

A SYSTEM FOR AUTOMATIC CLASSIFICATION OF SCIENTIFIC LITERATURE

EUGENE GARFIELD, *President*, MORTON V. MALIN,
Vice-President, HENRY SMALL, *Research Associate*

(*Institute for Scientific Information, 325, Chestnut Street, Philadelphia, Pa.*)

Received on December 22, 1974

ABSTRACT

A computer-based system for automatically classifying scientific articles is described. The unique feature of this system is that its structure is completely determined by citation patterns in the Science Citation Index. These citation patterns give rise to clusters or clumps of cited papers which correspond, in turn, to clusters of citing papers. Classification headings for each cluster are determined by examination of high frequency word pairs drawn from the titles of the citing papers. Classification of a new article is performed automatically by determining what clusters it cites and assigning appropriate (weighted) subject headings to it. The system will permit updating of the classification scheme on an annual basis, and the incorporation of new headings and deletion of old ones.

This paper describes an automatic classification system being developed at ISI, with the unique feature that its structure is completely determined by citation patterns derived from the *Science Citation Index*^R data base. I will also summarize some of the current research at ISI. ISI's research and development objectives include the development of new information products and services, and development of improved processing operations, methods, and systems. But we also conduct basic research in the area of information science; the last activity supports the first two. We believe that the more we learn about the characteristics of the scientific literature, and its relationship to science and research communication, the better we will be able to develop and provide services to the user. The automatic classification system discussed in this paper is an outgrowth of a basic research project now being conducted at ISI using the *Science Citation Index*^R data base.

This data base now consists of 13 years of back files containing nearly 3.4 million source articles and nearly 40 million citations. Thus, we have an unusual opportunity for conducting a broad program of research activities, using the file to study the characteristics of the literature and to

conduct citation behavior studies in the history and sociology of science. These studies are very productive both for ISI and the scientific community.

Before describing some of the research and results a brief description of citation indexing would be useful, because it is necessary to understand this data base in order to understand the research work we are doing. Since the literature already contains excellent detailed descriptions of citation indexing, [1, 2, 3]. I will not discourse in detail on the *SCI*^R, but only describe the concept and the data one has available in the printed Index and on computer tape.

In brief, a citation index is a cumulation of journal article references arranged so that one can determine what later or more current articles have cited any earlier article or book. The *Science Citation Index*^R arrangement is alphabetical by the cited author of each cited item. Under each cited item is listed all the later articles which have cited it during a specified time period, *e.g.*, three months, one year, or five years. ISI now processes about 2400 journals for the *SCI*^R, and all references in all of the articles in these journals are keyed into the data base, and eventually appear in the printed *SCI*^R cumulations. At the same time, a number of other data elements are keyed from the source articles. These include: all authors of a given article, author addresses, full title of the article, and journal, volume, page, and year. The number of references keyed in 1973 was roughly 5 million coming from approximately 400,000 source items. It is virtually the only data base, available, which includes the bibliographic, that is cited, references.

It was long ago pointed (1955) that these cited references are a unique and important group of indexing terms [4]. Salton [5], in particular, has confirmed the value of citation indexing for retrieval of information. Thus, bibliographic citations are important indicators of document content. Human indexers do not ordinarily think of citations as descriptors of the citing document, but they are, in fact alternate representations of the documents they identify [4, 6]. Were this not true, the automatic classification system described below would not be possible.

I will digress briefly from the main topic to describe a project which illustrates how we benefit from research with our data. This project is called the Journal Citation Index. To create this compilation, every citation from the source items processed during the last quarter of 1969 was extracted from the total year's file. Through a series of sorts, a new type of citation file was created which, instead of obtaining citations to articles, obtained citations to the Journals in which the cited articles were published.

Further programming then produced listings providing data showing for each cited journal which journals had been cited. Counts were made to show the frequency with which each journal was cited, and the year of the cited articles. The process of analysis continued until we were able to produce statistical indicators which would permit ranking of the journals based on factors other than just the frequency of citations. A description of this project can be obtained from my 1972 *Science* articles [7]. Indeed, this project illustrates well the classificatory power of citation analysis. What other means do we have available today for categorizing journal collections ?

The purpose of the JCI project was to test the hypothesis that citation frequency is a measure of impact of a journal. We believed that such data could help us and others in developing a core list of scientific journals as well as aid in journal evaluation and selection procedures. In 1973, the listing were published by ISI as a service for libraries under the title *Journal Citation Reports*.TM More recent data are now available covering the year 1972, and we plan to produce an updated *JCR*TM on a regular basis.

More relevant to the subject of this paper is the second research project which I will describe ; work being done by Dr. Henry Small of ISI's R and D staff, "Mapping of Scientific Specialties." The work is being supported by a grant from the National Science Foundation. Although primarily a project concerned with historical and sociological aspects of science, it has great relevance to information science, and to our automatic classification system. The objective of this research project was to test the hypothesis that citations to scientific articles could be used in identifying scientific specialties, in effect that citation data could be used for classificatory purposes. Thus, an understanding of the classification system described below requires initially an understanding of the research from which it stems.

It is appropriate also at this point to define automatic classification because clustering is an essential part of classification and of the specialty mapping research. Webster's New Collegiate dictionary defines classification as "systematic arrangement in groups or categories according to established criteria." An automatic classification system can thus be defined as a method for systematically arranging documents in groups or categories by a process that requires no human intervention, save the keying of the text. In this context, text may be full text titles, abstracts, or citations in a bibliography. A system for automatic classification is, therefore, one which satisfies the requirement of clustering or bringing like things together or as a process which groups objects resembling one another in terms of their properties.

