画数 | 汉字起笔至末笔笔形值 | 内码 |  |  |
| 四角号码法 | 四角号码 | 横笔数 | 总笔画数 | 汉字起笔至末笔笔形值 | 内码 |

## 4．规范检索

要解决标准正形汉字与繁体字与异体字相互连接问题。建立相互参见对照表，解决规范检索问题。

## 5．输入方法

要选择一种或几种适合古籍大字符集的输入方法，解决汉字输入问题。

## 6．汉字显示

古籍资源的应用是全球化的问题，需要解决古籍大字符集的显示问题。虽然，系统都支持UNICODE，但是没有扩A、扩B大字符集字库的支持也不能正确显示汉字。因此，需要解决支持大字符集的汉字字库问题。

总之，基于字符集的中文信息处理是当今数字图书馆古籍资源建设和应用的基础。呼吁各界重视基于字符集的中文信息处理的基础研究，再创我国中文信息处理领域的辉煌。

**参考文献：**

《中文信息处理技术的现状与进展》 通用中文代码国际联合会 1991年3月 Version 2.0

《试论中文文献的有序化》 翟喜奎《现代图书情报技术》1990年 第1期

# A Model for Scholarly Collaboration in the Development of On-line Reference Works: The Digital Dictionary of Buddhism

Charles Muller
Toyo Gakuen University
Kyoto University Institute of Humanities Presentation
Beijing January 22, 2004

## I. Technical Review

I began the compilation of the Digital Dictionary of Buddhism (DDB) and the CJKV-English Dictionary soon after my entry into graduate school in Buddhist Studies, upon my coming to awareness of the dearth of adequate lexicographical and other reference works in English language for the textual scholar of East Asian Buddhism in particular, and East Asian philosophy and religion in general. I decided, during my first Buddhist and Confucian/Daoist texts readings courses to save everything I looked up, and have continued that practice down to the present, through the course of studying scores of classical texts.

At the time that I began this process, I could not have dreamt of such a thing as the Internet, or even thought of the possibility of having this material available as a digital database—I was simply envisioning the eventual publication of a new, comprehensive printed work. But as developments in the IT world progressed, the newly appearing potentialities gradually began to dawn on me. Then, in 1995, I tasted the Internet, and once I figured out how to insert *<html>* tags at the beginning and end of a text file, I was on my way to preparing these dictionaries for web publication—the first version of which I uploaded in the summer of 1995. Soon after this, the dictionary was discovered on by

Christian Wittern, who promptly downloaded all the files, and applied a basic SGML structure, which is the ancestor of the XML markup system used today.

Due to the lack of widespread popular implementation of SGML, I did not make any special effort to develop this format for a few years. But after 2000, the popularity of XML began to suddenly increase, and so I began to take this format seriously. A major technical turning point in the history of the DDB came in January 2001, when I was contacted by Michael Beddow, a scholar of German Literature who was also an extremely accomplished XML programmer, and who had been using XML for some time already to develop his own only lexicographical project, the Anglo-Norman Dictionary (*http://www.anglo-norman.net*). Michael generated, based on the markup structure of the DDB, an array of indexes that used Xpointers to call up single-entry data units out of large files, each of which contained hundreds of entries. Michael also developed a CJK-Utf-8 search engine.[1]

## II. Content Development

In my first presentation of the DDB at the meeting of the Electronic Buddhist Text Initiative (EBTI) in 1996, the dictionary had 3,200 entries. Today, less than nine years later, the DDB now boasts 35,000 entries, making it by far the largest compilation of its type in the English language, and even larger than some of the best-known Japanese works, such as Oda's *Bukkyō daijiten*. An instrumental factor in this rate of growth is the aid received through JSPS research grants, which allowed us to hire graduate students to help digitize large amounts of data for input. But this stage ended in 2002, and we have entered a new phase, where we are finally receiving large contributions of data from

---

1 I have focused here on developments in the DDB, but please note that all of the same technological enhancements have been applied to the CJKV-E.

unselfish collaborators who understand the spirit of the project and its limitless potential for the future. Two of the largest recent individual contributions have come from Prof. KARASHIMA Seishi, who has contributed over 7,000 entries from his research on the *Lotus Sutra*, and from Dr. Stephen Hodge, who has contributed over 2,500 hundred terms from his translation work on the *Yogācārabhūmi-śāstra.* In addition to these unusually large contributions, we have recently been benefiting from a continuous stream of smaller contributions, amendments, and corrections, from an ever-increasing  number of scholars.[2]

