(3) Morgan, H. L., "The Generation of a Unique Machine Description for Chemical Structures," J. CHEM. Doc. 5, 107 (1965).

(4) Tate, F. A., H. L. Morgan, D. P. Leiter, and R. E. Stobaugh, "A Mechanized Registry of Chemical Compounds," presented at the 1965 Congress of International Federation for Documentation (FID), Washington, D. C., October 10-15, 1965 (unpub).

(5) "Text Searching," Chemical Abstracts Service, Columbus, Ohio, 1968.

(6) "Desktop Analysis Tool for the Common Data Base," Chemical Abstracts Service, June 1968, PB 179 900, Clearinghouse for Federal Scientific and Technical Information.

(7) "Preparation of Search Profiles," Chemical Abstracts Service, Columbus, Ohio, 1967.

(8) "Substructure Search—Background Information and Question Coding Instructions," Chemical Abstracts Service, Columbus, Ohio, 1968.

(9) Snedecor, G. W., "Statistical Methods," Iowa State College Press, Ames, Iowa, 1946.

(10) Kent, A. K., "United Kingdom Experiences in the Operation of a Retrieval and Dissemination Service Based on CAS Search Tapes," presented at the American Chemical Society Meeting, Atlantic City, N. J., September 1968.

(11) Unpublished experimental data, J. L. Carmon and M. K. Park, University of Georgia, Athens, Ga.

# Implementation and Evaluation of Two Computerized Information Retrieval Systems at the University of Pittsburgh

NEALE S. GRUNSTRA and K. JEFFREY JOHNSON

Pittsburgh Chemical Information Center, University of Pittsburgh, Pittsburgh, Pa. 15213

This article describes the Pittsburgh Chemical Information Center data processing group's implementation and evaluation of two information retrieval systems: TEXT-PAC, an information retrieval system developed by the International Business Machines Corporation, and a system developed by the Chemical Abstracts Service (CASCON). Both systems use the Chemical Abstracts Condensates (CA Condensates) tape as the input data base.

A group has been organized within the chemistry department of the University of Pittsburgh as an experimental station for the computerized dissemination of information.[1,2] One of the goals of this group, the Pittsburgh Chemical Information Center (PCIC), is to evaluate new chemistry data bases and the programs that search them. Thus, this study began shortly after Chemical Abstracts Service (CAS) began publishing Condensates in the fall of 1968.

CA Condensates is a machine-readable magnetic tape service of CAS. One tape is issued by CAS each week corresponding to the hard-copy issues of Chemical Abstracts (CA). Each record on the tape includes an abstract number, title, authors, bibliographic citation, and keywords which amplify the content of the article. The abstract and molecular formulas included in the hard-copy of CA are not included on the tape. Papers from approximately 12,000 journals are included. The odd-numbered issues of CA Condensates cover papers in biochemistry and organic chemistry. The even-numbered issues include chemical engineering, applied, physical and analytical chemistry, and macromolecular chemistry.

CA Condensates has three features that make it more appealing for current awareness than Chemical Titles (CT)[3], another CAS data base. First, since CT includes only 650 journals, CA Condensates offers significantly greater journal coverage as well as books and patents. Second, the keyword feature is absent in CT. And third, the CA Condensates user has the option of searching even or odd issues, or both. Thus, if he desires, the user may eliminate wide areas of chemistry, thereby reducing the number of irrelevant hits and search costs.

The Pittsburgh Chemical Information Center is currently providing current-awareness service to approximately 270 users. Chemical information specialists translate user interests into computer readable search strategies. These data are then submitted to a data processing group, which is responsible for both routine production processing as well as development of new information retrieval capabilities. The computer output is returned to the information specialists who maintain statistics on the processing costs and the number of citations per user. In addition, a group of behavioral scientists are in regular contact with Pittsburgh Chemical Information Center users to ascertain in-depth information about the ways in which chemists procure, use, and communicate scientific information. This information is obtained through a variety of sources including structured and unstructured interviews and feedback cards. The latter are used in the compilation of statistics concerning relevancy of the

retrieved information. Finally, there is a marketing group, which will be responsible for the sale of the information retrieval service in the future.

