

インターネット上の古文書画像データベースシステムの作成

岡部 建次 永田 大 広瀬 順皓
駿河台大学 文化情報学部

概要

インターネット上で古文書画像データベースを公開する。原文書コピーの画像データベースを、インターネット上で公開すれば利用は促進される。明治政治史料の個人文書では、多くのデータは書簡の形をしている。書簡は数ページから成り立っているため、1書簡を1Webとし、書簡を構成する画像を貼り付ければよい。こうして多量の1書簡1WebのWebファイル(HTMLファイル)から成り立つデータベースをインターネット上に作った。本稿は、インターネット上で利用できる古文書画像データベースとデータベース作成のノウハウの報告である。

Archives image database system on the Internet

Kenji Okabe Dai Nagata Yoshihiro Hirose
The Faculty of Cultural Information, Surugadai University

Abstract

We open an Archives image database system to the public through the Internet. The open through the Internet promotes the access by the users. Most private documents of Meiji political archives take the form of letters. The letter consists of several pages. We assigned one letter to one Web, and each page of the letter is pasted to each web as an image. Then the image database which consists of many Web pages is prepared on the Internet server computer. We report the way of producing Web database.

1. はじめに

原文書コピーの画像データベースを、インターネットで誰もがどこからでも手軽にアクセスできるようにすれば利用は促進される。インターネットホームページでは、Web に写真が貼り付けてあるものが多い。そこで我々は写真の代わりに古文書画像データを貼り付ければ、インターネット上で古文書画像データベースを公開できると考えた。古文書のうち、我々が対象としている明治政治史料の個人文書では、多くのデータは書簡の形をしている。書簡は数ページから成り立っている。従って1書簡を1Web とし、書簡を構成する各ページ画像を貼り付ければよい。こうして多量の1書簡1Web の Web ファイル (HTML ファイル) から成り立つデータベースをサーバー上に作った。データベースとデータベース作成のノウハウを報告する。

2. 問題点

2. 1 CD-ROM版画像データベースの場合

通常画像データベースはCD-ROMで提供される。CD-ROM版画像データベースは高画質な画像データを提供できる反面、更新ができない。インターネット経由の場合、最新の画像データを容易に提供できる。

2. 2 Webデータベース作成上の問題点

・画像貼付の方法

Web データベースでは、1書簡に相当するHTML ファイルに書簡の各ページ画像をリンクする。これはタグを用いて、図1のようなHTML 形式のファイルを作成する。図2にこのHTML ファイルをブラウザで表示した実際の画像を示す。

・大量の原文書画像

数枚(数ページ)から成る1書簡は図3のようにブラウザに表示される。これは図4のようなHTML ファイルを作ることで可能になる。個人文書データベースにおいては、画像ファイルは膨大な数になる。例えば、谷干城関係文書画像は約4000枚、書簡件数で約900に及ぶ。このような大量の画像ファイルをリンクした書簡 Web (谷干城関係文書の例では900) を手作業で1件ずつ作成するには、大変な時間と労力がかかる。

3. 解決法

大量の画像ファイル(各ページ)を書簡ごとにリンクする作業は、プログラムによって自動的に行う方法を考案した。

4. プログラムによるHTML化の方法

4. 1 原文書画像のファイル名と文書(書簡)の関係

原文書画像のファイル名は、文書(書簡)番号(5桁) + (文書(書簡)内での) 枝番号(3桁)の計8桁である。この8桁のファイル名に拡張子が3桁付与された形式になっている。このファイル名から、HTML形式によるリンクファイル(テキストファイル)をプログラムで作成する。

表1のように、原文書画像が存在するものとする。先頭のkは黒田清隆関連の文書であることを示す。次の4桁は、一つの書簡を表す数字である。一つの書簡は、複数の原文書画像から構成される。次の3桁は、一つの書簡内の枝番号となる。つまり、表1の例では、k0001001.gif から k0001007.gif までが一つの書簡であり、k0002001.gif から k0002003.gif までが、次の書簡を意味する。書簡によって、

構成する原文書画像の数はさまざまである。1つの画像のみ書簡もあれば、100を超える画像から構成されるものもある。

表1 原文書画像のファイル名 (サンプルであり実際のものとは、異なる)

| | | | | |
|--------------|--------------|--------------|--------------|--------------|
| k0001001.gif | k0001002.gif | k0001003.gif | k0001004.gif | k0001005.gif |
| k0001006.gif | k0001007.gif | k0002001.gif | k0002002.gif | k0002003.gif |

