

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

Financial support for this work has been received from the National Science Foundation under Grant GN 738. This work was initiated by E. M. Arnett, Professor of Chemistry and Director of the Pittsburgh Chemical Information Center, and Allen Kent, Director of the

Knowledge Availability Systems Center, University of Pittsburgh.

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

* Presented before the Chemical Documentation Section of the Fifth Middle Atlantic Regional Meeting, ACS, University of Delaware, Newark, Del., April 1, 1970.

Table I. *Chemical Titles*—Search Results

Profile Number	Computer Drops	Searcher ^c Located Items	Searcher + Computer ^b Located Items	Computer ^a Located Items	Total ^b Literature	Precision	Recall	False Drops
5	4	11	4	0	11	1.00	0.36	0
6	1	3	1	0	3	1.00	0.33	0
7	0	6	0	0	6	0
8	0	4	0	0	4	0
9	0	0	0	0	0	0
10	4	3	1	0	3	0.25	0.33	3
11	0	14	0	0	14	0
12	3	22	0	1	23	0.33	0.04	2
13	6	9	1	0	9	0.17	0.11	5
14	6	9	4	1	10	0.83	0.50	1
15	13	11	11	1	12	0.92	1.00	1
16	0	4	0	0	4	0
17	3	6	0	0	6	3
18	0	3	0	0	3	0
19	6	19	0	0	19	6
20	0	6	0	0	6	0
21	6	9	6	0	9	1.00	0.67	0
22	11	0	0	2	2	0.18	1.00	9
23	9	2	1	1	3	0.22	0.66	7
24	2	6	1	1	7	1.00	0.28	0
25	0	4	0	0	4	0
26	0	2	0	0	2	0
27	4	4	2	0	4	0.50	0.50	2
28	4	5	0	3	8	0.75	0.37	1
Sub Totals	82	162	32	10	172	0.51	0.24	40
29 ^e	613	28	22	14	42	0.06	0.86	577
Totals	695	190	54	24	214	0.11	0.36	617

^a This column gives the number of articles located by the searcher prior to computer searching. ^b This column gives the number of articles common to the individual's and computer search. ^c This column gives the number of articles located by the computer but missed by the searcher. ^d Searcher located items plus computer located, but searcher missed items. ^e This data is from a very broad profile on "Hydrocarbon Oxidation" and because of the large response has been separated statistically from the other profiles.

the title reproduced on the computer printout was of sufficient interest to merit investigating the original article. We feel the use of this criterion closely approximates the actual situation as the researcher scans journal tables of contents or current awareness tools such as *Chemical Titles* in a hectic effort to maintain some semblance of current awareness.

We evaluated our results in terms of the precision (the number of relevant documents selected over the total number of retrieved documents) and recall (the number of relevant and retrieved articles over the total literature pertinent to the question searched) parameters. The total literature pertinent to any research question is most likely a theoretical concept only as it is doubtful that the absolute number of articles published on any research question is ever truly known. Hence the denominator in the recall parameter contains an arbitrary element. As a first approximation we defined the literature as the total number of different pertinent articles found by both the computer and the individual researcher.

The results of our first attempt are shown in Table I. They represent the literature covered by three successive issues of *Chemical Titles*—i.e., the chemical literature that appeared in about 650 journals of interest to chemists over a six-week period. The researchers on the panel had interests in the areas of organic, analytical, physical, and catalytic chemistry. All of the profiles were prepared by the author; they contained on the average

2.5 parameters per profile with 4.2 terms per parameter.

The results varied from the extremes of two searches that had 100% recall (15 and 22) to 10 profiles that retrieved no pertinent literature although relevant articles had been identified by the individual searchers. Precision varied from 0 to 100%. Profile 29 concerned the very broad area of hydrocarbon oxidation, and because of the disproportionately large number of drops the results have been separated statistically from the rest of the data. Although the precision was quite low (6%), most of the literature (86%) deemed pertinent by the searcher was located. This search located 50% more articles of importance than located by the searcher himself and illustrates one of the greatest advantages of computerized searching techniques. Here also is an illustration of the problem of the amount of time the individual has available for reading the computer output versus the desired degree of literature coverage. Only three of the profiles captured 80% or more of the literature. The average value of precision was 0.51 and for recall it was 0.24.

