

ISR-9  
August, 1965

XXIV. INFORMATION SEARCH OPTIMIZATION AND  
ITERATIVE RETRIEVAL TECHNIQUES

J. J. Rocchio and G. Salton <sup>\*</sup>

ABSTRACT

The introduction of operationally effective time-shared computer systems, including suitable hardware for manual input, and display devices for computer-generated output, may be expected to produce important changes in the organization of automatic information retrieval systems. Specifically, the opportunity on the part of the user of such a system to interact with the system promises to lead to significant improvements in retrieval service and in the effectiveness of the search procedures.

In the present study several techniques are considered for exploiting a real-time communications link between the user and the retrieval system. Three principal methods are described in detail: the use of the system for the display of appropriate portions of the stored vocabulary to enable the requester to adjust the original query; the use of relevance judgments obtained from earlier partial searches to refine subsequent searches; and, finally, the use of a multiplicity of different information analysis methods. Examples are shown of each of the suggested processes, and evaluation coefficients are presented in an attempt to measure the amount of improvement obtainable with each procedure.

---

<sup>\*</sup> Computation Laboratory, Harvard University, Cambridge 38, Mass.  
This study was supported in part by the National Science Foundation under research grant GN-360.

## 1. Introduction

Automatic information retrieval systems must be designed to serve a multiplicity of users, each of whom may have different needs and may consequently require different kinds of service. Under these circumstances, it appears reasonable that the system should reflect this diversity of requirements by providing a role for the user in determining the search strategy. This is particularly important in automatic systems, where presently used one-shot (keyword) search procedures normally produce poor results.

In an automatic retrieval environment in which the user may be given access to the system - for example, by means of special input-output consoles - this can be achieved by two principal methods:

- (a) by providing automatic aids to the user in his attempt to formulate effective search requests;
- (b) by using several alternative analysis procedures, and a sequence of search steps to improve retrieval results.

In either case, the user can be made to control the retrieval process by asking him to furnish to the system information which subsequently determines, at least in part, the search strategy for a later pass.

Several methods may be employed to aid the user in formulating effective search requests. One of the simplest methods consists in providing some kind of automated dictionary which may be used to display certain pertinent parts of the stored information. Thus words, or

concepts, related to those originally included in a search request may be exhibited, and the user may be asked to choose from among these related terms in reformulating his request. The automated dictionary is then used as an aid in a manual reformulation of the request.

The iterative search process can also be mechanized more completely by leaving the search request largely unchanged, and by altering instead the information analysis process. In that case, the user furnishes to the system information concerning the adequacy of a preceding search operation, which is then used automatically to adjust the retrieval process for the next iteration.

In the present study, several alternative search optimization procedures are examined. In each case, the automatic SMART document retrieval process, presently operating on the IBM 7094 computer in a batch-processing mode, is used to simulate the real-time iterative search process.<sup>1,2,3</sup> The automatic evaluation procedures incorporated into SMART are utilized to measure the effectiveness of each process, and data are obtained which reflect the relative improvement of the iterative, user-controlled, process over and above the usual single-pass search procedure.

## 2. The Automatic Dictionary Process

In a conventional, batch-processing retrieval environment, the user normally relies on his intuition and experience - possibly aided by published references - in formulating an initial search request. Once the general

context has been established, the request must be normalized to a form suitable for use by the retrieval system. In a conventional coordinate indexing system, for example, this normalization would consist in a manual transformation of the original search request into an appropriate set of keywords. In certain automatic keyword search systems, a machine indexing process would generate the keywords, and stored synonym dictionaries might be used for normalization. After the analysis process, the normalized identifiers which specify the search request are matched with the identifiers attached to the documents, and correlation coefficients are obtained to measure the similarity between documents and search requests.

In the present section, a system is considered in which a communications link enables the user to affect the normalization process by making it possible for him to choose certain terms to be added and/or deleted from an original search formulation. Four main procedures appear to be of interest for this purpose:

- (a) a stored synonym dictionary, or thesaurus, may be used, given a set of thesaurus entries, to display all related entries appearing under the same concept category;<sup>1,2,3</sup>
- (b) a hierarchical arrangement of terms or concept classes may be available which, given a set of initial terms, can provide more general concepts by going "up" in the hierarchy, or more specific ones by going "down";<sup>1,2,3</sup>
- (c) a statistical term-term association matrix may be computed which can be used, given a set of terms, to find all those related terms which exhibit a tendency to co-occur in many documents of the collection with the terms originally specified;<sup>4</sup>



- (d) assuming the availability of a set of documents retrieved by an initial search operation, one may add to the terms originally specified in a search request, all those terms which occur in several of the retrieved documents but do not occur in the initial request.<sup>5</sup>

While it is potentially very useful to provide the user with a set of terms which may have been overlooked in formulating the original search request, it is probably even more important to furnish an indication of the usefulness in the retrieval process of each of the query terms. The most obvious indicator of potential usefulness is the density (or absolute number) of documents identified by each of the given index terms. The assumption to be made in this connection is that the usefulness of a term varies inversely with the frequency with which it is assigned to the documents of a collection.

