

テキスト語彙処理プログラム L E X の開発について：概要と論理

松 尾 雅 嗣

広島大学平和科学研究センター

ON THE DEVELOPMENT OF A PROGRAM PACKAGE FOR LEXICAL ANALYSES OF A TEXT: OUTLINE AND LOGIC

Masatsugu MATSUO

Institute for Peace Science, Hiroshima University

SUMMARY

LEX is a program package now being developed for various analyses of a text. LEX as a system has the following characteristics. First, it is intended for the use of researchers in social sciences and humanities who often deal with linguistic data but are not familiar with a computer-assisted data processing. Secondly, the system deals solely with characters and character strings, and not with numeric data. Thirdly, it is not a single-purpose system, but a general-purpose system, in that it is provided with a variety of ways of processing data.

The basic unit of processing by LEX is called “word”, which is defined as a string of characters. In addition, LEX can also deal with substrings of “words” and strings of “words” (called lexical strings). Moreover, LEX can perform data-selection operations at the level of substring, “word” and lexical string. These characteristics enable the user to obtain a rich variety of results.

The user of LEX communicates with the system by means of what might be called the LEX language. Its grammar is very limited and simple, and, roughly speaking, can be described by a few context-sensitive rewriting-rules.

The system itself consists of two components: one interpreting input LEX sentences, and the other performing data processings proper, such as compilation of frequency table, alphabetical list, cooccurrence table, concordance and so on.

0

社会科学，人文科学において，テキスト中の単語の処理，しかも相当量の単語の処理，が必要とされる場合が少なくない。本稿は，コンピューターによるテキスト，文書等の語彙処理を目的として現在開発中のプログラム・パッケージ，LEX (program package for lexical analyses of a text) の概要の報告である。

LEXは，テキストの構成要素である単語について，基礎的な諸統計，表，索引等を作成する機能をもった，一貫したシステムたることを目標とするが，このシステムの特徴として次の諸点を挙げることができよう。

まず第1に，LEXが想定する利用者は，一般の，数値計算を主目的とするコンピューター利用者ではなく，言語データを扱うことの多い，しかもデータのコンピューター処理にまったく馴染がないか不慣れな人文科学系，社会科学系の研究者である。このため，利用者の便宜を考えてファイルの取扱い，利用者の書くプログラムの記法等をできるかぎり単純にすることが意図されている。

第2に，LEXはその名称が示すように，もっぱら言語データ（あるいは非数値データ）を処理の対象とするシステムである。しかも，従来の言語データを処理するプログラムが，索引，コンコードダンス，頻度といった単一の出力を得るために作成された単一目的のものであったのに対し，LEXは，このような個々のプログラムをパッケージ化することにより，多様な出力が可能で，汎用性のあるものとなっている。更に処理をより多様なものにするために，LEXには，データの選別機能が備えられている。同様に処理単位に関しても選択機能を有する。す

なわち、L E Xでは単語連鎖と呼ばれる、単語より大きな単位、あるいは単語の構成要素である部分文字列、を単位としての処理も利用者の選択により可能である。このようなシステムの有する処理の多様性に加え、利用者は自分の作成したプログラム（実行用サブルーチン）を簡単にシステムに付け加えることもできる。

L E Xは現在広島大学計算センターのH I T A C 8700/8800のオペレーティング・システムのもとで稼動しているが、以下、この、言わばH I T A C版のL E X について概要を述べる。（ソース・プログラム（P L / I ），あるいはオブジェクト・プログラムを利用することにより、他機種、他機関での利用も十分に可能である。）

利用者の立場から見た場合、L E Xによるテキスト処理の流れは次の図1に示すようになる。

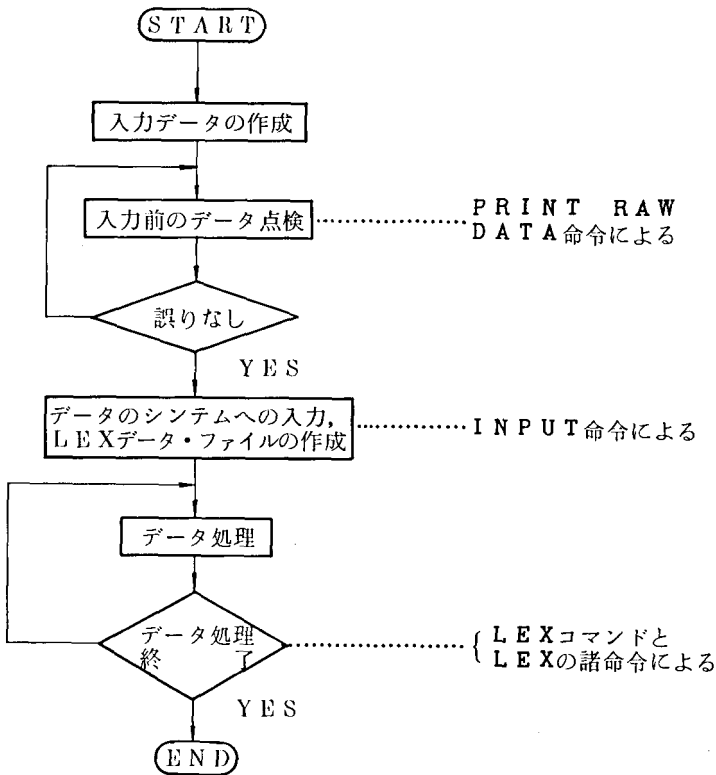


図1 利用者の立場から見たL E Xによるテキスト処理の流れ

利用者はまず現実のテキスト、即ち、利用者にとって有意味な単語の有意味な集合であるテキストをLEXに入力可能な形式に、例えばカードにパンチするなどの方法により、変換する。これがLEXに入力される素データ(テキスト)である。次のステップで素データ中のパンチミスなどの誤りを点検、修正する。

LEXにはこのために、`PRINT RAW DATA`という命令が用意されており、この命令によって、システムに入力前の素データをそのままの形で(即ちカード・イメージのまま)印刷することができる。(この場合の素データには後述するLEXのステートメントが含まれていてもよい。付録1参照。)データの点検が終了すれば、素データは`INPUT`命令を用いてシステムに入力される。システムは入力されたデータから後の処理に必要なデータ・ファイル(LEXデータ・ファイル)を作成する。LEX標準ファイルを用いれば後の処理では素データをその都度入力する必要はないし、非標準ファイルの場合でもごく簡単なファイル定義コマンド(OS7の場合DTFコマンド)を用意するだけでよい。

データが入力されれば、利用者は、マクロコマンドであるLEXコマンドに、必要なLEXの命令を付けてシステムにさまざまな処理を行わせる。言うまでもなく、データの作成、点検、入力はこのための準備段階であり、この部分が利用者にとっての最も(あるいは唯一の)重要なステップである。データの選別機能処理単位の選択機能を必要に応じて利用しつつ、このステップを繰返すことにより、多様な処理結果を得ることができる。

以下第1節では主として利用者の立場からLEXのデータ、ファイル、構文法といった概要について述べ、第2節では、システムがこれをどう解釈し実行するかという観点からシステムのロジックについて述べる。頻度表、索引といった個々の処理の実例を含めた利用法についてはスペースの関係で別稿に譲ることとし、ここでは触れない。