4. 2 「1書簡-1Webページ」によるリンク

1つのHTMLファイルから1つの書簡分の画像をリンクするようにする。つまり、「1書簡-1Webページ」である。「1書簡-1Webページ」にすることによって、文書目録データベースで検索した書簡をすばやく見ることができる。

「1書簡-1Webページ」であるから、同一の書簡の原文書画像が一つのHTMLファイルとしてプログラムで出力される。出力されるHTMLファイルのファイル名は、書簡の番号(原文書画像の先頭5桁)を付加した。表2・表3のように、HTMLファイルから1つの書簡分の画像がリンクされる。

表2 k0001.htm からリンクされる原文書画像 (書簡0001の画像)

| | | | | |
|--------------|--------------|--------------|--------------|--------------|
| k0001001.gif | k0001002.gif | k0001003.gif | k0001004.gif | k0001005.gif |
| k0001006.gif | k0001007.gif | | | |

表3 k0002.htm からリンクされる原文書画像 (書簡0002の画像)

| | | |
|--------------|--------------|--------------|
| k0002001.gif | k0002002.gif | k0002003.gif |
|--------------|--------------|--------------|

すなわちプログラムで文書1件ごとに画像ファイル名を読んで、解析して、図4のテキストファイル(HTMLファイル)を作成していく。図3に実際の画像を示す。

4. 3 「一定の画像数-1Webページ」にする必要性

「1書簡-1Webページ」によるリンクを行うと、書簡によっては、100を超える画像から構成されるものもある。このような書簡をすべて一つのHTMLファイルにしてしまうと、画像ファイルを読み込むだけで非常に時間がかかる。100以上の画像が並んでしまえば、かえって手紙が読みにくい。そこで、一定の画像数を超えた場合、1書簡の中でも別のファイルを分けることにした。例えば、表1の原文書画像の書簡0001において、5つ以上の原文書画像から構成される書簡を複数のWebページに分けた場合、表4・表5のようになる。プログラムを利用するので、何画像で別のHTMLファイルに分けるかの設定は、容易に変更できる。

表4 k0001.htm からリンクされる原文書画像 (書簡0001の前半の画像)

| | | | | |
|--------------|--------------|--------------|--------------|--------------|
| k0001001.gif | k0001002.gif | k0001003.gif | k0001004.gif | k0001005.gif |
|--------------|--------------|--------------|--------------|--------------|

表5 k0001_2.htm からリンクされる原文書画像 (書簡0001の後半の画像)

| | |
|--------------|--------------|
| k0001006.gif | k0001007.gif |
|--------------|--------------|

出力されるHTMLファイルは図5・図6のような形式である。

5. おわりに

1 Web が 1 文書になっているためインターネットで容易に利用できる。各書簡内の各ページ(1 画像)は 1 Web に収まる大きさで、画面で読める解像度、プリントして読める鮮明さを維持している。

既に谷干城文書の CD-ROM の公開配布をおこなっており、インターネット経由による web データベースの作成中である。本研究は私学振興財団学術研究奨励金、駿河台大学学内共同研究補助金、科学研究費(データベース促進)の助成を受けて行った。

図 1 HTML による画像のリンク

```
<HTML>
<HEAD>
<TITLE>k0001</TITLE>
</HEAD>
<BODY>
k0001001.gif<BR><IMG SRC=k0001001.gif><P>
</BODY>
</HTML>
```