At this point most of the profiles were revised. Those that had good recall values were changed little or not at all, while those that gave poor performance were extensively altered. In general, we reduced the number of parameters and expanded the number of terms within the remaining parameters. We hoped to increase the number of drops because if a profile is producing some output, even though it may not be relevant, one at least has

SEARCHING THE CURRENT CHEMICAL LITERATURE BY COMPUTER

Table II. *Chemical Titles*—Revised Profile Search Results

Profile Number	Computer Drops	Searcher Located Items	Searcher + Computer Located Items	Computer Located Items	Total Literature	Precision	Recall	False Drops
5	28	2	1	5	7	0.21	0.86	22
6	1	3	1	0	3	1.00	0.33	0
7	9	1	0	1	2	0.11	0.50	8
8	0	2	0	0	2	0
9	1	0	0	0	0	1
10	8	6	0	1	7	0.13	0.14	7
11	109	4	0	18	22	0.17	0.82	91
12	158	38	15	2	40	0.11	0.42	141
13	8	12	0	0	12	8
14	5	3	0	1	4	0.20	0.25	4
15	16	14	14	2	16	1.00	1.00	0
16	6	13	5	1	14	1.00	0.43	0
17	9	0	0	2	2	0.22	1.00	7
18	4	0	0	0	0	4
19	200	49	42	12	61	0.27	0.89	146
20	13	8	8	1	9	0.69	1.00	4
21	0	1	0	0	1	0
22	35	1	1	0	1	0.03	1.00	34
23	22	0	0	1	1	0.05	1.00	21
24	15	1	0	4	5	0.27	0.80	11
25	16	7	0	3	10	0.19	0.30	13
26	15	6	1	0	6	0.07	0.17	14
27	3	1	1	0	1	0.33	1.00	2
28	8	1	0	5	6	0.62	0.82	3
Sub Totals	689	173	89	59	232	0.215	0.64	541
29	484	26	21	24	50	0.09	0.90	439
Totals	1173	199	110	83	282	0.16	0.68	980

some idea how to improve the performance of the profile. If a profile generates no output whatever it is quite difficult to say how it might be improved.

Table II contains the results from searching three more issues of *Chemical Titles* using the revised profiles. This time the profiles averaged 2.3 parameters per profile containing 6.4 terms per parameter. The results again varied from six profiles with 100% recall to three that found no references although some pertinent literature was known to be present. Again precision varied from 0.0 to 100%. As expected the total number of drops increased almost eightfold (omitting profile 29 which had not been changed), and, as expected, the average value for the precision fell to 22%. However, the average value for recall more than doubled to 64%. This time 11 individual profiles found 80% or more of the literature. Profile 11, reduced from three to two parameters went from 0 to 109 drops, and the computer found four times as many pertinent articles as did the individual. Over all the profiles the computer found roughly 33% more articles than did the searchers themselves. In the case of profile 29, the computer located almost as many articles that the scientist had missed as he had found!

At this point in our study *CA Condensates* as a taped data base that covers over 7000 journals and is based upon a keyword index of *Chemical Abstracts* (CA) became available. Without any revision the profiles were scanned against *CA Condensates*. The data reported in Table III presents the results obtained from scanning six issues (1½ months of literature) of *CA Condensates* by the computer and the same six issues of CA by the individual researchers.

The researchers searched CA according to their own techniques and did not have their profiles available for guidance. We feel this procedure closely approximates the type of situation that might be expected to arise when computer current awareness services become widely utilized. Using this much broader input base we expected more drops and therefore a lower precision. However, insofar as a keyword index of an abstract is a better representation of the intellectual content of an article than is the title alone, we expected our recall to improve.

The results show that as expected the number of drops doubled and the precision dropped further to 18%. More importantly the average recall performance remained the same at 61%. This time 10 profiles found about 80% or more of the useful literature. Although profile 29 delivered over 1000 drops, it contained over 90% of the pertinent literature. Once again about 25% of the total literature was located by the computer and missed by the individuals.

Throughout this study we have attempted to maximize recall at the expense, if necessary, of precision, principally because we feel that the effort involved in reading a large number of false drop titles is small compared to the damage that could result from missing significant references.