While the DDB can certainly be viewed as a fairly successful model of the possibilities of online collaboration, it should be made clear that until we set up a mechanism to strongly encourage (perhaps "force" is the better term here) contribution, voluntary data submissions were few and far between. Initially we set up our password access/quota system to deal with hacking and data-theft problems. But we also discovered that we could take advantage of this same system to encourage contributions. Through this system, users who log onto the DDB web site to search for terms are able to freely look up ten items in a 24-hour period. After this, they are greeted by a message telling them that their quota is finished, but that they may gain an unlimited quota password by making a small data contribution.

In earlier days, when the content coverage of the DDB was still rather limited, this strategy did not generate that much response. But during the past year, with the expansion of the coverage to its present number, usage of the resource has also increased. The DDB has become a standard lookup tool for many Buddhist studies specialists—especially those who are doing intensive translation work. It is also used extensively in university classes in North America, and is a basic research tool listed on the syllabi of Buddhist

2   For a full list of contributors, see *http://www.acmuller.net/credits/credits-ddb.htm*.

Studies courses in such prestigious institutions as Harvard, Stanford, Princeton, Columbia, Berkeley, and other universities. As the DDB grows in both usefulness and in reputation as an essential reference tool for Buddhist studies research, scholars and students are increasingly coming to depend upon it, and thus eventually come to need unlimited access. Most serious scholars already have a large amount of specialized information on their hard drives that can easily be modified to become a DDB entry. All they need, it seems, is a small reason to make this effort, along with a little prodding.

For interested persons who do not have the specialized training to write or edit DDB entries, paid subscriptions are available. This approach was settled upon not with the expectation of making a lot of money, but simply to provide a recourse for persons who demanded full access in one way or another. As a by-product of this offering however, we decided to offer institutional subscriptions as well, and recently a number of major universities have decided to have their libraries subscribe to the DDB, including Columbia, Berkeley, Santa Barbara, and UCLA. While we are happy to gain a small amount of money to put back into the project, at this point, the greatest value of these subscriptions is in the recognition being accorded to the DDB as a primary reference tool. It is especially significant that this reference tool has been put together and produced, not by a major publishing company, but by a group of like-minded scholars.

## III. The Structure of a DDB Entry

As mentioned earlier, the DDB uses XML as its basic structural format. The DTD is based loosely on the recommendations of the Text Encoding Initiative, using many of the entities and attributes that are used for lexicons.[3] An entry is divided into three major

---

3  The reason that the DDB is not based more fully on TEI is simply that most of the structure was developed before I adequately understood the TEI model. I have thought from time to time about redoing the whole structure according to the TEI DTD, but the retooling of the stylesheets, as well as

sections: (1) A Pronunciation Section, wherein the readings of a Chinese term are provided in various East Asian languages and their romanization systems. (2) A Sense Section, which provides the translation of the term and other explanatory material, and (3) An External References Section, which provides references to the term in a variety of Buddhist Studies reference works. Each of these larger nodes has children nodes, and various other entities contained within. When a user selects a term either via hyperlink or by search engine lookup, and HTML page is generated. One sample page is given below:

---

numerous other aspects of the web implementation of the data set are simply too daunting for me to seriously consider at this point in time.

---

# 言說

[[Pronunciations](#)]

[py] yánshuō
[wg] yen-shuo
[hg] 언설
[mc] eonseol
[mr] ŏnsŏl
[kk] ゴンゼツ
[hb] gonzetsu

## Meanings

**[Basic Meaning:] verbal expression** [s.hodge]

Senses:

- (Skt. *vyavahāra*; Tib. *tha snyad*) [s.hodge]
- expresses, recounts; (Skt. *abhi-lap\**; Tib. *brjod par 'gyur ba*) [s.hodge]
- expressing, an expression, an utterance; (Skt. *abhilāpa*; Tib. *brjod pa*) [s.hodge]
- discourse; (Skt. *kathā*; Tib. *gtam*) [s.hodge]
- a figurative designation; (Skt. *upacāra*; Tib. *nye bar 'dogs pa*) [s.hodge]
- Language, speech (*vāc*), which is one of the three kinds of permeation of the store consciousness taught in the *Mahāyāna-saṃgraha*. 〔[攝大乘論](#) (T 1593.31.117c3)〕 [cmuller]
- The usage of language to teach the dharma (*deśanā*). [cmuller]
- Language as synergistic with the mental realm of phenomenal differentiation (*abhilāpa*). [cmuller]

## [Dictionary References]

Iwanami Bukkyō jiten 239, 293
Bukkyōgo daijiten (Nakamura)429b
Ding Fubao
Buddhist Chinese-Sanskrit Dictionary (Hirakawa)1072
Bukkyō daijiten (Mochizuki)(v.9-10)1043b
Bukkyō daijiten (Oda)582-1
Sanskrit-Tibetan Index for the Yogācārabhūmi-śāstra (Yokoyama and Hirosawa)

One distinctive feature that you will notice in this example, that one does not yet see in standard reference works, is that attribution is not simply given for the entire entry as a unit: responsibility is acknowledged for each segment (XML node) of the entry--and as often as possible, with the equivalent Sanskrit or Tibetan. Both characteristics are especially helpful for those who are doing research and translation. Of course, using XML like this, we can display more detailed information if we want to. But whether or not we decide to display it when we publish the HTML files, the users have the option of viewing the XML source data if they like. For the above-shown entry, the XML source data looks like this:

<entry ID="b8a00-8aaa" added_by="cmuller" add_date="1997-09-15" update="2003-10-11" rad="言" radval="07" radno="149" strokes="00">

<hdwd>言說</hdwd>
<pron_list>
<pron lang="zh" system="py" resp="cmuller">yánshuō</pron>
<pron lang="zh" system="wg" resp="cmuller">yen-shuo</pron>
<pron lang="ko" system="hg" resp="cmuller">언설</pron>
<pron lang="ko" system="mc" resp="cmuller">eonseol</pron>
<pron lang="ko" system="mr" resp="cmuller">ŏnsŏl</pron>
<pron lang="ja" system="kk" resp="cmuller">ゴンゼツ</pron>
<pron lang="ja" system="hb" resp="cmuller">gonzetsu</pron>
</pron_list>
<sense_area>
<trans resp="s.hodge">verbal expression</trans>
<sense resp="s.hodge">(Skt. <term lang="sa">vyavahāra</term>; Tib. <term lang="bo">tha snyad</term>)</sense>
<sense resp="s.hodge">
<trans resp="s.hodge">expresses, recounts</trans>; (Skt. <term lang="sa">abhi-lap*</term>; Tib. <term lang="bo">brjod par 'gyur ba</term>)</sense>
<sense resp="s.hodge">
<trans resp="s.hodge">expressing, an expression, an utterance</trans>; (Skt. <term lang="sa">abhilāpa</term>; Tib. <term lang="bo">brjod pa</term>)</sense>
<sense resp="s.hodge">
<trans resp="s.hodge">discourse</trans>; (Skt. <term lang="sa">kathā</term>; Tib. <term lang="bo">gtam</term>)</sense>
<sense resp="s.hodge">
<trans resp="s.hodge">a figurative designation</trans>; (Skt. <term lang="sa">upacāra</term>; Tib. <term lang="bo">nye bar 'dogs pa</term>)</sense>
<sense resp="cmuller">Language, <trans resp="cmuller"><term lang="en">speech</term></trans> (<term lang="sa">vāc</term>), which is one of the three kinds of permeation of the store consciousness taught in the <title>Mahāyāna-saṃgraha</title>. <bibl type="canoncite"><cit><xref idref="b651d-5927-4e58-8ad6">攝大乘論</xref> <biblScope source="T" div="num.vol.pg-col-line">T 1593.31.117c3</biblScope></cit>
</bibl></sense>
<sense resp="cmuller">The usage of language to teach the dharma (<term lang="sa">deśanā</term>). </sense>
<sense resp="cmuller">Language as synergistic with the mental realm of phenomenal differentiation (<term lang="sa">abhilāpa</term>). </sense>
</sense_area>
<dictref>
<dict><title>Iwanami Bukkyō jiten </title><page>239, 293</page></dict>
<dict><title>Bukkyōgo daijiten (Nakamura)</title><page>429b</page></dict>
<dict><title>Ding Fubao</title><page/></dict>
<dict><title>Buddhist Chinese-Sanskrit Dictionary (Hirakawa)
</title><page>1072</page></dict>
<dict><title>Bukkyō daijiten (Mochizuki)</title><page>(v.9-10)1043b</page></dict>
<dict><title>Bukkyō daijiten (Oda)</title><page>582-1</page></dict>
<dict><title>Sanskrit-Tibetan Index for the Yogācārabhūmi-śāstra (Yokoyama and Hirosawa)
</title><page/></dict>
</dictref>
</entry>

