## "Search Strategies Using the Science Citation Index"®

November 4, 1969

There is considerable literature concerning search strategies--how to search the literature so as to minimize the time and effort required. Many searchers are not aware that the methods one may employ to retrieve a few references in a "quick and dirty" search may be quite different than methods employed when the search must be comprehensive. Of course, no search is absolutely complete. One must always make economic decisions as to the time or money that can be invested. Since the working scientist is usually pressed for time it is important to learn how to select and use searching tools to turn up the maximum number of pertinent references in the minimum time. A few years ago Spencer (1) showed, in comparison to *Index Medicus* and *Chemical Abstracts,* that the first hours of searching the *SCI*® were far more productive. Subsequently, A.E. Cawkell (2) wrote a

paper on search strategies. His paper is reprinted in this issue of *CC*® to provide some useful examples. Reprints of this paper are available on request, as are the other items cited (3,4).

(1)   Spencer, C.C., "Subject Searching with *Science Citation Index;* Preparation of a Drug Bibliography Using *Chemical Abstracts, Index Medicus,* and *Science Citation Index* 1961 and 1964", *Am. Doc.* 18(2), 87-96, (1967).

(2)   Cawkell, A.E., "Search Strategies Using the *Science Citation Index,"* *Computer Based Information Retrieval Systems,* edited by Bernard Houghton, (Clive Bingley Ltd.), (1968).

(3)   Malin, M.V., "The *Science Citation Index:* A New Concept in Indexing", *Library Trends,* 16 (3), 374-387, (1968).

(4)   Garfield, E., "Citation Indexing: A Natural Science Literature Retrieval System for the Social Sciences", *Am. Behavioral Scientist,* 7 (10), 58-61, (1964).

# SEARCH STRATEGIES USING THE SCIENCE CITATION INDEX
## A.E. CAWKELL
UK and European Consultant to the
Institute for Scientific Information*

The systems to be described are derived from a common data base consisting of a magnetic tape composed each week from information gleaned from some 1,800 prime journals received from about forty five different countries, and covering all the major areas of science and technology; the information relates not only to articles, but also to editorials, corrections, letters to the editor, and so forth.

The manifest terms used to index these items include the journal title, volume, page and year, the words in the title, the names of the authors, the name of the place at which the authors work, and a complete list of any references included with articles.

The rate of processing, and the intellectual ability required, are of a different order to that required when texts are read and indexing terms are allocated by subject experts (in the hopes that others will use the same terms when seeking the article). However, the extraction of information from a huge number of journals with many different editorial practices, different formats, and in different languages is no mean task.

The task may be undertaken by persons having no subject know-ledge at a relatively fast rate. A very large number of articles may be processed with a short delay between first receipt of an article, and the appearance on tape of information about it; it becomes feasible not only to process a sufficient number of journals in any subject area in order to provide adequate coverage of that subject, but also to process journals covering many different subject areas. This has certain retrieval consequences best illustrated by an example.

Say a research engineer is investigating a new type of television device. He will be served with information about not only the ' litera-ture of television and electronics ', but also with information about the large and ill-defined ' literature of interest about television and electronics ', including, for instance, the important area of the subject reported upon in optical journals. This multi-disciplinary coverage without barriers is of obvious importance in many subjects.

With regard to indexing, the usefulness of references is less obvious than the other tags or terms which have been mentioned.

Given that we have, week by week, a magnetic tape carrying infor-mation about a large number of items, a computer loaded with that

*Reprinted from *Computer Based Information Retrieval Systems*, edited by Bernard Houghton (Clive Bingley Ltd.), (1968)

tape may be instructed to reorder and print out the stored information in some desired manner. For instance, the instruction may be ' print out a list of those current articles which have a particular reference in common, beneath that reference '.

In figure 1a the result of such an instruction is shown. Observe that two out of many thousands of recorded articles have been selected, because in each of these two there was a reference to the 1959 article by Kolin.

The computer may then be instructed to print out information about the citing articles in more detail (figure 1b).

If a man is engaged on the design of an electro-magnetic blood flowmeter and transducers, knows about Kolin's 1959 article, and wishes to be informed about current articles, he is led forward from the 1959 article, which he equated with the subject of interest, to two 1967 citing articles.

This is an example of the idea of retrieving information via the conceptual relationship between a cited and citing article; it is one of the methods used in certain end products and services to be described, derived from the magnetic tape; some answers will be provided to the questions ' how useful is this idea for information retrieval? Do authors habitually cite a proper selection of the prior art? Are the current articles so found likely to be relevant? How long does it take to carry out a search in this manner? '.

