

## WWWによるJAPAN/MARCの提供実験

安齋 宏幸 山本毅雄  
anzai@ulis.ac.jp, yamamoto@ulis.ac.jp  
図書館情報大学  
〒305 つくば市春日 1-2

JAPAN/MARC を UNIX-WS にコピーし提供実験を行った。本システムにおいて利用者は、WWW(World Wide Web) 環境によるクライアント・サーバー方式により任意の端末から書誌情報の検索を行うことができる。利用者に分かりづらい分かち書き語による書名検索には、複数ハッシュふるい分けを用いたキーワード切り出しシステムを活用することにより、柔軟な検索を行うことを可能にした。また、DBMS には ADABAS を用いることで検索速度の向上と DB 管理の省力化をはかっている。本研究では日本語検索に対する複数ハッシュふるい分け法の有効性、ならびに、書誌情報提供において WWW を利用することの価値とその限界について検証する。

## JAPAN/MARC EXPERIMENTAL SERVICE THROUGH WWW

Hiroyuki Anzai and Takeo Yamamoto  
anzai@ulis.ac.jp, yamamoto@ulis.ac.jp  
University of Library and Information Science  
1-2 Kasuga, Tsukuba-city, Ibaraki, 305 Japan

An experimental JAPAN/MARC bibliographic information retrieval service through WWW(World Wide Web) is described. Since Japanese words are not separated by obvious delimiters, ensuring same segmentation between the query and the database is a problem. The present system solves the problem by using the multiple-hash screening technique for processing both book titles and query strings, based on the same dictionary and using similar algorithms. Database management is handled by ADABAS, reducing management chores and response time. The effectiveness of multiple hashing screening technique for a Japanese text-based information system is examined, and the limitation of WWW's hyper-text environment for bibliographic information retrieval service is discussed.

## 1 はじめに

インターネットには無数のネットワーク情報資源が存在している。図書館などが提供している目録情報サーバーもまた、印刷出版物に関する情報を手に入れるための欠かせない情報資源となっている。しかしながら、その利用法は煩雑な場合が多く主題検索はおろか、目的とする書物のタイトルからでさえ満足に検索を行うことができないことがある。

書名検索を難しくしている原因のひとつが、分かち書きの存在である。日本語では1語という単位をはっきりと特定できない。そのため言葉をどこで区切って検索すればよいのか、果たして対象システムではトランケーションを使えるのか、キーワードはどういった基準で抽出されているのか、など、分かち書きに関連して利用者を悩ませる問題は多い。

今回我々が試作したシステムでは、中本が開発した複数ハッシュふるい分けによる日本語形態素抽出アルゴリズムを採用している。複数ハッシュふるい分けは、ICOT フリーソフトウェア No.33 形態素辞書をもとに作成するハッシュ表を使い、文字列中に存在している辞書該当語をすべて抽出する。データベースにおけるインデックス作成と利用者からの書名検索要求処理の双方に同抽出アルゴリズムを実装することで、日本語検索に伴う前述した欠点を克服することに成功した。

さらにWWWによるハイパーテキスト環境を広域ネットワーク（インターネット）への情報提供手段として用いることで、検索結果表示の著者名や件名、分類番号をクリックするだけで著者名・主題検索を行うことを可能にした。関心のある書誌を表示させ、その主題情報をもとに再び検索を行うことは一般によく用いられている。このような検索サイクルをハイパーテキスト環境で容易に行えるようにした。

## 2 システムの概要

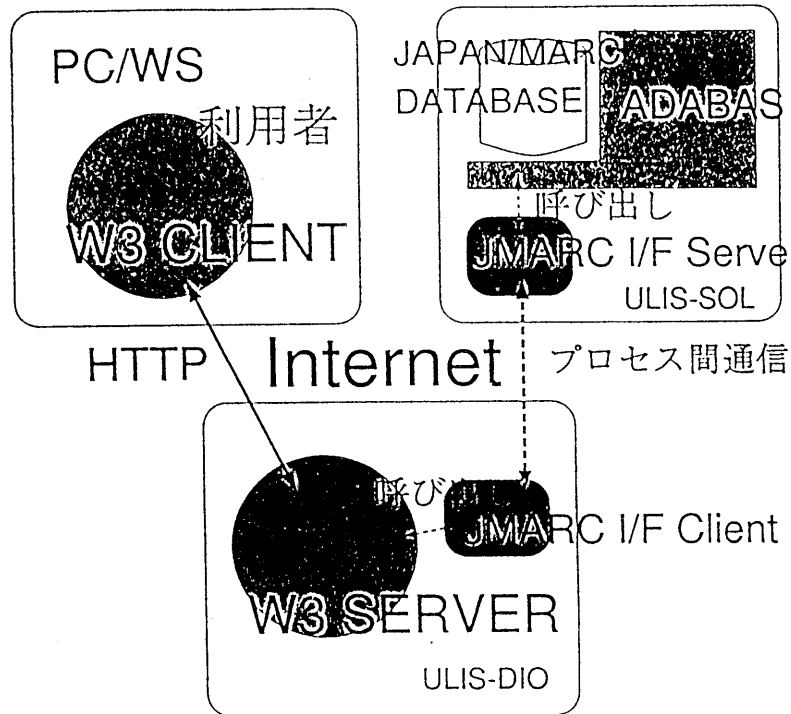
試験的に国立国会図書館の許可を得て提供したのは、JAPAN/MARC 1年分の約73,549レコード(1993 vol.13 - 1994 vol.12)である。本システムにおいて、JAPAN/MARC データベースはADABAS上に構築される。