## IV. Making Contributions

An important future enhancement for the DDB will the development of an input form for contributors, to allow them to readily add new entries, or modify presently existent ones. For the time being however, lacking a formal apparatus for the input of new materials, we have been adding material received in attached mail files, mostly in MS-Word format. As long as the contributors use a format with a uniform structure, and are able to submit their materials in Unicode, using Unicode-mapped diacritics and East Asian characters, there is not that much else that needs to be done, as we are able to do much of the main markup with macros and various scripts. We do, however, encourage users to submit their materials with XML markup to whatever extent they are able to handle it, ranging from a minimal type of markup, up to a fully marked-up document using our DTD. On our web site, we offer users the following options (from *http://www.acmuller.net/ddb/notes/ Basic_Formatting.html*):

-------------------------------------------------------------------------------------------

### Basic Formatting Suggestions for DDB Entries

## Topics:

A. Introduction
B. Basic DDB Entry Format
C. Basic DDB Entry Format (Simple XML Markup)
D. Basic DDB Entry Format (Fully Developed XML Markup)

Updated 2004.09.05

## A. Introduction

First and foremost, please understand well: **the usage of XML tagging is not necessary for contributing to the DDB. We will happily accept contributions in popular word processor file formats with no XML markup whatsoever.** If, however, you are interested in going a step or two beyond that, and would like to learn something about how we encode our materials, then please read on.

## B. Basic DDB Entry Format

Up to now, the basic organization of a DDB entry has been like this (with some abridgments for the sake of simplicity):

---

Headword: (Han characters)

Pronunciations:

Chinese (Pinyin):
Chinese (Wade-Giles):
Korean (Hangul):
Korean (Ministry of Education System):
Korean (McCuneReischauer):
Japanese (Katakana):
Japanese (Hepburn):
Translation: (Simple, short-phrase equivalent of the headword, if available)
Explanation: (Detailed explanation of the entry headword)

---

If you were adding a term, you would type the Chinese next to "Headword." You would then add the pronunciations for the languages you know. Someone else can supply the readings for the languages you can't handle. After the pronunciations, we usually make an attempt to offer one (or up to a few) common renderings of the term. If it were a person, place, temple, etc., we would just supply the commonly used name, such as "Zongmi," "Dongshan," "Jinglingsi," etc. If it were a concept, "middle way," etc. This is followed by a detailed explanation, which can have multiple nodes for multiple contributors, as necessary.

Let's look at example. This is an entry regarding the Korean monk Iryŏn. It is an entry for which I provided minimal information many years ago, and which badly needs to be expanded. But its present brevity makes it useful here:

---

**Headword: 一然**


**Pronunciations:**

**Chinese (Pinyin):** Yīrán

**Chinese (Wade-Giles):** I-jan

**Korean (Ministry of Education System):** Iryeon

**Korean (McCuneReischauer):** Iryŏn

**Japanese (Hepburn):** Ichinen

**Translation:** Iryeon

**Explanation:** (1206-1289) An important Goryeo monk. A prolific writer, who is most famous for his Samguk Yusa [Chinese title here], a collection of facts and anecdotes which is a basic text for the study of the history of Korean Buddhism.

---

**C. Basic DDB Entry Format: XML**

Now, for XML. Rather than starting off with an explanation of XML theory, I think it is simpler if I just re-present the above example using a simplified form of XML.