The computer used in this study was an IBM 360 Model 50 at the University of Pittsburgh Computer Center with 128k bytes of fast core storage, 1024k bytes of read-only storage, a 2314 disc drive, 5 magnetic tape drives, a card reader and punch, and a high speed printer. The programs were processed using the PCP (Primary Control Program) version of IBM's operating system.

During these studies neither multiprogramming nor spooling was available. Consequently, computer costs include tape set-up times, I/O processing, and searching at the flat rate of $200 per hour.

This article describes in detail two computer systems used to search *CA Condensates*—the CAS Condensates Search System (CASCON), and the IBM TEXT-PAC system.

## THE PROGRAMS

CAS supplies the CASCON search system without charge to *CA Condensates* tape subscribers upon request. CAS considers the CASCON system to be only a basic set of programs to be used with the *CA Condensates* tape. The source programs are in Basic Assembly Language (BAL) and contain approximately 7000 cards. The search program (plus Operating System) requires approximately 77k bytes for storage of the object program and will utilize all available core when loading search terms into memory.

A flow chart of the CASCON system is given in Figure 1. Input to the programs include an issue of *CA Condensates* and the interest profiles of the users. The program edits the profiles for validity, creates a profile table in the computer memory, and searches the entire tape, character-by-character, matching characters on the *CA Condensates* tape with entries in the profile table. Options must be specified in the search profile which permit the user to search author, coden, and/or title and key words. Records on the tape that match a given profile term are called alerts. The alerts for each profile are then sorted into sequence by question number, weight (i.e., significance of each alert as determined by the number of matches encountered between the document and

profile), and digest number. The sorted profiles are then printed and distributed to the user.

The TEXT-PAC programs were originally written by IBM[4] to search technical information on an in-house basis. They have versions that are compatible with either IBM 7090 or 360 series computers. TEXT-PAC represents a more general information retrieval system than CASCON. A conversion program is needed to translate the data base to be searched into the required TEXT-PAC input format. The source decks contain approximately 16,000 cards, and the programs require 128k of memory (the conversion program is the only program which requires 256k).

Flow charts of the TEXT-PAC system are shown in Figures 2, 3, and 4. For clarity, many of the TEXT-PAC options have been eliminated from the flow charts. However, these options have been itemized in the summary of the capabilities of both systems (Table I). The conversion of the *CA Condensates* data base into TEXT-PAC search format is illustrated in Figure 2. Two tapes are produced as output of the conversion subsystem— a searchable TEXT-PAC data base, where each word in the data base is sorted alphabetically by word length; and a condensed text tape, in readable context, used by
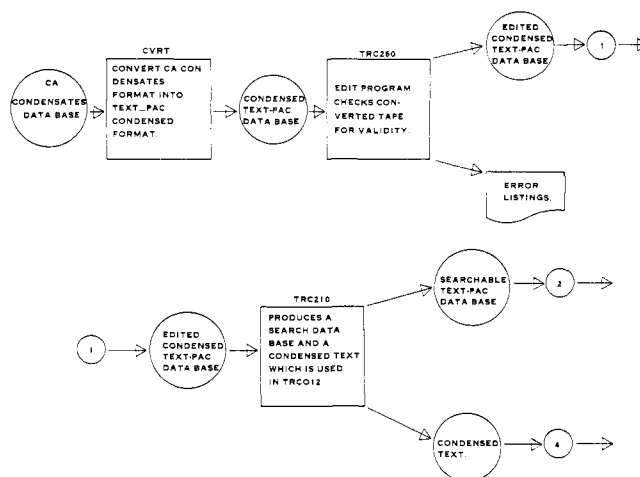