The specialty mapping project uses a computer based technique to be identify clusters of highly cited and co-cited scientific articles. Co-citation is defined as the number of times two publications are cited together in the literature. The clusters formed are found to correspond to scientific specialties. The technique employed begins with identification of highly cited papers in an annual file of the *SCI*. To initiate the experiment, the 1973 file was used, and all papers cited at least fifteen times were extracted. This reduced the total file of approximately 4,000,000 citations to a more manageable one containing 430,000 citations. Out of more than 2,000,000 unique cited items in 1973, only 16,000 items were selected. For each cited paper we extracted the list of citing papers and this new file was then resorted so that we could identify pairs of papers cited together, *i.e.*, co-cited. The number of identical pairs of cited papers were then counted to establish the co-citation strength of each pair of papers and a total of 710,000 distinct co-cited pairs were generated through this method.

The next step was to apply a clustering algorithm in order to group together the most highly co-cited documents. The clustering algorithm used is a single-link procedure. Briefly, to describe this procedure, a minimum linkage level is specified and the computer begins by selecting an initial document and retrieving all of its linkages to other documents which equal or exceed the minimum threshold. A cluster is complete when all documents have been identified which are linked together in a connected graph by linkages which satisfy the threshold criterion. At this point, the computer proceeds to the next unclustered document and generates another cluster. Clustering may be carried out at as many as four levels and the resulting clusters may be merged together to reveal the heirarchical or nested structure of the file.

Figures 1 through 3 show an example of a cluster obtained by this method. Figure 1 is a cluster as it emerges from the computer as a list of highly co-cited documents; Figure 2 is a cluster represented as a network with each document indicated by a circle and each line a co-citation linkage. Figure 3 is a list of the titles of the documents in the cluster which shows that the subject matter is quite narrowly focused on hormone releasing hormones.

At any single level, we may determine the linkages among clusters by counting the number of co-citations between documents in different clusters. This is called cluster co-citation. If, for example, we conduct a clustering run at level 11, the linkages among clusters are determined by co-citation links from one to ten. Using this information, we are able to draw as a

07/29/73

Figure 1
CLUSTER 33

PAGE 50

| | CITATION 1 | FREQ 1 | | CITATION 2 | FREQ 2 | STRENGTH |
|----|--------------|---------------------|-----------|-----------------|---------------------|-----------------|
| 1 | AMOS M | BIOCHEM BIOPHYS RES | 44 20571 | 19 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 13 |
| | AMOS M | BIOCHEM BIOPHYS RES | 44 20571 | 19 MONAHAN M | CR ACAD SCI PARIS | 273 50871 26 13 |
| 2 | BABA Y | BIOCHEM BIOPHYS RES | 44 45971 | 28 GEIGER R | BIOCHEM BIOPHYS RES | 45 76771 24 13 |
| | BABA Y | BIOCHEM BIOPHYS RES | 44 45971 | 28 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 22 |
| | BABA Y | BIOCHEM BIOPHYS RES | 44 45971 | 28 MATSUO H | BIOCHEM BIOPHYS RES | 45 82271 23 14 |
| | BABA Y | BIOCHEM BIOPHYS RES | 44 45971 | 28 SCHALLY AV | BIOCHEM BIOPHYS RES | 43 39371 42 15 |
| 3 | BURGUS R | CR ACAD SCI PARIS | 273161171 | 18 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 15 |
| | BURGUS R | CR ACAD SCI PARIS | 273161171 | 18 MONAHAN M | CR ACAD SCI PARIS | 273 50871 26 13 |
| 4 | GEIGER R | BIOCHEM BIOPHYS RES | 45 76771 | 24 BABA Y | BIOCHEM BIOPHYS RES | 44 45971 28 13 |
| | GEIGER R | BIOCHEM BIOPHYS RES | 45 76771 | 24 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 20 |
| | GEIGER R | BIOCHEM BIOPHYS RES | 45 76771 | 24 MATSUO H | BIOCHEM BIOPHYS RES | 45 82271 23 15 |
| | GEIGER R | BIOCHEM BIOPHYS RES | 45 76771 | 24 MONAHAN M | CR ACAD SCI PARIS | 273 50871 26 13 |
| 5 | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 AMOSS M | BIOCHEM BIOPHYS RES | 44 20571 19 15 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 BABA Y | BIOCHEM BIOPHYS RES | 44 45971 28 22 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 BURGUS R | CR ACAD SCI PARIS | 273161171 18 15 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 GEIGER R | BIOCHEM BIOPHYS RES | 45 76771 24 20 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 MATSUO H | BIOCHEM BIOPHYS RES | 45 82271 23 20 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 MONAHAN M | CR ACAD SCI PARIS | 273 50871 26 18 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 NISWENDER GD | P SOC EXP BIOL MED | 128 80768 71 18 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 RAMIREZ VD | ENDOCRINOLOGY | 73 19363 23 13 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 SCHALLY AV | BIOCHEM BIOPHYS RES | 43 39371 42 26 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 43133471 | 59 SCHALLY AV | SCIENCE | 173103671 44 20 |
| 6 | MATSUO H | BIOCHEM BIOPHYS RES | 45 82271 | 23 BABA Y | BIOCHEM BIOPHYS RES | 44 45971 28 14 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 45 82271 | 23 GEIGER R | BIOCHEM BIOPHYS RES | 45 76771 24 15 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 45 82271 | 23 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 20 |
| | MATSUO H | BIOCHEM BIOPHYS RES | 45 82271 | 23 MONAHAN M | CR ACAD SCI PARIS | 273 50871 26 13 |
| 7 | MONAHAN M | CR ACAD SCI PARIS | 273 50871 | 26 AMOSS M | BIOCHEM BIOPHYS RES | 44 20571 19 13 |
| | MONAHAN M | CR ACAD SCI PARIS | 273 50871 | 26 BURGUS R | CR ACAD SCI PARIS | 273161171 18 13 |
| | MONAHAN M | CR ACAD SCI PARIS | 273 50871 | 26 GEIGER R | BIOCHEM BIOPHYS RES | 45 76771 24 13 |
| | MONAHAN M | CR ACAD SCI PARIS | 273 50871 | 26 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 18 |
| | MONAHAN M | CR ACAD SCI PARIS | 273 50871 | 26 MATSUO H | BIOCHEM BIOPHYS RES | 45 82271 23 13 |
| 8 | NISWENDER GD | P SOC EXP BIOL MED | 128 80768 | 71 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 18 |
| | NISWENDER GD | P SOC EXP BIOL MED | 128 80768 | 71 RAMIREZ VD | ENDOCRINOLOGY | 73 19363 23 15 |
| | NISWENDER GD | P SOC EXP BIOL MED | 128 80768 | 71 SCHALLY AV | BIOCHEM BIOPHYS RES | 43 39371 42 16 |
| 9 | RAMIREZ VD | ENDOCRINOLOGY | 73 19363 | 23 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 13 |
| | RAMIREZ VD | ENDOCRINOLOGY | 73 19363 | 23 NISWENDER GD | P SOC EXP BIOL MED | 128 80768 71 15 |
| 10 | SCHALLY AV | BIOCHEM BIOPHYS RES | 43 39371 | 42 BABA Y | BIOCHEM BIOPHYS RES | 44 45971 28 15 |
| | SCHALLY AV | BIOCHEM BIOPHYS RES | 43 39371 | 42 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 26 |
| | SCHALLY AV | BIOCHEM BIOPHYS RES | 43 39371 | 42 NISWENDER GD | P SOC EXP BIOL MED | 128 80768 71 16 |
| | SCHALLY AV | BIOCHEM BIOPHYS RES | 43 39371 | 42 SCHALLY AV | SCIENCE | 173103671 44 15 |
| 11 | SCHALLY AV | SCIENCE | 173103671 | 44 MATSUO H | BIOCHEM BIOPHYS RES | 43133471 59 20 |
| | SCHALLY AV | SCIENCE | 173103671 | 44 SCHALLY AV | BIOCHEM BIOPHYS RES | 43 39371 42 15 |