利用者は任意のW3クライアントから本システムのURL(Uniform Resource Locator)を指定することで、検索フォームや検索例、データベースの情報などを受け取ることができる。W3サーバーとクライアント間は、NCSA httpdによるHTTP(HyperText Transfer Protocol)で通信を行う。検索フォームはWWWのFill-Out Form機能で実現されており、利用者が入力した検索式は送信ボタンを押すことによってW3サーバー(httpd)に送られる。

httpdはcgi(Common Gateway Interface)を通じてインタフェースクライアントプログラムを呼びだし、検索要求をプロセス間通信でインタフェースサーバーに連絡する。ADABASへの検索要求はC言語に埋め込まれたADABASコマンドによって行われる。JAPAN/MARCインタフェースサーバはクライアントから受け取った検索式の評価を行い、検索要求をADABASに送る。受け取った結果はHTML(Hyper-Text Markup Language)形式に変換され、クライアントに送信される。インタフェースサーバは書名検索文字列への複数ハッシュふるい分けや、ADABASコマンドに用意されていないトランケーション処理なども行う。

検索結果をあらわす HTML はそのまま W3 クライアントに送られ、利用者は自分の好みの W3 クライアントソフトウェアで結果をブラウジングすることができる。また、詳細な書誌情報の取り寄せや著者名、主題などからの再検索をクリック動作で行うことができる。システムの全体構成図を下に示す。

### システム全体構成図



### 2.1 JAPAN/MARC の利用

JAPAN/MARC MT を、HITAC M660 を経由して SUN WS 上にコピーし、これを標準的なテキストファイルに変換した。変換処理は国立国会図書館から出版されている JAPAN/MARC マニュアル [1] をもとに行った。

テキストファイルとして処理できるようになった MARC に対して、複数ハッシュふるい分けを用いた書名インデックス生成を行う [2]。

ここでいうふるい分けとは、あるキーが辞書中に存在しているかどうかを調べることである。複数のハッシング関数を用いて、辞書中に存在する単語を複数のハッシュ表の 1 ビットのフラグのエントリに変換する。ふるい分け対象のキーについてそのすべてのハッシュ表エントリに 1 ビットのフラグが認められれば、そのキーは辞書中存在しているとみなされる。

中本による研究では ICOT(Institute for New Generation Computer Technology) フリーソフトウェアの形態素辞書を使用している。本システムの作成においては中本による複数ハッシュふるい分けシステムを、書名検索におけるインデックス作成と照合処理の中核とした。

変換したレコード例、与えられた書名に対するインデックス例を下に示す。このようなインデックスは対象書誌のシリーズ名についても作成され、MARCレコードに追加される。

### 変換したレコード例

010A4-17-268005-4  
020AJP  
020B93065055  
100A19931117 1993 H1JPN 1312  
251A 地域経済レポート  
251D 平成5年  
251F 経済企画庁調査局／著  
270A 東京  
270B 大蔵省印刷局  
270D 1993. 5  
275A 262p  
275B 21cm  
350A 平成5年の副書名：地域経済の構造変化と調整局面  
360B 1800円  
551A チイキ ケイザイ レポート  
551X Tiiki keizai reporto  
551B 地域経済レポート  
551D 1993  
658A ニホン／ケイザイ  
658X Nihon／Keizai  
658B 日本／経済  
677A 332. 107  
685A DC55  
751A ケイザイ キカクチョウ チョウサキョク  
751X Keizai kikakutyootyosakyoku  
751B 経済企画庁調査局  
905A DC55-E510

