# Case study: Typesetting old documents of Japan

NAKANO Ken and KOBAYASHI Hajime

This is a report on the typesetting of *Komonjo books* (reprint books compiling old documents of Japan) published by Shiryo Hensan-jo, with pTeX (pLaTeX).

We would like to report that typesetting of *Komonjo books*, which are classified as difficult to typeset, can be made efficiently using TeX.

## About Shiryo Hensan-jo

From their web site (`http://www.hi.u-tokyo.ac.jp`): Shiryo Hensan-jo (the Historiographical Institute, HI), the University of Tokyo, has as its primary objective, rather than historiography in general, analysis, compilation, and publication of historical source materials concerning Japan.

The institute has become a major center of Japanese historical research, and makes historical sources available through its library, publications, and recently, databases.

## About us and TeX

Our company, Livretech, is a printing firm mainly producing textbooks and education-related books.

We've been using TeX since 1988, and have been concerned with some TeX-related topics. We developed a DVI driver program for the Shaken typesetter in 1989 (Shaken is a famous manufacturer in Japan making fonts and typesetting systems). We also typeset and output Japanese editions of *The TeXbook* (1992) and *The METAFONTbook* (1994).

## 1 Until TeX was chosen

Before TeX was chosen, *Komonjo books* were composed with hot metal types or phototypesetting systems developed for professional users. Several typesetting systems had been used, but commonly speaking, *Komonjo books* were difficult to typeset.

After typesetting was completed, contents of books were stored in the database named SHIPS (Shiryo hensanjo Historical Information Processing System), but another problem occurred.

In those days, Japanese typesetting systems worked with the creators' original character codes. For technical reasons and/or the creators' policies, typesetting data could not easily be re-used for another purpose. So, the final typesetting data would not automatically convert into database texts. To make the database text, typesetting contractor had to revise the final data manually, or in the worst case, input all texts again. And texts made in this way had to be proofread again.

To solve these problems, the professors of Shiryo Hensan-jo focused on TeX. This was in about 2001. They tried to typeset *Komonjo books* with TeX, and they found the possibility of typesetting and flexibility of text data.

And then, they contacted us and we started on this new challenge in 2002.

## 2 Typesetting problems

### 2.1 Problems with fonts

#### 2.1.1 Shortage of kanji characters

Massive numbers of kanji are used in *Komonjo books*. At that time, the Shaken system that was used to typeset *Komonjo books* had about 10,000 kanji characters. However, the fonts which could be used with pTeX had only about 6,500 kanji characters.

#### 2.1.2 Simplification of kanji shapes

In Japan, shapes of many frequently used kanji have been simplified to make them easier to read and write. Kanji of traditional shapes are called *Seiji* (正字), and those of simplified shapes are called *Ryakuji* (略字). Here are some examples of *Ryakuji* and *Seiji* character shapes:

| Ryakuji | Seiji | meaning |
|---------|-------|---------|
| 医 | 醫 | medical, doctor |
| 円 | 圓 | circle, yen |
| 塩 | 鹽 | salt |
| 害 | 害 | harm, damage |
| 学 | 學 | study, learn |
| 国 | 國圀 | country |

Like the letter "国," some *Ryakuji* characters have two or more *Seiji* characters.

In essence, *Komonjo books* are composed with *Seiji* characters. On the other hand, we use *Ryakuji* kanji in everyday life, and we rarely use *Seiji* characters. Because of this, *Seiji* characters have been put away in the dark recesses of fonts and kanji input methods. So it is difficult for us to input *Seiji* characters directly.

#### 2.1.3 *Hentai-gana*

変体仮名 (*Hentai-gana*, see ① in the figure on the next page) are Japanese old syllabary characters. *Hentai-gana* were made from kanji characters as well as modern *kana* characters. Here we show some *Hentai-gana* examples used in *Komonjo books*, with related kanji and kana characters.

| original kanji | 二 | 爾 | 而 | 八 | 者 | 者 | 者 | 江 | 三 | 茂 |
|----------------|----|----|----|----|----|----|----|----|----|----|
| *Hentai-gana* | ⼆ | ﾖ | ふ | ハ | ｾ | ﬞ | 志 | ね | ミ | 戈 |
| modern Kana | に | に | て | は | は | は | は | え | み | も |

Dainihon Ishin Shiryo, Ii-ke #27, p. 106

Dainihon Shiryo, No. 7 #32, p. 125

Dainihon Kinsei Shiryo, Hirohashi Diary #11, p. 194

Dainihon Shiryo, No. 7 #32, p. 227

Examples of *Komonjo books* (reduced 53%).