To this point we had not used NOT logic in our profiles. The output we received led us to believe that some increase in precision could be obtained by use of this technique to reject many of the false drops obviously not of interest to the researchers. Since this format takes precedence over all other logic, one must use it judiciously to avoid

Table III. CA Condensates—Search Results

Profile Number	Computer Drops	Searcher Located Items	Searcher + Computer Located Items	Computer Located Items	Total Literature	Precision	Recall	Fase Drops
1	62	17	3	1	18	0.06	0.22	58
2	67	31	11	17	48	0.42	0.79	39
3	9	6	2	0	6	0.22	0.33	7
4	16	0	0	3	3	0.19	1.00	13
5	76	5	4	7	12	0.14	0.92	65
6	0	6	0	0	6	0
7	4	2	0	0	2	4
8	1	4	0	1	5	1.00	0.20	0
9	0	0	0	0	0	0
10	18	7	1	1	8	0.11	0.25	16
11	166	6	0	4	10	0.02	0.40	162
12	329	45	25	7	52	0.10	0.62	297
13	16	50	1	0	50	0.06	0.02	15
14	4	3	3	0	3	0.75	1.00	1
15	18	16	16	2	18	1.00	1.00	0
16	4	9	1	2	11	0.75	0.27	1
17	71	8	7	5	13	0.17	0.92	59
18	11	2	2	0	2	0.22	1.00	9
19	375	69	58	30	99	0.23	0.89	287
20	48	25	25	4	29	0.60	1.00	19
21	11	13	9	2	15	1.00	0.73	0
22	64	1	0	0	1	64
23	61	3	1	3	6	0.07	0.67	57
24	40	0	0	10	10	0.25	1.00	30
25	17	5	1	4	9	0.29	0.56	12
26	14	6	0	2	8	0.14	0.25	12
27	9	0	0	0	0	9
28	14	11	3	2	13	0.36	0.38	9
Sub Totals	1525	350	173	107	457	0.184	0.61	1245
29	1047	38	32	36	74	0.06	0.92	979
Totals	2572	388	205	143	531	0.135	0.66	2224

any chance of rejection of useful articles. The profiles were once more revised and to each was added a parameter of NOT terms. Again we scanned six issues of *CA Condensates*, and the results are tabulated in Table IV. In spite of our efforts, our precision decreased further to 11%! We now realize that this occurred because of an error in preparing the profiles in which the NOT terms were collected into a single parameter added to the bottom of the profiles. Proper procedure requires that NOT parameters must not come last! Our more recent results, which have not yet been tabulated, indicate that precision can be improved by use of this logic, but the extent of the improvement remains unknown. The recall value remained at 62%. Nine profiles recovered about 80% of the literature. Profile 15 maintained its perfect record of locating all the applicable literature with almost perfect precision; the hydrocarbon oxidation profile (29) continued to find 90% of the pertinent literature.

At the conclusion of this study terminal interviews were held with each participant to get some measure of the degree of acceptance of this procedure. In general the degree of enthusiasm for computerized searching techniques varied directly with the degree of success of the individual's profiles. Those who experienced high recall values justifiably had a high degree of confidence in the technique. These individuals felt that their current searching time could be considerably reduced by this service, perhaps by as much as 75%. The majority of par-

ticipants felt that computerized searching for current awareness would be a very useful compliment to their current program of literature searching. If the service were to be provided, these people felt they could reduce their own searching time by 25 to 50% but would definitely not eliminate it entirely. A realistic operational procedure which was thought might evolve was that the average researcher would continue to scan perhaps the six most important journals in his area and rely upon computer techniques to cover the rest of the literature. Most of the participants felt that during the experiment they were getting tremendous literature coverage; however, the participants appeared to be improving their own ability to search the literature, as the percentage of computer-located, author-missed items (column 5) decreased from 29% in Table II to 27% in Table III and finally to 18% at the end of the study. Some individuals whose profiles were very unproductive still valued the computerized search since it gave them, perhaps unjustifiably, a very confident feeling that nothing of importance had been missed. Although these profiles did not drop out most of the existing literature, they still turned up useful articles missed by the scientist, and it is this ability of computer searching to locate remote articles, even though the profile itself is not very efficient nor does it recover most of the literature, that makes the service appealing to careful scientists.