1. LEXの概要

1.1 データの単位

LEXでは処理の基本単位を単語と称する。単語は、アルファベット、数字、片仮名、それに通常コンピューターで処理可能な特殊記号のうち、単語の区切り

及び単語連鎖の区切りとして用いられるものを除く特殊記号の、いずれかから成る16文字以内の記号列と定義される。従って、L E Xの単語は、日常的な意味あるいは言語学的な意味での単語である必要はなく、利用者にとって有意味な（少なくとも分析の対象とするに足る）記号列であればよい。L E Xで言うテキストはこの意味での単語の（少なくとも利用者にとっては有意味な）集合、正確には線状の連鎖と定義される。入力データ、素データ、データ・ファイル等の表現におけるデータは、このテキストの意味である。

この単語とテキストの概念的に中間の単位として、単語連鎖という単位も認められる。単語連鎖は、テキストの中の連続した単語から成る単位である。（但し、単語連鎖が唯一ひとつの単語から構成されることもありうる。）単語連鎖は最大4種類まで指定することが許されるが、この指定により、後の処理におけるデータの選別、処理単位の選択の幅が広がる。単語連鎖を指定したときには、境界標識を用いて個々の単語連鎖の境界を確定するとともに、個々の単語連鎖につける識別値を与えなければならない。例外は後に述べる行形式入力の場合だけである。識別値には数値型のS、F、Dの3つのモードと文字型のAモードの4つのモードがあり、個々の単語連鎖ごとに識別値を与える必要があるのはAモードの場合だけである。Sモードの場合は利用者の与える初期値（省略時は1）から始まる一連番号が識別値として個々の単語連鎖に与えられ、Fモードであれば、利用者が必要に応じて初期値を与えるたびにその初期値から始まる一連番号が識別値となる、Dモードであれば、当該の単語連鎖よりひとつ大きい単語連鎖が始まるたびに1から始まる一連番号を識別値とする。例えば、ページが変わるたびに行番号を新しく1から始めるような場合がそれである。このDモードの場合には、2種の単語連鎖間に、一方が他方を要素とする集合であるという関係がなくてはならないが、Dモードの場合を除きこのような制限は課せられない。またDモードの識別値から分るように、個々の単語連鎖の識別値は、個々の単語連鎖、あるいはそれに属する単語群をユニークに識別できるものである必要はない。また、複数の単語連鎖があるとき、その識別値の組合せがユニークである必要もない。個々の単語、あるいは単語連鎖はシステムがデータ・ファイル作成時に個々の単語（場合によってはこれに加えて行）に付ける一連の単語番号（WSQと言う）

(もしくは行番号(L S Q))によってユニークに識別でき、利用者はWS QあるいはL S Qを参照、引用できるからである。

以上のことからL E Xに入力されるデータはテキスト本体(データ本体)と単語連鎖の識別標識から構成されるということになる。L E Xでは、このほかデータ処理には何の関係もない註釈行の挿入が許される。註釈行は言うまでもなく、後のデータ処理や識別値、L S Q、WS Qなどの付与にはまったく関係しないが、行形式ファイルには記録されるので、入力後にも参照することができる。

1.2 入力データの処理

データ本体(テキスト)と識別標識と註釈行から成る入力データは、どのようなデータ・ファイルが作成されるかという意味での入力データの処理方式によって2種類の入力形式に分けられる。行形式入力と単語形式入力それぞれである。このふたつの入力形式の相違は、データ本体ではなくて、むしろ単語連鎖の境界標識や註釈行の処理と作成されるファイルの種類の相違である。

入力データは、すべてカード1枚を単位として(固定長80バイトのレコードとして)処理される。またデータ本体(テキスト)、識別標識、註釈行はそれぞれ別のカード(レコード)でなくてはならない。データ本体はひとつ以上の単語を(2つ以上であれば必要に応じてひとつ以上の単語の区切り記号によって区切られた単語を)パンチしたカード(もしくは記録したレコード)である。

行形式入力は、このカード1枚(1レコード)(行と称する)を1つの単語連鎖として処理することを要求する入力形式である。例えば印刷されたページの1行をそのままカード1枚にパンチして入力し、印刷された1行を単語連鎖として処理する方式である。この入力形式では、行と行番号と単語連鎖識別値を1レコードするデータ・ファイル(行ファイル)が作成され、次いで、入力された一行分を単語に分割し、単語と単語番号と行番号と識別値を1レコードするデータ・ファイル(単語ファイル)が作成される。行ファイルでは入力データがカンマ、ピリオド、空白などの区切りも含めて入力されたそのままの形で記録され、コンコーダンスのような行ファイルを処理する出力のときには、そのままのイメージで出力される。また註釈行も入力時のイメージのまま行ファイルに記録される。

これに対して単語形式入力では、カード1枚の入力データは単語に分割され、行形式入力と同じ形の単語ファイルが作成される。行ファイルは作成されず、註釈行は無視される。

実際のシステムの処理は、この2つ入力形式を統一的に扱う必要があることと単語連鎖の識別値を決める必要から、上に述べたのとは少し異なっている。前述の相違は利用者の立場から見た場合である。

L E Xシステムが作成するデータ・ファイルにはこのほかデータ情報ファイルがある。このファイルには、テキストのタイトル、単語連鎖の名称、識別モード単語の総数などの情報が格納されている。

L E Xデータ・ファイルは利用者が特に指定しないかぎり、計算センターの公用ボリューム上に恒久的ファイルとして作成される。これを標準ファイルと称する。但し厳密には、入力の際の標準ファイルはカードである。これはファイルの取扱いに馴染みの薄い利用者の便を考慮したためであるが、データが大量になれば磁気テープなどの補助記憶装置が必要になる。このような場合、利用者は非標準ファイルの指定をし、必要なファイル定義文を与えれば、L E Xデータ・ファイルは利用者の指定する媒体上に作成される。但し、データ情報ファイルだけは容量も少ないので(1 T R K)、常に公用ボリューム上に作成される。

データが入力された後には、L E Xコマンドと命令によってシステムに様々な出力を行わせることが可能となる。

1.3 L E X言語の統語法

データの入力であれ、データの処理であれ、システムに何らかの処理を行わせるためには、別の観点からすれば、システムと何らかのコミュニケーションを行うためには、利用者はシステムに何らかの指示、指命を与えなければならない。このコミュニケーションは、マクロ・コマンドであるL E Xコマンドとそれに続くステートメントによって行われる。L E Xコマンドはもっぱらファイルの指定に関するものであるので、第2節に譲り、ここではこのコマンド続に続けて与えるべき指示について述べる。

利用者がL E Xシステムに与える指示は、L E Xという言葉に属するひとつの

文と考えることができる。この意味でのL E Xの文はステートメントという単位から構成される。L E Xステートメントはカードと同一視しても差支えないが、必ずカードの第1カラムから始まる固定書式である。ステートメントは実行命令と従属ステートメントに大別される。書換え規則を用いれば次のように表わされよう。

L E Xの文→実行命令+(従属ステートメント)……①

実行命令はシステムが実行すべき処理を指定するものであり、独立に用いることができる。L E Xの文は実行命令で始まり(物理的には実行命令をパンチしたカードが最初のカードとなり)、実行命令の種類によってどのような従属ステートメントを与えうかが決まる。注意すべきは、どのような従属ステートメントを与えうかは実行命令によって決まるが、実際にその従属ステートメントを与えるか否かは原則として利用者の任意であるということである。これは、従属ステートメントは実行命令で指定した処理の詳細を規定するものであり、省略されればシステムの省略時解釈によって処理するからである。