NAKANO Ken and KOBAYASHI Hajime

*Hentai-gana* have different shapes with the same reading, while modern kana have one shape with one reading.

### 2.1.4 *Bonji*

梵字 (*Bonji*, see ②) characters are mystical letters used in Japanese Buddhism. *Bonji* consist of consonants, vowels and derivative letters, so many characters can be created. The following are some typical *Bonji* characters and their meanings:

| *Bonji* | sound | meaning |
|---|---|---|
| 𑖀 | a, ア | the most basic character |
| 𑖁 | āḥ, アーク | Dainichi Buddha, 大日如来 |
| 𑖮 | hrīḥ, キリク | Amitabha, 阿弥陀如来 |

## 2.2 Problems with notes

*Komonjo books* contain many notes to aid readers' understanding. Two kinds of notes occur frequently.

One is a "head note" (③), similar to `\marginpar` but a long note appearing at the end of a page has to be split over pages.

The other is an "interline note" (④) which can occur anywhere in the text. To avoid overlapping of the interline note and text or other materials, position adjustment is needed.

## 2.3 Problems with picture-like materials

Picture-like materials such as family tree diagrams (⑤) and handwriting-like lines (⑥) interrupt the flow of text. Handling them is troublesome in TeX.

## 3 Solutions to problems

### 3.1 Solutions for font problems

#### 3.1.1 New kanji fonts are available

The Adobe-Japan1-5 character collection was published in 2002. This collection has 20,316 characters total, including 12,668 kanji characters. The number of kanji characters approximately doubled from the font that we used till then. AJ1-5 fonts could mostly cover the range of *Komonjo books*.

After that, the UTF package was developed by SAITO Shuzaburo. With this package, we could now use all characters of AJ1-5 fonts with pTeX.

In addition, we made an OpenType font named `ut-gaiji` to gather the additional characters which AJ1-5 did not have ("ut" means "u-tokyo" and "gaiji (外字)" means "external character").

### 3.1.2 Input *Seiji* using *Ryakuji*

To avoid the difficulty of *Seiji* input, we decided to make *Seiji* texts from *Ryakuji* texts. Preparation in advance was as follows:

- Write a style file named `jitai.sty` ("jitai (字体)" denotes "character shape") to describe pairs of *Ryakuji* and *Seiji*, in order to specify *Seiji* kanji via *Ryakuji* kanji.

```
\def\seiji{\@ifnextchar[{\@seiji}{\@seiji[0]}}
\def\@seiji[#1]#2{\expandafter
 \ifx\csname @SJ#1#2 \endcsname\relax
  \typeout{!!! Seiji #2(#1) undef.}#2\relax
 \else \csname @SJ#1#2 \endcsname \fi }
\let\S\seiji
\def\DefSeiji[#1]#2#3{\expandafter
 \gdef\csname @SJ#1#2 \endcsname{#3}}
\DefSeiji[0]{亜}{亞}
\DefSeiji[0]{唖}{\CID{7633}}
\DefSeiji[0]{逢}{\CID{8266}}
\DefSeiji[0]{悪}{惡}
\DefSeiji[0]{医}{醫}
\DefSeiji[0]{学}{學}
\DefSeiji[0]{国}{國}
\DefSeiji[1]{国}{圀}
        :
```

  `jitai.sty` has about 1,000 such pairs. For example, after these definitions, the input `\S{学}` produces "學," and you can input `\S[1]{国}` to get optional *Seiji* character "圀."

  The `\CID` command in a kanji pair specifies the internal character code from AJ1-5 (Character ID) directly.

- Write a Ruby program which converts *Ryakuji* texts into *Seiji* texts, referencing the jitai file. When the Ruby program finds *Ryakuji* kanji in a text file, it inserts the *Seiji* command.

  For example, the kanji string "文学青年" ("Literary youth") is converted into "`\S{文}\S{学}\S{青}`年," and typeset as "文學青年."

