

MIRRF システムによる文献データの作成と利用について

有川 節夫 (九大理)

1. はじめに

MIRRF システム (Multistage Information Retrieval System Based upon Researcher Files) は 1974 年 ~ 75 年に、次の目的で開発した研究者用の文献情報検索システムである。[1]

- 1) 研究者の身の周りの文献データを蓄積し組織化すること。
- 2) データの加工を自動化すること。
- 3) 研究者の独自性を検索に積極的に反映させること。
- 4) 網羅性のあるデータ・ベースを利用すること。
- 5) 高次の知的検索を行えること。

こうした目的に従って、研究者ファイルの概念を導入し、文献を 2 種類の研究者ファイルの形に貯え、それを検索時に積極的に利用する方式にシステム全体をまとめた。ここにいう 2 種類の研究者ファイルとは、論文末尾にある参考文献の形の文献リストと、主として熟読された論文に対して研究者が自分の考えに従って付与する概念としての木構造をもつキーワードのファイルである。第 2 種の研究者ファイルを連関 (association) を用いて多段に利用するわけである。データ加工の自動化に関しては、文献のタイトルからキーワードを自動抽出する索引システム REKWEST、文献コード化、木構造付きキーワード・ファイルを部分ソーラスとして把握それを編集するシステム等を開発してきた。

最初の MIRRF システムはミニコン (F U-200/U-400) で実現して、実験的に使用し、またデータも約 20,000 件の文献データを作成してきた。このシステムでは上記の目的のうち 1) ~ 3) は或る程度満たされるものであった。4) に関しては、最近大型センターの FAIRS により INSPEC のデータを検索出来るようになった。5) に関しては、現在基礎的な研究を進めている最中であるが、近似的なシステムとして逐字サーチシステム TEXTIR が能力とスピードの両面から有用であるこ

とを見てきた。

一方、1978年までの特定研究「学術情報」M委員会において、学術情報用のデータ・ベース・システムをData Box, PRF (Private Researcher File), UDL (User Data Library), CDL (Center Data Library) といった各段階で把える考え方が提案され、その委員を通じて次第に定着しつつある。更に大型計算機センターの利用形態もTSS端末からの利用が一般的になってきた。

こうした経験や環境・状況の変化を正視して、研究者用の情報システムのあり方を再考してみると、我々が当初掲げた目標は現在でも正当であるように思われる。しかし、それだけでは十分でないことも明白である。このような情報システムの扱うべきデータは文献データに限定されたものではなく、論文作成や書類の整理、数値データの処理まで含めた日常のほとんどの研究活動を支援してくれる作業場としてのシステムであることも要求される。勿論外部のデータも自由にやり取りし、他の研究者と情報システムを通じて通信でき、更に、こうした情報システムを介した活動を基盤にして、4)でいうところの網羅性のあるデータ・ベースが構築できるものでなければならない。

このような観点に立って、現在九大大型計算機センターのM200に新たにMIRRFシステムを実現しつつある。最近、いわゆる研究者の作業場としてのシステムとしては一応の完成をみたので、ここにその概要を紹介する。全体のシステムは完成後センターを通じて公開する予定である。公開に伴う問題点とその解決法についても併せて述べてみる。なお、こうした研究者の作業場としての情報システムには便利で手ごわりのいい汎用のエディタが不可欠である。現在のMIRRFシステムにはTEDITシステムがある。これについては篠原[4]を参照されたい。

2. ファイルの構造

記憶域の有効利用とソフトウェアの作成効率を考慮して、2方向のポインタをもった逐次ファイルと各種ファイル名やキーワードの管理のためのB-treeをファイルの物理構造の基本としておいた。

2方向ポインタで連結された逐次ファイルは、主にゴールファイル用に使われるものである。現在のシステムでは、512文字(バイト)からなるブロック(ページ)を8のセクタに分割して図1のように構成した。このような何個かのセクタからなるファイルの消去は、このファイルをポインタもそのおまにしておき、消去されたファイルの末尾に追加することによって行う。従って、ファイルの消去は瞬時に行われる。ファイルの移動についても同様である。またファイルの作成に際しては、1) 使いかけのページにあるセクタ、2) 一度消去されたセクタ、3) 新しいセクタの順序で必要なセクタが発行されるので、記