従属ステートメントは大別して2種類に分けられる。書換え規則で図式的に示せば次のようになる。

従属ステートメント → (データ選別命令)+(明細指示)……②

前述のように、この規則は完全に文脈依存的である。データ選別命令は、実行すべき処理の対象となるテキストの部分集合を定義する。この命令が与えられなければL E Xデータ・ファイルに格納されたテキスト全体が処理の対象となる。これに対して明細指示は、実行命令ごとに異なっており、データの処理に関してより詳細な指定を与えるものである。

データ選別命令は選別、処理の単位によって、単語連鎖選別命令と単語選別命令にさらに下位区分される。

データ選別命令 → (単語連鎖選別命令) + (単語選別命令) ……③

この2種類の選別命令はそれぞれ次のように書換えられる。

語彙連鎖選別命令 →

$$\left\{ \begin{array}{l} \text{SELECT} \\ \text{REJECT} \end{array} \right\} + \text{選別単位名} + \text{識別値リスト} + * \text{END} \dots\dots ④$$

単語選別命令 →

$$\left\{ \begin{array}{l} \text{EXCLUDE} \\ \text{INCLUDE} \end{array} \right\} + (\text{SUBSTRING 指定}) + \text{文字列リスト} + * \text{END} \dots ⑤$$

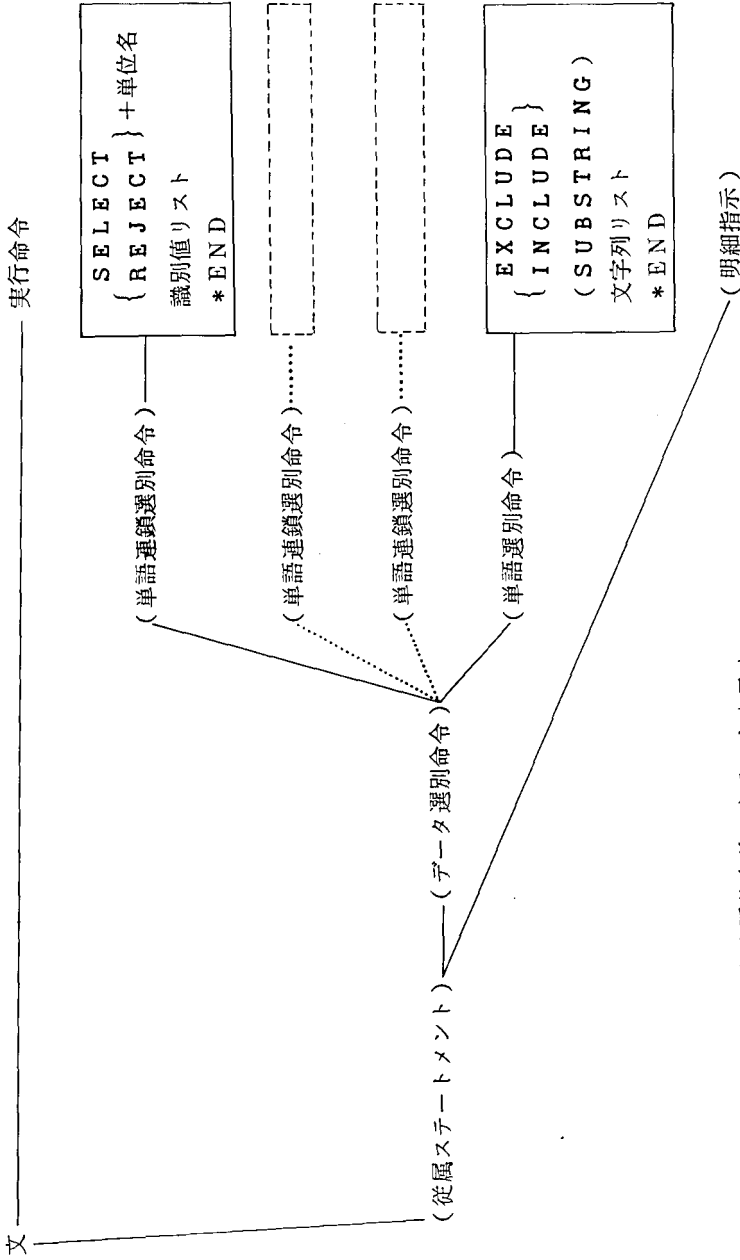
単語連鎖選別命令はSELECTまたはREJECTで始まる。SELECTは以下の識別値リストの値を識別値としてもつレコードを処理の対象とすることを意味する。REJECTはこの逆である。

これに続けて選別の対象となる単位名を与える。例えば、単語番号、行番号、単語連鎖名などである。次に選別すべき識別値のリストを与え、*ENDでリストを終える。

単語選別命令は、EXCLUDEまたはINCLUDEで始まる。EXCLUDE、INCLUDEの解釈はSELECT、REJECTの場合と同じである。これに続けて、必要があれば、SUBSTRING指定によって、与えられたリストと前方一致、中間一致等の一致方式を指定する。SUBSTRINGによる部分文字列の指定がなければ、完全一致が仮定される。

単語連鎖選別命令は実際には複数回用いることが許されており、最終的にデータ処理の対象となるのは、個々の単語連鎖選別命令によって決まるテキストの部分集合のそれぞれと単語選別命令によって決まる部分集合との論理積で定義されるテキストの部分集合である。

書換え規則①～⑤を使って、LEXの文を定義すれば、図式的には次の図2のようになる。



()は任意的であるか文脈依存であることを示す。

図 2 LEXの文の構造

1.4 主要な出力

L E Xのデータ処理の出力の詳細は別稿に譲るが以下にL E Xの主要な出力を列挙する。各々の出力はデータ選別機能を利用することにより更に多様な処理結果を得ることができる。

A 索引系の出力

- アルファベット（またはアイウエオ）順索引
- コンコーダンス（通常の形式と、K W I C形式の2通り）。処理単位としては単語だけでなく単語の集合、部分文字列あるいは部分文字列集合も可能。
- 脚韻語索引（幾通りかの押韻形式が可能である）

B 頻度表系の出力

- 頻度順リスト
- アルファベット順（アイウエオ順）リスト
- 脚韻語頻度順あるいはアルファベット順リスト

C 共起、連語に関する表

- 共起リスト（指定された単語と指定された範囲内で共起する単語のリスト）

D 素データ出力

（データ選別機能を用いれば簡単な情報検索にも利用できる。）

2. L E Xシステムの処理の流れ

2.1 L E Xコマンド

1.3で述べたように利用者がL E Xシステムに与える指示はL E XコマンドとL E Xの文から成る。L E Xの文が、システムの実行すべき処理とその明細を指定するのに対し、L E XコマンドはL E Xの文の実行に必要なファイルの割当てと、メイン・プログラム（正確にはプログラム・モジュール）の呼出しを主たる機能とする。そしてこのメイン・プログラムが、利用者の与えたL E Xの文を解釈し、実行する。

L E Xコマンドには、いくつかのオペランドがあるが、いずれもファイルの定義に関するものである。各オペランドの詳細はここでは触れないが、L E Xコマ

ンドの実行によってファイルがどのように割当てられるかについて大筋だけを述べておく。

まず、L E Xシステムで用いるファイルを列挙する。以下に示すのはすべてファイル定義名であり、ファイル名ではない。またソース・プログラムがP L N Iであるので、F O R T R A Nの場合のようにF T x x F 0 0 1といった形式でないことに注意されたい。

