

## 第3章 入力検索システム

### 3-1 システム設計の前提

入力検索システムNRDF2は、先に開発された同様のシステムNRDF1に基いて設計された。NRDF1からNRDF2への仕様の改訂において大きく変化した点は、

- ① システム内部でのデータの表現
- ② データ間の関連づけ

の2点である。いずれもNRDF2への移行において、基礎になる概念を単純化することによって多様な入力データへの対応をはかることを目指した。システム内部でのデータ表現に関しては、データの *compilation* と *interpretation* という概念を導入し、データのとりあつかい方によって異なる表現のしかたをする。これによってファイル内ではいくつかの異なる様相をもつ部分が現われる。

またデータ間の関連づけのために集合概念が利用された。このことによって、核データの多様性の一部を吸収すると同時にファイルの自己発展のための基礎を与えた。

#### 3-1-1 compiler から interpreter へ

NRDF1では *compiler* 方式をとっていた。(図3-1-1・A)

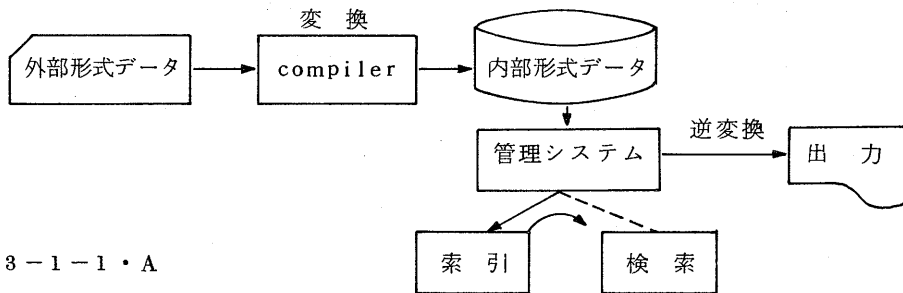


図3-1-1・A

*compiler* 方式とは、外部形式の入力データをすべて内部形式に変換し内部形式データの中からその一部を索引のための情報としてとり出して登録するものである。この場合、検索を行なって結果を出力するためには内部形式から外部形式への逆変換が必要である。*compiler* は外部形式、内部形式、Key項目、形式変換規則に関する情報をもつ必要がある。データを内部形式に変換すればシステムにとってデータ処理は容易となり、したがって検索も容易であるが、主たる検索の手がかりとしない情報をも一様に内部形式に変えることは、2次記憶域の使用効率と形式変換のプロセスを必要とする点でマイナスである。一般にある規則でcode化された情報を *explicit* なデータ構造の形に展開することは、個別のデータ構造の中に埋め込まれた規則情報の分だけ記憶域を多くとる。NRDF1はこの欠点を明らかにした。

これに対しNRDF2では interpreter方式をとることにした。

interpreter方式とは次のようなものである。(図3-1-1・B)

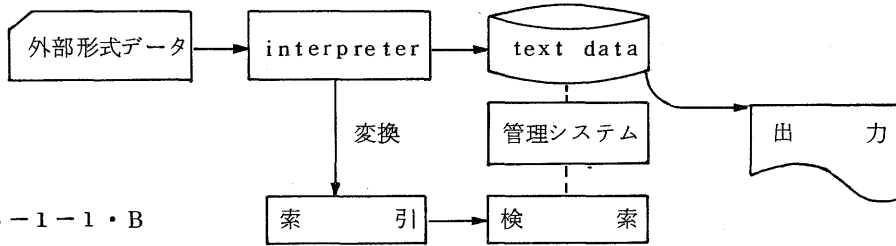


図3-1-1・B

一定の外部形式で与えられたデータを interpreterが読み込み、与えられた文法にしたがって内容を解釈し、指定されたkey項目に関する情報のみを抽出しこれを<key項目-key値>の対の形をした内部形式に変換して項目ごとの索引(inverted table)へ登録する。また読み込んだデータは文法チェックのあと変換なしに text dataとしてファイルに追加登録する。interpreterのもつ情報は外部形式(文法)とkey項目の種類及び索引のための内部形式である。