Thus, in a coordinate indexing system, in which the retrieval process is controlled by the number of matches between terms assigned to documents and terms assigned to the search requests, the indexing density provides a straightforward estimate of the number of documents likely to be retrieved in each particular operation. If a correlation function is used to compare keyword sets attached to documents and queries, the relation between number of retrieved documents and the size of each keyword set is less obvious. However, the general assumption that a query term with high indexing density will produce "broad" retrieval, whereas one with low indexing density produces "narrow" retrieval is still valid.

It seems reasonable under the circumstances, to require that each dictionary display provided to the user consist not only of the corresponding terms or concepts, but also of the frequency with which the various terms are assigned to the documents of the collection. The user can then utilize this information to refine the search request by promoting terms deemed important and demoting others which may be ambiguous or otherwise useless in the retrieval process.

Consider as an example the retrieval procedures illustrated for two different requests in Figs. 1 and 2, respectively. The original text of a request titled "IR Indexing" is shown in Fig. 1(a). Five concept numbers are obtained when this text is looked up in the "Harris III" thesaurus (version number 3 of the synonym dictionary) provided with the SMART system; these concept numbers, together with their frequencies in the document collection are shown in Fig. 1(b). The full thesaurus entries corresponding to these concept numbers are similarly presented in Fig. 1(c); finally, retrieval results obtained with the original search request are given in Fig. 1(d).

Under the assumption that the user examines the list of retrieved documents, and finds that the fifth and sixth documents (numbers 79 and 80) are useful to him, it is now possible to request that concepts attached to these documents, but not included in the original search request, be displayed. This is done in Fig. 1(e) for concepts jointly included in the relevant documents numbers 79 and 80.

It now becomes possible for the user to pick new terms from the list of Fig. 1(e) - for example, terms like "coordinate," "lookup," and

"Automatic Information Retrieval and Machine Indexing"
---

(a) Original query text for "IR Indexing"

Original Term Used in Request	Concept Number	Frequency of Concept (405 documents)
automatic	119	79
information	350	45
retrieval	26	6
machine	600	77
indexing	101	11

(b) Terms included in original request

Concept Number	Corresponding Thesaurus Entries
119	artificial, automatic, mechanical, machine-made, semi-automatic, semiautomatic
350	information
26	retrieval
600	machine
101	descriptor, flag, ID, index, key, keyword, label, subscript, tag, sub-script

(c) Thesaurus entries corresponding to original terms used

Processing of Request "IR Indexing"

Figure 1

Document Rank	Document Number	Correlation Coefficient	Relevant
1	167	0.41	no
2	166	0.38	no
3	129	0.33	no
4	314	0.33	no
5	79	0.33	yes
6	80	0.30	yes

(d) Retrieval results for original query  
(using version III of Harris thesaurus)

Concepts from Documents 79 and 80	Corresponding Thesaurus Entries
49	co-ordinate, <u>coordinate</u> , intercept, ordinate, pole, rectangular-to-polar
108	consult, look-up, look, <u>lookup</u> , scan, seek, search
114	<u>abstract</u> , article, auto-abstracting, bibliography, catalog, copy, etc.
170	noun, verb, sentence
497	science

(e) Concepts common to relevant documents number 79 and 80  
and not included in original request

Figure 1 (continued)

"Information Retrieval. Document Retrieval. Coordinate Indexing. Dictionary Look-up for Language Processing. Indexing and Abstracting of Texts."
---

(f) Modified query using terms  
from relevant documents

Retrieval Results Using Original Query			Retrieval Results Using Modified Query		
Ranks of Relevant Documents	Document Number	Correlation	Ranks of Relevant Documents	Document Number	Correlation
5	79	0.33	1	80	0.51
6	80	0.30	4	79	0.41
9	221	0.29	6	48	0.36
11	126	0.28	9	126	0.23
12	48	0.27	11	221	0.23
69	3	0.10	18	3	0.19

(g) Comparison of search results using original and modified queries

Figure 1 (continued)

"Can hand-sent Morse code be transcribed automatically into English? What programs exist to read Morse code?"

(a) Original query for "Morse Code"

Term Used in Request	Concept Number	Frequency of Concept (405 documents)
hand-sent	113	12
Morse	35	9 (LOW)
code	281	37
transcribed	570	25
automatically	119	70 (HIGH)
English	35	9
programs	608	104 (HIGH)
exist	234	55
read	569	25

(b) Terms included in original request

<p>Modification 1: "Can hand-sent Morse code be translated into English? Recognition of manual Morse code."</p> <hr/> <p>Modification 2: "Use original query and add 'Morse, Morse, Morse'."</p>
--

(c) Modified queries by deletion of common, high-frequency concepts and addition of important low-frequency concepts

Processing of Request "Morse Code"

Figure 2

Type of Query	Ranks of Relevant Documents	Document Number	Correlation
Original Query "Morse Code"	7	394	0.29
	30	305	0.13
Modification 1	4	394	0.33
	8	305	0.26
Modification 2	4	394	0.30
	16	305	0.13

(d) Comparison of search results using original and modified queries

Figure 2 (continued)

"abstract" - and to use them to rephrase the search request as shown in Fig. 1(f). A comparison of the retrieval effectiveness for both original and modified queries is shown in Fig. 1(g). It may be noticed by examining the ranked document output produced in the SMART system, that the relevant documents have much lower rank, and correspondingly higher correlation coefficients for the modified search request than for the original. The lowest relevant document, in fact, places only 18th out of a total of 405 documents when the modified query is used, whereas it originally ranks 69th.