However before dealing with systems and their effectiveness it is of interest to ask ' how can a particular information retrieval system be evaluated and compared with others? '. One answer to this question has been given recently by Cooper,[1] which seems to have considerable merit in a real situation.

Cooper suggests that a useful criterion is the amount of time saved when using a particular retrieval system, as compared with the time spent during a random search of a document collection. He also suggests that an often made assumption that a collection may be divided into two categories—those documents which are relevant, and those which are not—is not very realistic. In reality any collection may be ranked in many levels of relevance and interest. The capability of a retrieval system in drawing attention to some arbitrary number of upper levels should be considered. He suggests assessment of the average expected search time per desired relevant document or, alternatively, the time taken to screen out unwanted documents for each relevant document found. Some results expressed in this way will be given for the Institute for Scientific Information's (ISI) systems.

The first end product to be discussed is the Automatic Subject Citation Alert (ASCA) SDI service operated by the ISI. A weekly printout is provided about current articles according to the standing instructions given in a subscriber's profile. The magnetic tape is searched, item by item, for those articles tagged with the terms listed

in the profile. Figure 2 shows a profile circumscribing the subject 'Design of electromagnetic blood flowmeters', and also an example of a printout showing information about articles retrieved via the profile.

This profile includes various types of terms to exemplify usage; the printout example lists information about actual current articles retrieved. The number in brackets indicates relevance in the opinion of a subscriber, scaled from 1 to 5 (1 being highly relevant). It usually costs about £40 to run a profile for one year, but the cost will depend on the number of terms necessary to circumscribe adequately the subject of interest.

To assist in the selection of words in title, ISI publish a high frequency word list, part of which is shown in figure 3; the information content and frequency of occurrence of words in titles may be deduced from it. A low priced or unlisted word, when used as a 'word in title term' on its own, will retrieve, on average, a very small number of articles each week. A more costly high frequency word will retrieve a large number of articles, many of which may be of no interest because the word is not specific. The word costs in the separate 'combo' column reflect the lower probability of occurrence and higher information content of word pairs.

Three years of ASCA experience shows that the wide journal coverage, lack of inter-disciplinary barriers, and availability of several different types of term, combine to provide an effective, timely SDI system. A subscriber's co-operation is essential if the profile is to be effective; unless he knows something about literary practice in his discipline, is prepared to devote care to profile compilation, and is willing to circumscribe his subject properly and adequately, then his weekly printout is likely to exclude relevant articles (profile too specific) or to include excessive noise (profile too broad). Often a profile can be improved upon in the light of results; feedback leading to a modification of the profile may be necessary, and again this is up to the subscriber and is provided for in the system.

Some tests with ASCA on chemical subjects have recently been reported upon by Abbot.[2] ASCA performed well, particularly when the subject had inter-disciplinary aspects, but it should be noted that neither cited article terms, word stems, nor floating word stems were included in the results.[3]

To conclude this brief discussion of ASCA, some mention must be made of the significance of noise, as this is not always seen in its true perspective when the so-called 'precision' of information systems is described. It may be necessary to accept a high proportion of noise, in the knowledge that the profile broadening which brings it about is also bringing to light a high proportion of the published 'relevant',

and ' of interest ' articles. An ASCA printout can be read and assessed at a rate of about ten article entries per minute. Hence the possibly substantial benefit of being informed about, say, one additional relevant article per week, may be accompanied by the trivial penalty of a list of a few unwanted articles which may be quickly ignored.

ASCA is for keeping up.[4] The second end product derived from the magnetic tape file is for catching up. Reference is made to figure 4, which shows a part of the Citation index, Source index, and Permuterm Subject index, collectively called the *Science citation index* (SCI).[5]

The Source index section is an author-ordered list of articles for a particular calendar year.

The Citation index section is an ordered list of the references from these articles.

The articles from the Source index which have a reference in common are listed beneath that reference in the Citation index.

The cited references, spreading backwards across the centuries, are those published scientific works which today's scientists consider to be worth citing. Because of the accelerating growth of science a high proportion of the prior art, cited in a current year, lies within the preceding five years.