Figure 2. Flow chart of *CA Condensates* to TEXT-PAC base conversion
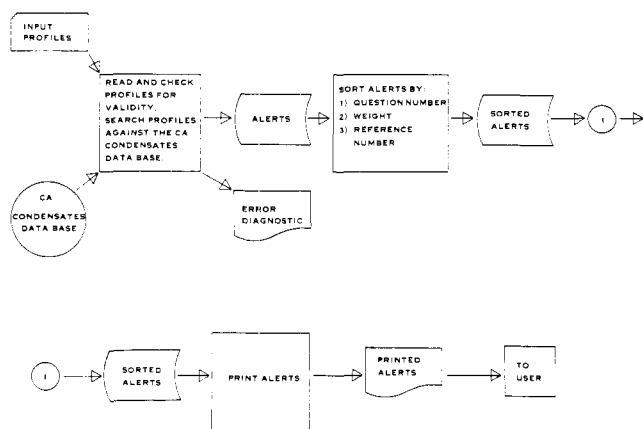


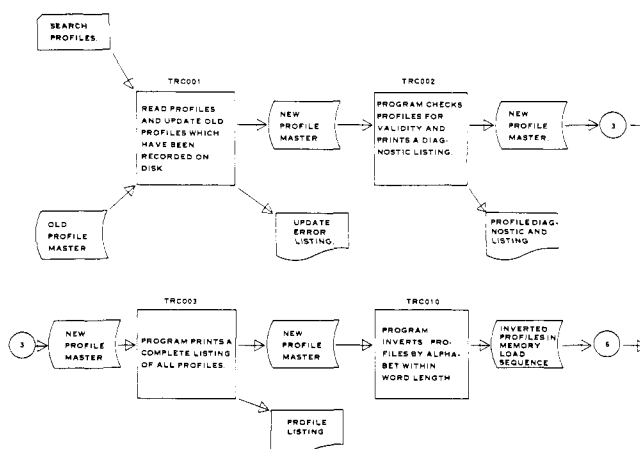Figure 1. Flow chart of CASCON processing



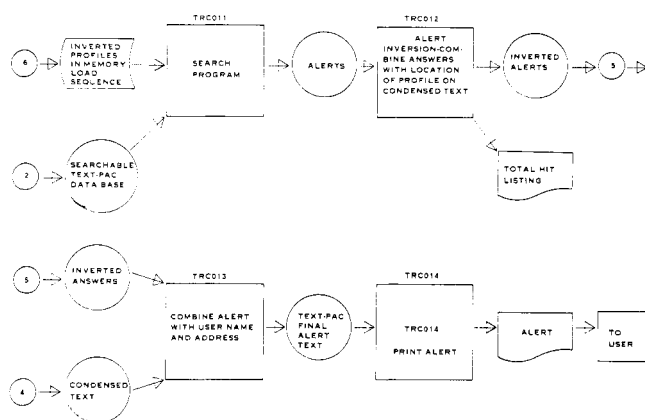Figure 3. Flow chart of TEXT-PAC profile preparation

Figure 4. Flow chart of TEXT-PAC search processing

the program which reformats the TEXT-PAC output (TRC013).

Figure 3 describes the handling of profiles. The profile update program (TRC001) allows the user to store profiles on tape or disk and to add, change, or delete selected items as desired. The profiles are edited, listed, and sorted alphabetically within word length.

As shown in Figure 4, these profiles and the converted *CA Condensates* data base are the input to the search program (TRC011). The resulting alerts are sorted according to profile number, merged with the condensed text data, printed, and sent to the user.

Table I summarizes the capabilities of both the CASCON and TEXT-PAC search systems. In general, the TEXT-PAC system seems to provide the user with more options, from the standpoint of profile logic as well as program capabilities. TEXT-PAC allows the user to specify nested logic. TEXT-PAC can also selectively search a wider variety of records on the data base including searches by date, location, and security classification. TEXT-PAC also provides greater variety in the formatting of output, allowing for 3 × 5 cards, bibliographic citations, and 80 column continuous form cards, printed two cards across.

Auxiliary reports are also available in the TEXT-PAC system, such as a word frequency report which shows the total number of occurrences of each word in the data base as well as a count of the number of documents in which each word occurs.

A dictionary compare program allows the TEXT-PAC user to check the correctness of spelling of words in the data base.