A second example, and a different dictionary feedback process, is illustrated in Fig. 2 for the request "Morse Code." The original text for the request is given in Fig. 2(a), and the corresponding thesaurus entries and their frequencies appear in Fig. 2(b). The user who examines the output of Fig. 2(b) may notice that concepts 119 (obtained from "automatically") and 608 (from "programs") appear with an excessively high frequency and that it may therefore be useful to remove them from the request statement.

Similarly, the crucial concept 35 ("Morse") appears with very low frequency. The reformulations of Fig. 2(c) reflect the corresponding deletions and additions.

The success of the request alteration may be evaluated by examining the ranks of the two relevant documents (numbers 305 and 394) as shown in Fig. 2(d). It may be seen that retrieval results are improved for both modifications 1 and 2 over the original, but that the better result is obtained for the first modification where the relevant documents are ranked fourth and eighth, respectively.

### 3. Request Optimization Using Relevance Feedback

The vocabulary feedback process illustrated in the preceding section appears to be both easy to implement and effective in improving search results. It does, however, put considerable demands upon the user who controls not only what is displayed by the system, but also what is returned in the way of modified information. A variety of search optimization methods should, therefore, be considered which place a much larger burden on the system, and a correspondingly smaller one on the user. One such procedure is the relevance feedback method.

In essence, the process consists in effecting an initial search, and in presenting to the user a certain amount of retrieved information. The user then examines some of the retrieved documents and identifies each as being either relevant (R) or not relevant (N) to his purpose. These relevance judgments are then returned to the system, and are used automatically to adjust the initial search request in such a way that query terms



or concepts present in the relevant documents are promoted (by increasing their weight), whereas terms occurring in the documents designated as nonrelevant are similarly demoted.

The amount of improvement to be obtained from the feedback process depends critically on the manner in which the search request is altered as a function of the user's relevance judgment. The following process which has been used experimentally with the SMART system appears to be optimal in this connection. Consider a retrieval system in which the matching function between queries and documents (or between query and document identifiers) induces a metric, or a monotonic function of a metric, on the space of query and document images (e.g., on the space of keyword vectors).<sup>6</sup> In such a case, it is possible to produce an ordering of the documents with respect to the input query in such a way that increasing distance between document and query images reflects increasing dissimilarity between them.

Let  $D_R$  be the nonempty subset of relevant documents from the source collection  $D$ , relevance being defined subjectively and outside the context of the system. An optimal query can now be defined as that query which maximizes the difference between average distances from the query to the relevant document set, and from the query to the nonrelevant set. In other words, the optimal query is the one which provides the maximum discrimination of the subset  $D_R$  from the rest of the collection ( $D - D_R$ ). More formally, let  $\delta(q, d)$  be the distance function used in the matching process between query  $q$  and document  $d$ . The optimal query  $q_0$  may then be defined as that query which maximizes the function

$$C = \widetilde{\int}_{d \notin D_R} (q, d) - \widetilde{\int}_{d \in D_R} (q, d), \quad (1)$$

where  $\tilde{\rho}$  is the average distance function, and decreasing distance implies stronger query-document association.

Clearly, equation (1) is of no practical use, even under the assumption that the optimal query  $q_0$  can be determined as a function of  $D$  and  $D_R$ , since knowledge of the set  $D_R$  (the relevant document subset) obviates the need for retrieval. However, if instead of producing the optimal query  $q_0$ , the relation (1) is used to produce a sequence of approximations to  $q_0$ , starting with some initial query which identifies a part of the set  $D_R$ , then a method for automatically generating useful query modifications becomes available. The system can, in fact, produce the optimal query to differentiate the partial set of relevant documents, identified by the user, from the remaining documents; the resultant query can then be resubmitted, and the process may be iterated, as more complete sets of relevant documents become available through subsequent retrieval operations. One may hope that only a few iterations will suffice for the average user; in any case, the rate of convergence will be reflected in the stability of the retrieved set.

In the SMART automatic document retrieval system, the query-document matching function normally used is the cosine correlation of the query vector with the set of document vectors, defined as

$$\rho(\bar{q}, \bar{d}) = \frac{\bar{q} \cdot \bar{d}}{|\bar{q}| |\bar{d}|} = \cos \theta_{\bar{q}, \bar{d}}, \quad (2)$$

where  $\bar{q}$  and  $\bar{d}$  are the vector images of query  $q$  and document  $d$ , respectively. Since the vector images are limited to nonnegative components, the range for the correlation is  $0 \leq \rho \leq 1$ , corresponding to an angular separation of

$q_0 \leq \theta \leq 0$ . Under these conditions, the correlation coefficient is a monotonic function of the angular distance metric. Furthermore, since the correlation decreases with increasing distance, relation (1) may be rewritten as:

$$C = \widetilde{\rho}(\bar{q}, \bar{d}) - \widetilde{\rho}(\bar{q}, \bar{d}), \quad (3)$$

$\bar{d} \in D_R$                        $\bar{d} \notin D_R$

where  $\widetilde{\rho}$  is the average cosine function  $\rho$ . It can be shown,<sup>7</sup> that in this case C is maximized for

$$\bar{q}_0 = \frac{1}{n_0} \sum_{\bar{d}_i \in D_R} \frac{\bar{d}_i}{|\bar{d}_i|} - \frac{1}{N-n_0} \sum_{\bar{d}_i \notin D_R} \frac{\bar{d}_i}{|\bar{d}_i|}, \quad (4)$$

where  $n_0 = n(D_R)$ , the number of elements in the set  $D_R$ , and  $N = n(D)$ , the number of elements in the collection.