### 書名インデックス例

公務員試験労働法  
公/公務/公務員/務/員/試/試験/験/労/労働/働/法/  
最新データ比較政治ハンドブック  
最/最新/新/データ/比較/政/政治/治/ハンド/ハンドブック/ブック/  
精選不動産登記法商業登記法先例書式60問  
精/精選/選/不/不動/不動産/動/動産/産/登/登記/記/記法/法/商/商業/業/登/登記/記/記法/法/先/先例/例/書/書式/式/問/  
お願い！かなえて!!トランプ恋占い  
お/願/願い/い/か/かな/かなえ/な/なえ/え/えて/て/ト/トランプ/ランプ/恋/占/占い/い/  
コンピュータでビジネスチャンスをつかむ本  
コンピュータ/で/ビジネス/チャンス/を/つ/つか/つかむ/か/かむ/む/本/  
オブジェクト指向分析手法入門  
オブジェクト/オブジェクト指向/指/指向/向/分析/手/手法/法/入/入門/門/  
ファイリング革命最前線  
ファイリング/革/革命/命/最/最前/最前線/前/前線/線/

## 2.2 検索文字列へのキーワード抽出

実際に利用者が入力した書名検索語は、同様に複数ハッシュユふるい分けを用いて形態素が抽出される。抽出されたキーワードは検索での処理効率を上げるために、他の文字列の一部として存在している部分文字列の除去を行う。

検索文字列の処理の一例

1. 公務員試験労働法
2. 公/公務/公務員/務/員/試/試験/験/労/労働/働/法/
3. 公務員/試験/労働/法/

## 2.3 書名による検索

同じ辞書を用いてふり分けられた書名インデックスとキーワードは、もとの文字列が同じ場合には完全に検索を行うことができる。一見当たり前のように感じられる話ではあるが、利用者にとって分かち書きという存在がこの容易な検索を難しくしてきたという事実は前述のとおりである。さらに、一般に行われる「経済」「入門」などといったワード単位での検索の論理積作業が、事実上、一度の検索で自動的に行われるなどのメリットも存在する。直感的に入力された書名文字列は、その文字列に含まれる部分ワードの論理積検索に分解されることにより、さらに再現率の高い検索へと変化する。

一方、固有名詞や比較的新しい現代用語など、そもそも辞書に登録されていない単語は当然のことながら抽出されない。しかし、それら辞書に登録されていない未知語は、登録されているもっと短い文字列の単語の集合としてインデックス処理されることが多い。特に漢字やひらがなは1文字で意味をなす形態素として辞書に登録されている場合が多く、これらの漢字やひらがなから構成される未知語に対しても本システムにおける検索はかなりの柔軟性を示している<sup>1</sup>。ただし、検索速度と精度のさらなる向上のためには将来的に改善を要求される部分であり、辞書そのものの管理を行うツール群の開発が必要である。

## 2.4 ADABAS でのデータベース構築

加工した MARC レコードをデータベースとして運用するために ADABAS[3] を用いた。ADABAS はインバーテッド型 DBMS である。ディスクリプタ指定をおこなったフィールド値からの高速検索が可能で、検索集合同士のブール演算なども高速に行うことができる。同 DBMS を選択したことによって、検索から検索結果の抽出までの時間を十分満足できる範囲内に押さえ、複雑な検索処理も効率よく行うことができた。マシンの負荷状態や質問内容にもよるが、検索処理自体は数秒で行われる。検索要求と検索結果のネットワークでの転送にかかる時間や、検索フォームを W3 クライアントが描写するのにかかる時間のほうが一般に長くなる。

## 2.5 WWW での提供

インタフェースプログラムに対して WWW サーバーである NCSA httpd が Common Gateway Interface(CGI)[4] を使って検索を行えるようにした。インタフェースプログラムは検索要求を ADABAS に伝え、得られた検索結果を HTML 形式に変換して WWW サーバーに返す。主題検索を実現するためのアンカーや、検索結果が 0 件の場合の情報ページへのリンクも、この際に HTML のなかに折り込まれる。

<sup>1</sup>一方、カタカナの新語はその部分文字列も含めてまったく抽出できないことが多い。幸いにもカタカナ語は前後のひらがなや漢字との区切りが明確であるため、複数ハッシュふり分けとは別にこういったカタカナの新語を抽出する仕組みを実装している。

## 2.6 海外への対応

本システムには書名、著者名、分類記号など英数文字で記述が用意されている情報を海外に提供するための半角英数文字による検索・表示モードが用意されている。外国からの利用者が、日本語を使えない環境でも日本における書誌情報を手に入れることのできるための配慮である。インターネットでの情報提供における標準語が英語である以上、情報提供における非日本語環境への対応は無視できない。