## COMPARISON OF THE CASCON/TEXT-PAC CURRENT-AWARENESS SYSTEMS

Both IBM and CAS provided the search programs on a program tape in source and object format, accompanied by documentation providing instructions on how to implement the system.

The documentation supplied with the TEXT-PAC system (although in draft form) was more complete than the documentation supplied with the CASCON system. TEXT-PAC was received as a pre-release of a type-three program from IBM. The availability of the TEXT-PAC system to a prospective user under the new IBM

Table I. Summary of the Capabilities of the TEXT-PAC and CASCON Systems

| Description of Features | Available with CASCON System | Available with IBM TEXT-PAC System |
|---|---|---|
| Profile update program available | | X |
| User statistical information | | |
|   Word frequency distribution | | X |
|   Alerts per profile | X | X |
|   Total number of alerts | X | X |
|   Total records scanned | | X |
|   Flag terms which caused alert | | X[a] |
|   KWOC index | | X |
| Edit program capabilities[b] | | |
|   Check profile logic | | X |
|   Check profile content | X | X |
|   Print profiles | | X |
| Output capabilities | | |
|   Bibliographic citation | X | X |
|   80 column cards, 2-up | | X |
|   3 × 5 cards | | X |
| Program available for feedback of user information | | X |
| Profile composition capabilities | | |
|   Left truncation | X | |
|   Right truncation | X | X |
|   AND/OR logic | X | X |
|   Absolute | | X |
|   NOT logic | X | X |
|   WITH logic | | X |
|   Option to search by | | |
|     Author | X | X |
|     Title and key words | X | X |
|     Dates | | X |
|     Corporate author | | X |
|     Locations | | X |
|     Journal coden | X | X |
|   Weighting capability | X | |
|   Nested logic capability | | X |
|   Adjacent logic | X | X |
|   Security check feature (for confidential information) | | X |
|   Upper or lower case capability | | X |

[a] Available in retrospective search only. [b] The CASCON edit program is part of the search program. This means that all acceptable profiles *must* be searched. The TEXT-PAC system, by contrast, allows the edit programs to be run independently of the search.

"unbundling" policy is still to be resolved. Both the CASCON and the TEXT-PAC systems required two to three months to implement. This time, however, could have been reduced considerably had it not been necessary to work in an over-the-counter environment in the University of Pittsburgh Computer Center.

Data set names used by the organizations who wrote the original programs had to be changed to correspond to those familiar to the Operating System (OS) of the University of Pittsburgh. This presented a potential hazard in that, if a data set name is mispunched during the transition, no OS diagnostic is given, yet the program will not run correctly, and considerable time can be devoted to locating this problem.

An effort was made to improve processing efficiency

of the CASCON and TEXT-PAC search programs in an effort to reduce costs and extend the number of user options. To achieve this the following features were added: The CASCON print program was modified to record the citations that are normally printed on to a magnetic tape (7 track, 556 BPI). This tape was then printed off-line on a small IBM 1401, available at a substantially reduced charge, rather than printing the tape on the 360/50. This has reduced processing costs substantially (15–25%). A similar modification will be made shortly for the TEXT-PAC system.

Further, this program was written to print the output on continuous form cards, two cards across, rather than on stock paper. Identical information is printed on both of the two cards, which were both distributed to the user. The user is requested to return one of these cards after completing some questions, pre-printed on the card, concerning relevancy. This information will be read back into the computer for feedback studies. This type of output should prove useful to the members of the Chemical Information Center involved in user feedback studies and also provide the user with a more easily handled (and filed) citation. Also, additional coding has been added to the CASCON programs to provide the user with key words as part of the printed output. This capability was available with the TEXT-PAC system upon receipt of the system.

Both TEXT-PAC and CASCON have been modified to cause all processing messages to print on the printer rather than on the computer console. This permits all processing messages to be printed on one device in a chronological manner, which is helpful in debugging. This is particularly important in an over-the-counter processing environment where the user does not have hands-on access to the computer or to the console log. Service to users has been improved by separating the user profiles according to odd-even tape preference. Prior to this all profiles were combined and processed against both the odd and even issues of *Condensates*.