The query modification algorithm employed may now be written in the form:

$$\bar{q}_{i+1} = n_1 n_2 \bar{q}_i + n_2 \sum_{i=1}^{n_1} \frac{\bar{r}_i}{|\bar{r}_i|} - n_1 \sum_{i=1}^{n_2} \frac{\bar{s}_i}{|\bar{s}_i|}, \quad (5)$$

where  $\bar{q}_i$  is the  $i$ th query of a sequence, and  $R = \{\bar{r}_1, \bar{r}_2, \dots, \bar{r}_{n_1}\}$  is the set of relevant document vectors retrieved in response to query  $\bar{q}_i$ , and  $S = \{\bar{s}_1, \bar{s}_2, \dots, \bar{s}_{n_2}\}$  is the set of nonrelevant document vectors retrieved in response to  $\bar{q}_i$  (the specification of the sets R and S constitute the feedback from the user after the  $i$ th iteration of the process).

Document Rank	Document Number	Correlation	User Feedback
1	351	.65	relevant
2	353	.42	relevant
3	350	.41	relevant
4	163	.36	-
5	82	.35	-
6	1	.32	-
7	208	.27	not relevant
8	225	.25	not relevant
9	54	.24	-
10	335	.21	not relevant

(a) Retrieval results using original query for  
 "Pattern Recognition" (version II  
 of Harris thesaurus)

Retrieval Results Using Original Query			Results Using Query Modified by User Feedback		
Ranks of Relevant Documents	Document Number	Correlation	Ranks of Relevant Documents	Document Number	Correlation
1	351	.65	1	351	.66
2	353	.42	2	350	.60
3	350	.41	3	353	.55
4	163	.36	5	163	.37
6	1	.32	6	1	.32
9	54	.24	7	54	.29
26	205	.17	11	314	.23
27	224	.17	16	205	.19
33	314	.16	17	39	.19
34	39	.12	30	224	.16
Recall .972 Precision .864			Recall .989 Precision .923		

(b) Comparison of search results using original and modified queries

Query Processing Using Relevance Feedback

Figure 3

Search Request	Original Query		Modified Query	
	Recall	Precision	Recall	Precision
IR Indexing	0.976	0.728	0.991	0.928
Pattern Recognition	0.972	0.864	0.989	0.923
Analog-Digital	0.984	0.870	0.983	0.918
Average over 17 Search Requests	0.970	0.876	0.975	0.918

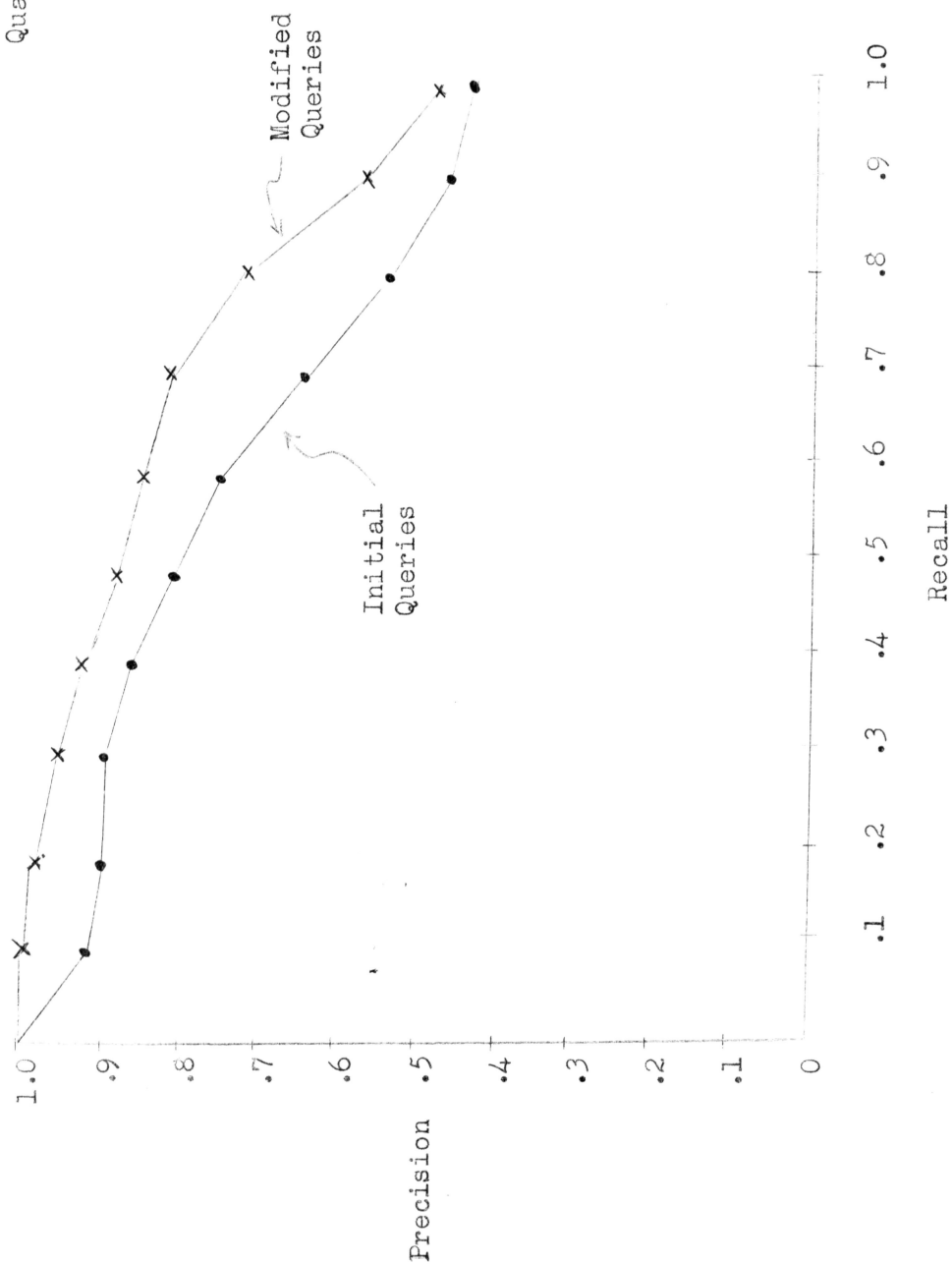
Evaluation Figures Showing Effect of One-step  
Relevance Feedback