I D A T A……データ情報ファイル。

L D A T A……行ファイル。

W D A T A……単語ファイル。

この3つのファイルがL E Xデータ・ファイルであり、いずれもデータ入力時には出力ファイル、その他の場合には入力ファイルとして用いられる。この3つのファイルが標準ファイルであれば、公用ボリューム上に恒久的ファイルとして作成され、非標準ファイルであれば利用者のファイル定義コマンドに従って他の記憶媒体上に作成される。ただし、データ情報ファイルだけは常に標準ファイルとして作成される。

L E Xではこの3つのデータ・ファイルの他に、作業用の一時ファイルを用いる。L E Xで用いる一時ファイルには、行形式データのための作業用ファイル、単語形式データのための作業用ファイル、データ選別のリスト用のファイルの3種類がある。

このようなファイルがL E Xコマンドによってどのように割当てられるかを次の図3に示す。

L E Xコマンドの実行が開始されると、まずファイルを必要とするかどうか調べられる。これは、システムに入力前の素データの印刷など、ファイルを必要としない場合があるからである。ファイルが必要なければ直ちにメイン・プログラム(L X M A I N)が実行される。他方ファイルが必要であれば、データの入力、即ち、L E Xデータ・ファイルの作成ジョブかどうかを判定する。データの入力であれば、I D A T Aを出力ファイルとして定義し、次にL E Xデータ・ファイルの出力先が標準ファイルか非標準ファイルかを判定する。非標準ファイルであれば何もせず、標準ファイルであればW D A T Aを出力ファイルとして定

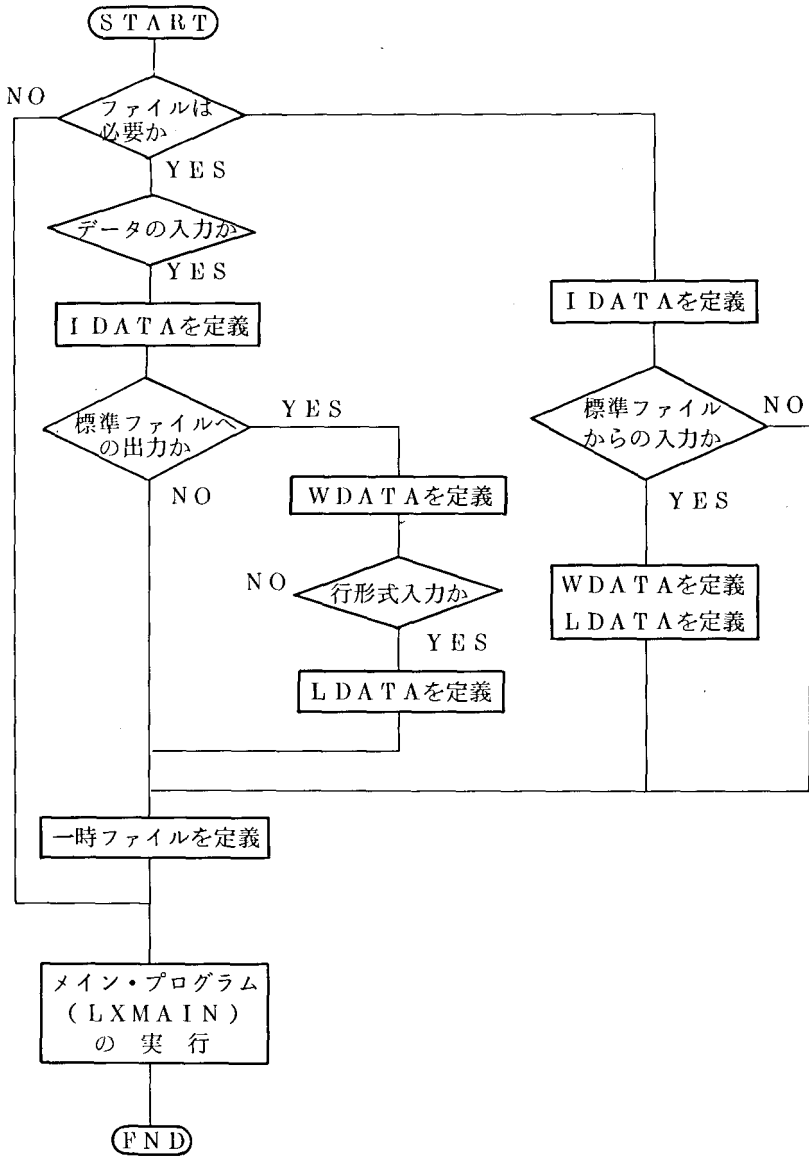


図3 LEXコマンドの実行

義し、必要に応じてLDATAを定義する。続けて作業用の一時ファイルを定義して、LXMAINを実行する。

これに対して、データの入力でない場合には、入力ファイルとしてまずIDATAを定義し、次に入力が標準ファイルであれば、WDATAとLDATAを入力ファイルとして定義して、一時ファイルの定義の後、メイン・プログラムを実行する。

このようにしてファイルの割当てが終了するとシステムはメイン・プログラムLXMAINの実行に移るわけであるが、この段階では利用者がLEXコマンドのオペランドやファイル定義コマンドで指定したファイルの割当てが正しいかどうかはシステムは判定しない。エラーがあった場合それが検出されるのはメイン・プログラムが実行され、実行ルーチンが呼び出される時である。しかもこれはLEXシステムではなく計算機システムによって行われる。

2.2 メイン・プログラムと実行用ルーチン

LEXのメイン・プログラムはLXMAINと呼ばれる。LXMAINは形式上はメイン・プログラムであるが、実際のLEX命令の実行はLXMAINによって呼び出される個々の実行用ルーチンによって行われる。LXMAINの機能は図4に示すように、LEX命令を解釈して必要な実行用ルーチンと呼ぶことと、LEXの文にエラーがあったとき、即ち実行用ルーチンが異常終了したときエラーメッセージを印刷することである。従って利用者が自分のプログラムをLEXに付け加えたければ、実行用命令をLXMAINに登録し、それに対応する事実上のメイン・プログラムである実行用ルーチンを付け加えればよい。

図4では、実行用ルーチンの実行終了後、正常終了かどうか判定されるようになっているが、実際には、後の図5にも示すように、実行用ルーチンがLEX文のシンタックスを点検中にエラーが検出されれば、制御は直ちにLXMAINに返される。データ選別命令の場合は、図5に示すように実行用ルーチンからデータ選別ルーチンに制御が移り、このサブルーチンで統語法がチェックされるが、ここでエラーが検出されるとデータ選別ルーチンは直ちに実行用ルーチンに制御を戻し、実行用ルーチンは直ちにLXMAINに制御を戻す。このことから