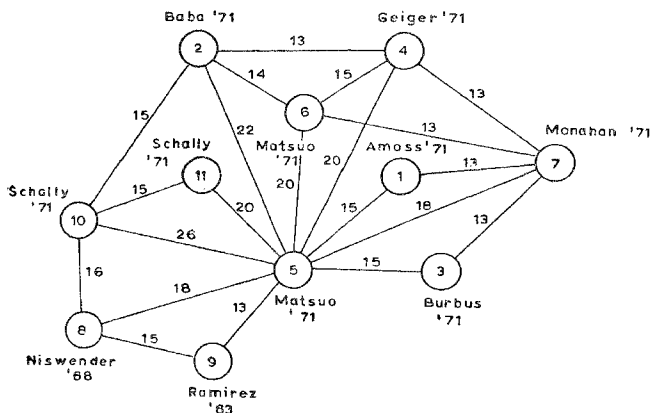


FIG. 2

FSH and LH—Releasing Hormone

- | | | | |
|-----|---|---------------------|-----------|
| 1. | Amoss M | Biochem Biophys Res | 44 20571 |
| | Purification, Amino-Acid Composition and N-Terminus of Hypothalamic Luteinizing Hormone Releasing Factor (LRF) of Ovine Origin | | |
| 2. | Baba Y | Biochem Biophys Res | 44 45971 |
| | Structure of Porcine LH-Releasing and FSH-Releasing Hormone 2. Confirmation of Proposed Structure By Conventional Sequential-Analysis | | |
| 3. | Burgus R | Cr Acad Sci Paris | 273161171 |
| | Molecular Structure of Hypothalamic Factor of Ovine Origin Controlling secretion of Hypophysial Gonadotropic Luteinizing-Hormone (LH) | | |
| 4. | Geiger R | Biochem Biophys Res | 45 76771 |
| | Synthesis and Characterization of a Decapeptide Having LH-RH/FsH-RH Activity | | |
| 5. | Matsuo H | Biochem Biophys Res | 43133471 |
| | Structure of Porcine LH-Releasing and FH-Releasing Hormone 1. Proposed Amino-Acid Sequence | | |
| 6. | Matsuo H | Biochem Biophys Res | 45 82271 |
| | Synthesis of Porcine LH-Releasing and FSH-Releasing Hormone By Solid-Phase Method | | |
| 7. | Monahan, M | Cr Acad Sci Paris | 273 50871 |
| | (FR) Total Synthesis By Solid-Phase of Decapeptide Stimulating secretion of Hypophysial Gonadotropin LH and FSH | | |
| 8. | Niswender G. D. | P Soc Exp Biol Med | 128 80768 |
| | Radioimmunoassay For Rat Luteinizing Hormone with Antiovine LH Serum and Ovine LH-1311 | | |
| 9. | Ramirez V D | Endocrinology | 73 19363 |
| | A Highly Sensitive Test for LH-Releasing Activity-Ovariectomized, Estrogen Progesterone-Blocked Rat | | |
| 10. | Schally A V | Biochem Biophys Res | 43 39371 |
| | Isolation and Properties of FSH and LH-Releasing Hormone | | |
| 11. | Schally A V | Science | 173103671 |
| | Gonadotropin-Releasing Hormone-One Polypeptide Regulates, Secretion of Luteinizing and Follicle-Stimulating Hormones | | |

FIG 3

graph, the network of most active specialties in science. By drawing such a map for each year, over a period of years, we can study how the links between specialties have changed, and where new specialties have emerged or old ones declined.

