

1966a

**ASLIB**  
**CRANFIELD RESEARCH**  
**PROJECT**

**FACTORS DETERMINING THE PERFORMANCE**  
**OF INDEXING SYSTEMS**  
**VOLUME 2. TEST RESULTS**

by

**Cyril Cleverdon and Michael Keen**

**An investigation supported by a grant to ASLIB**  
**by the National Science Foundation**

ASLIB CRANFIELD RESEARCH PROJECT

FACTORS DETERMINING THE PERFORMANCE  
OF INDEXING SYSTEMS  
VOLUME 2

Cyril Cleverdon and Michael Keen

An investigation supported by a grant to Aslib  
by the National Science Foundation

Cranfield

1966

Copyright © 1966

Cyril Cleverdon and Michael Keen

Published by  
Cyril Cleverdon  
Wharley End  
Bedford

REPORT ON THE TESTING AND ANALYSIS OF  
AN INVESTIGATION INTO THE COMPARATIVE  
EFFICIENCY OF INDEXING SYSTEMS

ERRATA

- p.18, bottom line "facet schedules as used in the project test were particularly adaptable"
- p.22, Table 3.2. Heading to columns should be added:  
16 mins.      12 mins.      8 mins.      4 mins.      2 mins.
- p.25, Table 3.8. Final line change "619" to read "69".
- p.44, Section (e) First line of third paragraph should read:  
"73 documents are included in this total of  
88 failures".
- p.55, line 12 After "collection" insert before bracket:  
"which were retrieved".
- p.105, Table 10.3. The terms "RECALL" and "RELEVANCE" should be placed against the vertical and horizontal co-ordinates respectively.



## SUMMARY

The test results are presented for a number of different index languages using various devices which affect recall or precision. Within the environment of this test, it is shown that the best performance was obtained with the group of eight index languages which used single terms. The group of fifteen index languages which were based on concepts gave the worst performance, while a group of six index languages based on the Thesaurus of Engineering Terms of the Engineers Joint Council were intermediary. Of the single term index languages, the only method of improving performance was to group synonyms and word forms, and any broader groupings of terms depressed performance. The use of precision devices such as links gave no advantage as compared to the basic device of simple coordination.

All results have to be considered within the context of the experimental environment, but they can be said to substantiate or clarify many of the findings of Cranfield I. It is conclusively shown that an inverse relationship exists between recall and precision, whatever the variable may be that is being changed. The two factors which appear most likely to affect performance are the level of exhaustivity of indexing and the level of specificity of the terms in the index language. For any given operational situation, the optimum levels cannot be categorically stated in advance, but can only be determined by an evaluation of the system, the main consideration probably being the subject field.

It would be unusual if the characteristics of the subject field used for this test were such as to make it unique, so the high performance obtained with the single terms in natural language can be considered to be of some importance in regard to the use of natural language text as input to mechanised systems.



## PREFACE

It was intended that this should be the final volume of the Report on Cranfield II. This may still be the case, but as the results were being prepared for publication, we were continually aware of the gaps that needed to be filled. The delay in the appearance of this volume is partly due to attempts to obtain some of the missing data, but a great deal still remains to be done. The detailed analysis of the reasons for failure to retrieve relevant documents or for the retrieval of non-relevant documents was an important part of Cranfield I, but so far in this project it has only been attempted in a superficial manner. It is most desirable that such analysis should be done, the more so because of the completely unexpected test results.

To the acknowledgements already made in Volume I, I would also wish to include Professor Salton and Professor Wilkins. It was a most happy chance that the National Science Foundation should fund two projects which, when they started, appeared to have quite different objectives, but which, as they progressed, were seen to be closely related. Cooperation with Professor Salton has enabled us to go much further than would otherwise have been the case, and I am most grateful for his continued assistance over the past two years and for providing data on his work which we have been able to present briefly in this report.

It was on the recommendation of a colleague that I read 'Social Deviance' by Professor L.T. Wilkins. Although dealing with a different subject, the author's opinions on what should be done in his particular field were entirely in line with my views on what is required in the field of documentation. For his agreement to allow me to include a series of short extracts from his book, I am most grateful to Professor Wilkins.

Cranfield  
December 1966

Cyril Cleverdon



## CONTENTS

	PAGE	
Chapter 1	Introduction	1
Chapter 2	Test Environment	34
Chapter 3	Methods for presentation of results	31
Chapter 4	Main test results	78
Chapter 5	Simulated ranking and document output cut-off	192
Chapter 6	Supplementary tests and results	221
Chapter 7	Citation indexing and bibliographic coupling	243
Chapter 8	Conclusions	252
References		264
Appendix 3A	Tables of generality number, and fallout, recall and precision ratios	265
Appendix 4A	Set of individual test results	288
Appendix 5A	Formula for document ranking based on probability considerations, by G.H. Stearman	296
INDEX		299



## CHAPTER 1

### Introduction

If two scientists disagree on any issue, and the issue is within the ambit of science, then it must be possible for them to agree on a procedure which they can both accept as a critical test of their points of difference. For reasons of personality they may not be able to get together to work out such a test procedure, but it must exist as a possibility. If such a critical test cannot be imagined as possible, then the issue between them is not a scientific issue. The scientific method does not vary with the subject-matter, but is the same irrespective of its results and basically the same in all the sciences.

L. T. Wilkins: Social Deviance, page 4.

In reviewing the final report on Cranfield I, N.D. Stevens (Ref. 1) described it as being 'extremely complex'; even after 'several careful readings' he found parts of it 'still bewildering', and said that 'there are so many side issues that the author neglects the clear and detailed presentation of the main headings; the reader finds himself sidetracked by these, or other interesting diversions'. Since this reviewer was by no means the only person who made such comments, it has been our particular endeavour in this report to make quite clear what has been done, how it has been done and what has been the outcome, even though at times this has led to what some people may consider undue verbosity and repetition.

In one respect, this project is easier to report, for, being in a more concentrated field, it does not raise many of the side-issues - such as indexing times, indexer qualifications, etc. - which came up in Cranfield I and which were sufficiently interesting to sidetrack the reader. On the other hand, to those who have been involved in Cranfield II, the earlier project seems to have been child's play to what has now been attempted, and the complexity of the present work is inevitably reflected in what has to be reported. To those readers who, like ourselves, tend to view with dismay the many papers on information retrieval which consist substantially of some twenty pages of mathematics, we can only apologise that it has become necessary to introduce a number of equations into this volume. However, it is certain that there is no mathematics

which should be beyond the comprehension of a schoolboy of average intelligence, but, even so, it is suggested that Chapter 3 might be skipped by those who are not closely concerned with the particular problem of performance measurement. Important though the work in this Chapter is felt to be, yet the arguments may well be of marginal interest to many readers. At the beginning of Chapter 4, which presents the main set of test results, full information is given concerning the performance measures which are actually used; Chapter 3 explains in some detail why those measures were selected in preference to other possible measures.

To a lesser extent, the same is true of Chapter 2 which discusses at length the variables which were being investigated and the environment in which the test was carried out. Again, we have tried to make Chapter 4 complete in itself in that such matters are briefly recorded therein. Only if the reader is puzzled as to why such seemingly unnecessarily tortuous actions have been taken, need he refer to Chapter 2 to find the possible justification.