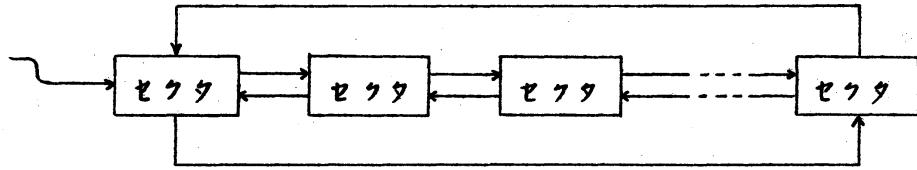


図1. 逐次ファイルの構造

憶域は常に効率よく使用されることになる。

B-treeのレコード構造は図2のようなものである。図において KLP の部分に実際のキーの長さが記入される。KLP の位置は各B-treeごとに指定できる。

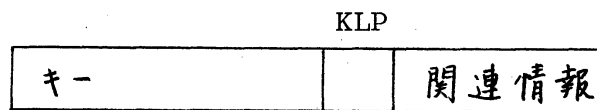


図2. B-treeにおけるレコード

一定の長さの場合には $KLP < 0$ であり $|KLP|$ でその長さを指定する。このようなレコード構造を採用し、B-treeにおけるキーの検索、挿入、削除のルーティンでは関連情報の部分を全く扱わないようにしたことによって、B-treeをファイル名やキーの管理だけでなく、語の頻度調査などデータの整列化を必要とする作業にもそのまま適用することができる。

3. MIRRF システムにおけるファイルの構成

現在の MIRRF システムのファイル構成は図3に示す通りである。TEDIT ファイルと勿論外部データ・セットを除いて全て512バイトのレコード長のダイレクト・アクセス・ファイル上に構成されている。

MIRRF ファイルは図1, 2で示した逐次ファイルとB-treeからなる。

B-treeはファイル名を管理する。論理的には最大10段階までの階層構造をもつものである。ここでのファイルはpass numberによって保護されるが、不用意な消去から護るため及びファイル実体の共有のためにrenameの機能がある。これはファイルの実体に新しい名前を追加するだけであって、一方の名前のもとに変更が生じたときにはじめて実体のコピーが作られる。

作業用ファイルは利用者がMIRRFシステムを使って作業をするとき扱われるもので、図1の構造をもつものである。このシステムでは端末機器等も全てファイルと見做される。そうしたファイルや他の各種ファイルへの1からのデータの転送は全てsave/loadで行われる。最後にload又はpull upされたフ