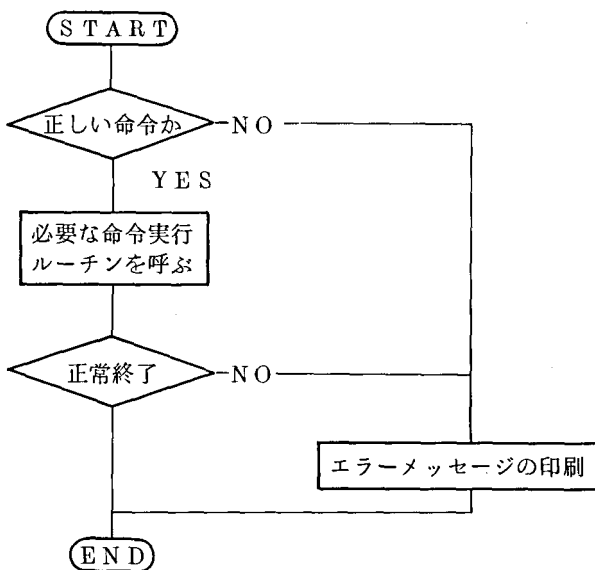


図4 LXMAINの流れ

も明らかなように、利用者の与えたLEX文に誤りがあれば、システムは直ちに実行を中止し、エラー・メッセージを印刷して（利用者の立場から見れば）異常終了する。この意味で、LEXにはエラーのレベルは設定されておらず、すべてのエラーは実行の中止に至る。エラーのレベルを設定するなり、エラー発生時に標準値を設定して処理を継続するという方式も考えられないではないが、当面現在の方式を変更する予定はない。

LXMAINによって呼び出される実行用ルーチンは、前述のように事実上は個々の命令を実行するためのメイン・プログラムと見做してもよい。実行用ルーチンは、図5に示すように、原則として、LEX文を（正確には、LEX文の従属ステートメントを）解釈、実行する部分と、それに従って命令を実行する部分のふたつの部分から成っている。

一般の実行用ルーチンではまず従属ステートメントを処理する。最初に単語連鎖選別命令がチェックされ、あれば単語連鎖選別ルーチンと呼ぶ。単位やリスト

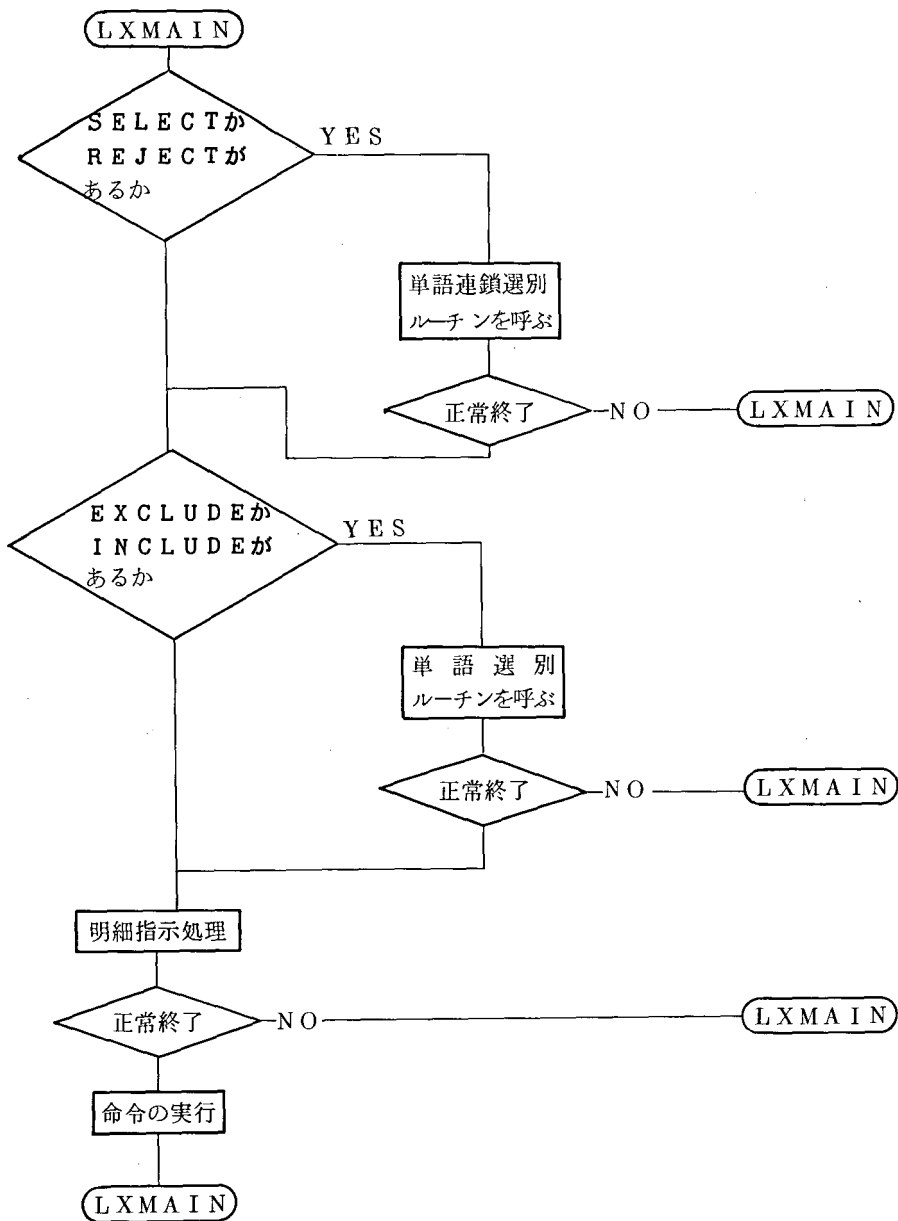


図5 実行用ルーチンの概念図

がここで処理されエラーがなければ選別処理を行う。エラーを検出した場合の L X M A I N への戻り方は前出の通りである。次に単語選別命令の有無がチェックされ、あれば単語選別ルーチンと呼ぶ。ここでも同様のやり方でリストが処理され選別処理が行われる。データの選別が終ると、各実行ルーチンに特有の明細指示の処理が行われ、エラーがなければ、従属ステートメントの処理は正常に終了する。

従属ステートメントの処理が終わってはじめて、命令が実行される。この部分が各実行ルーチンの中核部であることは言うまでもない。命令の実行が終了すると制御は再び L X M A I N に戻り、ジョブまたはジョブ・ステップの終了となる。

3. 入出力の実例

3.1 入力データの実例と L E X の文

次の図 6, 図 7 に単語形式, 行形式で入力されるデータの実例を I N P U T 命令とともに示す。図 6, 7 ともに P R I N T R A W D A T A 命令の出力例であり, 左端の一連番号は, システムがつけた入力順のカード番号である。

次の図 6 が単語形式の入力例である。

①

I N P U T 命令

②～⑦

I N P U T 命令に対する明細指示。

T I T L E = テクストの名称。F O R M A T = 入力の形式。L S *i* は, 第 *i* 単語連鎖の存在を宣言し, 必要に応じて名称 (N A M E), 識別モード (I D M O D E), 初期値 (S T A R T) を与える。明細指示は省略されれば標準値が仮定される。⑦の R E A D D A T A は L E X 文の終り, あるいは入力テキストの始まりを示す。