Figure 4 is a map of biomedical clusters derived from the 1972 *SCI*^R file. Only clusters containing three or more documents have been included and only if they have been linked with another cluster on the diagram by a cluster co-citation threshold of 100. The map shows four major areas of biomedicine. In the upper left hand corner are chromosomes and RNA viruses, and in the upper right is work on immunology. Attached to immunology in the lower right is research related to biological membranes. To the left of this, in the lower left hand corner is work related to cyclic AMP. The pattern of specialties and linkages changes from year to year, and we can observe the evolution of this network over time.

The purpose of the mapping of the science project is to increase our understanding of the processes of growth and change in science, and to apply this understanding in the area of science policy. The important finding of our mapping work is that the basic unit of science appears to be the scientific specialty, not the discipline or the isolated researcher. Further, we have found that growth and change in specialties can be extremely rapid. These findings have important implications for information retrieval. First, they indicate that we must gear information services and classification schemes to the specialty scale, because this scale is probably most relevant to the working scientist, and is the one at which he generates and utilizes information.

Second, a classification scheme, if it is to be effective, must be capable of changing very rapidly. Probably an annual update is needed to respond to new developments and growth in some new specialties, but this will vary for different specialties. Some have a lifetime of as little as one year—others ten years or more. The precise life expectancy of a specialty is a question of considerable interest, and one which we should be able to answer using the ISI data base.

The application of our clustering work to classification is, therefore a highly natural one. Only one important modification is necessary to adapt an essentially science policy oriented system, where the criterion is the level of activity, to information science oriented system where criterion is to generate as many classification categories and classify as many articles as possible. The change consists in adopting a normalized linkage measure,