The test results presented in Chapter 4 make up the main bulk of the report. Some may cavil at the way in which, at the slightest provocation, we include plots of the results. Undoubtedly these add to the bulk, but we can only hope that they will allow readers more quickly to get a general idea of what has been happening. The following chapter presents substantially the same set of results in a simpler but probably more controversial manner. In Chapter 6, extracts have been taken from the main test results and presented in such a way as to illustrate different aspects of the investigation.

Subsidiary to the main test was an attempt to make a comparative evaluation of citation indexing and bibliographic coupling. While there should be no serious problems in making such an evaluation under operational conditions, the value of testing this form of index in an artificial environment appears dubious. However, with considerable reservations the results are given in Chapter 7.

Up to this stage the results have been presented without any attempt being made to draw conclusions. All such have been relegated to the final chapter of this volume, in which an attempt is also made to relate the results to other investigations in this field.

There is one general apology that should be made and that is for the introduction into this report of yet more jargon. Many terms first used in reporting Cranfield I now appear to have gained general acceptance, but it is unlikely that such phrases as 'maximum starting term coordination level method' or 'proportional coordination level method' will crop up very frequently in the literature - and we certainly hope they won't - but it has been necessary to find terms to describe certain procedures so that, in reports of other tests, one has a chance of knowing which of several

possible methods has been used. We can only plead that we have not - as some apparently delight to do - concocted new terms to describe measures or methods when existing terminology has already appeared in the literature. If we have offended in this way, it was unintentional and we hope that our attention will be drawn to any lapse. The only case of which we know is where a term has been changed from that used in Volume I. The term 'generality ratio' has been dropped in favour of 'generality number'. It is hoped that the argument in Chapter 3 will provide the reasons for this change.

## CHAPTER 2

### Test Environment

Communication is the means which enables society to adjust itself to alterations of technology and education and other social changes. The scientific method can offer no grand vision, no global strategy, no panacea. It will never be possible to demonstrate that anything is absolutely right or even completely scientifically true.

L.T. Wilkins: Social Deviance, page 28.

In the first volume were considered the general plan of the test design, the variables that were to be investigated and the methods to be used. In the course of the project, changes were made regarding certain details, and this chapter presents the environment in which the testing was actually done.

While an information retrieval system may be defined in its scope as 'all stages from the receipt of a document within a system, to the making of that document (or a representation of it) available to an enquirer', not all these stages have been included in the investigations in the present project. The central concern was the effect of index language devices on the operational performance, but in addition a number of other variables or factors have been included for various reasons. In order to clarify later discussions, a breakdown of an indexing system into four main groups is suggested, namely environmental factors, software factors, operational factors and hardware factors (see Fig. 2.1).

The environmental factors relate to the environment or conditions in which a given system has to operate. Four general factors are given, and, in the case of an operational system, they are all determined to a great extent by the needs of the user group which the system exists to serve. The subject field, the questions asked and the relevance needs directly depend on the users, while the collection size will be determined by the management largely in relation to user needs. However, for an

ENVIRONMENTAL FACTORS

SUBJECT FIELD; precision of terminology, overlap of terminology with other fields.

COLLECTION SIZE; Number of recognisable subject fields.

QUESTIONS ASKED; Broad survey, specific request, etc.

RELEVANCE NEEDS; graded decisions.

SOFTWARE FACTORS

CONCEPT INDEXING; level of exhaustivity.

INDEX LANGUAGE; hospitality for specificity and provision of devices.

SEARCH STRATEGY; flexibility to vary exhaustivity and specificity.

OPERATIONAL FACTORS

SUBJECT COVERAGE

TIME; indexing and searching.

EFFORT; intellectual and physical.

PERSONNEL; indexers and searchers, qualifications and performance.

CLERICAL ROUTINES

RETRIEVAL PERFORMANCE

HARDWARE FACTORS

TYPE OF STORE

INPUT

EXPANSION CAPACITY

UPDATING ABILITY

PHYSICAL FORM OF OUTPUT

FIGURE 2.1 MAIN FACTORS IN AN INFORMATION RETRIEVAL SYSTEM

experimental test a set of environmental conditions has to be created, and some of those are inevitably, to a greater or lesser extent artificial.

The software factors relate to the intellectual design of the storage and retrieval parts of an indexing system. The three main software factors are all the subject of management decisions in a given situation, and such decisions are always centred (although often unconsciously) on the twin parameters of exhaustivity and specificity (defined and discussed in Ref. 2).

The operational factors are concerned with the routine operation of a system, i.e. all the processes required to make documents available to enquirers when the system has been set up. The factors in Fig. 2.1 are not intended to be an exhaustive list, but are given to illustrate the range of operations involved. Any basic evaluation of such factors is complicated by an infinite number of possible compromises between the least effort and the best quality, with both effort and quality being subjective notions notoriously difficult to measure.

The hardware factors refer to the purely physical aspects of system operation that involve man-made entities. A brief and incomplete listing of five items is given.

If one considers Cranfield II within this framework, it can be seen basically to have investigated the software factors, in the context of a laboratory situation in which the environmental factors and operational factors have been strictly controlled. Hardware factors have been ignored because, in this investigation, the measurements are being made on those software factors which are quite unaffected by changes in hardware. The operational factor of retrieval performance is the main measurement made, and details of how this is done are given in Chapter 3. In the artificial environment created for the test it was found that a limited set of changes could be investigated; these included several sets of questions picked by different criteria, relevance judgements made in four different grades, collections of three different sizes and tests in two related but different subject fields.

#### Software factors

The software factors examined in the tests will be described and discussed first. A simplified table (Fig. 2.2) shows the variables that have been examined, listed under the three major factors of indexing, index languages and searching.

### CONCEPT INDEXING

1. Manual indexing, at three levels of exhaustivity
2. Natural language abstracts and titles

### INDEX LANGUAGES

1. Single terms
2. Simple concepts
3. Controlled terms
4. Abstracts and titles
5. Recall devices
  - a. Single term indexing, eight languages
  - b. Simple concept indexing, fifteen languages
  - c. Controlled term indexing, six languages
  - d. Abstracts and titles, four languages.
6. Precision devices
  - a. Single term indexing, four types
  - b. Simple concept indexing, one type
  - c. Controlled term indexing, two types
  - d. Abstracts and titles, one type

### SEARCH RULES

1. Coordination levels, all possible levels
2. Combination rules, six types.

### FIGURE 2.2 SOFTWARE FACTORS EXAMINED IN TEST

#### Concept-Indexing

The manual indexing carried out on the document collection is described in Chapter 4 of Volume 1, and this constituted the main body of data tested; of particular importance was the fact that three levels of exhaustivity of indexing were distinguished. The results of this variation in exhaustivity have been evaluated on the single term languages, but not on the simple concept or controlled term languages. In addition, Professor Salton prepared (with the SMART programme) a KWIC type index of the titles and abstracts of 200 documents (subset 1); in this connection abstracts and titles can be considered as variant forms of concept indexing, and the test searches which were made enabled direct comparison to be made with the manual indexing carried out by the project staff.