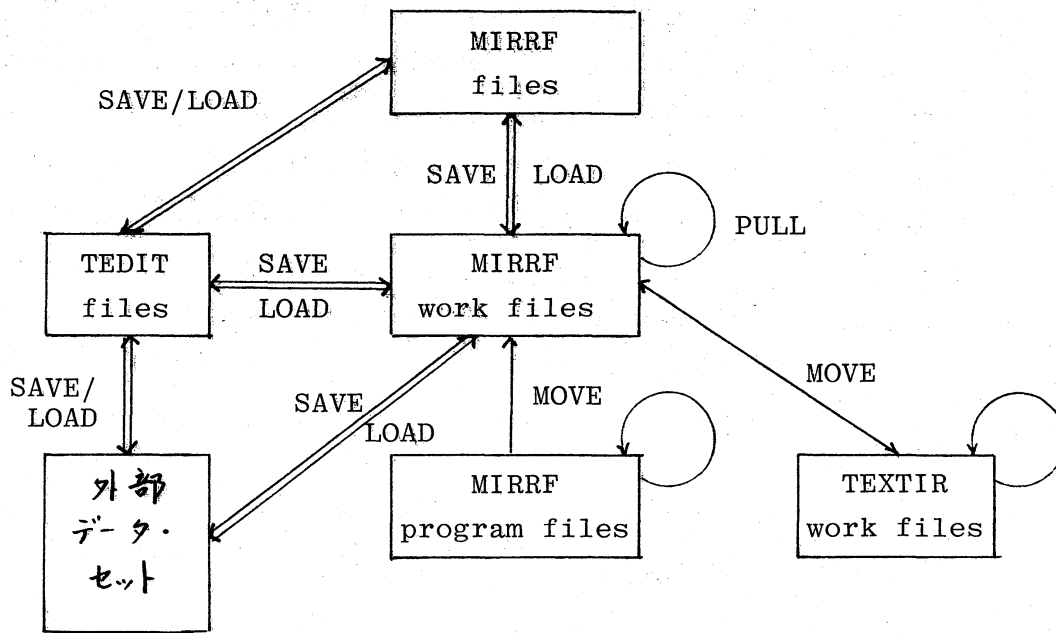


図3. MIRRF システムにおけるファイル

イルがトップにあるようになってくる。従って古いファイルは新しいファイルがもってこられる度に push down されることになる。作業用ファイルの領域が物理的にいっぱいになるときは、作業中に最近の10個のファイルを残して古いものから順次消去されて作業域が確保される。古いファイルをトップにもってくる (pull) こともできるので、作業中のファイルを不注意による消滅から守ることが出来る。

図1に示すように、作業域は TEDIT ファイル、MIRRF プログラム・ファイル、TEXTIR 作業用ファイルだけでなく外部データ・セットと連絡できる。

MIRRF プログラムファイルは MIRRF システムにおけるキーボードからの全てのコマンド、データを記録するためのものである。構造は作業用ファイルと同じである。

TEXTIR 作業用ファイルは MIRRF システムのサブシステム TEXTIR の検索結果をファイル機番、レコード番号、セクタ番号、開始文字位置、長さ、ヒットした質問番号の組んで記録するためのものである。構造はプログラム・ファイルと全く同じである。

〈注意〉 図3は、MIRRF 作業用ファイル又は TEDIT ファイルから見たものである。⇒ はファイル実体のコピーを示し、→ はポインタのみの移動を示す。

4. MIRRF システムの機能

MIRRF システムの主要な機能について例を挙げて説明する。MIRRF システムには MIRRF 状態と DO 状態とがある。MIRRF 状態はこのシステムの初期状態であり、この状態でコマンドが投入されるとプログラムファイルが新規に作成される。DO 状態ではコマンドやデータの記録は現在使用中のプログラムファイルにとられる。両者はその他の点では全く同一である。状態の遷移は図 4 に示すようにコマンドによって行われる。図において READY は TSS における READY 状態であり、COMMAND は以下に説明するような END 以外の MIRRF システムのコマンドである。MIRRF 状態で開かれたプログラムファイルは END コマンドで再び MIRRF 状態に戻るときに閉じられて、プログラムファイル領域のトップに置かれる。

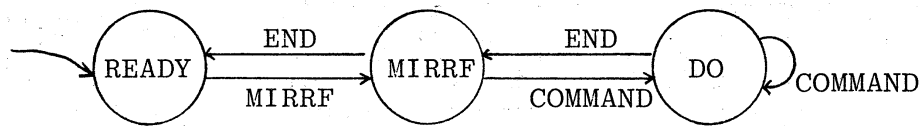


図 4. MIRRF システムにおける状態遷移

END MIRRF システムから離れるためのコマンドである。DO 状態にあるときは、上図に見られるように END を 2 回投入しなければならない。

KEYIN キーボードから入力されるデータを作業域のトップに作業用ファイルとして作成する。プロンプトとして ">" が出される。入力の終了は CR ; ; ; で指示する。

LOAD MIRRF ファイルや TEDIT ファイル、外部データ・セットを MIRRF のトップの作業用ファイルに持つてくるためのコマンドである。FILE:= に続いて必要なファイル名を指定する。指定されたファイルがパス・ナンバーをもっているときには、PUSS NUMBER:= に従ってそれを与えなければならない。なお、キーボードは KB という名前のファイルと見做される。従って、FILE:=KB とすると、KEYIN コマンドと同じ働きをすることになる。

SAVE 作業域のトップにある作業用ファイルの内容を MIRRF ファイルや TEDIT ファイル、外部ファイルへセーブするためのコマンドである。ファイル名やデータ・セットの指定は LOAD の場合と同様である。端末の画面は VD という名前のファイルと見做される。

LOOK トップにある作業用ファイルの初めの 20 行の表示を行う。 CR ; ; ; に次に 20 行ずつの表示がくり返される。終了は E CR で指定する。

LIST MIRRF システムがもっているファイル名やファイルの内容をリストするためのコマンドである。ファイル名は $w_1 \cdot w_2 \cdots w_n$ の形で階層的に定義できる。完全なファイル名が指定されたときには、そのファイルの内容がリストされる。不完全な指定の場合には、その指定された名前を親とするようなファ