Automatic Classification of Scientific Literature

Fig. 5

| | | | | | | | | | |
|-------------------------|---------------|------|------|------|------------------------|-------------|-------------|------|------|
| GRAH BE | J PALEONTOLOG | 46 | 233 | 72 | ROELANDT J | ES ISSAT J | 34 | 1232 | 72 |
| 64 MODIFIED NOMIC AID B | | | | | TICZON AR | CIRCULATION | 47 | 443 | 72 |
| MELLIARD J | BIO-MED ENG | 7 | 518 | 72 | EDMONDS M | | | | |
| ARMSTRONG JA | | | | | 56 CANCER RES | 16 | 222 | | |
| TO BE PUBLISHED | | | | | NAKAZATO H | J BIOL CHEM | 246 | 1472 | 73 |
| LANDAUER R | J APPL PHYS | 44 | 1156 | 73 | GIRON ML | BIOC BIOP A | 287 | 638 | 72 |
| 50 J ANAT | | | | | MCLAUGHLIN CS | J BIOL CHEM | 248 | 1475 | 73 |
| DAVYDOVA IV | TSITOLOGIA | 15 | 22 | 73 | NAKAZATO H | BIOC BIOP A | 248 | 1472 | 73 |
| 55 INTERPRETATION INJUR | 1135 | | | | ORBIEN SJ | NATURE-BIOL | 242 | 632 | 73 |
| KREFFI S | J RECHTSMED N | 71 | 251 | 73 | BARAKAS HJ | BIOCHEM | 12 | 920 | 73 |
| 51 VIROLOGY | 14 | 27 | | | SHINESNESS D | NATURE-BIOL | L | 241 | 245 |
| SILVA RP | J COMP PATH | 83 | 161 | 73 | MCLAUGHLIN CM | BIOC BIOP B | 59 | 737 | 71 |
| 52 PHYS REV | 127 | 1918 | | | 62 J BIOL CHEM | 237 | 4636 | | |
| ASPMES DS | PHYS REV | 6 | 4648 | 72 | BARAKAS HJ | BIOCHEM | 12 | 920 | 73 |
| CHEMLA US | ANN PEELECOM | 37 | 417 | 72 | 63 J BIOL CHEM | 238 | 3185 | | |
| FUJITAMA T | B CHEM S J | 46 | 87 | 73 | SHINESNESS D | NATURE-BIOL | L | 241 | 245 |
| IEU H | J APPL PHYS | 32 | 41 | 73 | 65 NOBLEIC ACIDS | STRUCT | 1 | | |
| HO K | J PHYS ZAP | 14 | 138 | 73 | 73 PAPILL GM | BIOC BIOP R | 69 | 737 | 71 |
| KARPLYUK NS | PLASMA PHYS | 15 | 113 | 73 | 69 J BIOL CHEM | 244 | 1514 | | |
| SOROKIN PP | ILUC J O EL | QR | 9 | 27 | 74 FAULT CH | BIOCHEM | 12 | 925 | 73 |
| 54 APPL PHYS LETT | 8 | 196 | | | MCLAUGHLIN CS | J BIOL CHEM | 248 | 1465 | 73 |
| LAK M | PHYS REV A | 7 | 750 | 73 | MOLLOY CR | P NAS US | 61 | 3084 | 72 |
| 54 APPL MICROBIOL | 12 | 132 | | | SHINESNESS D | NATURE-BIOL | L | 241 | 245 |
| HEROME K | ARCH G VIR | 39 | 353 | 72 | SLATER I | P NAS US | 70 | 966 | 73 |
| HELDENBA PK | INFECTION | 7 | 255 | 73 | 70 BIOCHEM BIOPHYS RES | 41 | 1531 | | |
| 67 APPL PHYS LETT | 10 | 16 | | | 71 P NAT ACAD SCI USA | 58 | 1336 | | |
| DEMISTER DN | CHEM P LETT | 13 | 469 | 73 | CHRISTMA JK | BIOC BIOP A | 254 | 153 | 73 |
| CVANDRILA JL | ZH EKSP TSIC | 64 | 465 | 73 | 71 P NAT ACAD SCI USA | 58 | 1336 | | |
| KRYUKOV NG | ILUC J O EL | QR | 62 | 2636 | 72 | DEKAREVI A | PERS LETTER | 29 | 164 |
| SOROKIN PP | IEED J O EL | QR | 9 | 27 | 73 | BIADRUS S | J VIROLOGY | 10 | 1166 |
| 68 ACTA VIROL | 12 | 15 | | | ELBOREL G | P NAS US | 70 | 354 | |
| HELDENBA PK | INFECTION | 7 | 265 | 73 | COOPER HL | TRANSPLAN P | 11 | 3 | |
| 71 APPL MICROBIOL | 21 | 723 | | | CORNUELL L | BIOC BIOP A | 294 | 543 | |
| DAHL H | ACT PAT S B | B | 80 | 863 | 72 | DELCARCO J | BIOC BIOP R | 59 | 465 |
| STEWART WE | J VIROLOGY | 10 | 895 | 72 | EATON BT | VIROLOGY | 50 | 855 | |
| TAN TH | J EXP MED | 137 | 317 | 83 | FAUST CH | BIOCHEM | 12 | 925 | |
| 71 J EXPERIMENTAL MED | 134 | 713 | | | GIRON ML | BIOC BIOP A | 287 | 638 | |
| EVANS MS | INFECTION | 7 | 76 | 72 | 73 | | | | |
| KAMA K | J GEN VIROL | 25 | 297 | 72 | GRENNER JR | | 287 | 69 | |
| NERMUT MV | J GEN VIROL | 57 | 317 | 72 | HUNT JA | BIOCHEM J | 131 | 315 | |
| 72 SCIENCE | 176 | 526 | | | LINDBERG U | J VIROLOGY | 16 | 99 | |
| BIADRUS S | J VIROLOGY | 10 | 1126 | 73 | LUZZATI D | BIOCHEMIE | 54 | 1157 | |
| CORNUELL L | BIOC BIOP A | 294 | 543 | 73 | MCLAUGHLIN CS | BIOC BIOP A | 248 | 1465 | |
| DELCARCO J | BIOC BIOP R | 59 | 486 | 73 | MILLER BL | J GEN VIROL | N | 17 | |
| EATON BT | VIROLOGY | 50 | 855 | 72 | MOLLOY CR | P NAS US | 69 | 3684 | |
| FAUST CH | BIOCHEM | 50 | 855 | 72 | MONIER P | BIOCHEMIE | 54 | 1157 | |
| FRENDS JD | J GEN CHEM | 12 | 955 | | MURPHY H | P NAS US | 70 | 118 | |
| GILLESPIE | SCIENCE | 179 | 1328 | 73 | NAKAZATO H | J BIOL CHEM | 246 | 1472 | |
| GRENNER JR | BIOC BIOP A | 287 | 361 | 72 | ORBIEN SJ | NATURE-BIOL | 242 | 32 | |
| MILLER BL | J GEN VIROL | 25 | 349 | 72 | PERLMAN S | P NAS US | 70 | 359 | |
| NAKAZATO H | J BIOL CHEM | 248 | 1472 | 73 | BARAKAS HJ | BIOCHEM | 12 | 920 | |
| PERLMAN S | P NAS US | 70 | 359 | 73 | ROSEMOND N | J VIROLOGY | 41 | 399 | |
| SARKAR PK | BIOC BIOP R | 50 | 308 | 73 | SARKAR PK | BIOC BIOP R | 50 | 228 | |
| SCHLOM J | SCIENCE | 179 | 466 | 73 | SCHLOM J | SCIENCE | 179 | 466 | |
| SLATER I | P NAS US | 70 | 406 | 73 | SHINESNESS D | NATURE-BIOL | L | 241 | |
| TAYLOR JM | BIOCHEM | 12 | 465 | 73 | SLATER I | P NAS US | 70 | 406 | |
| ARMSTRONG JB | 363 | | | | WHYATT PL | BIOC PERM | 22 | 229 | |
| 67 GENETICS | 56 | | | | 72 CAN MED ASSOC J | 107 | 534 | | |
| ARONSTRON JB | CAN J MICRO | 18 | 1695 | 72 | FLOCKICE E | CAN MED A J | L | 168 | |
| SILVERMAN M | J BACT | 113 | 105 | 73 | VANDERME WF | CALIF MED | 118 | 28 | |
| | | | | | 72 J GEN VIROL | 25 | 227 | | |
| | | | | | NAKAZATO H | J BIOL CHEM | 248 | 1472 | |
| | | | | | EDMONDS MW | | | | |
| | | | | | 70 J CLIN END | 31 | 468 | | |
| | | | | | LAMBE L | J CLIN END | 36 | 358 | |

Armstrong JA Science 176 526 72 Edmonds M P Nat Acad Sci 68 1336 71

cited 15 times

cited 29 times

co-cited 12 times

$$\text{coefficient of Jaccard-Sneath} = \frac{12}{15 + 25 - 12} = .37$$

rather than an absolute measure. Earlier, I described the procedure for determining the frequency of co-citation between two highly cited documents. It is a simple step to convert this absolute frequency into a percentage overlap. In clustering terminology, this is known as a Jaccard-Sneath matching coefficient. For example, if paper A is cited fifteen times and paper B is cited twenty times, and together they are co-cited five times, the matching coefficient or percentage overlap is .16 per cent $[5/(15+20-5)]$.