Data concerning the usage of terms in the single term language is given in Fig. 5.1 of Volume 1; some additional information on term usage is given in Fig. 2.3 in relation to the simple concept and controlled term languages, the average postings per document being 18 and 24 respectively. Fig. 2.4 gives similar data for the abstracts, with the average postings of key terms being 74. This latter figure is not strictly comparable, since the same word may be 'posted' several times for the same document.

SIMPLE CONCEPTS

Collection size	200 documents
Total terms in vocabulary	2,798
Average posting per document	18

CONTROLLED TERMS

Collection size	200 documents	350 documents
Total terms in vocabulary	816	985
Terms in E.J.C. Thesaurus	694	827
Additional terms	122	158
Added lead-in vocabulary terms	1,285	1,514
Average postings per document	24	24

FIGURE 2.3 DATA CONCERNING USAGE OF TERMS IN SIMPLE  
CONCEPT AND CONTROLLED TERM INDEX LANGUAGES

COLLECTION SIZE	200 abstracts
Total postings of all words	33,042
Total postings of words less those on restriction list	14,783
Distinct words on restriction list	204
Distinct words not on restriction list	3,123
Average postings of all words per document	165
Average postings of words not on restriction list per docu- ment	74

First ten terms ranked by usage

FLOW  
NUMBER  
MACH  
PRESSURE  
RESULTS  
WING  
EFFECTS  
SHOCK  
BOUNDARY  
LAYER

FIGURE 2.4 DATA CONCERNING USAGE OF WORDS IN ABSTRACTS

Index languages

As described in Vol. I, Chapter 5, the languages tested fall into three main groups:

- I Single Terms, with the base being the natural language concept indexing split into unit terms,
- II Simple Concepts, with the base also being the natural language concept indexing, with some of the more complex pre-coordinated concepts split into simple concepts,
- III Controlled Terms, with the base being the controlled vocabulary derived from the E.J.C. Thesaurus, and indexing performed by translating the natural language concepts into the controlled vocabulary.

In defining any particular index language, these three main types will be denoted by the Roman numerals I, II and III; the various sets of recall devices tested are denoted by Arabic numerals and the precision devices by lower case letters.

Recall devices

The starting point of each series of tests is the use of the basic terms as indexed. From this base, various recall and precision devices are added, both separately and in different aggregates. In the single term languages, four different recall devices were tested, namely control of synonyms, confounding of word forms, control of quasi-synonyms and control of clusters of terms by means of reduced vocabularies based on hierarchies. A total of eight aggregates was tested, and a tree diagram giving details of the eight languages is given in Fig. 2.5.