図 2 HTML によるリンクの画像

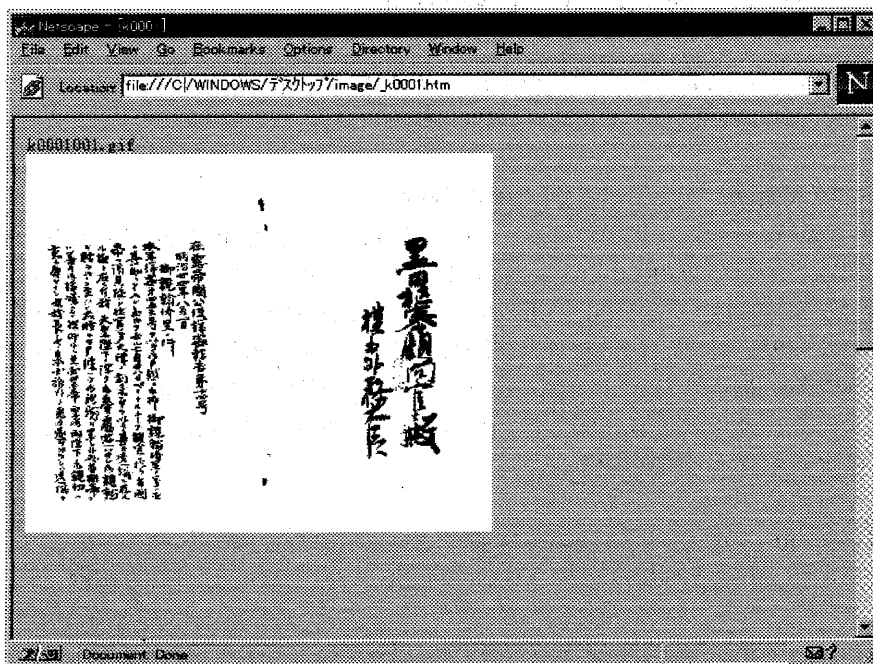


図3 「1書簡-1Web ページ」によるリンクの画像

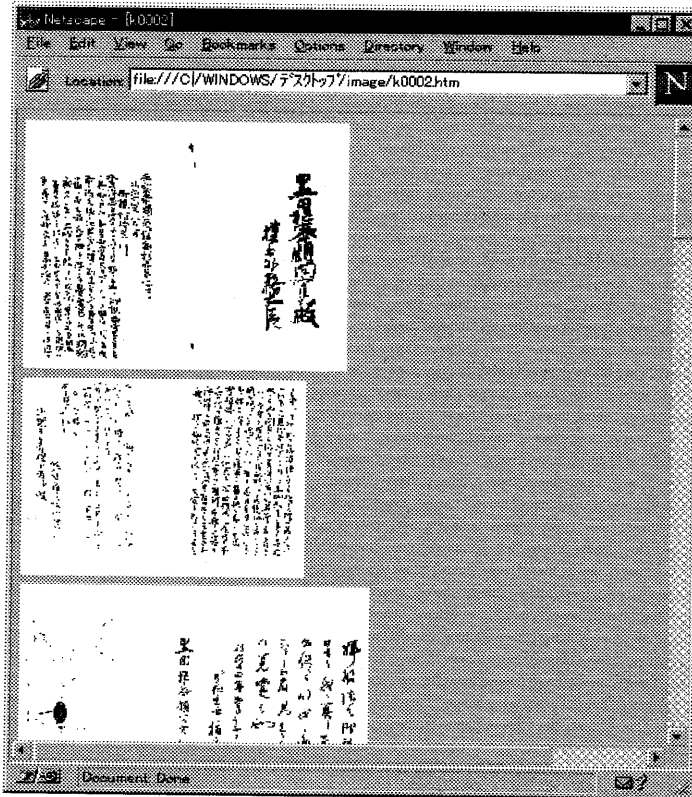


図4 「1書簡-1Web ページ」によるリンクのHTMLファイル (ファイル名 k0001.htm)

```
<HTML>
<HEAD>
<TITLE>k0001</TITLE>
</HEAD>
<BODY>
k0001001.gif<BR><IMG SRC=k0001001.gif><P>
k0001002.gif<BR><IMG SRC=k0001002.gif><P>
k0001003.gif<BR><IMG SRC=k0001003.gif><P>
k0001004.gif<BR><IMG SRC=k0001004.gif><P>
k0001005.gif<BR><IMG SRC=k0001005.gif><P>
k0001006.gif<BR><IMG SRC=k0001006.gif><P>
k0001007.gif<BR><IMG SRC=k0001007.gif><P>
</BODY>
</HTML>
```

図5 「一定の画像数-1 Web ページ」によるリンクのHTMLファイル (ファイル名 k0001.htm)

```
<HTML>
<HEAD>
<TITLE>k0001</TITLE>
</HEAD>
<BODY>
k0001001.gif<BR><IMG SRC=k0001001.gif><P>
k0001002.gif<BR><IMG SRC=k0001002.gif><P>
k0001003.gif<BR><IMG SRC=k0001003.gif><P>
k0001004.gif<BR><IMG SRC=k0001004.gif><P>
k0001005.gif<BR><IMG SRC=k0001005.gif><P>
</BODY>
</HTML>
```

図6 「一定の画像数-1 Web ページ」によるリンクのHTMLファイル (ファイル名 k0001_2.htm)

```
<HTML>
<HEAD>
<TITLE>k0001</TITLE>
</HEAD>
<BODY>
k0001006.gif<BR><IMG SRC=k0001006.gif><P>
k0001007.gif<BR><IMG SRC=k0001007.gif><P>
</BODY>
</HTML>
```