This technique is illustrated in Figure 5 where we have calculated the Jaccard-Sneath coefficient for documents by Armstrong and Edmonds.

The results of our analysis of the 1973 *SCI^R* illustrates the normalized clustering method. An initial citation frequency threshold of fifteen was selected, and a file consisting of 15,923 cited documents obtained. Citations among these documents were determined and the Jaccard-Sneath coefficients were calculated for each pair of cited documents. Clusters were formed at level .16 (16 per cent), and a total of 16,001 clusters were formed, the largest cluster consisting of 117 cited documents.

These clusters were then used to retrieve 1973 source items processed by ISI for the *SCI^R* in 1973. About 25 per cent of the source items were retrieved. A higher fraction of the source items would be classifiable using a lower initial citation frequency threshold.

Automatic indexing and classification is a goal which may never be completely attained, and creation of a system for classifying new documents may require some intervention of human judgement. The system described here is not entirely automatic because it requires human judgement to assign "headings" or labels to the groupings. This judgement is made on the basis of scanning titles with the aid of work pair frequency counts. This is, however, the only point at which human intervention is necessary.

Figure 6 shows a portion of the citing titles obtained for the clusters on hormone releasing hormones, and Figure 7 is a list of word pairs derived from these titles. The naming of the cluster can be readily done using, both the titles and word pairs. The goal of an automatic classification system is to classify new source documents and, therefore, the test of this system is whenever the clusters obtained from the 1973 file are capable of classifying articles published in 1974. Our research is still proceeding, and in the following, I will outline the procedure which is being developed.

The complete list of cluster names and identifying numbers are maintained on one disc file. A second disc contains the cluster number and all the cited references contained in that cluster. As a new document is being entered into the ISI data base, it is possible to match the cited references in the document against the file containing clustered documents and associated cluster numbers. If a new source document contains one or more references which match the cluster file, one or more classification headings can be assigned to the source document.

Suppose, for example, that a particular source document contains five references, three of which cite documents in one cluster, and two which cite

Figure 6

| 12/22/73 | SOURCE TITLES ASSOCIATED WITH CITATION CLUSTERS | | | | 1972 ANNUAL | LEVEL 11 | PAGE 250 | |
|----------|---|--------|----------------|-----|-------------|----------|----------|--|
| CLUSTER | SOURCE | AUTHOR | SOURCE JOURNAL | VOL | PAGE | CODE | FREQ | *** TITLE *** |
| 60 | ABE K | | ENDOCR JAP | 19 | 77 | | 3 | EFFECTS OF SYNTHETIC LUTEINIZING HORMONE-RELEASING HORMONE ON PLASMA LEVELS OF LUTEINIZING-HORMONE AND FOLLICE-STIMULATING HORMONE IN MAN |
| 60 | ARANDE EO | | LANCET | 2 | 112 | | 2 | EFFECT OF SYNTHETIC GONADOTROPIN-RELEASING HORMONE IN SECONDARY AMENORRHEA |
| 60 | AKANDE DO | | LANCET | 11 | 112 | | 3 | EFFECT OF SYNTHETIC GONADOTROPIN-RELEASING HORMONE IN SECONDARY AMENORRHEA |
| 60 | AMOS M | | J CLIN ENDO | 54 | 454 | N | 6 | STIMULATION OF OVULATION OF RABBIT TRIGGERED BY SYNTHETIC LRF |
| 60 | AMOS M | | J CLIN ENDO | 35 | 175 | N | 4 | RELEASE OF GONADOTROPINS BY ORAL ADMINISTRATION OF SYNTHETIC LRF OR A TRI-PEPTIDE FRAGMENT OF LRF |
| 60 | ARIMURA A | | ENDOCRINOL | 90 | 163 | | 5 | STIMULATION OF RELEASE OF LH BY SYNTHETIC LH-RH IN-VIVO. I. COMPARATIVE STUDY OF NATURAL AND SYNTHETIC HORMONES |
| 60 | ARIMURA A | | ENDOCRINOL | 91 | 529 | | 7 | STIMULATION OF FSH RELEASE IN-VIVO BY PROLONGED INFUSION OF SYNTHETIC LH-RH |
| 60 | ARIMURA A | | P SOC EXP M | 139 | 851 | | 3 | RELEASE OF LUTEINIZING-HORMONE BY SYNTHETIC LH-RELEASING HORMONE IN EWE AND RAM |
| 60 | BESSER GM | | BR MED J | 3 | 267 | | 3 | HORMONAL RESPONSES TO SYNTHETIC LUTEINIZING-HORMONE AND FOLLICLE STIMULATING HORMONE-RELEASING HORMONE IN MAN |
| 60 | BEYERMAN H C | | REC TR CHIM | 91 | 1239 | | 8 | SYNTHESIS OF DECAPEPTIDE SEQUENCE PROPOSED FOR LH-RELEASING AND FSH-RELEASING HORMONE |
| 60 | BISHOP W | | ENDOCRINOL | 91 | 643 | | 2 | ACUTE AND CHRONIC EFFECTS OF HYPOTHALMIC-LESIONS ON RELEASE OF FSH, LH AND PROLACTIN IN INTACT AND CASTRATED RATS |
| 60 | BOGDANOV EM | | MED C VIRG | 8 | 5 | | 2 | CURRENT KNOWLEDGE OF GONADOTROPIN RELEASING FACTOR%\$ |
| 60 | BORGEAT P | | P NAS US | 69 | 2677 | | 10 | STIMULATION OF ADENOSINE 3',5'-CYCLIC MONOPHOSPHATE ACCUMULATION IN ANTERIOR-PITUITARY GLAND IN-VITRO BY SYNTHETIC LUTEINIZING HORMONE-RELEASING HORMONE |
| 60 | BORVENDE J | | J ENDOCR | 55 | 207 | N | 2 | OVULATION INDUCED BY SYNTHETIC LUTEINIZING-HORMONE RELEASING FACTOR IN ANDROGEN-STERILIZED FEMALE RATS |
| 60 | BRETON B | | CR AC SCI D | 274 | 2530 | | 6 | FR RECIPROCAL ACTIVITY OF HYPOTHALMIC FACTORS OF RAMS %OVIS-APIES AND TELEOSTEAN FISH ON SECRETION IN-VITRO OF GONADOTROPIC HORMONES C-HG AND LH RESPECTIVELY BY HYPOPHYSIS OF CARP AND RAMS |
| 60 | BURGER HG | | J ENDOCR | 54 | 227 | | 4 | LUTEINIZING-HORMONE RELEASING FACTOR IN ULTRAFILTRATES OF BLOOD COLLECTED FROM PITUITARY STALK OF OVARIECTOMIZED RATS AND RATS SUBJECTED TO ELECTRICAL STIMULATION OF PREEPTIC AREA |
| 60 | BURCUS R | | P NAS US | 69 | 278 | | 10 | PRIMARY STRUCTURE OF OVINE HYPOTHALMIC LUTEINIZING HORMONE-RELEASING FACTOR %LRF |
| 60 | CHANG JK | | BIOC BIOP R | 47 | 727 | | 7 | HYPOTHALMIC HORMONES . 37. LUTEINIZING RELEASING HORMONE. SYNTHESIS AND ARGS-ANALOGS. AND CONFORMATION-SEQUENCE-ACTIVITY RELATIONSHIPS |
| 60 | CHANG JK | | BIOC BIOP R | 47 | 1256 | | 7 | STUDIES ON ANALOGS OF LUTEINIZING RELEASING HORMONE TOWARDS ELUCIDATION OF RELEASE MECHANISM |
| 60 | CHANG JK | | J MED CHEM | 15 | 623 | | 3 | HYPOTHALMIC HORMONES . 36. SYNTHESIS OF LUTEINIZING-RELEASING HORMONE OF HYPOTHALMUS AND 8-LYSINE ANALOG |
| 60 | DEBELJUK L | | ENDOCRINOL | 90 | 585 | N | 3 | STUDIES ON PITUITARY RESPONSIVENESS TO LUTEINIZING HORMONE RELEASING HORMONE %LH-RH IN INTACT MALE RATS OF DIFFERENT AGES |