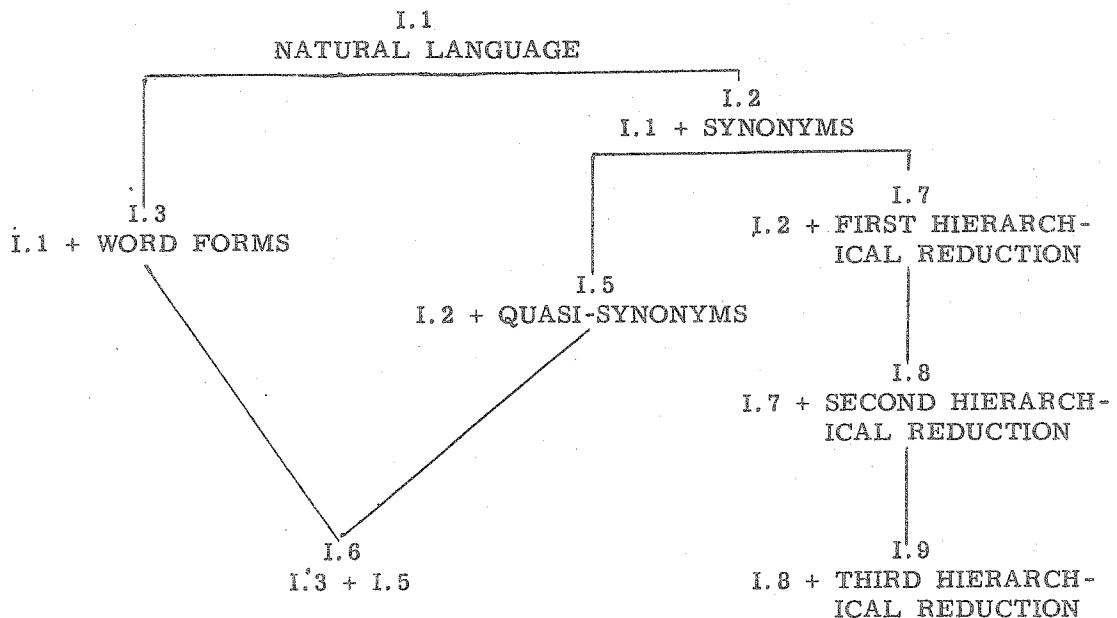


FIGURE 2.5 SINGLE TERM INDEX LANGUAGES

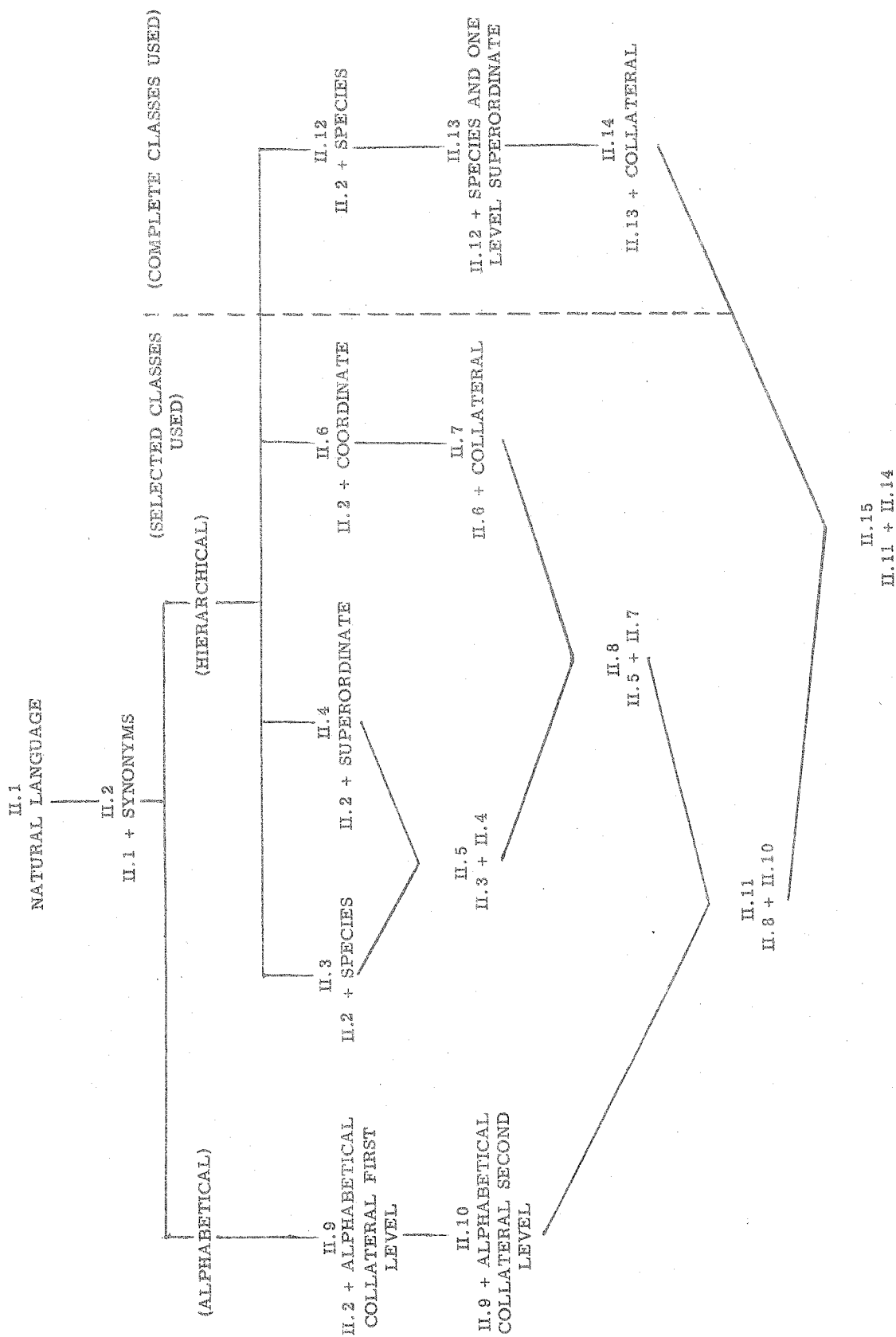


FIGURE 2.6 SIMPLE CONCEPT INDEX LANGUAGES

From this it can be seen that quasi-synonyms were tested together with synonyms, and that synonym control was also the base from which the three levels of reduction by hierarchy were tested.

The recall devices tested with the series of simple concept languages were the most comprehensive investigated. They involved one alphabetical and seven hierarchical devices, in fifteen different aggregates as shown in Fig. 2.6 (discussion on the hierarchies used and the rotated alphabetical list of concepts was given in Vol. 1, pages 74-83). It should be noted that recall devices 12, 13, and 14 of Fig. 2.6 involved the use of the complete classes of terms in the various hierarchical reductions, but, with the other languages, selections, based on intellectual decisions, were made from the various classes.

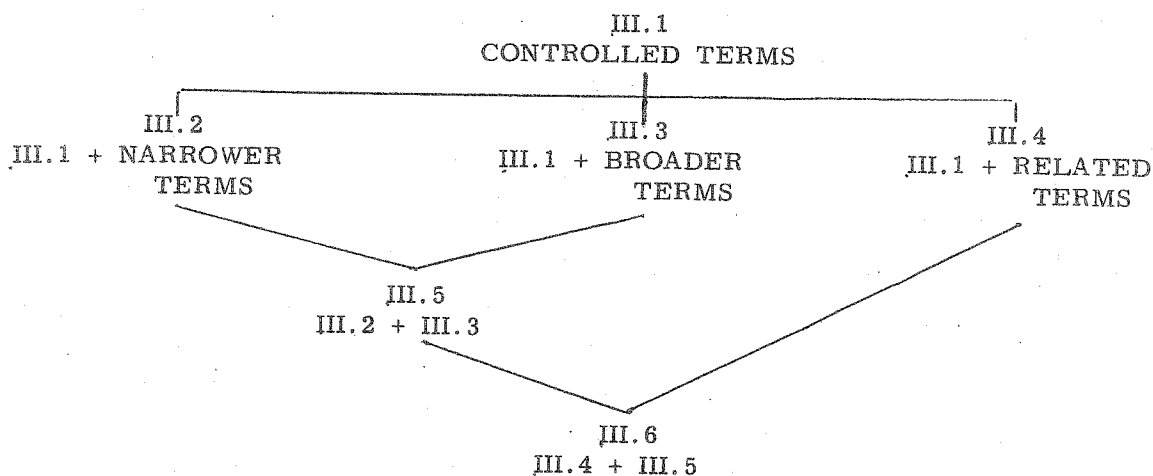


FIGURE 2.7 CONTROLLED TERM INDEX LANGUAGES

With the controlled terms, six index languages were tested. These consisted first of the basic terms, followed by the three classes of related terms as used in the E.J.C. Thesaurus (i.e. broader terms, narrower terms and related terms). In addition, two aggregates were tested; the six languages are listed in Fig. 2.7.

#### Precision devices

All the languages mentioned were tested for recall without any precision devices; this involved searches which accepted any one single term in the question. The fundamental precision device of coordination was also investigated in every test made, and all the basic tables of results in Chapter 4 show the coordination level in the rows of the tables. Two

additional precision devices were tested on the single term languages, namely partitioning and interfixing, as shown in Fig. 2.8.

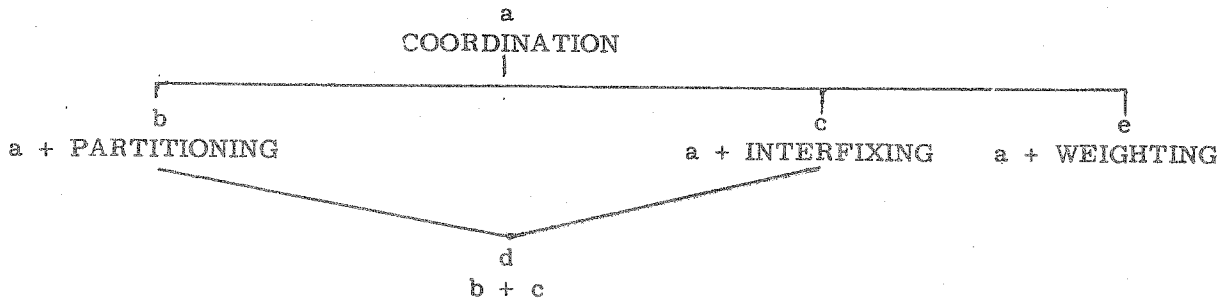


FIGURE 2.8. PRECISION DEVICES

No precision devices other than coordination were tested on the simple concept languages. The device of weighting was tested on the controlled terms. In this weights are assigned to the search term and a match sought with the weights assigned to the terms in indexing.

All the index languages tested may now be specified; for example II.2.a represents Simple Concept Index Language (II), with the recall device of Synonyms controlled (2), and coordination (a) as the precision device. The code for Single Term Index Language, with the recall device of Quasi-synonyms and the precision devices of partitioning and interfixing would be I.5.d.

#### Search Rules

In the search programmes for the questions tested an exhaustive extraction of all the possible notions contained in each question was made in the natural language of the questions as they were received. All these notions were included in the search prescription initially prepared for the three main index languages. After the basic question terms had been recorded, all the additional terms included in a logical sum relationship were pre-formulated by the very structure of the various languages already described. For example in Question 61 'Are there any papers dealing with acoustic wave propagation in reacting gases'. The terms underlined made up the search prescription, and these terms, as they are, were used for Index Language I.1. For Index Language I.2, Synonyms controlled, reference to Appendix 5.2 of Vol. I shows that the term Sound is now combined with Acoustic. For Index Language I.3 Word endings, the term Acoustically is combined with Acoustic; Waviness and Wavy are combined with Wave and there

are similar groupings for the other terms. For Index Language I.5, Quasi-synonyms, the term Sonic is combined with Acoustic, and Reaction now forms a group which includes the quasi-synonyms Energy, Force, Action, Behaviour, Kinetic, Response. With Index Language I.7, I.8 and I.9, the groups for each starting term are determined by the decisions taken in the compilation of the single term hierarchies as given in Appendix 5.3 of Volume I. There is nothing to add regarding the search prescription, for it was the search rules that were capable of variation; this could be achieved by varying the coordination level or by selecting acceptable combinations of the search terms.