⑧, ⑨

註釈行。

⑩, ⑬, ⑯, ⑲, ⑳, ㉓

*** RAW DATA TO BE INPUT ***

1 INPUT
2 TITLE=FINAL DOCUMENT ON DISARMAMENT SESSION
3 FORMAT=WORM
4 LS1 NAME=SENT,
5 LS2 NAME=PARGRH,
6 LS3 NAME=SLCTN,IUMODE=A,
7 READ DATA
8 / TEXT=FINAL DOCUMENT OF THE TENTH SPECIAL SESSION OF THE GENERAL ASSEMBLY
9 / A-RES-S-10-2 13 JULY 1978
10 * * * INTR
11 ATTAINMENT OF THE OBJECTIVE OF SECURITY, WHICH IS AN INSEPARABLE ELEMENT OF
12 PEACE, HAS ALWAYS BEEN ONE OF THE MOST PROFOUND ASPIRATIONS OF HUMANITY.
13 *
14 STATES HAVE FOR A LONG TIME SOUGHT TO MAINTAIN THEIR SECURITY THROUGH THE
15 POSSESSION OF ARMS.
16 *
17 ADMITTEDLY, THEIR SURVIVAL HAS, IN CERTAIN CASES, EFFECTIVELY DEPENDD ON
18 WHETHER THEY COULD COUNT ON APPROPRIATE MEANS OF DEFENCE.
19 *
20 YET THE ACCUMULATION OF WEAPONS, PARTICULARLY NUCLEAR WEAPONS, TODAY CONSTITUTES
21 MUCH MORE A THREAT THAN A PROTECTION FOR THE FUTURE OF MANKIND.
22 *
23 THE TIME HAS THEREFORE COME TO PUT AN END TO THIS SITUATION, TO ABANDON THE USE
24 OF FORCE IN INTERNATIONAL RELATIONS AND TO SEEK SECURITY IN DISARMAMENT, THAT IS
25 TO SAY, THROUGH A GRADUAL BUT EFFECTIVE PROCESS BEGINNING WITH A REDUCTION IN
26 THE PRESENT LEVEL OF ARMAMENTS.
27 *
28 THE ENDING OF THE ARMS-RACE AND THE ACHIEVING OF REAL DISARMAMENT ARE TASKS

図 6 単語形式入力例

単語連鎖の境界識別標識。より正確には、個々の第 i 単語連鎖の開始を示す。
* は第 1 単語連鎖 (LS 1) の識別標識 (ここでは、SENT (ence) という名称が与えられている)。このデータでは、⑪の“ATTAINMENT”から⑫の“HUMANITY”までがひとつのSENTを成し (識別値は 1), ⑭の“STATES”から⑮の“ARMS”までがまたひとつのSENTである (識別値は 2)。#, % はそれぞれ第 2 単語連鎖 (LS 2, 名称は PARGRH), 第 3 単語連鎖 (LS 3, 名称は SECTN) の識別標識である。第 3 単語連鎖は識別モードが A モードなので、⑩で、“INTR” という識別値が与えられている。次に、“%識別値” が現われるまですべての単語に第 3 単語連鎖 SECTN の識別値として“INTR” が与えられる。

第 1 単語連鎖, 第 2 単語連鎖の識別モードはいずれも省略時解釈により S モードなので, 省略時解釈による初期値 1 から始まる一連番号が識別値として与えられることになる。

⑪, ⑫, ⑭, ⑮, ⑰, ⑱, ⑳, ㉑, ㉓-㉔, ㉖

実際のデータ, 即ちテキスト本体である。

次の図 7 は, 行形式入力の例で, 開発を主題とする文献のタイトルを 1 レコードとして入力する例である。

①-⑤

INPUT 命令と明細指示。

FORMAT による指定がないので, 省略時解釈により, 行形式入力と仮定される。

⑧

註釈行。

⑦, ⑪, ㉒, ㉔

単語連鎖の境界識別標識。LS 1, LS 2, LS 3 とともに A モードなので識別値が与えられている。行形式入力の場合の第 1 単語連鎖 (LS 1) については, カードが区切りと仮定されるので通常は識別標識 * を使わない。しかし A モードの時には“*識別値”は第 1 単語連鎖の区切りとは見做されず, 次の“*識別値”

*** RAW DATA TO BE INPUT ***

```
1 INPUT
2 TITLE=LITERATURE ON INTERNATIONAL DEVELOPMENT
3 LS1 NAME=YEAR, IDMODE=A,
4 LS2 NAME=PUBMOD, IDMODE=A,
5 LS3 NAME=VOL-NO, IDMODE=A,
6 READ DATA
7 #1973 #X #9-1
8 / RAW DATA TAKEN FROM THE BIBLIOGRAPHY OF THE INTERNATIONAL DEVELOPMENT JOURNAL
9 POLITICS OF DEVELOPMENT ADMINISTRATION: SOME HYPOTHESES
10 COMPARATIVE PUBLIC ADMINISTRATION: TOWARDS ECOLOGICAL DEVELOPMENTAL ORIENTATION
11 #1974
12 DISEASE & LABOR PRODUCTIVITY
13 PRESENT STATE OF MARINE SCIENCE & OCEANOGRAPHY IN LESS DEVELOPED COUNTRIES
14 FACTOR PROPORTIONS & URBAN EMPLOYMENT IN DEVELOPING COUNTRIES
15 TRADE PATTERNS & ECONOMIC EVOLUTION
16 PUBLIC FINANCE & INEQUALITY
17 TARIFF REVENUE & OPTIMAL CAPITAL ACCUMULATION IN LESS DEVELOPED COUNTRIES
18 CULTURAL MYTHS & REALITIES OF PROBLEM-SOLVING
19 MODIFYING BUREAUCRATIC SYSTEMS IN DEVELOPING WORLD
20 TOWARD THEORY OF POLITICAL MOBILIZATION
21 PRIORITIES IN ENVIRONMENTAL HEALTH
22 #BOOK
23 CHOICE & ADAPTATION OF TECHNOLOGY IN DEVELOPING COUNTRIES
24 #X
25 DESIGN FOR MARKET ECONOMY
26 IMPACT OF AGRICULTURE'S DOMESTIC TERMS OF TRADE
27 ENVIRONMENTAL CONTROLS & ECONOMIC GROWTH
28 NEEDED . GLOBAL STRATEGY OF DEVELOPMENT
```

图 7 行形式入力例

が現われるまで個々の第1単語連鎖に与えるべき識別値と見做される。このような場合、単語連鎖という概念よりも、図8、図9からも明らかのように、変数という概念のほうがふさわしいと言える。識別値、例えば、“1973”は、個々の行レコードあるいは単語レコードがとる、YEARという変数の変数値と見做すべきであろう。このような両義性はある意味ではLEXの欠陥と見做しうるかもしれないが、むしろ長所と考えるべきである。というのは、一方でテキスト中の連続した単語の連鎖の場合、即ち厳密な意味での単語連鎖であり、同一の識別値をもつことはテキストの有意味な構成単位であることを意味し、その単位によってデータ処理を制御する可能性と、他方で不連続であっても同一識別値であることによってデータ処理を制御する可能性というふたつの可能性が開けるからである。

⑨, ⑩, ⑫-⑰, ⑳, ㉕-㉞

データ本体である。行形式入力であるから各行が(各カード1枚が)ひとつの第1単語連鎖(即ち行)となる。

3.2 LEXデータ・ファイル

図7に示した入力データがどのような形でLEXデータ・ファイルに記録されているかをPRINT WORD FILE命令, PRINT LINE FILE命令の出力によって図8, 9に示す。

図8は図7のデータから作成された単語ファイルの内容を示すが、実際のファイルでは、LSQ, WSQは、個々の単語の後にWSQ, LSQの順に置かれている。入力形式が行形式なのでLSQには行番号が入っている。単語形式であれば何も記録されない。