It follows that if a searcher enters the SCI at some item of the prior art identified with a particular subject, then he will be led forward from that item to the current citing articles listed beneath it. This is the simplest form of search strategy using the Citation index

If the citing articles are obtained, selected references from them may be used as new citation index entry points. This process, known as ' cycling ', again brings the searcher forward in time to a crop of current citing articles. A network of articles inter-connected by references is built up.

The Permuterm subject index section is an index to the words in the titles of all the articles in the Source index. All possible pairs of words are permuted and ordered in the manner shown in figure 4c. Certain very high frequency words such as ' is ' and ' and ' are excluded.

A search will be described (see figure 5) using these three sections of the SCI.

The subject is ' The design of electromagnetic blood flowmeters '; the starting points used are the words ' flow ' and ' square-wave ' (this word refers to a function in a class of instruments of interest). The search initiator also suggested one reference which he identified with the subject—an article by Kolin in the *Proceedings of the Society of Experimental Biology and Medicine,* 1941 v 46 page 235.

This article by Kolin, number 1 on the diagram, was used as an entry point in the 1967 SCI, but nobody cites it in 1967. The article was obtained, and from the references given in it, articles were selected

```
                KOLIN A    59    P NAT ACAD SCI   45   1312
Figure 1a       BENFIELD JR  DIS CHEST    67   52    321
                KHOURI   EM  J APPL PHYSL 67   23    395
```

(Ordering of citing current articles beneath the common
cited item)

---

```
   BENFIELD JR  COON R   CREE EM
      DIS CHEST   52   321  67   18R   N3   99580
         Current methods in canine pulmonary research in-
         cluding description of improved bronchiospirometry
         tube.

   KHOURI EM   GREGG DE
      J APPL PHYSL   23   395   67   9R   N3   99592
         An inflatable cuff for zero determination in
         blood flow studies.
```

Figure 1b   (Details of current citing articles)

---

```
   KOLIN A                 SOURCE AUTHOR
   KOLIN A                 REFERENCE AUTHOR
   WYATT DG                SOURCE AUTHOR
   WYATT DG                REFERENCE AUTHOR
   MILLS CJ                SOURCE AUTHOR
   VIRGINIA MASON RESEARCH CENTRE   (ORGANISATION)
   SHERCLIFF JA            THEORY OF ELECTROMAGNETIC
                              FLOW MEASUREMENTS 1962
   WETTERER E              ZF BIOL   98   26   1937
   DENISON AB              CIRC RES   3   39   1955
   SPENCER MP              IRE TRANS MED   ME-6   220   1959
   BLOOD FLOW
      MEASUREMENT/         (WORD TYPE 1)
   FLOWMETER/              (WORD TYPE 2)
   BLOOD                     "    "    "
   MAGNETIC                  "    "    "
   SQUARE WAVE               "    "    "
   ELECTROMAGNETIC/          "    "    "
   CATHETER TIP              "    "    "
   INDUCTION                 "    "    "
```

(The meaning of the word questions is - inform me
about any article which has the word phrase BLOOD
FLOW MEASUREMENT/ in the title, and about any article
which has any two words in the TYPE 2 list in the
title regardless of order).

Figure 2a (Profile circumscribing the subject THE
DESIGN OF ELECTROMAGNETIC BLOOD FLOWMETERS, (incom-
plete).)

```
SHERCLIFF JA    62    THEORY OF ELECTROMAGNETIC =
   CITED BY   JAGENEAU   AH    SCHAPER WKA              (3)
      ACT PHYSL N   14   346   67  M  3R  N3   95284
      Flow and pressure measurements in unrestrained dog
```

| Journal· | Page | Category | (M= proc. of meeting) |
|---|---|---|---|
| Volume | Year | | |

```
SPENCER MP     IRE TRANS MED   ME-6   220   59
DENISON AB     CIRC RES   3   39   55                    (1)
   CITED BY BOND RF
      J APP PHYSL   22   358   67  4R  N2   89513
(TERM) In vivo method for calibrating electromagnetic
      flowmeter probe
```

| Number of ref | ISI ref. no. |
|---|---|
| Journal issue no. | |

```
KOLIN A     P NAT ACAD   45 1312   59
SPENCER MP     IRE TRANS MED ME-6   220   59            (4)
   CITED BY DAGGETT WM   AUSTEN WG
      AM J SURG   114   139   67  M  98R  N1   95285
      Biomedical engineering applications to
      cardiovascular surgery.


SHERCLIFF JA    62    THEORY OF ELECTROMAGNETIC =
   CITED BY BROUILET. EC   LYKOUDIS PS                (5)
      PHYS FLUIDS   10   995   67  18R  N5   94718
      Magneto-fluid-mechanic channel flow 1. Experiment


SOURCE AUTHOR     KOLIN A
KOLIN A   J APPL PHYS   15   150   44
KOLIN A   REV SCI INST   16   109   45
KOLIN A   SCIENCE   130   1088   59
WYATT DG   MED BIOL ENG   4   17   66                    (1)
   CITED BY KOLIN A
      P NAS US   57 1331   67  8R  N5   93597
(TERM)  An electromagnetic intravascular blood
        flow sensor
```