As has been noted, all possible levels of coordination (logical product) were investigated at every stage, and therefore the effect of any rules that might be postulated concerning a minimum coordination level that would be acceptable can be seen from the tables of results. For instance, if a question had six terms, then the results would have been recorded for a search made with all six terms, then for a search with five terms, then with four terms and so on down to a single term search. No test was made in which the searches of a set of questions either commenced or were terminated by a subjective decision that varied from question to question.

The main variations introduced as search rules concerned the combinations of terms that were accepted. The six variations tested are given in Fig. 2.9.

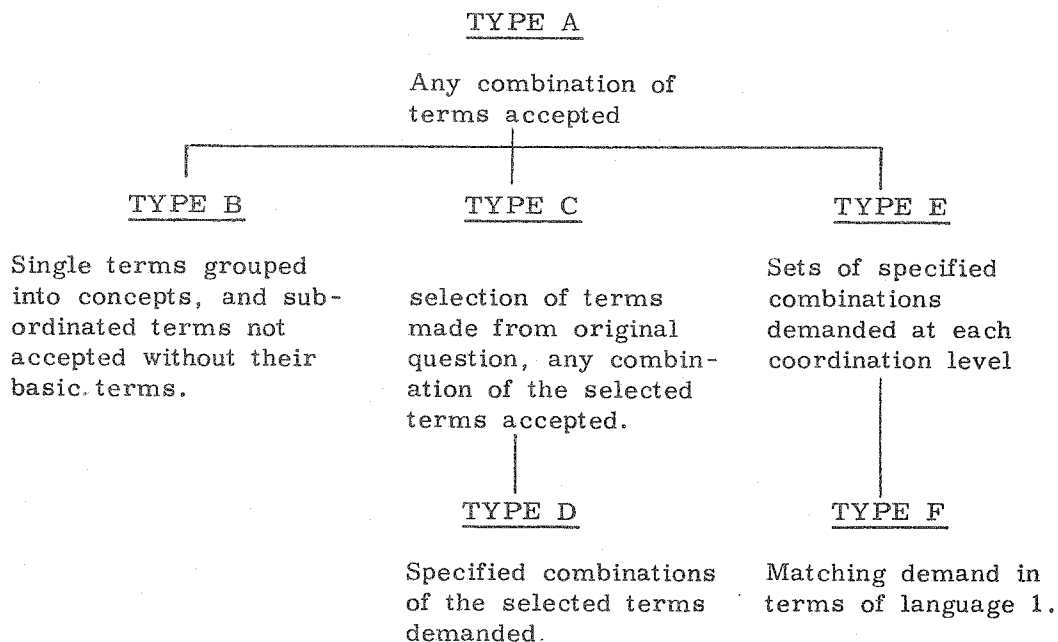


FIGURE 2.9

SUMMARY OF SEARCH RULES

In search Type A, any combination of terms was always accepted, without regard for the cases where some combinations accepted might be meaningless. For example, consider a question with the search terms Methods, Testing, Analysis, Investigating, Static, Dynamic, Stability, Characteristics, Re-entry, Body, Free, Flight, Tests. At, say, a coordination level of four, any combination of these search terms would be accepted, such as Methods, Static, Re-entry, Free. This is only one of many non-sensical combinations of the search terms at this level of coordination. The use of this search rule for investigating nearly every other variable was adopted, since it could be applied with equal consistency to all the different languages, with the exception of the tests of the precision devices of partitioning and inter-fixing on the single term languages. For these tests it was felt that a certain amount of intellect should be put into the search rules, and this consisted of a rule (Type B) which did not permit 'subordinate' terms to be accepted unless the associated 'basic' terms was present. The distinction between basic and subordinate terms became apparent when the single terms of the search questions were grouped into concepts, prior to the test of inter-fixing. In the example mentioned there are certain concepts that would emerge, such as Static stability characteristics, Re-entry body, Methods of testing. Basic terms in these concepts might be Stability, Body and Testing, for these terms are meaningful on their own in the context of the search question. Therefore Search Rule B would require, for instance, that Re-entry would not be accepted unless Body was also present, nor would Static be accepted unless Stability was present. The importance of adopting this rule before making a test of interfixing is that at least two terms from a concept must be present for interfixing to be tested. If, in the indexing of a document, the two single and separate terms Static and Stability appeared, and the demand for interfixing was added, if they were not interfixed then only one of the single terms could be accepted (which would have to be Stability to accord with Search Rule B). Without Search Rule B the single term Static could be accepted in this case, or in a case where Stability did not occur at all.

Searches C and D were carried out on the single term index languages, and represented an attempt to discover the effect of including more intelligence in searching. The first stage, Search C, involved making a selection of the original starting terms taken from the search question. This was to eliminate from the search prescriptions certain terms such as Problem, Applied, Variation, Influence, Solution, Comparison, Determination, Effect, etc. This search rule was tested on a set of twenty questions, all of which originally had seven starting terms; the selections made resulted in a range of from two to six of the terms, with the average being 4.1. In using these selected sets of search terms, any combination of these was still accepted, as in Search A.

Search D used the selected search terms of search C, and made strict

restrictions concerning the actual combination of terms that would be accepted at every coordination level, so as to eliminate the non-sensical combinations.

The most satisfactory and carefully applied search rules were applied to the controlled language tests, since it was thought that intelligence in searching would be best tested on an index language that also had an average degree of intelligence used in its formulation. This was Search E, where all the combinations of acceptable terms were individually selected for each coordination level. It was usual to accept a number of such combinations, with the object of retaining as many of the relevant documents as possible. This search rule was applied to the controlled term index languages (III.1 - III.6) both with and without the precision device of weighting. The sets of acceptable combinations were formulated on the basis of the starting terms of the question, and thus the use of Search E in testing languages other than III.1 (Basic terms) may have resulted in a poorer performance for the languages than is theoretically possible; the reason for this is that the grouping of a number of terms in the later languages might result in non-sensical combinations of terms.

One further additional rule designed to be used with the various recall languages was tried. This was Search Type F, also carried out on the controlled term index languages III.2 to III.6. The reasoning behind this search was that in all previous rules tested, the terms that actually made a match between a document and search prescription were all treated 'equally'. For example, if two documents had a match of five terms with a question using the controlled term index language III.5a (related terms), no distinction would be made between a document which actually had four starting terms, *and only one* related term, and a second document which was matched only by related terms, without a single starting term. The first document clearly represents a closer match with the search prescription, and it might generally be assumed that a starting term match is more desirable than any related term match. In Search F, a record was made of the number of starting terms that came up in a given match, and was done with the rules of Search E in use. This was used to make up sets of results with a given minimum match demanded, and results will be given for controlled term languages III.5 and III.6.