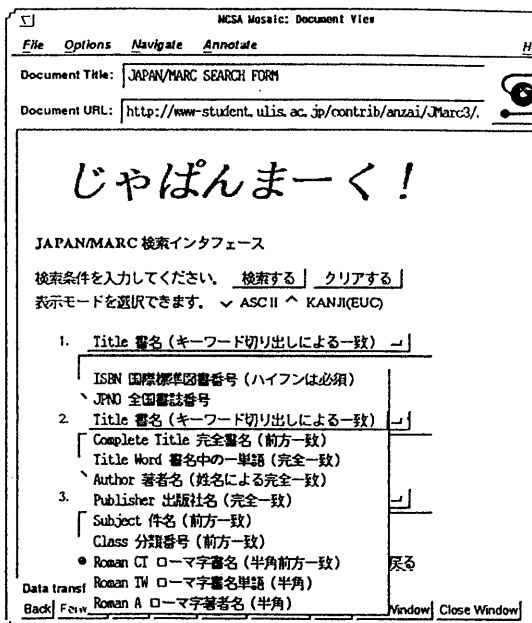
## 3 システム実行例

実際の検索例を以下に示す。

1. ホームページからはリンクをたどることによって、望むページを見ることができる。
2. 検索フォームでは、プルダウンメニューの形式で検索対象のフィールドを指定できる。結果表示は漢字と ASCII から選択することができる。



1



2

3. 検索結果は件数に応じて表示形態が変化する (簡略表示例)。書名をクリックすれば、書誌の詳細表示を見ることが出来る。
4. 書誌の詳細表示。著者名、件名、分類番号から再検索が可能。

Document Title: JAPAN/MARC SEARCH RESULT

Document URL: http://www-student.ulis.ac.jp/cgi-bin/anzai/jmarc3/

書名検索文字列(T)は次のように評価されました。  
経/経済/経済入門/入門  
経済入門

該当レコードは 21 件です。

## Number 1

全国書誌番号:93017245  
経済英語入門 / 石塚理彦著  
東京: 日本経済新聞社, 1992.12  
4-532-10466-1 1992.12

## Number 2

全国書誌番号:93018715  
経済刑法入門 / 中山研一(ほか) 編著  
東京: 成文堂, 1992.11  
4-7923-1298-1 1992.11

## Number 3

全国書誌番号:93018730  
ベーシック/アジア経済入門 / 日本経済新聞社編  
東京: 日本経済新聞社, 1993.1  
4-532-10650-8 1993.1

## Number 4

全国書誌番号:93024357

Back Forward Home Reload Open... Save As... Clone New Window Close Window

3

5. 検索ではブール演算が可能。

6. 半角英数での結果表示例。利用者が指定した「漫画」という書名単語が、複数ハッシュふるい分けによるインデックス対象となっているかを調べている。

Document Title: JAPAN/MARC SEARCH RESULT

Document URL: http://www-student.ulis.ac.jp/cgi-bin/anzai/jmarc3/

全国書誌番号 93018730  
ISBN 4-532-10650-8  
書名 1 ベーシック/アジア経済入門  
著者 1 日本経済新聞社編  
出版地 東京  
出版者 日本経済新聞社  
出版年 1993.1  
ページまたは冊数 154P  
大きさ 18CM  
登録名 日経文庫  
注記 1 900円

## データエレメント数

記号数: 1  
双書名数: 1  
注記数: 1

書名ディスクリプタ数: 10  
ベーシック/アジア/ケイザイ/ニューモン/経済/経済入門/入門  
著者名ディスクリプタ数: 4  
ニホン/ケイザイ/シンブン/シヤ/日本経済新聞社  
件名ディスクリプタ数: 1  
アジア-経済  
分類番号ディスクリプタ数: 2  
332.2 DC141

Back Forward Home Reload Open... Save As... Clone New Window Close Window

4

Document Title: JAPAN/MARC SEARCH EXAMPLE

Document URL: http://www-student.ulis.ac.jp/contrib/anzai/jmarc3/

ラジオボタンで出力の漢字/半角切り替えが出来ます。また、各検索条件の論理関係はAND, OR から選択できます。ORは完全一致の同じタグ同士を接続できるものと、なんでも接続できるものがあります。

結合強度はOR(同), AND, OR(異)の順です。

前方一致を実現しているタグ、キーワード切りだしを利用しているタグは、内部的に論理演算を行うことで表現しています。これらのタグをAND, OR で結合することはできません。検索結果については保証されません。

表示モードを選択できます。 √ ASC II ^ KANJI(EUC)

1. Title Word 書名中の一単語 (完全一致) 
  - コミック
  - ^ OR(同じタグ) √ AND √ OR(異なるタグ)
2. Title Word 書名中の一単語 (完全一致) 
  - 漫画
  - √ OR(同じタグ) ^ AND √ OR(異なるタグ)
3. Author 著者名 (姓名による完全一致) 
  - アイカフ

Data transfer complete.