interpreter方式は、必要最小限の情報のみを「compile」して検索のための内部形式に変えるため、フォーマットの異なるデータの受理に際して、key情報の抽出さえできればよく、compiler方式のようにすべてのデータを内部形式に変換するという強い条件がつかず抵抗が少なくなる。また必要に応じて部分的に他の内部形式へ変換していく自由度を残している。これは実行時に text dataの解釈をしながら判断したりデータ操作をしたりすることを意味し、実行速度の低下を無視すれば、interpreterの媒介によって、すでに内部形式で情報が存在している compiler方式の場合のような処理手順と表面的には同じ手続きを実行することになる。

interpreter方式と compiler方式をまとめて比較すると次のようになる。

	interpreter方式	compiler方式
使用する2次記憶域	少ない	多い
情報操作	遅い	早い
key項目以外の検索	遅い	早い
処理プログラムの大きさ	小さい	大きい
異種のフォーマットに対して	抵抗が小さい	抵抗が大きい

結局NRDF2で格納されている情報は、comment文として補足的に使う free text(自然語)と、syntaxとcodeで規制された coded informationおよびその一部をexplicitなデータ構造に展開した compiled informationの3段階の様相をもっている。

free textはman-orientedな情報で、現実世界のできごとを記述するための強力な表現力をもっている。しかしまた反面、一意性や情報密度が犠牲にされている。coded informationは人工語であるという点で人間とのよいinterfaceをもち、同時に定義が厳密であるという点で機械にもよくなじむ。compiled informationはcoded informationにおいて文法によって一次元的なひろがりの中に圧縮されている「意味」をtableやlistなどのデータ構造の形に変換したもので、機械によるデータ操作のための最適化が可能である。これはmachine-orientedな情報である。

free text型の情報(あるいは学術論文)からcoded informationへの変換は、研究者が入力データをいかに作成するかという人間の側の技術的問題であり、coded informationの型式すなわち入力フォーマットを定める場合にはこの変換による情報洩れや人工語としての表現力の問題が重要となる。これに対してcoded informationからcompiled informationへの変換はシステム開発の技術上の核心である。coded informationからどの範囲の、どの位の深さの内容をどのようなデータ構造に変えるのかということが、システムの性格を決定する。interpreter方式とは、この問題に対してcoded informationの全部ではなく、keyとなるべき一部の情報を索引構造の中に転写しようという一つの方針をとることを示すものである。したがってkey項目の選択、索引構造の詳細については別途これを考慮する必要がある。compiled informationの型式の設定は、データ構造としてシステム内での情報検索戦略の裏付けとなるべきもので、質問文からはじまって目的データに至るまでの情報から情報への経路を準備することである。

### 3-1-2 データ間の関連とdata set

核反応データ、とりわけ荷電粒子核反応データは、反応様式の多様性を反映して多くの種類のデータをその中に含む。したがってこれを収容するファイルには適応性と記述能力が特に要求される。また、多種多様なデータでありながら、これを計算機独自の表現によらずに、しかも簡単にとりあつかえる形で保持するということは、研究活動の自然な延長としての情報活動を考える時、学術情報にとって基本的な要求である。しかし核反応データの場合、それが多様であっても、単位となる情報はいつもひとつのデータテーブルであって文献情報などと同じように比較的是っきりしている。異なる点は、実際上の単位情報はこのデータテーブルに各種の(非数値的な)情報を加えたもので、この付加情報とデータテーブルの、場合によってはネットワーク的な、結びつきが一種の情報構造を与えているということである。この点を糸口としてNRDF2のファイル構成が行なわれる。

一回の核反応実験から得られるデータは、いくつかのデータテーブルの形にまとめられる。これらのデータテーブルおよび関連する情報の集まりをdata streamと呼ぶ。data stream

はいくつかの data set の集合である。data set とはひとつのデータテーブルとそれに関連する情報の集合であって、その要素は section と呼ばれる。section には、データテーブル本体の DATA section、書誌情報をもつ BIB section および実験方法に関する情報をもつ EXP section の 3 種類がある。data stream の中で複数個の data set が section を共有してもよい。たとえば、BIB section は通常 1 data stream 中にひとつしかなく、data stream 中のすべての data set に共有される。入力にあたって各 section はその種別と帰属する data set の識別番号を先頭に表示する。section はさらに sub section の集合である。sub section には description sub section と data sub section の 2 種類があり、各々その記述の format を異にしている。

description sub section は、

<項目>=<値>; 又は <項目>=( <値>, <値>, …………… );