Figure 2b  (Example of a weekly print-out (composite)
showing information obtained about articles via a
profile shown in figure 2a)

| DOLLARS | | WORD TERM |
| ALONE | COMBO | |
| --- | --- | --- |
| 7 | 4 | Head |
| 45 | 5 | Health |
| 71 | 8 | Heart |
| 76 | 9 | Heat |
| 26 | 4 | Helium |
| 22 | 4 | Hemoglobin |
| 7 | 4 | Hemolytic |
| 9 | 4 | Hemorrhage |
| 19 | 4 | Hepatic |
| 8 | 4 | Heterogeneous |
| 207 | 22 | High |
| 23 | 4 | Higher |
| 8 | 4 | Histamine |
| 33 | 4 | Histochemical |
| 22 | 4 | History |
| 14 | 4 | Hormone |

Figure 3 (Part of the high-frequency word term list with prices)

```
KOLIN A ------------- 36 ------P SOC EXP BIOL MED   35  53
     DOUTHEIL U              PFLUG ARCH      66  287 111
     FONTAINE JL            PATH BIOL       66   14 332
     GAULT JH               CIRCULATION     66   34 833
     HIRSCH HH              Z KREISLAUF     66   55 765
          ----------- 41 ------ P SOC EXP BIOL MED   46 235
     HILAL SK               AM J ROENTG     66   96 986
     RYAN DP                REV SCI INST    66   37 486
          ----------- 45 ----- REV SCI INST        16 109
     ATTINGER ED            CIRCUL RES      66   19 230
   ┌FERGUSON DJ             CIRCUL RES      66   19 917
   │       •
   │       •
   │       •
   │     ----------- 65 ----PROTIDES BIOL FLUIDS  12 410
   │ MCDOUGAL EI            BR MED BULL     66   22 115
   │     ----------- 65 ----- Z ALLERGEIFORSCH  128 117
   │ KOLIN A                J IMMUNOL       66   97 261
   │
   │
   └──────────────────────▶ (to figure 4b)
```

Figure 4a

(Part of the 1966 Science Citation Index)

FERGUSON DJ   LANDAHL HD
    CIRCUL RES   19 917  66  10R  N5  85131
    Magnetic meters - effects of electrical resist-
    ance in tissues on flow measurements and an
    improved calibration for square-wave circuits.