Figure 4

The foregoing method of automatic query optimization was tested by performing a single iteration using a set of 17 search requests, previously available in the SMART system. Figure 3 shows the results for a request on "Pattern Recognition." The original retrieval results, using version 2 of the "Harris" thesaurus, are given in Fig. 3(a). The user identifies documents 351, 353, and 350 as relevant, and 208, 225, and 335 as nonrelevant. The query is then automatically modified, in accordance with the expression of equation (5), and retrieval performance is compared in Fig. 3(b). It may be seen that drastic improvements are obtained both in the ranks of the relevant documents and in the magnitude of the correlation coefficients. The "recall" and "precision" measures, shown in Fig. 3(b), are the normalized evaluation measures incorporated into the SMART system.<sup>8,9</sup>

Figure 4 shows individual retrieval results for three queries, and normalized recall and precision values averaged over 17 different requests.

Thesaurus: Harris II  
Cosine Correlation  
Quasi-Cleverdon Graph



Precision Versus Recall for Initial Queries  
and Queries Modified by Relevance Feedback  
(averaged over 17 search requests)

Figure 5

The feedback process may be seen to improve both recall and precision in almost every case. This is also reflected in the precision versus recall plot of Fig. 5. The recall and precision values in Fig. 5 are the standard evaluation measures, originally introduced by Cleverdon,<sup>10</sup> and not the normalized measures previously shown in Figs. 3 and 4. The method of construction of the precision-recall graph of Fig. 5 has previously been described in detail.<sup>9</sup> The positive effect of the relevance feedback is again obvious from the figure.

#### 4. Automatic Modification of the Analysis Process

The last search optimization process to be described depends, like its predecessor, on feedback provided by the user, and results in selective changes in the document and request analysis process. However, instead of furnishing relevance judgments based on the output of a previous retrieval operation, the user makes a qualitative assessment of the effectiveness of an initial search operation. For example, he may find that the documents obtained from the system show that his request was interpreted too narrowly (since all retrieved documents belong to some small subfield of the larger area which he expected to cover), or too broadly, or too literally, or too freely.

Depending on the type of interpretation furnished by the user, the system now proceeds to initiate a new search operation under altered analysis procedures. If the user's verdict was "too narrow," a hierarchical subject arrangement similar to the one previously mentioned in part 2 might be consulted, and each original query term could be replaced by a broader one; if,

on the other hand, the initial search was "too broad," more specific terms might be obtained from the hierarchy. If the interpretation was too literal, the use of a synonym dictionary might provide more reasonable results; and so on.

Automatic retrieval systems are particularly attractive in such a situation, because these systems make it possible to provide at relatively little extra cost, a variety of indexing procedures which may be called upon as needed. The SMART system, in particular, provides a large variety of indexing methods including the following:

- (a) a null thesaurus procedure which uses the word stems originally included in documents and search requests for content identification;
- (b) a synonym dictionary ("Harris" thesaurus) which replaces original word stems by synonym classes or concept numbers;
- (c) a hierarchical arrangement of concept numbers which can be used, given a set of concepts to obtain more general ones ("hierarchy up"), or more specific ones ("hierarchy down");
- (d) a statistical phrase procedure which is used to replace pairs or triples of co-occurring (related) concepts by a single "phrase" concept (e.g., the concepts "program" and "language" might be combined into "programming language");
- (e) a syntactic phrase process which generates phrases only if the components in fact exhibit an appropriate grammatical relationship;
- (f) a variety of so-called merged methods,<sup>9</sup> in which the system proceeds iteratively through two or three simple processes and combines the output.



Obviously, the ability to generate a multiplicity of distinct index images for each document does not necessarily imply that each modification in the analysis process results in large-scale improvements in the search effectiveness. Experiments conducted with the SMART system have, however, shown that in many cases a considerable increase in retrieval effectiveness is obtainable when changes in the analysis are adapted to the aims of each particular user.