という形式で非数値的情報を記録している。この形式を statement という。いくつかの statement を (と) でくくるとことや、associator (logical pointer) を用いて相互に関連づけをすることができる。また、/\*と\*/ではさまれた free text を comment として記録することもできる。

data sub section は column free な数値データテーブルで、各数値は可変長であって、空白で区切られていけばよい。

このように、データファイルとは表(データテーブル)の集まりであって、検索はその中から希望する一部の表を取り出すことにほかならない。(図 3-1-2・A) しかし厳密には数値だけからなる表では、その表の意味内容が何であるのかを示すという自己説明性に欠けるため、表に、書誌事項などの付加情報を合わせてひとつの集合としたものを data set としてとりあつかいの単位にする。したがって同じ条件の下で実験が行なわれた複数個の表は、実験条件の section、EXP section を互いに共有する複数の data set となる。これは、intersection をもちうる集合の集まりとしてデータファイルをとらえるものである。このようなデータファイルでは、異なる表の間の関連は、それらの表を核とする data set がその外縁に何らかの共通部分(intersection)をもつという形で表現される。(図 3-1-2・B)

図 3-1-2・A

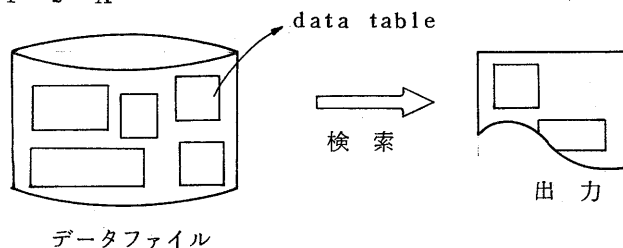
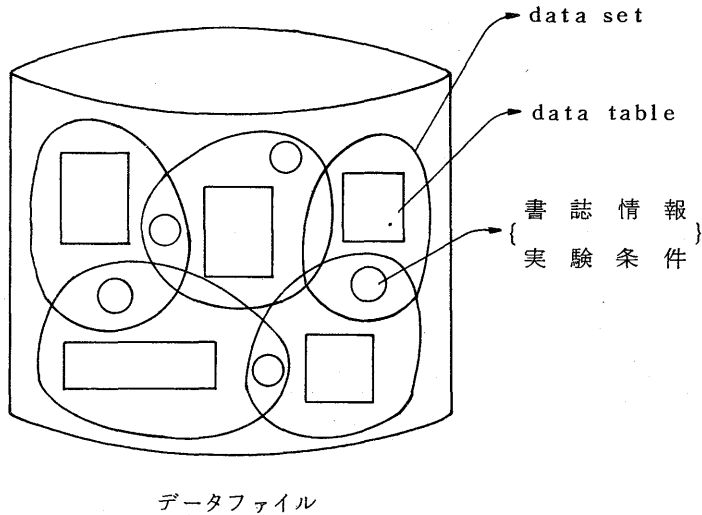