#### Document relevance

Before demonstrating the form of the results obtained when these variables are tested, a single environmental variable will be mentioned. This is the variation made in document relevance, resulting from the scale of four grades of relevance that was followed by the questioners in assessing the relevant documents (see Vol I, p.21). In finding the effect on retrieval performance of these decisions, four sets of results

were obtained, comparing first a set of questions when only relevance 1 documents were accepted as relevant, then with documents of relevance 1 or 2, next with documents of relevance 1 or 2 or 3, and finally with documents of relevance 1 or 2 or 3 or 4. Apart from the particular test to measure this variable, the broadest relevance decision, namely 1 - 4, was always used in other tests.

### The Composite Table

Some idea of the volume, variety and complexity of the tests carried out can be seen from the composite table, (Fig. 2.10) which gives results for various combinations of six variables tested on the single term index languages I.1 to I.6. The basic set of questions used is subset 1, which has 35 questions, each having seven starting terms, but some of the results are based on two selections of these, namely 19 questions of subset 4 and 20 questions of subset 6. Four of the variables are listed at the head of the table, and the other two at the left side; the table divisions consist of the following factors:-

1. The coordination level varies from 1 to 7, which would result in seven main sections of the table. However, due to problems of presentation in this report, the table is truncated by the omission of the figures relating to the first three levels, so that it only presents four main sections covering the coordination levels of 4, 5, 6 and 7.
2. Four search rules (A,B, C and D) are next varied, and are applied in order of increasing intelligence within each coordination level.
3. The precision devices (a, b, c and d) are recorded next, with most results using no linking devices, apart from the three columns near the centre of each section.
4. The final factor at the head of the table is document relevance, with the three higher grades listed first, followed by the lowest grade used for all subsequent combinations (1, 1-2, 1-3, and 1-4).
5. The rows are first divided into five, representing the index languages I.1, I.2, I.3, I.5 and I.6.
6. The final variable is indexing exhaustivity, the three levels being repeated as divisions of each index language in turn.

The meaning of the codes used in this table has already been described earlier in this chapter.

The search results are shown as percentages for recall and precision.

Thus each set of recall and precision devices can be understood by examining the columns above, and the row to the left of a set of ratios, and then reading off the particular combination of variables being tested. For example, if the first section of the table as printed is examined

Co-ordination		4+																		
Search Rules		A				B				C	D									
Precision Device		a				b	c	d	a											
Document Relevance		1	1-2		1-3		1-4													
Recall Device	Exhaustivity	R	P	R	P	R	P	R	P	R	P	R	P	R	P					
I.1	1	28	2	25	9	17	17	19	24											
	2	44	1	35	5	28	10	30	15											
	3	44	1	38	4	33	10	33	14	28	23	20	26	19	32	12	31	28	29	21
I.2	1	28	2	25	8	17	16	19	23											
	2	50	1	37	5	30	10	31	15											
	3	50	1	39	4	34	10	35	13	29	21	21	23	19	32	12	31	29	30	
I.3	1	39	3	30	9	20	16	21	23											
	2	56	1	41	5	33	10	33	13											
	3	56	1	43	4	37	8	36	11	33	15	24	24	22	24	15	29	33	24	
I.5	1	39	1	27	4	23	9	25	14											
	2	56	1	41	2	36	5	38	7											
	3	56	1	44	2	42	4	44	6	40	8	27	12	26	11	18	15	36	14	
I.6	1	44	1	32	4	24	8	26	12											
	2	56	1	42	2	37	4	40	6											
	3	56	1	47	1	44	4	45	5	44	7	30	11	29	11	21	15	40	11	27

FIGURE 2.10a THE COMPOSITE TABLE. COORDINATION LEVEL 4+  
R = RECALL RATIO, P = PRECISION RATIO  
(Performance figures are expressed as percentages)

Co-ordination		5+											
Search Rules		A				B				C	D		
Precision Device		a				b	c	d	a				
Document Relevance		1	1-2	1-3	1-4								
Recall Device	Exhaustivity	R	P	R	P	R	P	R	P	R	P	R	P
I.1	1	11	5	9	17	7	37	8	54				
	2	28	4	19	11	14	22	15	31				
	3	28	3	20	9	16	18	16	26	12	64	7	64
I.2	1	17	7	9	16	7	35	8	51				
	2	33	4	22	12	15	22	16	32				
	3	33	3	23	9	17	18	18	25	13	65	7	64
I.3	1	17	6	10	17	8	36	8	51				
	2	39	4	23	11	17	22	17	29				
	3	39	3	25	8	19	17	19	23	16	51	9	55
I.5	1	11	2	11	9	8	16	10	27				
	2	39	2	25	5	17	10	19	15				
	3	44	2	30	5	21	8	23	12	21	23	11	48
I.6	1	22	4	15	11	10	21	11	30				
	2	44	2	27	5	19	9	21	13				
	3	50	1	30	4	24	8	25	11	22	21	13	36
										11	36	7	45
												28	27
													23
													63

FIGURE 2.10b. THE COMPOSITE TABLE. COORDINATION LEVEL

5+

Coordination		6+																		
Search Rules		A				B				C		D								
Precision Device		a				b		c		d		a								
Document Relevance		1	1-2		1-3		1-4													
Recall Device	Exhaustivity	R	P	R	P	R	P	R	P	R	P	R	P							
I.1	1	11	17	5	33	4	75	4	83											
	2	11	5	11	21	8	40	8	44											
	3	11	3	11	15	8	28	8	38	8	100	5	100	3	100	3	100	5	33	5
I.2	1	11	17	5	33	4	75	4	83											
	2	11	4	11	20	8	37	8	50											
	3	11	3	11	15	8	27	8	37	8	100	5	100	3	100	3	100	5	33	
I.3	1	11	17	5	33	4	75	4	83											
	2	17	6	13	21	9	40	9	52											
	3	17	4	13	14	9	28	9	36	8	100	5	100	3	100	3	100	5	33	
I.5	1	11	13	5	25	5	63	4	75											
	2	28	6	19	18	12	30	12	43											
	3	28	4	19	12	12	28	12	28	12	44	5	58	3	44	3	44	10	50	
I.6	1	11	12	5	24	5	59	4	71											
	2	33	8	19	15	12	26	13	37											
	3	33	4	22	12	13	19	13	26	13	37	6	62	3	44	3	44	10	50	5

FIGURE 2.10c. THE COMPOSITE TABLE. COORDINATION LEVEL 6+

Co-ordination		7+							
Search Rules		A				B			
Precision Device		a				b	c	d	
Document Relevance		1	1-2	1-3	1 - 4				
Recall Device	Exhaustivity	R P	R P	R P	R P	R P	R P	R P	R P
I.1	1	11 29	4 43	2 71	2 71				
	2	11 13	4 19	3 38	3 50				
	3	11 13	4 19	3 38	3 50	3 100	3 100	2 100	2 100
I.2	1	11 29	4 43	2 71	2 71				
	2	11 13	4 19	3 38	3 50				
	3	11 13	4 19	3 38	3 50	3 100	3 100	2 100	2 100
I.3	1	11 29	4 43	2 71	2 71				
	2	11 13	4 19	3 38	3 50				
	3	11 13	4 19	3 38	3 50	3 100	3 100	2 100	2 100
I.5	1	11 25	4 38	3 75	2 75				
	2	11 9	4 14	3 32	4 45				
	3	11 7	5 14	4 28	4 38	5 100	3 100	2 100	2 100
I.6	1	11 25	4 38	3 75	2 75				
	2	17 12	6 20	4 36	4 48				
	3	17 9	8 19	5 31	5 41	6 53	4 100	2 100	2 100