<entry>

<hdwd>一然</hdwd>

<pron_list>

<pron>Yiran</pron>

<pron>I-jan</pron>

<pron>Iryeon</pron>

<pron>Iryŏn</pron>

<pron>Ichinen</pron>

</pron_list>

<trans>Iryeon</trans>

<sense> (1206-1289) An important Goryeo monk. A prolific writer, who is most famous for his <title>Samguk Yusa</title> 三 國遺事, a collection of facts and anecdotes which is a basic text for the study of the history of Korean Buddhism.</sense>

</entry>

If you look at this for a minute, you will see that there is not much difference between the first example and the XML-tagged example. The basic difference is that here we are using opening and closing tags to delimit information. You will notice that inside the <sense> tags, the title of Iryeon's text, *Samguk Yusa*, is enclosed with the tags <title></title>, indicating that this is the name of written work. We also use similar tags for <term>technical terms</term>, <foreign>foreign words</foreign> and other elements. When this entry is published as HTML, these words will automatically be italicized. We can also use these tags to build indexes.

If can cooperate by using this simple level of XML structuring, it would be greatly appreciated. But once again, it is not absolutely necessary for the task.

**D. Basic DDB Entry Format (Fully Developed XML Markup)**

The above example shows the barest XML framework—what are called ELEMENT tags. The tags <entry>, <pron>, <title>, etc. are all known as "elements" in XML parlance. But elements can also be enhanced by a very useful secondary layer of information, which is known as ATTRIBUTE information. Please see the same entry, again presented in a manner much closer to the way it is actually contained in our data set:

<entry added_by="cmuller" add_date="1990-09-21" update="">

<hdwd>一然</hdwd>

<pron_list>

<pron lang="zh" system="py" resp="c.wittern">Yīrán</pron>

<pron lang="zh" system="wg" resp="cmuller">I-jan</pron>

<pron lang="ko" system="mc" resp="cmuller">Iryeon</pron>

<pron lang="ko" system="mr" resp="cmuller">Iryŏn</pron>

<pron lang="ja" system="kk" resp="cmuller">イチネン</pron>

<pron lang="ja" system="hb" resp="cmuller">Ichinen</pron>

</pron_list>

<sense_area>

<trans resp="cmuller"><person_entry loc="ko">Iryeon</person_entry> </trans>

<sense resp="cmuller"> (1206-1289) An important Goryeo monk. A prolific writer, who is most famous for his <title lang="ko">Samguk Yusa</title> 三國遺事, a collection of facts and anecdotes which is a basic text for the study of the history of Korean Buddhism.</sense>

</sense_area>

&lt;/entry&gt;

---

I believe that the point of most of the attributes should be obvious, but one of the most important that I would like to draw your attention to is that of "resp", which means "responsibility"—thus, "accreditation." Far distinguished from paper publishing counterparts, the usage of XML in a digital reference work allows us to give credit to the person responsible for every small part of the &lt;entry&gt;. Thus, if someone wanted to add another &lt;sense&gt; element (or "node") to this entry, it could easily be done, giving that person credit in the "resp" attribute.

Also commonly used in the DDB is the "lang" attribute, which tells us the language of the text or foreign word that will be italicized. For texts, we also have a "prov" (provenance) attribute. For temples and geographical entries, we have a "loc" (location) attribute. There are a number of others as well.

Using attributes allows for all kinds of programming possibilities, including various font transformations on presentation, creation of detailed indexes, and so forth.

However, once again, for those for whom this is a headache, it is fine if you want to terminate your exposure to XML here. Ensuing discussions will go into a bit more detail on XML for those who are interested, so you may ignore these if you wish.

## V. Near-Term Future Prospects for the DDB

During the past year, we have reached a distinctive new stage in the development of the DDB, wherein suddenly a large number of recognized scholarly experts in Buddhist Studies have begun to contribute data, and major university libraries have decided to subscribe. We presently have some 5,000 entries on the queue awaiting input, with new contacts from interested scholars coming weekly. Thus, we appear to be on the verge of being able to declare the DDB project a major success. We still eventually need to figure out a way to round out the balance of the coverage, so that there is more equal representation in terms of sects and cultural traditions, but this can probably be solved by the attainment of another significant grant or two. But if we can continue to grown at the rate of 5-10,000 terms a year for the next five years or so, it will probably be the right time to begin to turn our full attention to the proper completion of the sister project of the

DDB--the CJKV-E dictionary.