Listing of word pairs for clusters 1972 Level 11 Fred 2

| | | | | | | | | |
|---------------------|---------------------|----|----|----|---|---|----|-----|
| Releasing | Factor | 60 | 7 | 0 | 0 | 0 | 7 | 1-0 |
| Synthetic | Hormone-releasing | 60 | 0 | 5 | 0 | 1 | 7 | 2-8 |
| Synthetic | Luteinizing-Hormone | 60 | 5 | 0 | 0 | 1 | 7 | 2-1 |
| LH | Hormone | 60 | 0 | 0 | 3 | 2 | 9 | 5-4 |
| Synthesis | Hormone | 60 | 0 | 1 | 0 | 1 | 9 | 7-1 |
| LH-Releasing | Hormone | 60 | 9 | 1 | 0 | 0 | 10 | 1-1 |
| Luteinizing-Hormone | Hormone | 60 | 0 | 5 | 1 | 2 | 10 | 4-2 |
| Synthetic | Releasing | 60 | 1 | 6 | 2 | 1 | 11 | 2-5 |
| Hormone | LH-RH | 60 | 10 | 1 | 1 | 0 | 12 | 1-2 |
| Hormone-Releasing | Hormone | 60 | 10 | 0 | 0 | 0 | 12 | 1-7 |
| LH | FSH | 60 | 9 | 1 | 0 | 1 | 13 | 2-5 |
| Luteinizing | Hormone-Releasing | 60 | 12 | 0 | 0 | 0 | 13 | 1-3 |
| Luteinizing-Hormone | Releasing | 60 | 11 | 0 | 0 | 0 | 13 | 2-3 |
| Luteinizing | Hormone | 60 | 1 | 12 | 1 | 0 | 16 | 2-5 |
| Releasing | Hormone | 60 | 16 | 0 | 0 | 1 | 17 | 1-1 |
| Synthetic | Hormone | 60 | 0 | 7 | 7 | 2 | 19 | 3-3 |
| <hr/> | | | | | | | | |
| Aromatic | Compounds | 61 | 0 | 2 | 0 | 0 | 2 | 2-0 |
| Electron-spin | Reactions | 61 | 0 | 0 | 1 | 0 | 2 | 4-3 |
| Electron-spin | Resonance | 61 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Electron-spin | Studies | 61 | 0 | 1 | 0 | 1 | 2 | 3-0 |
| Formation | Decay | 61 | 1 | 1 | 0 | 0 | 2 | 1-5 |
| Ions | Aqueous-solutions | 61 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Ketyl | Radicals | 61 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Pulse | Ions | 61 | 0 | 0 | 1 | 1 | 2 | 3-5 |
| Radiolysis | Ions | 61 | 0 | 1 | 1 | 0 | 2 | 2-5 |
| Resonance | Reactions | 61 | 0 | 1 | 0 | 1 | 2 | 3-0 |
| Resonance | Studies | 61 | 1 | 0 | 1 | 0 | 2 | 2-0 |
| Studies | Reactions | 61 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Pulse | Aqueous-solutions | 61 | 0 | 1 | 0 | 1 | 3 | 3-6 |
| Radiolysis | Aqueous-solutions | 61 | 1 | 0 | 1 | 1 | 3 | 2-6 |
| Pulse | Radiolysis | 61 | 5 | 0 | 0 | 0 | 5 | 1-0 |
| <hr/> | | | | | | | | |
| Amphetamine-induced | Behavior | 62 | 0 | 1 | 1 | 0 | 2 | 2-5 |
| Amphetamine-induced | Stereotyped | 62 | 1 | 1 | 0 | 0 | 2 | 1-5 |
| Apomorphine | L-Dopa | 62 | 0 | 2 | 0 | 0 | 2 | 2-0 |
| Apomorphine | Rats | 62 | 0 | 2 | 0 | 0 | 2 | 2-0 |
| Behavioral | Lesions | 62 | 0 | 0 | 0 | 0 | 2 | 6-0 |
| Behavioral | Rat | 62 | 0 | 0 | 0 | 2 | 2 | 4-0 |
| Central | Action | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Central | Dopamine | 62 | 0 | 1 | 1 | 0 | 2 | 2-5 |
| Central | Dopaminergic | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Central | Effect | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Central | Mondamine | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Central | Rats | 62 | 0 | 0 | 1 | 0 | 2 | 4-5 |
| Dopamine | Receptor | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Dopamine | Receptors | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Effect | Activity | 62 | 0 | 0 | 1 | 1 | 2 | 3-5 |
| Effect | Apomorphine | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Effect | Locomotor | 62 | 0 | 1 | 1 | 0 | 2 | 2-5 |
| Effects | Central | 62 | 0 | 1 | 0 | 1 | 2 | 3-0 |
| Evidence | Dopamine | 62 | 0 | 1 | 1 | 0 | 2 | 2-5 |
| Evidence | Receptors | 62 | 0 | 0 | 1 | 1 | 2 | 3-5 |
| Induced | Rat | 62 | 0 | 2 | 0 | 0 | 2 | 2-0 |
| Locomotor | Activity | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Model | Tardive | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Mondamine | Neurons | 62 | 2 | 0 | 0 | 0 | 2 | 1-0 |
| Rat | Lesions | 62 | 1 | 0 | 1 | 0 | 2 | 2-0 |
| Amphetamine | Fat | 62 | 1 | 1 | 0 | 1 | 3 | 2-3 |
| Central | Neurons | 62 | 0 | 2 | 0 | 1 | 3 | 2-6 |
| Effect | Rats | 62 | 0 | 0 | 2 | 0 | 3 | 3-6 |