FIGURE 2.10d. THE COMPOSITE TABLE. COORDINATION LEVEL 7+

(coordination level 4), and the ratios at the top left corner examined (28% recall, 2% precision), the following variables are shown to have produced that result: a search at coordination level of four terms; search rule A (any combination); precision device 'a' (no linking in the index language); relevant documents graded 1 only accepted; recall language 1 (natural language terms); and indexing exhaustivity 1 (low exhaustivity). After this, a move across this section of the table to the right will first alter the document relevance grades, then introduce a search rule, then include the three precision devices and finally test three more search rules. A move into the next section will increase the coordination level of the search, and in any section a move down the table will increase the indexing exhaustivity before a new recall language is brought in.

The position of these variables in the table is of no significance; the table could, for instance, first have been divided into the five recall languages, with the seven coordination levels repeated at each stage, etc. and hundreds of variations are possible. The actual combinations of different variables for which results have been presented in the complete composite table total 609, which is a choice of the most useful combinations out of the theoretical total of 6720 combinations possible.

Each set of recall and precision ratios is an average of results from the set of 35 questions and it is estimated that the composite table represents more than 16,000 individual results. When it is considered that the scope of the whole project extends to 221 questions, that there are some 28 other index languages which are not included in this table and that there are a number of other new variables, the individual results available are estimated to exceed 200,000.

### Environmental Factors

The main environmental factors involved in the testing are listed in Fig. 2.11. For various reasons, as the test proceeded, different sets of questions and collections of different sizes were used. To consider first the sets of questions. Although 279 questions were available for use, the largest set for which results are presented numbers 221. The balance of 58 were multi-themed questions, that is they really consisted of more than one question, e.g. Question 3 'How can one describe the aerodynamic forces and the heating rates acting on high speed aircraft'. Four of these were used in some of the smaller question sets only. The first series of tests, on the recall devices of the single-term index languages, were made of the complete collection of 221 single-theme questions. The major problem that then arose was to find a satisfactory method of totalling the results of searches based on different numbers of starting terms (this matter is considered at length in Chapter 3). For this reason, we investigated the results on a set of 35 questions each of which had seven starting terms. The tests on interfixing and partitioning were particularly difficult to do, because of the painstaking clerical work necessary. These were therefore done on

two subsets which had 19 questions with 7 starting terms and 17 questions with 11 starting terms.

### QUESTIONS

1. Relevance assessments, 4 grades
2. Differing number of starting terms and retrieving terms
3. Differing totals of relevant documents
4. Two sources of questions, 'basic' and 'supplementary'
5. Question sets of different sizes, picked according to different criteria, searched on collections of varying sizes.

### COLLECTION SIZE

1. 1400 documents
2. 350 documents from the 1400 documents.
3. 200 documents from the 350 document subset.

### SUBJECT TERMINOLOGY

1. Aerodynamics
2. Aircraft Structures.

## FIGURE 2.11 SUMMARY OF MAIN ENVIRONMENTAL FACTORS

By the time we came to investigate the simple concept languages and the controlled term languages, the clerical effort involved in carrying out searches precluded the use of the full sets of questions, and accordingly a set of 42 questions was prepared, consisting entirely of questions in the field of aerodynamics. It is this set which is used for presenting the majority of the test results in Chapter 4. At a later stage, this subset was extended to 77 questions in the field of aerodynamics; finally an additional set of 42 questions in the field of structures was compiled for purposes of comparison, with the aerodynamic question set of similar sizes. The subsets of questions are all numbered, and details of these appear in Fig. 2.12. Lists of the question numbers for subsets 1, 2 and 3 were given in Vol. I, Appendix 3E; the remaining subsets are shown in Appendix 3.2 of this volume.

Reduced collection sizes were also used for reasons of the effort involved in testing. This was not only the clerical effort involved in the searching, but also the intellectual effort involved in compiling word lists for the various index languages. When it was decided to test simple concepts, a set of 200 documents was chosen, and the initial task involved re-formulating the indexed concepts from the original

Question Subsets	No. of Questions	No. of Documents in Collection Tested	No. of Relevant Documents	Generality Number
<u>Subset 1</u> All have seven starting terms in single term languages, covering aerodynamics and structures.	35	1400	287	5.9
<u>Subset 2</u> Starting terms vary, all questions aerodynamics only. 38 are drawn from Subset 3, and 4 from the 58 questions not used. The number of relevant in the 1400 collection is actually 201, but the three documents concerned (1329 in Q119, 2289 in Q145 and Q146) are deleted from the collection in results of this subset since they did not appear in collection Subset 1.	42	1400	198	3.4
	42	200 (subset 1)	198	23.6
	42	350 (subset 2)	198	13.5
<u>Subset 3</u> Starting terms vary, covering aerodynamics and structures. The largest set of questions available, all single theme in single term languages.	221	1400	1590	5.1
<u>Subset 4</u> Part of Subset 1, all having seven starting terms.	19	1400	131	4.9
<u>Subset 5</u> All having eleven starting terms in single term languages, covering aerodynamics and structures.	17	1400	109	4.6
<u>Subset 6</u> Part of Subset 1, all having seven starting terms.	20	1400	147	5.3
<u>Subset 7</u> Includes all the questions in Subset 2, aerodynamics only.	77	350 (subset 2)	454	16.8
<u>Subset 8</u> Structures only.	42	1400	255	4.3

FIGURE 2.12

indexing records. The choice of the subset of 200 documents extracted from the 1400 was governed by:-

1. Use of a set of aerodynamic questions and documents.
2. Choice of the largest set of questions that could be tested on a subset of 200 documents.
3. Choice of questions restricted to those not having any relevant documents in the range of Numbers 1001 - 1299 (because of the different weighting method used at that stage of the indexing).

The third rule restricted the choice of questions quite considerably, and the second rule was modified by not allowing 'similar' questions - mainly two or more questions having an overlapping set of identical relevant documents asked by the same questioner. The 42 questions finally used had 198 relevant documents in the subset of 200 documents. None of the base documents to the 42 questions were included.

The second subset chosen was one of 350 documents (subset 2), and this subset included the 200 documents of subset 1. Subset 2 also consisted entirely of aerodynamic documents, with an additional 35 questions having all their relevant documents in the subset. These, together with the 42 questions for subset 1, resulted in the 77-question subset (subset 7) which was used for the tests on the controlled vocabulary.

In presenting the test results, the majority of results are based on these smaller subsets of documents and questions. The first tests were made with 221 questions on the 1400 collection, and these tests were repeated on smaller documents and question subsets in order to validate the use of such subsets. It will be shown in the next chapter how the difference in performance can be adequately accounted for, and the use of smaller subsets does not, we believe, impare either the validity or, to any appreciable extent, the accuracy of the results and findings.