Some explanation is required to describe the variables relating to the analysis of the two systems. The input profile format of each search system differs. As a consequence, a program was written to automatically translate the *Condensates* profiles into TEXT-PAC profile format. This assures as nearly as possible that the same profiles (and consequently similar search logic) will be used by both search systems. Since TEXT-PAC does not handle left-hand truncation, this option was removed from both the TEXT-PAC and CASCON input.

Wherever possible, programs which are not available in both systems have been eliminated in making this comparison. For example, one of the TEXT-PAC edit programs prints a listing of the profiles. This program is not available in the CASCON system and therefore was eliminated from the timings altogether. In addition, the print timings for both systems have been eliminated for three reasons: to reduce test costs, because the print costs would be fairly constant for both systems, and most of the printing for production runs would be handled by an off-line printer.

To establish a constant data base for both systems with the same number of input records, the same issue of *CA Condensates* (volume 70, issue 16) was used as input for all tests. The *CA Condensates* issue contains 5413 records. The reformatted TEXT-PAC issue contains 4992 records. The difference in the number of records is due to the fact that TEXT-PAC converts the *CA Condensates* file into a more efficient TEXT-PAC search format. During this process, TEXT-PAC eliminates any records which are not in acceptable TEXT-PAC format. This procedure is currently eliminating approximately 10% of the original file which, of course, is a serious factor as far as relevancy is concerned. This difference in the number of abstracts available to the search programs will also affect the processing times in favor of TEXT-PAC. This factor has been normalized in all timings included in this document. Improvements have been made subsequent to the timings presented in this study which will restore records discarded by the conversion program.

The CASCON search program uses a character-by-character search. Generally, this type of search has proven to be extremely slow in processing. This is particularly true when one considers that the data base being searched is serial in nature. The profile edit (checking) program in the CASCON system is built into the search program. This means that it is not possible to edit the profiles without going directly into the search phase. With such an arrangement, profiles which are in error could not be corrected prior to the search phase. Under the CASCON system, erroneous profiles are merely ignored, or if serious discrepancies arise, the search is cancelled. This proves to be a costly arrangement from the point of view of machine time used.

The TEXT-PAC system has done away with many of the disadvantages outlined above. For example, TEXT-PAC searches on full words rather than on a character-by-character basis. This is done through the *CA Condensates*–TEXT-PAC conversion programs that put the *CA Condensates* data base into TEXT-PAC format and at the same time inverts the terms in each document in the data base alphabetically within word length. This procedure speeds up the search significantly as the search program merely looks through the list of words in each document of the same length as the term specified in the profile, rather than searching all terms in the document. At the same time, TEXT-PAC does search all of the data in the document rather than searching only selective fields. However, the structure of the inverted file precludes the availability of left-hand truncation in the search profile (right-hand truncation is available). Another disadvantage of the TEXT-PAC system is the absence of weighting of the search terms.

Table II shows the processing costs for various quantities of input profiles for each system. These results have been obtained using an average of 12 search terms per input profile. The term "cost per profile" is calculated by dividing the total processing costs by the number of input profiles. "Costs per alert" is computed by dividing the total processing costs by the number of alerts received as output from each set of programs.

It is interesting to note that, while costs per profile decrease for larger input volumes, costs per alert have increased slightly for both the CASCON and TEXT-PAC systems. It is felt that this slight increase in cost per alert is due to the reduced number of alerts produced compared to the increase in the number of input profiles.
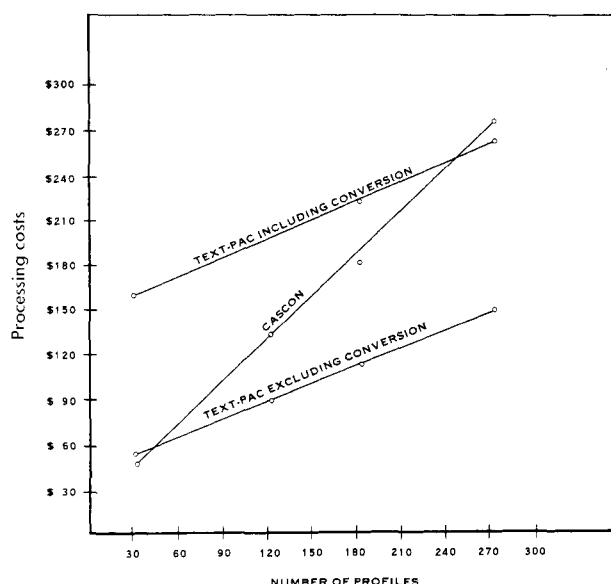
Figure 5. Comparison of CASCON-TEXT-PAC processing costs