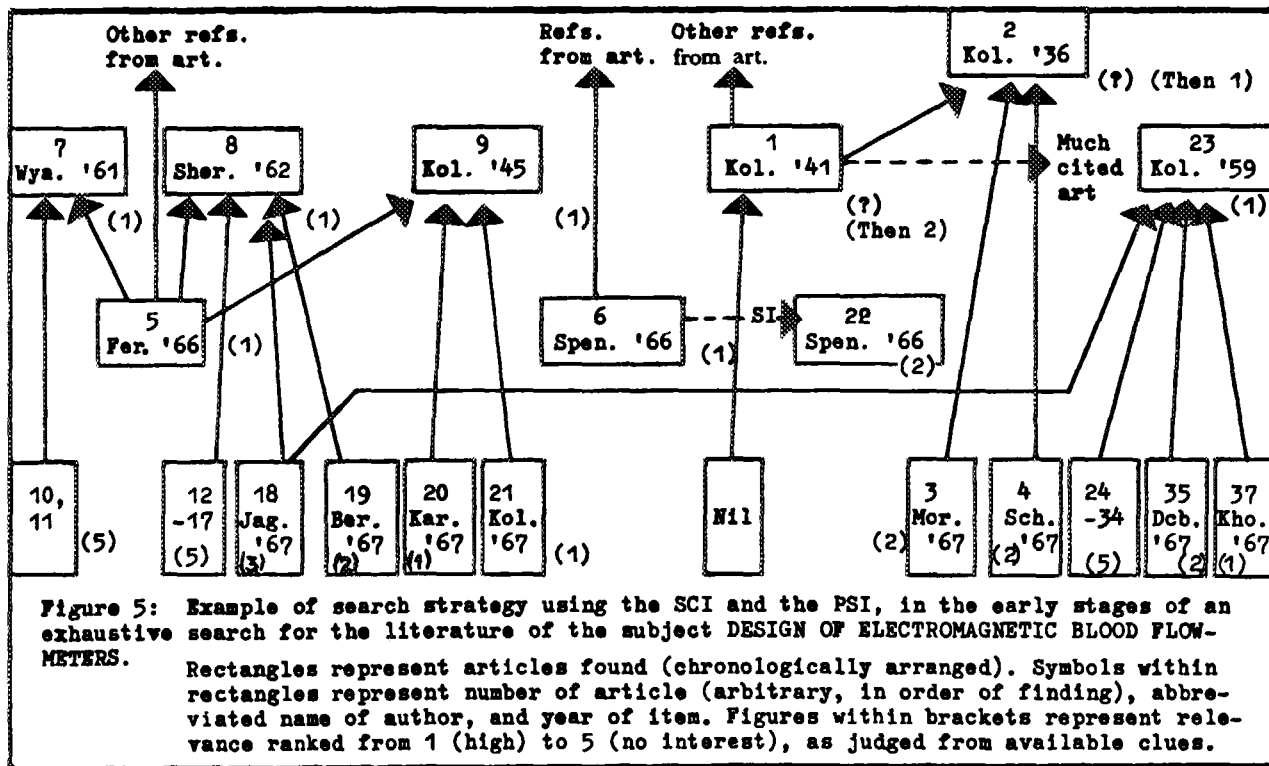    Figure 4b (Part of the 1966 Source Index)

| PRIMARY TERM | FIRST | ISI NO. |
| CO-TERM | AUTHOR | AND CODE |
| FLOW | | |
| SPOOLS | RISCH R | 72758 P |
| SQUALUS | MILLEN JE | 79980 |
| | MURDAUGH HV | 82037 L |
| SQUARE | LARSEN PS | 79843 M |
| SQUARE-EDG. | ANNAND WJD | 86257 N |
| SQUARE-WAVE | FERGUSON DJ | 85131 |
| | SPENCER MP | 74188 M |
| SQUIRE | BAUER J | 80971 |
| . | | |
| . | | |
| . | | |
| TISSUES | FERGUSON DJ | 85131 |
| TITANIUM | BROOKES CA | 78196 |
| . | | |
| . | | |
| . | | |

Tech. Note

Letter

Proc. of
meeting

Figure 4c

(Part of the 1966 Permuterm Subject Index)

56

Other refs.
from art.

Refs.
from art.

Other refs.
from art.

2
Kol. '36

(?) (Then 1)

7
Wya. '61

8
Sher. '62

9
Kol. '45

1
Kol. '41

Much
cited
art

23
Kol. '59

(1)

(1)

(1)

(?)
(Then 2)

(1)

5
Fer. '66

(1)

6
Spen. '66

SI

22
Spen. '66

(1)

(2)

10,
11

(5)

12
-17

(5)

18
Jag.
'67
(3)

19
Ber.
'67
(2)

20
Kar.
'67
(1)

21
Kol.
'67

(1)

Nil

(2)

3
Mor.
'67

4
Sch.
'67
(2)

24
-34

(5)

35
Dcb.
'67
(2)

37
Kho.
'67
(1)

Figure 5: Example of search strategy using the SCI and the PSI, in the early stages of an
exhaustive search for the literature of the subject DESIGN OF ELECTROMAGNETIC BLOOD FLOW-
METERS.

Rectangles represent articles found (chronologically arranged). Symbols within
rectangles represent number of article (arbitrary, in order of finding), abbre-
viated name of author, and year of item. Figures within brackets represent rele-
vance ranked from 1 (high) to 5 (no interest), as judged from available clues.

as new entry points. In this particular article only the briefest biblio-graphic references are provided, so without being a subject expert it is difficult to decide which articles would be of interest. Firstly a reference to an article by the same author, Kolin, was selected—article no 2 on the diagram—and upon entering the 1967 SCI, two citations to this article were found. Other references given in article no 1 yield further 1967 citing articles (excluded from figure 5 for clarity).