Consider, in this connection, the evaluation output for a variety of analysis methods produced by the SMART system, reproduced in Figs. 6 and 7. Figure 6 contains output for the search request titled "Automata Phrases," with nine relevant documents. Six simple analysis methods are shown: null thesaurus, "Harris Two" (version 2 of the regular synonym dictionary), statistical phrases, syntax phrases, hierarchy up, and hierarchy down. Thirteen "merged" methods, each including two simple components, are also included in Fig. 6, as well as nine triple merges. For each method, the output is presented in two parts: the left part includes the document numbers of the first 15 documents retrieved by that method, whereas the right-hand side consists of only the relevant document numbers and their ranks in decreasing correlation order with the request. Below the lists of document numbers, a variety of recall and precision measures are provided for each analysis procedure, to reflect the effectiveness of the corresponding process.

An examination of Fig. 6 reveals, for example, that for the request on "Automata Phrases," improved retrieval is obtained by switching from the word stem procedure to the synonym recognition process using the regular thesaurus (labeled ① in Fig. 6). This is reflected both by the magnitude

AUTOMATA PHR		9 RELEVANT		HARRIS TWO		NULL THES		STAT PHRASE		SYNTAX PHR		HIERARCHY UP		HIER DOWN		HARRIS TWO	
STAT PHRASE		SYNTAX PHR		HARRIS TWO		NULL THES		STAT PHRASE		SYNTAX PHR		HIERARCHY UP		HIER DOWN		HARRIS TWO	
TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT	
1 316	1 316	1 129	1 129	1 371	1 371	1 371	1 371	1 371	1 371	1 371	1 371	1 264	1 264	1 316	1 316	1 316	1 316
2 129	2 129	2 316	2 316	2 316	2 316	2 316	2 316	2 316	2 316	2 316	2 316	2 316	2 316	2 129	2 129	2 129	2 129
3 313	3 313	3 167	3 167	3 372	3 372	3 372	3 372	3 372	3 372	3 372	3 372	3 212	3 212	3 173	3 173	3 313	3 313
4 176	4 176	4 173	8 313	4 212	6 313	4 002	5 313	4 316	5 313	4 316	5 313	4 218	22 313	4 238	5 176	4 167	5 176
5 167	7 371	5 166	10 241	5 002	8 176	5 313	7 241	5 313	6 176	5 129	50 176	5 129	50 176	5 176	7 371	5 176	10 371
6 249	20 372	6 213	12 372	6 313	9 241	6 135	6 176	6 176	7 241	6 127	68 371	6 305	33 315	6 305	33 315	6 173	14 241
7 371	36 241	7 371	16 264	7 139	10 129	7 241	5 129	7 241	8 129	7 361	116 372	7 371	30 372	7 371	30 372	7 166	17 372
8 166	42 315	8 313	56 315	8 176	42 315	8 176	39 315	8 129	36 315	8 263	133 241	8 249	43 241	8 249	43 241	8 249	25 264
9 385	74 264	9 263	104 176	9 241	84 264	9 129	64 264	9 245	58 264	9 177	196 315	9 167	92 264	9 167	92 264	9 213	70 315
10 045		10 241		10 129		10 085		10 249		10 340		10 161		10 371		10 371	
11 173		11 259		11 245		11 245		11 283		11 123		11 166		11 385		11 385	
12 101		12 372		12 089		12 245		12 167		12 128		12 204		12 263		12 263	
13 204		13 060		13 249		13 212		13 166		13 045		13 045		13 045		13 045	
14 117		14 106		14 283		14 283		14 173		14 171		14 209		14 241		14 241	
15 238		15 067		15 167		15 167		15 101		15 213		15 348		15 259		15 259	
RNK REC= 0.2250	RNK REC= 0.2250	RNK REC= 0.2083	RNK REC= 0.2083	RNK REC= 0.3061	RNK REC= 0.3061	RNK REC= 0.3261	RNK REC= 0.3261	RNK REC= 0.3543	RNK REC= 0.3543	RNK REC= 0.4059	RNK REC= 0.4059	RNK REC= 0.2000	RNK REC= 0.2000	RNK REC= 0.3061	RNK REC= 0.3061	RNK REC= 0.3061	RNK REC= 0.3061
LOG PRE= 0.6347	LOG PRE= 0.6347	LOG PRE= 0.6111	LOG PRE= 0.6111	LOG PRE= 0.7077	LOG PRE= 0.7077	LOG PRE= 0.7338	LOG PRE= 0.7338	LOG PRE= 0.7465	LOG PRE= 0.7465	LOG PRE= 0.7465	LOG PRE= 0.7465	LOG PRE= 0.6158	LOG PRE= 0.6158	LOG PRE= 0.6867	LOG PRE= 0.6867	LOG PRE= 0.6867	LOG PRE= 0.6867
NCR REC= 0.9565	NCR REC= 0.9565	NCR REC= 0.9520	NCR REC= 0.9520	NCR REC= 0.9714	NCR REC= 0.9714	NCR REC= 0.9739	NCR REC= 0.9739	NCR REC= 0.9770	NCR REC= 0.9770	NCR REC= 0.9770	NCR REC= 0.9770	NCR REC= 0.9405	NCR REC= 0.9405	NCR REC= 0.9471	NCR REC= 0.9471	NCR REC= 0.9471	NCR REC= 0.9471
ACR PRE= 0.8208	ACR PRE= 0.8208	ACR PRE= 0.8019	ACR PRE= 0.8019	ACR PRE= 0.8714	ACR PRE= 0.8714	ACR PRE= 0.8771	ACR PRE= 0.8771	ACR PRE= 0.8943	ACR PRE= 0.8943	ACR PRE= 0.8943	ACR PRE= 0.8943	ACR PRE= 0.8057	ACR PRE= 0.8057	ACR PRE= 0.8575	ACR PRE= 0.8575	ACR PRE= 0.8575	ACR PRE= 0.8575
OVERALL= 0.8597	OVERALL= 0.8597	OVERALL= 0.8194	OVERALL= 0.8194	OVERALL= 1.0118	OVERALL= 1.0118	OVERALL= 1.0559	OVERALL= 1.0559	OVERALL= 1.1008	OVERALL= 1.1008	OVERALL= 1.1008	OVERALL= 1.1008	OVERALL= 0.9259	OVERALL= 0.9259	OVERALL= 0.9158	OVERALL= 0.9158	OVERALL= 0.9158	OVERALL= 0.9158
NOR OVR= 1.6034	NOR OVR= 1.6034	NOR OVR= 1.5620	NOR OVR= 1.5620	NOR OVR= 1.7283	NOR OVR= 1.7283	NOR OVR= 1.7566	NOR OVR= 1.7566	NOR OVR= 1.7792	NOR OVR= 1.7792	NOR OVR= 1.7792	NOR OVR= 1.7792	NOR OVR= 1.5532	NOR OVR= 1.5532	NOR OVR= 1.7148	NOR OVR= 1.7148	NOR OVR= 1.7148	NOR OVR= 1.7148