Back Forward Home Reload Open... Save As... Clone New Window Close Window

5

Document Title: JAPAN/MARC SEARCH RESULT

Document URL: http://www-student.ulis.ac.jp/cgi-bin/anzai/jmarc3/

Title Word 漫画 is a Kanji Index.

## Bibliography

JPN0 93067695  
ISBN 4-408-30151-5  
Title 1 NATUKASI MANGA00  
Author AIKAWA NOBUHIKO

## Data Elements: 1

Title Roman Keywords: 2  
NATUKASI MANGA00  
Roman Author Keywords: 2  
AIKAWA NOBUHIKO  
Roman Subject Heading Keywords: 1  
MANGA  
Classification codes: 2  
726.1 KC486

Copyright 1993-1994, Japan National Diet Library.  
Copyright 1995, Hiroyuki Anzai.  
All Rights Reserved.  
● Back to search form

anzai@ulie.ac.jp

Back Forward Home Reload Open... Save As... Clone New Window Close Window

6

## 4 システムの評価

この提供実験において、本システムは1日に100件前後の検索要求を受け続けている。情報システムである以上、蓄積されたデータがその価値を大きく左右する。残念ながら1年という収録年月では、利用者に本格的な検索のニーズを生じさせることができない。このため明らかに「通りすがり」的な利用のみが目立っており、検索要求の履歴からは本システムの利用者の明確な反応は得ることができなかった。

ハイパーテキストを用いた検索結果のブラウジングなどは、すでにたくさんのWWWを基盤に据えた情報システムで実現されている。しかし、次の点で本システムは評価できる。

- 複数ハッシュふるい分けを用いて日本語書名に対する柔軟な検索を実現した
- ADABASで検索処理速度の向上と豊富なアクセスポイントを実現した
- ハイパーテキスト(W3)を用いた主題や著者からの検索フィードバックを可能にした

本システムでは、検索結果集合を保存してブール演算を行うことができない。この課題は本システムで採用したWWW環境が、コネクションレス指向のクライアント・サーバー方式の情報システムであることから派生している。利用者には要求のたびごとに結果集合の書誌情報すべてが提供されるため、件数を絞り込んでから表示したり、結果集合を順に組み立てながら検索を行うことができない。

Z39.50[5][6]が、こういった書誌情報の提供には本来望ましい。Z39.50はANSI(American National Standards Institute)が定めた分散環境下におけるコンピュータ間のコネクション指向情報検索プロトコルである。このプロトコルを用いることによって、結果集合をサーバー側に作成し組み立てて検索を行うことができる。また、検索結果のソートや検索語のスキヤニング、必要とするエレメントの絞り込みや検索結果上限などの設定等をサーバーに連絡することを可能にしている。

## 5 これからの展望

本システムではJAPAN/MARCを対象書誌データとして運用した。我々は同様の手法を取り入れて、図書館情報大学附属図書館OPACをインターネット上に提供する準備を始めた。この新しいシステムにはインターフェースプログラムにZ39.50のサーバーとしての特徴を取り入れることを目標にしている。

### 参考文献

- [1] 国立国会図書館編。JAPAN/MARC マニュアル 図書編。東京、国立国会図書館、1992,129p.
- [2] 中本賢一、山本毅雄、長谷部紀元。複数ハッシュふるい分け法の日本語情報システムへの応用。情報システム 48-7。(1994.3)
- [3] ソフトウェア・エージェンシー。ADABAS(UNIX) マニュアル。東京、ソフトウェア・エージェンシー。
- [4] National Center for Supercomputing Applications. Common Gateway Interface.  
<http://hoohoo.ncsa.uiuc.edu/cgi/>
- [5] Z39.50-1995 Maintenance Agency Text.  
<http://lcweb.loc.gov/z3950/agency/1995doce.html>
- [6] William Moen. ANSI/NISO Z39.50 Protocol : Information Retrieval in the Information Infrastructure.  
<http://www.cni.org/pub/NISO/docs/Z39.50-1992/www/50.brochure.toc.html>