Fig. 7

documents in another cluster. The source item would then be assigned two classification headings, one with a weight of three and the other with a weight of two. A test of the effectiveness of this method must involve a comparison of the results of this automatic classification procedure with manual indexing procedure performed on a sample of source documents. The system must also be tested in user studies, since a great deal will depend on how well we have identified and named the subject of each of the clusters. As with any system, we cannot hope to please every user, but rather to develop a system which will satisfy the needs of a maximum number of users for the minimum cost. The advantage of the automatic procedure described in this paper is that all manipulations, save the naming of the clusters, are totally automatic and require no human judgement.

Since the theme of this conference is the ordering of global information networks, it is appropriate that we discuss the connections between our citation clustering experiment and the need for a global classification scheme. The application of citation data in the creation of a classification scheme has the advantage of being closely geared to the international activity of the scientific community which have established these citation patterns through their publications. Since scientific specialties do not have national boundaries, we believe the citation approach is a fair procedure for identifying subject areas which are of interest to many different countries. Secondly, bibliographic citations themselves are an international language. Clusters of citations may, therefore, be named in any language, but the content remains defined by the cited documents. Hence, it is possible to envision truly international classification scheme based on the *Science Citation Index*^R with subject experts in every country naming clusters according to that country's scientific usage. This may not really be necessary if English becomes the international language of science, but even if this does not occur, citation indexing still remains an indexing language which is essentially free of semantic or linguistic problems. Our main problem is in dealing with the variety of alphabets and symbol systems in Japanese, Chinese, Russian, etc. Such a system would go a long way towards improving worldwide exchange of information to the benefit of all countries involved.

REFERENCES

- [1] Garfield, E. .. *Science Citation Index—A New Dimension in Indexing, Science*, 1964, 144 (3619) 649-654.
- [2] Malin, M. V. .. *The Science Citation Index—A New Concept in Indexing, Library Trends*, 1968, 16 (3), 374-387.

- (3) Weinstock, M. .. Citation Indexes. *Encyclopedia of Library and Information Science*, 1971, 5, 16-40.
- [9] Garfield, E. .. Citation Indexes for Science, *Science*, 1955, 122, (3150) 108-111.
- [5] Salton, G. .. Associative Document Retrieval Techniques Using Bibliographic Information, *Journal of Association of Computing Machinery*, 1963, 10, 4, 440.
- [6] Garfield, E. .. Can Citation Indexing Be Automated? in *Statistical Association Methods for Mechanized Documentation: Symposium Proceedings* (M. E. Stevens, V. E. Giuliano and L. B. Heilprin, eds.), National Bureau of Standards, Washington, D.C., Miscellaneous Publication, 1965, 269, 189-192.
- [7] Garfield, E. .. Citation Analysis as a Tool in Journal Evaluation, *Science*, 1972, 178 (4060), 471-479.
- [8] Sparck-Jones, K. .. Some Thoughts on Classification for Retrieval, *Journal of Documentation*, 1970, 26, (2)
- [9] Small, H., and Griffith, B.C. The structure of Scientific Literatures I: Identifying and Graphing Specialties, *Science Studies*, 1974, 4, 17-40.