図 3-1-2・B

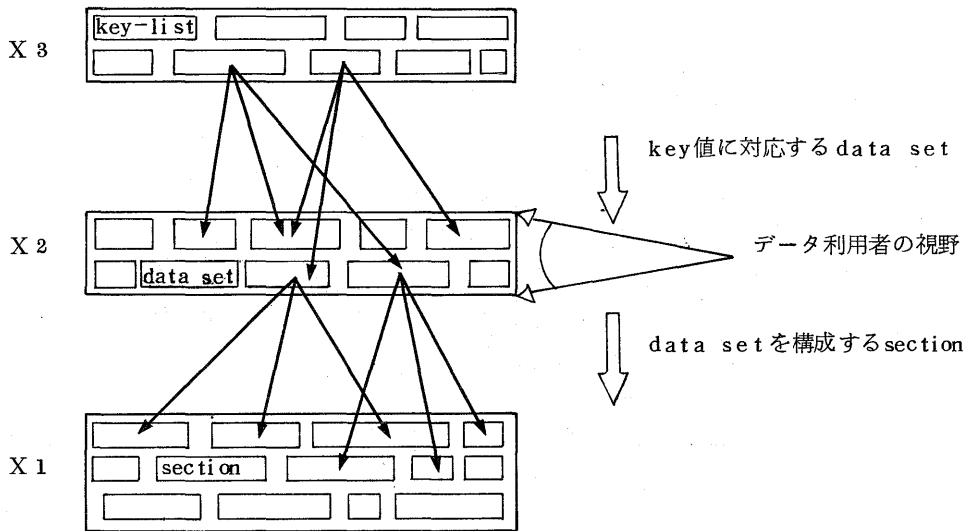


### 3-2 ファイルの構成

NRDF2では、X1、X2、X3という3種のファイルをもっている。X1は、入力された情報を section 単位で、input image のまま格納したもので、可変長レコードである。X2は、data set の集合である。実際には、各 data set が section の集合であるため、対応する X1 の要素 ( section ) への pointer の組でひとつの data set を論理的に表わしている。data set は検索の target となるものであって、データ利用者は、data set の全集合であるこの X2 ファイルのみを意識すればよい。X3はX2を検索するための各種索引である。索引は、data set 中の非数値的情報の中から、システムの定める項目 ( key 項目 ) についての対応する値にもとづいて作られる。X3を構成する要素は、key 値とそれに対応する X2 の部分集合である。この部分集合は、実際には X2 への pointer の組で表わされている。X3はX2をその内容にもとづいて分類した情報という意味で、X2に対する部分転置ファイルである。同様に X2 は X1 に対する部分転置ファイルとしてとらえることができる。X1、X2、X3 の関係を図 3-2 に示す。

すなわち、section の集合 X1 を基底集合とした部分集合族が X2 であり、X2 を基底集合とした部分集合族が索引 X3 である。データ利用者の論理的な視野を規定する X2 を中心にとらえれば、X1 はデータの実体を収容しており、X3 はデータの分類を与えているのである。このようなファイル構成に対してシステムは、データ入力時には、入力されたデータを適宜 X1、X2、X3 に振り分ける。またデータ検索時には、データ利用者の X2 上での探索のために、必要な時点でシステムが自動的に X1、X3 を参照して、要求されたデータを取り出す。

図 3-2



### 3-3 検 索

N R D F 2において検索操作とは、希望する条件を満たすデータテーブルを探し、とり出すことである。ここでデータテーブルとは実際には、表とその表に関する非数値的情報をワンセットにした data set であって、この data set の集合であるファイル X 2 が検索の主な舞台となる。

検索の操作はしたがって X 2 からひとつの部分集合をとりだすことであり、このとき望みの部分集合を得るために利用されるのが X 3 である。また最終的にとり出すべきデータ本体は X 1 に格納されている。X 3 は検索する者にとってはファイルの分解能を与える役割をもっている。すなわち、大雑把な分類しかもたない index では粗っぽい検索しかできず、詳細にわたる情報にもとづく index は精密な検索操作を保障する。

ファイルに対する質問は：

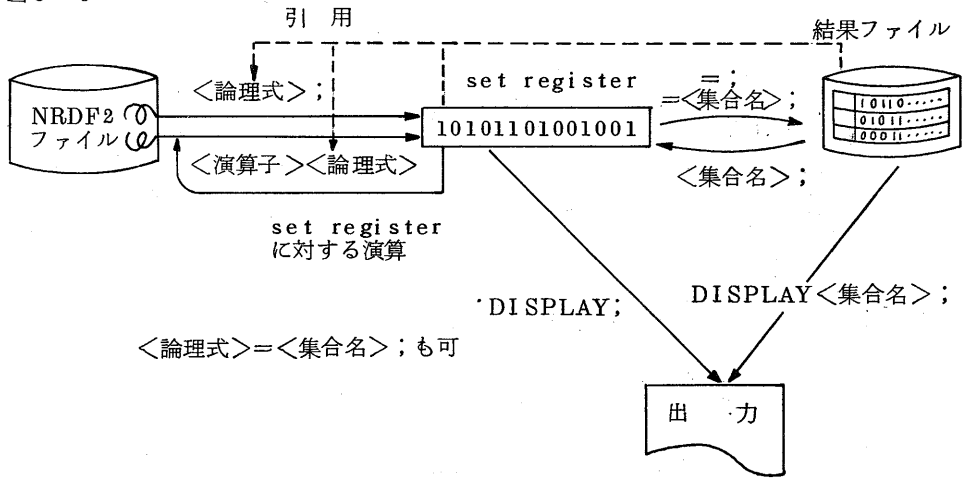
<論理式>=<集合名>;