Although *CA Condensates* provides an unparalleled view

SEARCHING THE CURRENT CHEMICAL LITERATURE BY COMPUTER

Table IV. CA Condensates—Revised Profile Search Results

Profile Number	Computer Drops	Searcher Located Items	Searcher + Computer Located Items	Computer Located Items	Total Literature	Precision	Recall	False Drops
1	22	15	4	4	19	0.36	0.42	14
2	81	24	11	13	37	0.30	0.65	57
3	6	2	1	0	2	0.17	0.50	5
4	23	2	1	3	5	0.17	0.80	19
5	52	10	5	1	11	0.12	0.55	46
6	1	2	1	0	2	1.00	0.50	...
7	0	2	0	0	2
8	0	3	0	1	4	...	0.25	...
9	0	0	0	0	0
10	11	10	2	2	12	0.36	0.33	7
11	116	8	1	1	9	0.02	0.22	114
12	976	29	23	12	41	0.03	0.85	941
13	27	19	1	0	19	0.04	0.05	26
14	30	2	0	0	2	30
15	19	17	17	1	18	0.95	1.00	1
16	6	10	4	1	11	0.83	0.45	1
17	26	6	0	1	7	0.04	0.14	25
18	11	0	0	0	0	11
19	390	78	57	10	88	0.17	0.76	323
20	47	22	20	2	24	0.47	0.92	25
21	26	15	5	1	16	0.23	0.37	20
22	41	2	0	0	2	41
23	53	1	0	4	5	0.08	0.80	49
24	46	2	1	6	8	0.15	0.87	39
25	19	11	5	2	13	0.37	0.54	12
26	27	0	0	1	1	0.04	1.00	26
27	14	0	0	0	0	14
28	21	5	0	0	5	21
Sub Totals	2091	297	159	66	363	0.108	0.62	1866
29	970	55	48	11	66	0.06	0.89	911
Totals	3061	352	207	77	429	0.092	0.66	2777

of the total published literature, it has one major drawback in that it is not too current. The compilation of CA takes considerable time and, therefore, the data obtained from CA Condensates may be over a year old. This shortcoming is minimized when the individual does some current awareness searching by himself—i.e., the half dozen most important journals in his area. In contrast, the very wide coverage of the world's journals by CA gives the user a very confident feeling that no significant articles have been missed if he has developed an efficient profile. Computerized scanning using *Chemical Titles* as an input base assures very up-to-date results with a good profile since the journal often has access to tables of contents of the issues it covers prior to their publication. Of course with this input one is searching a smaller, more select group of journals. However, the fact that our recall remained the same for both *Chemical Titles* and CA Condensates searching implies that the additional keywords obtained from the abstracts in the latter case disclose very little additional conceptual material beyond that already contained in the title.

The essential requirement to use computer searching techniques for current awareness remains the construction of a proficient profile. That this can be done no longer remains in doubt. That this can *always* be done does remain in doubt. The techniques of proper profile construction can stand considerable refinement, but when a good profile is developed the results obtained from current

awareness searching are excellent. Another searching technique expected to give very efficient and accurate results but one with which we have had minor experience, is the use of author scanning. If the individual researcher keeps an author file of the important investigators in his field, their publications can be pinpointed with ease and efficiency by computer scanning techniques.

It seems quite reasonable to suppose that a professional literature scientist could easily increase the proportion of satisfactory profiles with more refined techniques such as the use of weighted terms. We therefore feel that computerized current awareness services will provide the practicing scientist with a very powerful weapon in his battle against the continually expanding technical literature.

REFERENCES

- (1) Kent, A. B., "United Kingdom Experiences in the Operation of a Retrieval and Dissemination Service Based on CAS Search Tapes," paper presented before the Chemical Literature Division, 156th Meeting, ACS, Atlantic City, N. J., September 1968.
- (2) Lancaster, F. W., "MEDLARS: Report on the Evaluation of Its Operating Efficiency," *Amer. Doc.* 21, 119-42 (1969).
- (3) Savage, T. R., "The Interpretation of SDI Data," *Amer. Doc.* 18, 242-46 (1967).
- (4) Wagner, R. H., "A Selective Current-Awareness System Using Engineering Indexes Plastics Data Base. II. Performance," *J. Chem. Soc.* 9, 85-8 (1969).