図9は図7のデータから作成された行ファイルの内容を示す。LSQは実際のファイルでは、行の後に置かれている。註釈行にはLSQの値として0が与えられている。

(LSQ)	MSG	YEAR	PUBMOD	VOL-NO	(LSQ)	MSG	YEAR	PUBMOD	VOL-NO
1	1	1973	X	9-1	1	2	1973	X	9-1
1	3	1973	X	9-1	1	4	1973	X	9-1
1	5	1973	X	9-1	1	6	1973	X	9-1
2	7	1973	X	9-1	2	8	1973	X	9-1
2	9	1973	X	9-1	2	10	1973	X	9-1
2	11	1973	X	9-1	2	12	1973	X	9-1
2	13	1973	X	9-1	3	14	1974	X	9-1
3	15	1974	X	9-1	3	16	1974	X	9-1
3	17	1974	X	9-1	4	18	1974	X	9-1
4	19	1974	X	9-1	4	20	1974	X	9-1
4	21	1974	X	9-1	4	22	1974	X	9-1
4	23	1974	X	9-1	4	24	1974	X	9-1
4	25	1974	X	9-1	4	26	1974	X	9-1
4	27	1974	X	9-1	4	28	1974	X	9-1
5	29	1974	X	9-1	5	30	1974	X	9-1
5	31	1974	X	9-1	5	32	1974	X	9-1
5	33	1974	X	9-1	5	34	1974	X	9-1
5	35	1974	X	9-1	5	36	1974	X	9-1
6	37	1974	X	9-1	6	40	1974	X	9-1
6	39	1974	X	9-1	7	42	1974	X	9-1
7	41	1974	X	9-1	7	44	1974	X	9-1
7	43	1974	X	9-1	7	46	1974	X	9-1
7	45	1974	X	9-1	8	48	1974	X	9-1
8	47	1974	X	9-1	8	50	1974	X	9-1
8	49	1974	X	9-1	8	52	1974	X	9-1
8	51	1974	X	9-1	8	54	1974	X	9-1
8	53	1974	X	9-1	8	56	1974	X	9-1
9	55	1974	X	9-1	9	58	1974	X	9-1
9	57	1974	X	9-1	9	60	1974	X	9-1
9	59	1974	X	9-1	10	62	1974	X	9-1
9	61	1974	X	9-1	10	64	1974	X	9-1
10	63	1974	X	9-1	10	66	1974	X	9-1
10	65	1974	X	9-1	11	68	1974	X	9-1
10	67	1974	X	9-1	11	70	1974	X	9-1
11	69	1974	X	9-1	11	72	1974	X	9-1
11	71	1974	X	9-1	12	74	1974	X	9-1
12	73	1974	X	9-1	12	76	1974	X	9-1
13	75	1974	X	9-1	13	78	1974	X	9-1
13	77	1974	BOOK	9-1	13	80	1974	BOOK	9-1
13	79	1974	BOOK	9-1	13	82	1974	BOOK	9-1
13	81	1974	BOOK	9-1	13	84	1974	BOOK	9-1
13	83	1974	BOOK	9-1	14	86	1974	X	9-1
14	85	1974	X	9-1	14	88	1974	X	9-1
14	87	1974	X	9-1	14	90	1974	X	9-1
15	89	1974	X	9-1	15	92	1974	X	9-1
15	91	1974	X	9-1	15	94	1974	X	9-1
15	93	1974	X	9-1	16	96	1974	X	9-1
16	95	1974	X	9-1	16	98	1974	X	9-1
16	97	1974	X	9-1	16	100	1974	X	9-1
16	99	1974	X	9-1	17	102	1974	X	9-1
17	101	1974	X	9-1	17	104	1974	X	9-1
17	103	1974	X	9-1	18	106	1974	X	9-1
17	105	1974	X	9-1	18	108	1974	X	9-1

図 8 単語ファイル

	YEAR	PUBMOD	VOL-NO
0 / RAW DATA TAKEN FROM THE BIBLIOGRAPHY OF THE INTERNATIONAL DEVELOPMENT JOURNAL	*****	*****	*****
1 POLITICS OF DEVELOPMENT ADMINISTRATION: SOME HYPOTHESES	1973	X	9-1
2 COMPARATIVE PUBLIC ADMINISTRATION: TOWARDS ECOLOGICAL DEVELOPMENTAL ORIENTATION	1973	X	9-1
3 DISEASE & LAND PRODUCTIVITY	1974	X	9-1
4 PRESENT STATE OF MARINE SCIENCE & OCEANOGRAPHY IN LESS DEVELOPED COUNTRIES	1974	X	9-1
5 FACTOR PROPORTIONS & URBAN EMPLOYMENT IN DEVELOPING COUNTRIES	1974	X	9-1
6 TRADE PATTERNS & ECONOMIC EVOLUTION	1974	X	9-1
7 PUBLIC FINANCE & INEQUALITY	1974	X	9-1
8 TARIFF REVENUE & OPTIMAL CAPITAL ACCUMULATION IN LESS DEVELOPED COUNTRIES	1974	X	9-1
9 CULTURAL MYTHS & REALITIES OF PROBLEM-SOLVING	1974	X	9-1
10 MODIFYING BUREAUCRATIC SYSTEMS IN DEVELOPING WORLD	1974	X	9-1
11 TOWARD THEORY OF POLITICAL MOBILIZATION	1974	X	9-1
12 PRIORITIES I: ENVIRONMENTAL HEALTH	1974	X	9-1
13 CHOICE & ADAPTATION OF TECHNOLOGY IN DEVELOPING COUNTRIES	1974	BOOK	9-1
14 DESIGN FOR MARKET ECONOMY	1974	X	9-1
15 IMPACT OF AGRICULTURE'S DOMESTIC TERMS OF TRADE	1974	X	9-1
16 ENVIRONMENTAL CONTROLS & ECONOMIC GROWTH	1974	X	9-1
17 NEEDED - GLOBAL STRATEGY OF DEVELOPMENT	1974	X	9-1
18 DEVELOPED & DEVELOPING	1974	X	9-1
19 SYSTEMS OF HEALTH CARE DELIVERY	1974	X	9-1
20 BANKS & REGIONAL ECONOMIC DEVELOPMENT	1974	X	9-1
21 TRADITIONAL PATRIMONIALISM & MODERN NEOPATRIMONIALISM	1973	BOOK	9-1
22 ENERGY & WORLD AGRICULTURE	1974	X	9-1
23 SIX PALUSTRIE THESES ABOUT PEASANTS' PERSPECTIVES IN DEVELOPING WORLD	1974	X	9-1
24 FOOD SCIENCE IN DEVELOPING COUNTRIES: SELECTION OF UNSOLVED PROBLEMS	1974	BOOK	9-1
25 POLITICS OF DISTRUST: FIELD PROBLEMS IN COMPARATIVE RESEARCH	1974	X	9-1
26 MINERAL RESOURCES & ECONOMIC GROWTH	1974	X	9-1

