# Computer Documentation System for Small- and Medium-Sized Information Collections

P. J. LEWI and W. W. BRAET

Scientific Data Processing Department, Janssen Pharmaceutica,
Research Laboratoria, Beerse, Belgium

Received November 10, 1969

A computer service was designed for use in scientific and engineering communities of a pharmaceutical company. Alphabetic information is divided into main subjects, entered in free-format on punched cards, and transferred to magnetic tape. The master tape can be searched on any subject by strings of word fragments and Boolean operators. The program was written in Fortran IV for an IBM-1800 computer installation.

This paper describes how information processing problems in a Research and Development division of a pharmaceutical company have been solved by computer in a simple and straightforward manner.

The system designer, exploring the needs for information and documentation in a scientific and technical community, is usually confronted with a great variety of requirements.

As Altmann[1] has recognized, this is mainly due to R & D people working in small teams on highly specialized subjects for varying periods of time. Some groups operating separately, will sometimes be served best by an individual and manually operated tool, such as a catalog file, peak-a-boo, or edge-notched cards, the usefulness of which is discussed in detail by Kent[4].

In other cases, a personalized computer system, using numerically coded descriptors and a few levels of subordination, such as described by Gillis[2], can be very helpful.

We selected English sentences as input, using the free-format recording method of Korein[5]. The retrieval language consists of English words and Boolean operators, as in the method employed by Heaps[3].

Our system however can be used equally well by those who are more familiar with hierarchical classification, tagged descriptors, or subject headings. Although a natural language system with random ordering has been rated low by Meadow[6], it can be used advantageously in a diversified R & D community, provided that the size of the collections can be kept between 5000 and 10,000 records. The method described here stands halfway between a keyword extraction system and the processing of languages with semantic, syntactic, and statistical analyses, discussed by Salton[7] and Simmons[8].

Our principal objectives were to encourage user responsibility for the performance of the system and to ensure easy formulation of requests and output specifications. Furthermore, we were required to prepare indexes by various subjects, such as those reported by Teal[9].

In a generalized information system, performance is often measured in terms of precision and recall. In a personalized service, these measures are not meaningful, as the precision will generally be very high. Also the recall ratios can be anything between 0 and 1 depending on the effort and endurance of the individual who assumed responsibility for the system. Therefore, we propose a measure that could be termed the utility ratio. If a bibliography for a scientific paper were compiled, the utility would be the number of references derived from the computer, divided by the total number of references (including those obtained from colleagues, personal library, or other services). The utility reflects the user's satisfaction or dissatisfaction and also his willingness to continue or discontinue the documentation service.

## DESIGN OF THE SYSTEM

The system assumes that information derived from a document can be divided into a limited number of independent subjects. In a chemical literature retrieval system, we would have the subjects: author names, title, source, abstract, and chemical names.

The number of alphanumerical entries—e.g., names, keywords, phrases—covered by a subject is only limited by the field length of the input record.

For practical reasons (computer memory), we limited the number of subjects to five and assigned a maximum of 10 fields with 74 characters to each subject.

Since 80-column tab cards are used as primary input medium, each field of 74 characters was made a separate transaction. Information is entered on the cards in free format, and a transaction number (columns 79-80) defines
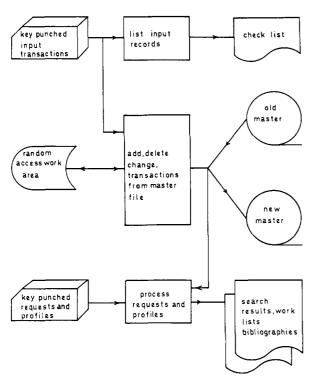
Figure 3. Block diagram of computer procedures

from core or tape into a single field of 25 × 74 characters. Access to a given subject of the record is obtained indirectly via a branch table which contains the transaction numbers. This technique has the advantage that the entire field can be scanned in a continuous way, without being hindered by possible word segmentations at the end of a transaction field.