TOP 15 DOCUMENTS

RELEVANT DOCUMENTS

RANKS OF RELEVANT DOCUMENTS

HARRIS TWO		HARRIS TWO		HARRIS TWO		HARRIS TWO		STAT PHRASE		STAT PHRASE		SYNTAX PHR		SYNTAX PHR	
STAT PHRASE		SYNTAX PHR		HIERARCHY UP		HIER DOWN		HIERARCHY UP		HIER DOWN		HIERARCHY UP		HIER DOWN	
TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT		TOP 15 RELEVANT	
1 316	1 316	1 316	1 316	1 316	1 316	1 316	1 316	1 371	1 371	1 371	1 371	1 371	1 371	1 371	1 371
2 371	2 371	2 371	2 371	2 371	2 371	2 371	2 371	2 371	2 371	2 371	2 371	2 371	2 371	2 371	2 371
3 129	3 129	3 129	3 129	3 129	3 129	3 129	3 129	3 313	3 313	3 313	3 313	3 313	3 313	3 313	3 313
4 313	4 313	4 372	4 372	4 313	4 313	4 176	4 176	4 372	4 372	4 372	4 372	4 316	4 316	4 129	4 129
5 372	5 372	5 313	5 313	5 173	5 173	5 238	5 238	5 371	5 371	5 371	5 371	5 313	5 313	5 212	6 313
6 176	6 176	6 212	7 176	6 176	11 371	6 167	34 372	6 002	9 129	6 002	9 129	6 173	9 129	6 313	8 176
7 002	11 241	7 176	10 241	7 218	50 372	7 245	39 315	7 218	12 241	7 238	11 241	7 218	10 176	7 238	10 241
8 167	46 315	8 167	43 315	8 167	65 241	8 385	43 241	8 313	14 176	8 176	42 315	8 313	12 241	8 176	42 315
9 249	73 264	9 249	68 264	9 249	72 315	9 371	82 264	9 129	68 315	9 139	74 264	9 129	64 315	9 385	69 264
10 139		10 241		10 127		10 166		10 139		10 385		10 176		10 241	
11 241		11 166		11 371		11 345		11 127		11 241		11 127		11 249	
12 166		12 385		12 361		12 161		12 241		12 249		12 241		12 245	
13 385		13 245		13 166		13 173		13 361		13 167		13 361		13 167	
14 045		14 045		14 263		14 161		14 176		14 089		14 263		14 161	
15 089		15 173		15 385		15 264		15 263		15 161		15 245		15 283	
RNK REC= 0.2980	RNK REC= 0.2980	RNK REC= 0.3147	RNK REC= 0.3147	RNK REC= 0.2103	RNK REC= 0.2103	RNK REC= 0.2074	RNK REC= 0.2074	RNK REC= 0.3719	RNK REC= 0.3719	RNK REC= 0.3000	RNK REC= 0.3000	RNK REC= 0.3982	RNK REC= 0.3982	RNK REC= 0.3103	RNK REC= 0.3103
LOG PRE= 0.7488	LOG PRE= 0.7488	LOG PRE= 0.7523	LOG PRE= 0.7523	LOG PRE= 0.6488	LOG PRE= 0.6488	LOG PRE= 0.6175	LOG PRE= 0.6175	LOG PRE= 0.7621	LOG PRE= 0.7621	LOG PRE= 0.7397	LOG PRE= 0.7397	LOG PRE= 0.7405	LOG PRE= 0.7405	LOG PRE= 0.7390	LOG PRE= 0.7390
NOR REC= 0.9703	NOR REC= 0.9703	NOR REC= 0.9725	NOR REC= 0.9725	NOR REC= 0.9526	NOR REC= 0.9526	NOR REC= 0.9517	NOR REC= 0.9517	NOR REC= 0.9787	NOR REC= 0.9787	NOR REC= 0.9705	NOR REC= 0.9705	NOR REC= 0.9809	NOR REC= 0.9809	NOR REC= 0.9719	NOR REC= 0.9719
ACR PRE= 0.8956	ACR PRE= 0.8956	ACR PRE= 0.8975	ACR PRE= 0.8975	ACR PRE= 0.8315	ACR PRE= 0.8315	ACR PRE= 0.8071	ACR PRE= 0.8071	ACR PRE= 0.9028	ACR PRE= 0.9028	ACR PRE= 0.8905	ACR PRE= 0.8905	ACR PRE= 0.9125	ACR PRE= 0.9125	ACR PRE= 0.8900	ACR PRE= 0.8900
OVERALL= 1.0468	OVERALL= 1.0468	OVERALL= 1.0670	OVERALL= 1.0670	OVERALL= 0.6591	OVERALL= 0.6591	OVERALL= 0.6248	OVERALL= 0.6248	OVERALL= 1.1340	OVERALL= 1.1340	OVERALL= 1.0397	OVERALL= 1.0397	OVERALL= 1.1788	OVERALL= 1.1788	OVERALL= 1.0493	OVERALL= 1.0493
NOR OVR= 1.7469	NOR OVR= 1.7469	NOR OVR= 1.7600	NOR OVR= 1.7600	NOR OVR= 1.5944	NOR OVR= 1.5944	NOR OVR= 1.5658	NOR OVR= 1.5658	NOR OVR= 1.7962	NOR OVR= 1.7962	NOR OVR= 1.7431	NOR OVR= 1.7431	NOR OVR= 1.8171	NOR OVR= 1.8171	NOR OVR= 1.7497	NOR OVR= 1.7497