3.3 その他の出力例

図10はFREQUENCY命令の出力例である。データは図7に示した開発を主題とする文献のキーワードである。

TOTAL NUMBER OF WORDS PROCESSED=		375		
TOTAL KIND OF WORDS PROCESSED=		205		
WORD	FREQ	PCT	ACC	
			FREQ	
			PCT	
&	29	7.7	29	7.7
OF	25	6.7	54	14.4
IN	22	5.9	76	20.3
COUNTRIES	14	3.7	90	24.0
DEVELOPING	14	3.7	104	27.7
DEVELOPMENT	10	2.7	114	30.4
ECONOMIC	9	2.4	123	32.8
DEVELOPED	4	1.1	127	33.9
FOOD	4	1.1	131	34.9
PROBLEMS	4	1.1	135	36.0
WORLD	4	1.1	139	37.1
ADMINISTRATION	3	0.8	142	37.9
FDR	3	0.8	145	38.7
INTERNATIONAL	3	0.8	148	39.5
LESS	3	0.8	151	40.3
PUBLIC	3	0.8	154	41.1
REGIONAL	3	0.8	157	41.9
RESOURCES	3	0.8	160	42.7
TECHNOLOGY	3	0.8	163	43.5
TRADE	3	0.8	166	44.3
CASE	2	0.5	168	44.8
COMPARATIVE	2	0.5	170	45.3
CONCEPT	2	0.5	172	45.9
CONTROLS	2	0.5	174	46.4
DISEASE	2	0.5	176	46.9
ECONOMY	2	0.5	178	47.5
ENERGY	2	0.5	180	48.0
ENVIRONMENTAL	2	0.5	182	48.5
EXPORT	2	0.5	184	49.1
GROWTH	2	0.5	186	49.6
HEALTH	2	0.5	188	50.1
LABOR	2	0.5	190	50.7
LOCS	2	0.5	192	51.2
NEED	2	0.5	194	51.7
OPTIMAL	2	0.5	196	52.3
POLICY	2	0.5	198	52.8
POLITICAL	2	0.5	200	53.3
POLITICS	2	0.5	202	53.9
RESEARCH	2	0.5	204	54.4
SCIENCE	2	0.5	206	54.9
SOME	2	0.5	208	55.5
SYSTEMS	2	0.5	210	56.0
THEORY	2	0.5	212	56.5
URBAN	2	0.5	214	57.1
ABOUT	1	0.3	215	57.3
ABSORPTION	1	0.3	216	57.6
ACCUMULATION	1	0.3	217	57.9
ADAPTATION	1	0.3	218	58.1
ADDENDA	1	0.3	219	58.4
ADVERSE	1	0.3	220	58.7
AGRICULTURE	1	0.3	221	58.9

図10 FREQUENCY命令出力例

図11はALPHA命令の出力例である。

TOTAL NUMBER OF WORDS PROCESSED= 375
 TOTAL KIND OF WORDS PROCESSED= 205

WORD	FREQ	PCT
&	29	7.7
ABOUT	1	0.3
ABSORPTION	1	0.3
ACCUMULATION	1	0.3
ADAPTATION	1	0.3
ADDENDA	1	0.3
ADMINISTRATION	3	0.8
ADVERSE	1	0.3
AGRICULTURE	1	0.3
AGRICULTURE'S	1	0.3
AID	1	0.3
APPROACH	1	0.3
AS	1	0.3
ASPECTS	1	0.3
ASSESSMENT	1	0.3
AT	1	0.3
BANKING	1	0.3
BANKS	1	0.3
BETWEEN	1	0.3
BEYOND	1	0.3
BIG	1	0.3
BOTTLE-FEEDING	1	0.3
BUREAUCRATIC	1	0.3
BY	1	0.3
CAPITAL	1	0.3
CARE	1	0.3
CASE	2	0.5
CHANGE	1	0.3
CHOICE	1	0.3
CITIES	1	0.3
COMMUNICATION	1	0.3
COMPARATIVE	2	0.5
CONCEPT	2	0.5
CONTRACEPTION	1	0.3
CONTROLS	2	0.5
COUNTRIES	14	3.7
COUNTRIES'	1	0.3
CREDIT	1	0.3
CRIME	1	0.3
CRIMINAL	1	0.3
CRISIS	1	0.3
CULTURAL	1	0.3
DELINEATION	1	0.3
DELIVERY	1	0.3
DEPENDENCE	1	0.3
DESIGN	1	0.3
DESIGNATION	1	0.3
DETERMINANTS	1	0.3
DEVELOPED	4	1.1
DEVELOPING	14	3.7
DEVELOPMENT	10	2.7
DEVELOPMENTAL	1	0.3

図11 ALPHA命令出力例(部分)

図12は COOCCURRENCE 命令の出力例で，“ギジュツシンボ”という単語と同一単語連鎖内（ここでは文献内）で共起する単語の共起頻度リストである。データは，技術移転に関する文献のキーワードである。

SPECIFIED ITEM=ギ*シ*ユツシンボ*
OCCURS 27 TIMES

COOCCURRS WITH	TIMES	PCT	CR
ケンキユウカイハツ	4	14.8	14.8
チホウシ*チタイ	3	11.1	11.1
ギ*シ*ユツカイハツ	3	11.1	11.1
アメリカ	2	7.4	7.4
セイフ	1	3.7	3.7
ギ*シ*ユツヨソク	1	3.7	3.7
ヤンキ*ヨウシヤカイカク	1	3.7	3.7
ケイサ*イシスウ	1	3.7	3.7
ケイカク	1	3.7	3.7
トツキヨセイト*	1	3.7	3.7
コウキヨウシ*キ*ヨウ	1	3.7	3.7
トシカイハツ	1	3.7	3.7
ユツカセイサク	1	3.7	3.7
トツキヨ	1	3.7	3.7
ライセンス	1	3.7	3.7
カンリシヤ	1	3.7	3.7
カンキヨウホセ*ン	1	3.7	3.7
ハツチントシ*ヨウコク	1	3.7	3.7
ヒユ_マツフアクタ	1	3.7	3.7

TOTAL= 27

図12 COOCCURRENCE 命令出力例

図13はKWIC INDEX命令の出力例で、この例は“DEVELOPMENT”に関するKWICインデックスである。

DEVELOPMENT (10)				X	9-1
	POLITICS OF DEVELOPMENT ADMINISTRATION: SOME HYPOTHES	1973		X	9-1
	NEEDED - GLOBAL STRATEGY OF DEVELOPMENT	1974		X	9-1
	BANKS & REGIONAL ECONOMIC DEVELOPMENT	1974		X	9-1
	REGIONAL DELINEATION: DESIGNATION OF DEVELOPMENT	1974		X	9-1
	DEVELOPMENT & INTERNATIONAL ECONOMIC ORDE	1974		X	9-1
	POLICY OPTIONS FOR RURAL DEVELOPMENT	1973		X	9-1
	OPTIMAL DEVELOPMENT PROGRAMME UNDER UNCERTAINTY:U	1974		X	9-1
	CREDIT CONTROLS AS INSTRUMENTS OF DEVELOPMENT POLICY IN LIGHT OF ECONOMIC T	1974		X	9-1
	COMMUNICATION IN DEVELOPMENT	1974		X	9-1
	REGIONAL DEVELOPMENT IN LOGS	1974		X	9-1

図13 KWIC INDEX出力例

付 記

原稿提出後，新たな命令の追加等の修正，拡張のほか，いくつかの拡張を施した。主なものは次のふたつである。

- ① 会話ジョブでも使用可能にした。

- ② LEXステートメントの書式を所謂自由書式にした。