Once the limits between which the input field is to be scanned are determined, a program segment retrieves the next term (if any) from this subfield. Comparison between the retrieved term and the terms of the request statements is performed after elimination of all special characters.

Depending upon the subject chosen, the matching routine uses suffixing, prefixing, or a combination of both rules and the result is entered into a truth/false table.

Request processing, tape update, and list options are grouped into a single main program. Each option is selected by appropriate control cards. The total length of the program is about 8K words (16 bits), not including the disk and tape utility programs. Grouping of program segments can result in core overlays of less than 4K, so that the system can be made operational on machines of limited core memory size.

The program was initially coded and executed in FORTRAN IV. The most critical segments (word matching and elimination of special characters) were rewritten in assembly language.

## DISCUSSION

The main advantages of the system described here were the free format, entirely alphabetic, data input and request formulation. Experienced keypunch operators had little difficulty in preparing and verifying the input cards. The one-to-one correspondence between data records on card, tape, and printer together with the single file concept simplifies maintenance and updating procedures of the master file. Although the present system does not pretend to compete in speed and versatility with larger sophisticated programs, the authors have found that the present system can solve a large number of needs arising in a scientific community. Practical usage of the system is limited, however, to smaller data collections, with an upper limit of 10,000 records of an average size of 1000 characters each.

## SUMMARY

A computer system is described, designed for data collections of up to 10,000 records and to be run on a computer with limited available core storage and peripheral equipment.

The design of data input, output, and request formulation was generalized to an extent that the system can be tailored to individual needs. Free format, alphabetic fields, single record layout, and alphabetic word matching combined with Boolean relations are its main characteristics.

The system can be used with advantage as an intermediate stage between manually-operated and large-scale document retrieval systems.

## ACKNOWLEDGMENT

## LITERATURE CITED

1 ALTMANN,B.
A NATURAL LANGUAGE STORAGE AND RETRIEVAL (ABC) METHOD ITS RATIONALE, OPERATION AND FURTHER DEVELOPMENT PROGRAM.
J. CHEM. DOC., 6, 154-157 (1966).

2 GILLIS,C.N.
BIOMEDICAL INFORMATION RETRIEVAL.
A COMPUTER-BASED SYSTEM FOR INDIVIDUAL USE.
J. CHEM. DOC., 7, 98-100 (1967).

3 HEAPS,H.S.
BOOLEAN, FRACTIONAL AND ASSOCIATIVE SEARCHES ON TRUNCATED TITLE WORDS.
PROC. AM. SOC. INFORM. SC., 5, 179-184 (1968).

4 KENT,A.
CHAPTER 7, P. 172-182, IN 'SPECIALIZED INFORMATION CENTERS', SPARTAN BOOKS, WASH. D.C., 1965.

5 KOREIN,J., GOODGOLD,A.L. AND RANDT,C.T.
COMPUTER PROCESSING OF MEDICAL DATA BY VARIABLE-FIELD-LENGTH FORMAT.
J. AM. MED. ASS., 186, 132 (1963).

6 MAEDOW,C.T.
CHAPTER 5, P. 169-173, IN 'THE ANALYSIS OF INFORMATION SYSTEMS', JOHN WILEY AND SONS, INC., N.Y., 1967.

7 SALTON,G.
AUTOMATIC LANGUAGE PROCESSING.
TECHN. REP. 68-6, DEPT COMP. SC., CORNELL UNIV., ITHACA, N.Y., 1968.

8 SIMMONS,R.
ANSWERING ENGLISH QUESTIONS BY COMPUTER, A SURVEY.
PART III8, IN 'THE GROWTH OF KNOWLEDGE', M. KOCHEN, ED., JOHN WILEY AND SONS, INC., N.Y., 1967.

9 TEAL,T.W. AND GREENBERG,S.M.
MANAGING LITERATURE IN THE PHARMACEUTICAL INDUSTRY.
DRUG INFORM. ASS. BULL., 2, 136-143 (1968).