Then, work steps are as follows:

1. Input text file with *Ryakuji* kanji.
2. Convert *Ryakuji* texts into *Seiji* texts.
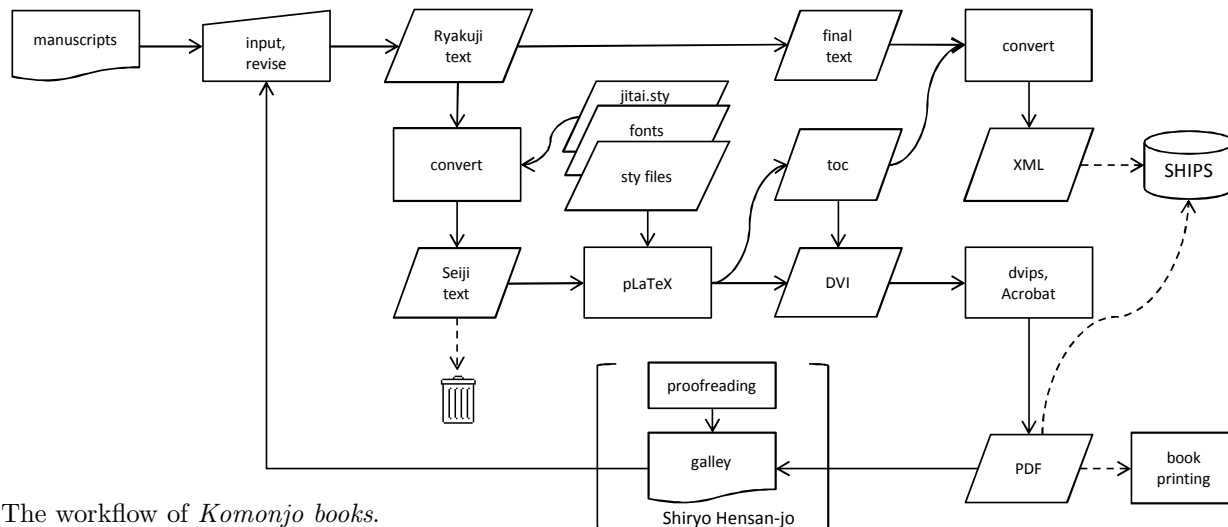3. Typeset *Komonjo books* with *Seiji* texts.

We revise only *Ryakuji* text files. We discard *Seiji* text files after typesetting has completed, since *Seiji* texts are hard to read due to the automatic insertion of *Seiji* commands.

### 3.1.3 *Hentai-gana*

We made *Hentai-gana* as external characters and put them into `ut-gaiji`.

### 3.1.4 *Bonji*

First, we used the free *Bonji* font, but it had only typical characters, so we often had to create additional characters. Therefore, now we use the 今昔文字鏡 (Konjaku-Mojikyo) font package which has about 2,000 *Bonji* characters. And we also use the `mojikyo` macro package developed by HONDA Tomoaki, to specify *Bonji* characters.

The workflow of *Komonjo books*.

## 3.2 Solutions for notes

### 3.2.1 Customize \marginpar for head note

A head note must be split when it "falls off" the end of a page. To implement this, we customized the definition of \marginpar. The original maintains a \@marbox to store a page of marginal notes. We just \vsplit \@marbox to \textheight, and put the remaining note text at the top of the \@marbox of the next page.

### 3.2.2 Interline note

An interline note is placed at the right or left side of the text. We defined \rnote and \lnote respectively, and an optional position adjustment, as follows:

> \rnote[*v-adjust*][*h-adjust*]{*note text*}.

## 3.3 Picture-like materials

We made the complicated family tree diagram as an integral number of line units, so that it can divide at the end of page.

For handwriting-like lines, we defined a \Curve macro, taking three points of a curve. Our \Curve is essentially the \bezier function with \unitlength of 1 zw (= width of a kanji). The second curve of ⑥ ( ∫ ) was drawn with 3 lines as follows:

```
\Curve(.4,-.5)(.6,-.5)(.6,.2)
\Curve(.6,.2)(.6,.5)(.6,.5)
\Curve(.6,.5)(.6,1.1)(1.1,1.1)
```

## 3.4 Generate database texts

Another purpose of typesetting with TeX was to generate database texts from final TeX files. We wrote a Ruby program for this, and it works as follows: delete unnecessary information except text,

unify *Seiji* in *Ryakuji*, convert character code from Shift-JIS to Unicode, and convert from TeX file to the XML format of the database.

## 3.5 Result

We show the whole workflow of *Komonjo books* in the figure above. Using our methods, more than half of the *Komonjo books* of Shiryo Hensan-jo have been typeset with TeX now.

As a result, we think that the two purposes to be expected in TeX (efficiency of typesetting and re-using of input for the database texts) have been accomplished.

⋄ NAKANO Ken
  k-nakano (at) livretech dot co dot jp

⋄ KOBAYASHI Hajime
  koba (at) livretech dot co dot jp