イル名がリストされる。従って、ファイル名を空語 (CR だけ) で指定すると MIRRF システムが現在もっているトップ・レベルのファイル名がリストされる。

RENAME MIRRF ファイルの実体に新しい名前を追加するためのコマンドである。ファイル名やパス・ナンバーの指定は上記のコマンドと同様である。以下ファイル名等の指定についての説明は省略する。

APPEND 指定された MIRRF ファイルにトップにある作業用ファイルを追加するためのコマンドである。

PULL 作業用ファイルやプログラム・ファイル, TEXTIR の作業用ファイルの指定されたものをトップにもってくるためのコマンドである。

PULL

FILE:=WF

FILE POSITION:=3

のように使用する。ファイル名は WF (作業用ファイル) の他に PF (プログラム・ファイル), TF (TEXTIR 作業用ファイル) を指定できる。上の例では、MIRRF の作業用ファイルの 4 番目のものがトップにもってこられる。

FILE POSITION:=0 とすると、トップにあるファイルが末尾におかれる。作業域がいっぱいになったときに先ず最初に消去されることになる。

WFMV 作業用ファイルのトップにあるものを、TEXTIR の作業用ファイルのトップに移動させる。

PRMV プログラム領域のトップにあるファイルを作業域のトップに移動させる。

TXMV TEXTIR の作業域のトップのファイルを MIRRF の作業域のトップに移動させる。

LENGTH 指定されたファイルの長さを表示する。ファイル名 WF でトップにある作業用ファイルの長さが分る。

REST MIRRF のファイル領域や作業用ファイルの使用可能なセクタ数を表示する。

REPLACE トップにある作業用ファイルを対象にして、指定された記号列と他の指定された記号列ですべて置き換えたファイルを作り、トップに置く。以下のように使用する：

REPLACE

NEW LINE:=!

A1:=x1

B1:=y1

A2:=x2

B2:=y2

⋮

⋮

$\underline{A_n} := x_n$
 $\underline{B_n} := y_n$
 $\underline{A_{n+1}} := \textcircled{CR}$

この場合、トップにある作業用ファイルに含まれているすべての x_i が y_i で置き換えられたファイルが作られる。なお $\{x_1, \dots, x_n\}$ をキーワードをみて、パターン・マッチング・マシンを構成し、それがファイル上を左から右へ走ることになる。

TEXTIR 一方向の逐字サーチシステムである。詳細については、[3]を参照していただくことにして、ここでは使用例だけを挙げておく。下の例において、:=以下の部分は利用者による入力である。区切語 (delimiters) は文字列であり、キーワード (keyword) は文字列又は $x \dots y \dots z$ のようなトリプルドット...を含む文字列である。論理式 (logical formula) はキーワード変数

DO:TEXTIR
 <<TEXTIR/E>>

0...RESET, 1...NORESET:=
 NEW LINE:=!
 0...LIST, 1...NOLIST OF TABLE:=0

OUTPUT DELIMITERS
 D1:=#
 D2:=

INPUT DELIMITERS
 D1:=

ADDITIONAL INPUT DELIMITERS
 D1:=

KEYWORDS
 A1:=TURING
 A2:=COMPLEXITIES
 A3:=ONE-WAY
 A4:=

LOGICAL FORMULAE
 Z1:=A1.A2
 Z2:=A3
 Z3:=
 FILE:=MIYANO.PAPERS

RETRIEVED TEXTS
 TOTAL = 6
 QUESTION 1 (Z 1) = 0
 QUESTION 2 (Z 2) = 6

FILE:=
 0...LIST, 1...NOLIST OF RESULTS:=0

(A1, A2, ...), 式変数(Z1, Z2, ...), 整数, ε-変数(E), か, こ [,] , (or), ・ (and), ^ (not), +, -, *, /, <, >, =, <> (≠) を用いて作られるものである。TEXTIR による検索結果は、先に述べたように、ヒットした質問番号をコード化した整数, ファイルの機番, レコードの番号, セクタの番号, 開始文字位置, 長さの組の形で、TEXTIR 作業用域のトップにあるファイルに記録される。

TXEP これは TEXTIR による検索結果を編集して新しいファイルを作るためのコマンドで、下の例のように使う。QUESTION:= は TEXTIR での質問の番号

```
DO:TXED
QUESTION:=2
(QUESTION 2=      6 TEXTS)
NEW LINE:=!
OUTPUT DELIMITER:=NO.
NUMBERING (Y OR N):=Y
SORT BY:=AY
```