For example, the average number of alerts retrieved (TEXT-PAC) when processing 120 profiles was 23.4 per profile while the number of alerts retrieved at the 180 profile level was only 17.2, and 12.3 for 270 profiles. The average number of alerts retrieved per profile decreases in the same manner for the CASCON system.

It should be noted that these timings are for an even issue of *CA Condensates*. It is expected that the cost per profile would decrease similarly for an odd issue, but the cost per alert behavior may change due to the difference in the subject matter.

Timings were conducted for 30, 120, 180, and 270 input profiles for both the CASCON and TEXT-PAC systems. These results are presented in Figure 5. While multiple timings were conducted for each input level, the lowest values have been presented here, using the rationale that, since the processing was performed in an over-the-counter environment, the only variation between timings at a given volume of input profiles would be the result of set-up time (since the content of the profiles and data base remained constant).

Two sets of costs have been used for TEXT-PAC— one including the processing costs required to convert the *CA Condensates* tape into TEXT-PAC search format and the other excluding conversion costs. Both sets of figures have been included because the conversion processing generates TEXT-PAC searchable tapes which are also usable in the TEXT-PAC retrospective search system. Hence, the costs for conversion processing could be added to the current-awareness costs if no retrospective search processing is anticipated by a prospective user of the system, included with retrospective costs, or included on a pro rata basis in the processing costs of both the current-awareness and retrospective systems.

## CONCLUDING REMARKS

Obviously, cost justification of TEXT-PAC compared to CASCON depends upon the distribution of TEXT-PAC conversion costs. This can be seen from Figure 5. Excluding the conversion costs, the break-even point comparing TEXT-PAC to CASCON is approximately 45 profiles. Including conversion costs, the break-even point is in the area of 240 profiles. Splitting conversion costs between both TEXT-PAC current-awareness and retrospective systems places the break-even point at approximately 140 profiles. Costs per input profile as well as costs per unit of output (alert) are indicated in Table II.

TEXT-PAC looks more attractive cost-wise either where a user anticipates continually large volumes of input, or where conversion costs can be amortized through retrospective processing.

Costs, however, do not tell the complete story. TEXT-PAC seems to have a wider variety of user options available in the system which, to some extent, may serve to offset small differences in processing costs. In the TEXT-PAC system, for example, an optional word frequency report is available which indicates both the total number of occurrences of a word in the data base as well as the total number of documents in which each word occurs. This type of information can prove useful in the preparation of input profiles. TEXT-PAC also provides greater variety in the formatting of printed output.

The emphasis in this study has been on a cost comparison between the two search systems. The comparative

Table II. TEXT-PAC and CASCON Comparative Processing Costs

| | TEXT-PAC | | | | | | CASCON | | |
|---|---|---|---|---|---|---|---|---|---|
| | EXCLUDING TAPE CONVERSION COSTS | | | INCLUDING TAPE CONVERSION COSTS | | | | | |
| Number Of Profiles | Cost Per Profile | Number Of Alerts Retrieved | Cost Per Alert | Cost Per Profile | Number Of Alerts Retrieved | Cost Per Alert | Cost Per Profile | Number Of Alerts Retrieved | Cost Per Alert |
| 30 | $ 1.703 | 85 | $ .601 | $ 5.336 | 85 | $ 1.883 | $ 1.623 | 866 | $ .056 |
| 120 | .745 | 2868 | .031 | 1.653 | 2868 | .069 | 1.110 | N/A | N/A |
| 180 | .625 | 3094 | .036 | 1.231 | 3094 | .072 | 1.004 | 3963 | .046 |
| 270 | .556 | 3342 | .045 | .960 | 3342 | .078 | 1.029 | 4899 | .057 |