の形式で与えられ、<論理式>は ( A T H = A . B C D ) のように key 項目とその指定する値を関係演算子 = , > = etc でつないだものを要素として任意の論理条件を組合わせることができる。

<集合名>は<論理式>で指定された条件にもとづいて検索を行ない、その結果として得られる data set のひとつの集合に付ける名前、この名前の下に一時記憶域に結果が格納される。= や <集合名>は省略可能で各々省略時解釈が適用される。

この質問から得られる結果は一時記憶域とは別に常に、set register と呼ばれる特別の作業領域に保持され、これを媒介にデータファイルと検索結果ファイルとの間の情報の受けわたしが行なわれる。( 図 3-3 )

図 3-3

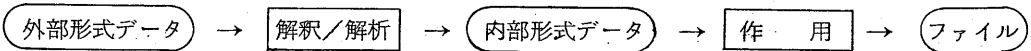


なお検索操作の間、論理式などによってとり出されたデータの集合は常にデータの address の集合であってデータそのものではない。このデータ本体は、最終的な出力命令 DISPLAY が発せられた時、はじめて X2 と X3 を参照してとり出すことができる。また、論理式中で与えられる AND、OR 等の演算はそのつど集合演算としてとりあつかわれる。すなわち、set register の上ではもっぱら集合演算が行なわれその結果として集合が残る。

### 3-4 システムの構成

NRDF2 におけるシステムの基本動作とは、図 3-4・A のように、人間側から与えられた外部形式データを受けとり、これを解釈・解析し内部形式データに変換し、この内部形式データによってファイルへの各種の作用を及ぼすというものである。