This simple sequence of events might be all that is necessary, should the requirement be for a quick search in order to reveal a small selection of current articles about a given subject. The fact that the starting article provided in this case did not describe any aspect of the subject of great importance, and was not cited (at least not in 1967), did not preclude its use as a starting point.

Referring now to the starting words provided, two articles having both of the words of interest were found in the 1966 Permuterm subject index. The articles are numbered 5 and 6 on the diagram. Being of recent origin, they have not yet been cited (as at September 1967), so the articles were obtained, and their references observed. In the case of article no 5 by Ferguson, the references included in it are explicit, and three were selected as points of entry—articles num-bered 7, 8 and 9. These three articles are cited by twelve 1967 articles, four out of the twelve being relevant or of interest.

Many possibilities now exist for continuing the search, and two are given by way of example.

As the Source index is author-ordered, it may be entered to see whether any of the known authors have written any other articles about the same subject. For instance in the 1966 Source index, it was observed that the author Spencer has written a second article about the same subject—number 22 on the diagram. This kind of procedure can be tried for other authors and for other years of the Source index.

Another technique which can be very rewarding is to observe a list of the cited references under a known author's name in the Citation index, paying particular attention to those references which have been much cited.

Referring to the cited author Kolin in the 1967 Citation index, a particular article by him is much cited—the article published in 1959 and numbered 23 on the diagram. It is not known at the outset whether the author is writing about the same subject in this particular cited article, as only an abbreviated reference is given, but this is soon discovered, because detailed information about the citing articles can be obtained from the Source index, and the relevance of the cited article thereby deduced. Kolin's 1959 article was cited fourteen times in 1967, two of the citing articles being of great interest.

So far in this search, slightly more than one hour has been spent, by far the greater proportion of it in making notes about the information found. This information relates to eighteen articles of interest, half of them published during 1967, and it may be considered that cut-off time has been reached—that is, that no further time is justified for continuing the search.

Let us consider what has been achieved: during the hour eighteen articles ranging from ' highly relevant ' down to ' of some interest ' have been found, and thirty nine articles of no interest have been ignored. In all of these thirty nine articles there are good reasons for citing the earlier articles, but the conceptual relationship is by reason of the application of an electromagnetic flowmeter during an operation or surgical experiment. It will be remembered that the searcher is interested only in instrument design, so these thirty nine articles are, rather strictly, classified as being of no interest. (They could well be of interest to an associate.)

During the search hour the average time taken to find relevant or of interest articles was $2\frac{1}{2}$ minutes per article (inclusive of writing time), and the average time taken for an unwanted article was twenty three seconds.

The relevance ranking given is based on a combination of the clues, the title, the author, and the journal and the activity of the author.

If the search is being carried out by the research engineer, his relevance assessment of the articles may be slightly revised when he has obtained and read the articles. However if someone who is not a subject expert is carrying out the search on behalf of the engineer, then he may considerably revise the assessment.

If the engineer does the searching, his knowledge of the subject will enable him to filter out noise before it breeds further noise during an extended search, thereby saving search time.

The purpose of a literature search of this kind may be to find the maximum number of relevant articles in some justifiable time or may be to find relevant articles—let us say six—in the minimum possible time. There may be other requirements, particularly if the searcher is interested in historical aspects of the subject; he may even wish to know about the outstanding or ' milestone ' articles which have been published.

Searches have been carried out which provide larger numbers of relevant articles in a shorter time than in the search just described. On the other hand some searches have been less rewarding and have taken longer. At least in the early stages, any search will depend on the effectiveness of the starting points chosen. If it so happens that a subject expert chooses a starting paper which he knows is quite

outstanding in the subject area, then one look-up may be sufficient, as a high proportion of current articles may cite that well known paper.

Spencer* has reported a comprehensive search carried out using SCI, in which search times and articles found compared favourably with search times and articles found using well known abstracting journals. In one search, articles found from an abstracting journal were used as starting points, and then articles were found using SCI at the rate of 1·1 per minute for the first 5·6 hours.

The test brings to mind a further application of SCI—the up-dating of an existing bibliography. Presumably the bibliography will be subject-oriented, the compiler having made a selection of the literature for some specific purpose. The bibliography may be regularly up-dated simply by entering the current edition of SCI (quarterly or annually) at each reference to see who has cited it. It could be up-dated on a weekly basis by running an ASCA profile consisting of the entire bibliography.