number of retrievals has been given, but the relevancy of the retrievals has not been indicated. A group within the Pittsburgh Chemical Information Center is analyzing user feedback. They have found that for CASCON users the retrievals are approximately one-third relevant. The corresponding figure for TEXT-PAC users and data relating to the effect of no left truncation and no term weighting in the TEXT-PAC system will be reported at a later date.

## ACKNOWLEDGMENT

## LITERATURE CITED

(1) Arnett, E. M., "A Chemical Information Test Station," Chemistry, 42 (3), 16 (1969).
(2) Bloemeke, M. J., and S. Treu, "Searching Chemical Titles in the Pittsburgh Time-Sharing System," J. CHEM. DOC. 9, 155 (1969).
(3) Freeman, R., J. Godfrey, R. Mainzell, R. Rice, and W. Shepard, "Automatic Preparation of Selected Title Lists for Current Awareness Services and Annual Summaries," J. CHEM. Doc. 4,107 (1964).
(4) Esposito, A. V., R. Fleischer, S. D. Friedman, S. Kaufman, S. Rogers, S. Skye, M. Shotkin, "TEXT-PAC, S/360 Normal Text Information Processing, Retrieval, and Current Information Selection Systems 360d-06. 7. 020," IBM, Armonk, N. Y., December 1968.

# Searching the Current Chemical Literature by Computer*

HOWARD P. ANGSTADT
Sun Oil Co., Marcus Hook, Pa.   19061

Received July 1, 1970

**This paper presents the results obtained from a study of the efficiency and reliability of two currently available computerized current awareness services. Once proficient profiles were developed, virtually no difference was seen in the precision and recall values obtained from using either Chemical Titles or CA Condensates as data bases.**

There is no longer any doubt that the increasingly complex nature and rapid growth of the technical literature necessitates a new approach to the problem of technical awareness. Historically the scientist is confronted with two distinct problems concerning the literature—namely, obtaining the results of the previous research effort in a given area (information retrieval) and being aware of current progress in his chosen specialty (current awareness). This paper reports the results we obtained from an experiment designed to evaluate the relative efficacy of computerized searching of Chemical Titles and CA Condensates as current awareness tools.

Our approach to the problem of evaluating computerized current awareness instruments consisted of having several researchers scan by their normal procedures the same material that was computer searched and then compare the two outputs. More specifically, several individuals were asked to formulate precisely in single questions several of their research interests. Each statement was carefully keyword coded into sets of words called parameters that in turn were linked together by AND, OR, or NOT logic according to the techniques given in "Preparation of Search Profiles," a publication of the Chemical Abstracts Service, published by the American Chemical Society. Prior to searching, the entire profile was checked with

the author of the question to ascertain its fidelity. After being coded into the necessary computer format, each profile was compared to the data base for the time period under study, and the titles of all matching articles were printed out by the computer. Simultaneously, but without the benefit of the profiles, each researcher was asked to evaluate the same data base with respect to the questions he posed and subsequently to do the same for the computer printout.

The information we sought was the number of articles uncovered by the researcher scanning the data base in his normal way that were judged by him to be pertinent to the specified research area; the number of these articles also located by the computer; the number of titles located by the computer but not found by the researcher yet judged by him to be interesting enough to look up the original article; and the total number of articles located by the computer.

The problem of deciding just what constitutes a pertinent article relative to a specific question is quite difficult, and several approaches have already been suggested.[1-4] We relied solely upon the judgment of the researcher submitting the question being fully aware that certain problems are inherent in this approach. For our purposes it was not realistic to evaluate the pertinence of an article by means of a panel, nor did we distinguish between degrees of relevancy. Our sole criterion of a pertinent article was the individual's on the spot judgment that