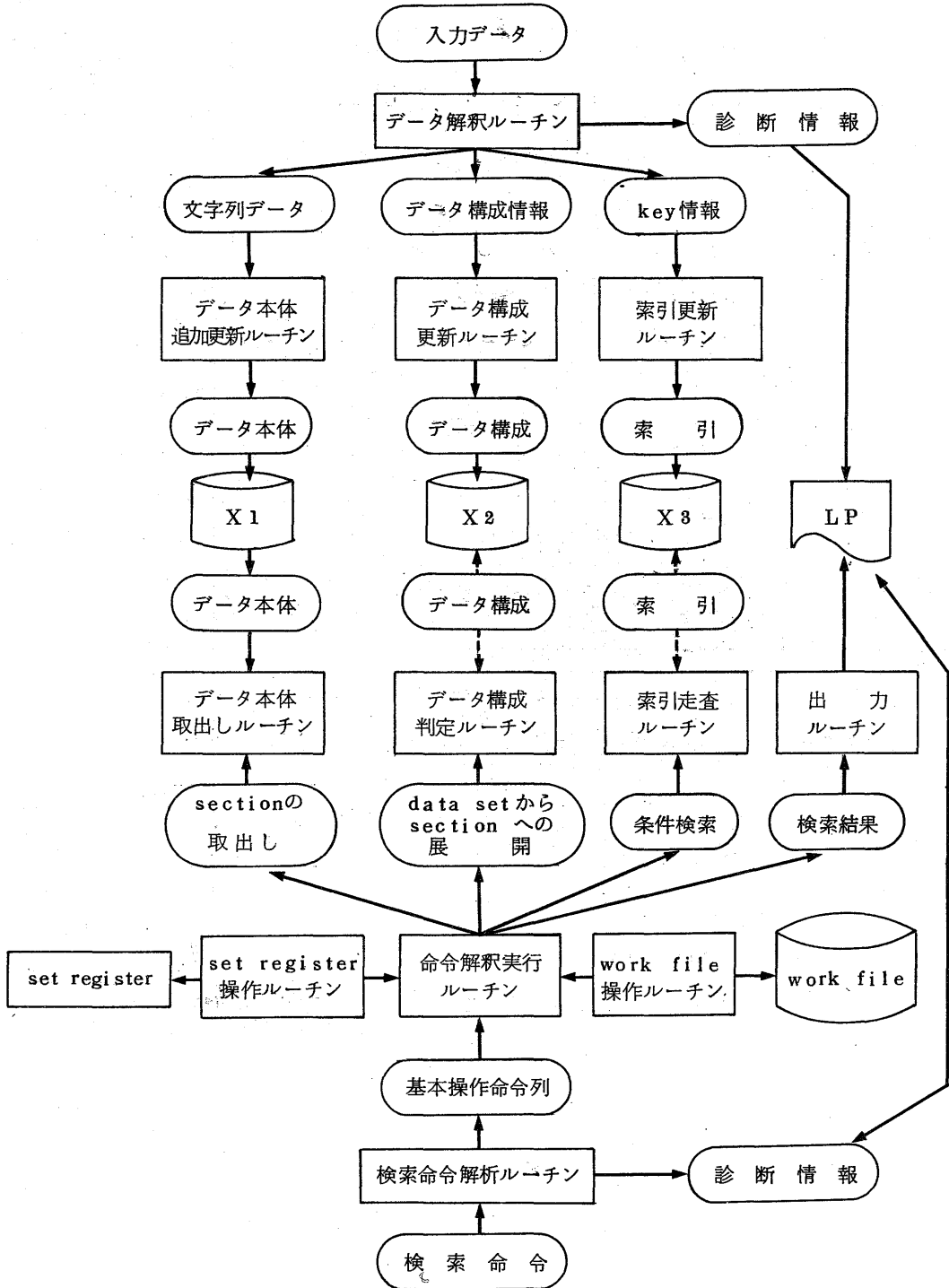
図 3-4・A



ここで、データ入力時においてはこの基本的動作は、入力データの解釈と分類、各種内部データの生成、さらに X1、X2、X3 の各種ファイルの更新となる。また、データの検索時には、検索命令の解釈と基本命令への変換、そしてその基本命令の実行によるファイルの参照・とり出しとなる。

図 3-4・B にデータ入力およびデータ検索のためのシステム構成と、主なデータの流れを示す。

図 3-4 \* B システム構成



### 3-4-1 データ入力

データ入力時にはまず、データ解釈ルーチンが起動され、computer readable formatで書かれた入力データを読み込む。この時に、文法上のcheckを行ない、もしデータに文法エラーがあれば、それについての診断情報を出力する。また、システム内で指定されたkey項目リストを参照しながら、key項目についての値の記述があれば、それを抽出してkey情報リストを生成する。さらに、入力データの各sectionごとに表わされたdata setの構成についての情報を抜き出して、sectionとdata setとの間のn:mの(多対多の)対応関係情報リストを作る。最後にシステムは、入力データをsectionごとにちぎって分解し、データ本体の登録にそなえる。

次に解釈ルーチンが生成した各種の情報にもとづいて各ファイルの更新に入る。

まず、抽出されたkey情報のリストにもとづいてこのリストの分だけ索引に追加すべく、索引更新ルーチンを起動して索引ファイルX3の更新を行なう。データ構成情報は各data setがどのsectionから構成されているかを表わしており、これにもとづいて、ファイルX2を更新する。さらに、sectionごとにちぎられた、文字列のままのデータ本体が、ファイルX1に追加される。

### 3-4-2 データ検索

検索命令が与えられた時には、検索命令解析ルーチンによって、命令の形式のcheckのあと、ひとつの検索命令をいくつかの基本的な操作命令の列に変換する。この時の基本操作命令には、key値指定に対する索引走査、work file操作、set register操作、data setの展開、および出力のためのデータ本体参照操作がある。

これらの基本操作命令列を、命令解釈実行ルーチンが、順次読み込んで実行する。索引走査の際にはファイルX3が調べられ、指定された条件にかなうdata setの集合(addressの集合)がとり出される。work file、set registerは、X3からとり出されたaddress集合の一時保存のための場所として使われ、各集合は名前をつけて参照される。

検索結果の最終的な出力を行なうために、各data setがどのsectionから構成されているのかを調べる必要がある。このために、データ構成情報の集まりであるファイルX2が参照され、data setごとに構成sectionのaddressが得られる。さらに、このaddressをもとにして、データ本体が文字のまま格納されているファイルX1の中から各sectionをとり出して、出力にそなえる。こうして、X3→X2→X1と順次たどることによって、指定した条件を満たすdata setをLine-Printer上に文字の形で出力することができるようになる。またwork file、set registerは、最終出力に至る途中でのさまざまな途中結果をメモがわりにとっておくために利用される。