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 15 RELEVANT

TOP 1



of the evaluation coefficients, and by the ranks of the last relevant document (104th out of 405 for the null thesaurus, and 74th for "Harris Two"). An improvement is also obtained by switching from the "Harris" thesaurus to the phrase procedures, and from statistical phrases to syntax phrases (labeled ② ). The third example from Fig. 6 shows that the merged procedure which combines the statistical phrases with the hierarchy results in an increase in performance over and above each of the component methods. A further improvement is obtained by adding the regular "Harris Two" thesaurus process (example ④ of Fig. 6).

Figure 7 shows evaluation output obtained for the same 28 analysis methods previously shown in Fig. 6, but averaged over 17 different search requests. The output of Fig. 7 is presented in the form of precision versus recall graphs, similar to that shown in Fig. 5 (the actual graphs are not drawn but tables are presented instead). The five examples specifically indicated in Fig. 7 again confirm the earlier results that improvements are obtainable from method to method.

Each of the three search optimization procedures described in this study appears to be useful as a means for improving the retrieval effectiveness of real-time, user-controlled search systems. Additional experimentation with larger document collections and with an actual user population may be indicated before incorporating these procedures in an operational environment. Iterative, user-controlled search procedures appear, however, to present an interesting possibility, and a major hope, for the eventual usefulness of large-scale automatic information retrieval systems.

## REFERENCES

1. Salton, G., "A Document Retrieval System for Man-machine Interaction," Proceedings of the ACM 19th National Conference, Philadelphia (1964).
2. Salton, G. and Lesk, M. E., "The SMART Automatic Retrieval System - An Illustration," Comm. of the ACM, Vol. 8, No. 6 (June 1965).
3. Salton, G. et al., "Information Storage and Retrieval," Reports No. ISR-7 and ISR-8 to the National Science Foundation, Computation Laboratory, Harvard University (June 1964 and December 1964).
4. Giuliano, V. E. and Jones, P. E., "Linear Associative Information Retrieval," Vistas in Information Handling, P. Howerton ed., Spartan Books, Washington (1963).
5. Curtice, R. M. and Rosenberg, V., "Optimizing Retrieval Results with Man-machine Interaction," Center for the Information Sciences, Lehigh University, Bethlehem, Pennsylvania (1965).
6. Rial, J. F., "A Pseudo-metric for Document Retrieval Systems," Working Paper W-4595, Mitre Corporation, Bedford, Massachusetts (1962).
7. Rocchio, J. J., Relevance Feedback in Information Retrieval, Report No. ISR-9 to the National Science Foundation, The Computation Laboratory of Harvard University, to appear August 1965.
8. Rocchio, J. J., "Performance Indices for Document Retrieval Systems," Report No. ISR-8 to the National Science Foundation, Sec. 3, The Computation Laboratory of Harvard University (December 1964).
9. Salton, G., "The Evaluation of Automatic Retrieval Procedures - Selected Test Results Using the SMART System," Report No. ISR-8 to the National Science Foundation, The Computation Laboratory of Harvard University (December 1964).
10. Cleverdon, C. W., "The Testing of Index Language Devices," ASLIB Proceedings, Vol. 15, No. 4 (April 1963).