Finally let us consider the task of providing a comprehensive bibliography from scratch, including the significant prior art and current important articles about a specified subject. Let us assume that the compiling is to be done by a person who knows little about the subject, and the purpose of the exercise is to provide a scientist, who is about to start on a new area of research, with a list of important background reading. This has been done for the flowmeter subject already discussed. The principle employed was to collect information as quickly as possible up to some arbitrary number of articles if necessary using more than one edition of SCI. The number must be large enough for certain trends to be obvious, primarily to answer the question ' is there a manageably small number of articles, one or more of which will probably be cited by any relevant current article? '. This is really the same question that should be asked by anyone compiling an ASCA profile. However the ASCA subscriber will usually be in a position to provide the answer; in this case we have to deduce it.

A convenient way of arranging information is to list the arbitrary number of articles alphabetically by author, and to mark against each the number of ' mentions ', distinguishing between self-citations, citations (references) observed in one of the obtained articles, citations found in SCI, or information found by author search in the Source index. Suitable weight must be given to the fact that articles will not be equally eligible for citing—an article published in 1963 of some repute may have accumulated more citations by 1966 than an article of equal repute published in 1965.

From this list, some articles will be mentioned once, others several times. It may be helpful to draw a citation diagram of the type shown in figure 5 for some manageable number of preferred articles (if more than about forty articles are included on such a diagram the picture becomes confused). The purpose of the drawing is to assist in establishing or confirming relevance by observing common references in a group of articles, perhaps to articles of known relevance (bibliographic coupling).

By these procedures, sufficient clues will accumulate to indicate that certain articles are in the 'milestone' category—they may be much cited and probably describe some notable advance in the subject; or they may be review articles or even controversial articles. Twelve articles of this kind were identified in the flowmeter subject. These were used to find current articles as each 1967 SCI quarterly was published, and would of course be used as key terms in an ASCA profile. Three specific items of this type are included in the profile of figure 2a—the book by Shercliff and the articles by Denison and Spencer.

These methods have been tried out with some success with subjects as diverse as ' holography ', ' explosive welding ', and ' foot and mouth disease virus replication '. One difficulty during assessment is to find out subsequently what articles have been missed. In the flowmeter subject two articles published in 1967 would not have been retrieved because they appeared in journals not covered by ISI. Presumably they will appear in the system in due course when sufficient time has elapsed for them to be cited. Although a careful lookout has been kept, it is by no means certain that there are not other missed articles.

The subject 'holography' is of some interest; 125 articles published during 1967 had at least one of six articles in common. It took only fifteen minutes to appreciate the significance of the six articles which are by two authors. The 125 articles are found by entering the Citation index at two places—that is at the two cited authors. However in a small group of articles (1968) about vibration analysis and temporal recall, the citations are to very recent work only, not to the six classic articles.

These references reflect a new trend, and rather exceptionally, do not couple up with the prior art on which they are based. They would have been retrieved by words in title through ASCA or the Permuterm index, but for retrieval through citations, only a man closely following developments would have looked up the appropriate starting article in SCI.

From the search described, from other searches which have been carried out, and from experience obtained during the operation of a number of ASCA profiles, it may be concluded that the systems described are practical and effective. As in all information systems there are advantages and disadvantages; the disadvantages may be minimised if the general principles of the possible search strategies are clearly understood. The multi-disciplinary coverage, and the facilities provided for retrieval through the association of concepts and ideas are unique for a system of this magnitude. The results obtained in practice from the proper employment of these special attributes indicate that the system compares very favourably with other information retrieval systems.

REFERENCES

1 Cooper W S ' Expected search length; a single measure of retrieval effectiveness based on the weak ordering of retrieval systems ' *American Documentation* 19 30 1968.
2 Abbot M T J, Hunter I S, Simkins M A ' Current awareness searches on *CT, CBAC, and ASCA* ' *Aslib proceedings* 20 129 1968.
3 Cawkell A E Letter to the editor *Aslib proceedings* 20 233 1968.
4 Garfield E, Sher I H ' ISI's experiences with ASCA—a selective dissemination system ' *Journal of chemical documentation* 7 147 1967.
5 Garfield E ' Primordial concepts, citation indexing and historio-bibliography ' *Journal of library history* 2 235 1967.
6 Spencer C C ' Subject searching using the *Science citation index*. Preparation of a drug bibliography using *Chemical abstracts, Index medicus* and SCI, 1961 and 1964 ' *American documentation* 18 88 1967.