を指定する。次の () 内はその質問に対してヒットしたテキストの個数である。NEW LINE:= は次にある出力区切り語 (output delimiter) の中で return キーの代わりに使う文字である。出力区切り語は任意の文字列である。番号付けを行いたいときは、NUMBERING (Y OR N):= で Y を選ぶ。最後の SORT BY:= は TEXTIR で使う標準的な文献テキストに対して有効なもので、C (文献コード), A (著者), Y (年代), J (誌名) の中から 0 個以上選ぶ指定する。上の例では、著者順にソートして、同一著者に関して (論文の) 発行年順にソートされることになる。こうした標準的な文献以外のテキストに対しては何も指定しなければよい。なお、ソーティングにはファイル名等の管理に使った B-tree を使っている。

? Help がある。MIRRF システムが入力要求をした任意の時点で使用できる。このシステムでは help に対する案内文は 1 つの MIRRF ファイルとして扱われる。従って、利用者は必要に応じて内容の追加・変更を行うことができる。

TEDIT MIRRF のテキスト・エディタにはいる。([4] を参照)

以上が現在の MIRRF システムの基本的なコマンドである。他に MIRRF ファイルや作業用ファイルを初期化するコマンド等があるが、ここでは省略する。

5. 交信記録とその利用

先に述べたように、MIRRF システムでは、MIRRF 状態から次の MIRRF 状態までのキーボードからの入力はコマンドやデータも全てプログラムファイルに記録される。これは通常のファイルと全く同様に扱われるだけでなく、一種のコマンド・プロシジャとして利用できる。MIRRF ファイルにセーブされたこ

した記録は、入力要求が行われた行意の時点で、

/FILE NAME

//FILE NAME

の形で呼び出して使用できる。前者に対しては、全ての作業がキーボードから直接指示した形で実行される。即ち、コマンドやデータも全て表示され、実行される。後者に対しては、こうした表示は行われぬ。

ファイル名としてキーボード (KB) を指定することもでき、ファイルの内容に再び上記の形の指示を含むことができるので、利用者は簡単に自分用のシステムを MIRRF の中に持つことができよう。

このような通信の記録は当然エディティングの対象になる。したがって、それをプログラムとして使う場合には、一種の構文のチェックが必要である。このシステムでは、こうした記録をプロンプト付きでとっているので、構文のチェックは極めて簡単に (FORTRAN の 1 ステートメント) 実行中に行われる。

6. おわりに

エディタを除いて、現在の MIRRF システムの機能について簡単に説明した。今後、いわゆる PRF マネージメント・システムとして整備し、九大センターを通じて一般の利用者に公開していくつもりである。

安全性の問題 作業用ファイルや TEDIT ファイルに関しては、ファイルの安全性について或る程度注意を払った。しかし、MIRRF ファイルに関しては、直接アクセス・ファイルを使っているために、計算機の crash 等に対して現在は無防備である。この問題もデータの共有の問題と併せて解決したので、本稿で述べた範囲内の MIRRF システムを SIGMA (E) という名前で近く九大大型計算機センターから公開することになっている。

PRF → UDL → CDL の構図 公開される SIGMA システムは、少なくとも文献データに関しては、PRF から UDL までカバーできるものと考えている。CDL としては、現在 FAIRS によって INSPEC がアクセスできるので、そちらの充実を期待したい。また MIRRF システムも CDL のデータの蓄積に寄与できるものにしていきたい。

参考文献

1. Arikawa, S. and Kitagawa, T.: Multistage Information Retrieval System Based upon Researcher Files, Res. Rept. Res. Inst. Fund. Inform. Sci., Kyushu Univ., No. 51 (1975), 1-34.

2. Arikawa, S., Kanō, S., Kitagawa, T. and Takeya, S.: Organization and Use of Private Researcher Files in Scientific Research Works, Research on Scientific Information System in Japan (1980).
3. Arikawa, S.: A One-way Sequential Search System and Its Applicability to Medical Information Processing, Proc. HISCC-13 (1980).
4. 篠原武: MIRRF システムにおけるエディタについて, 「知識の表現とそれを利用する情報検索システムの研究」報告書 (1981), 247-258.