The use of these subsets enabled the various environmental factors involved to be investigated. For example, the effect of the change in collection size from 1400 to 200 documents with a fixed set of questions was investigated. Comparison was also possible between the 350 and 200 collections.

In the case of the questions, different subsets were made up and results obtained when environmental factors such as those listed in Fig. 2.11 were being investigated. The four grades of document relevance are included in the main test results (Chapter 4) and the effect of the other factors that are listed is considered in Chapter 6.

An attempt was made to compare the two distinct subject fields that existed in the 1400 document collection. Many of the question sets contained both aerodynamics and structures questions, but the direct

see fig. 2.12 page 287

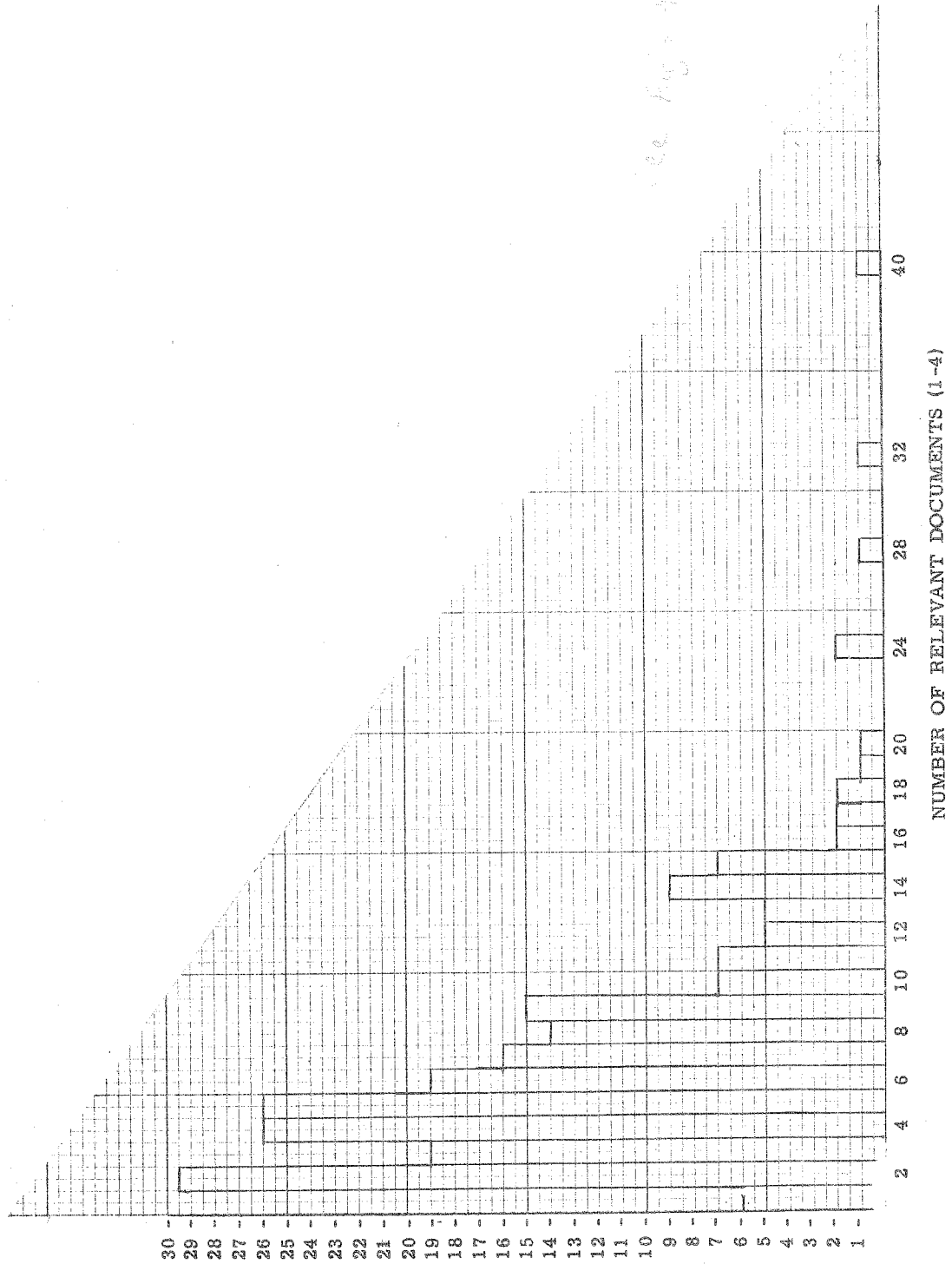


FIGURE 2.13 DISTRIBUTION OF THE RELEVANT DOCUMENTS IN THE 221 QUESTIONS

comparison was made of 42 questions on aerodynamics and 42 questions on structures.

Fig. 2.13 is a chart showing the distribution of the number of relevant documents throughout the 221 questions, from which it can be seen that the range is from six questions each having only one relevant document to one question which had forty relevant documents.

#### Sample Precision Results

At low precision ratios, the clerical work involved in obtaining correct figures was so great that in some cases it did not appear to be justified. This was due to the large number of non-relevant documents which would be retrieved and therefore had to be recorded. With index language I.1a, results were obtained down to the single term level but with other index languages the decision was taken that, with the searches in the 1400-document collection, no attempt would be made to obtain precision figures below 5%. This, however, introduced a variation between questions, since for a question having six starting terms, a precision figure lower than 5% might not be reached until the coordination level was down to two terms. However, with a ten starting term question, this figure might be reached by the coordination of four terms.

In the presentation of the test results, note has been taken of this point and also the additional point regarding the number of questions capable of giving results, this being dependent on the number of starting terms which each question had. This can be best illustrated by referring to Fig. 2.14, which presents condensed results for 221 questions on the 1400-document collection with Index Language I.2a. The column headed 'z' presents the figures for the number of questions that were potentially capable of giving results, and it can be seen that at a coordination level of 2, every question came in this category. However, at a coordination level of 3, the total has dropped to 220, this indicating that there is one question which had only two starting terms. At a coordination level of 4, the total drops to 212, showing that there are eight questions with only three starting terms. As the coordination level rises, so the number of questions drops until, at a level of 15, it is seen that only three questions have this number of starting terms.

The column headed 'y' shows the number of questions which actually contributed figures for the calculation of the precision and fallout ratios - not, it should be noted, for the recall ratio which was always checked down to single term level. In Fig. 2.14, y is equal to z from a coordination level of 15 down to a coordination level of 7, and therefore the precision and fallout ratios can be based on complete data. However, at a coordination level of 6, only 161 questions were searched, and the precision and fallout ratios have been calculated on the basis of the non-relevant documents retrieved in these 161 searches. To indicate this, an asterisk

Index Language I. 2. a (S. T. Synonyms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 221 (Subset 3)

Number of Relevant Documents 1,590

Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,514	(-)	95.2%	(-)	(-)	221	0	221
2	1,313	59,734*	82.6%	2.2%*	19,406%*	221	44*	221
3	981	23,654*	61.7%	4.0%*	7,680%*	216	109*	220
4	644	8,850*	40.5%	6.8%*	2,873%*	192	142*	212
5	355	2,946*	22.3%	10.4%*	0,957%*	139	177*	197
6	169	928*	10.6%	15.4%*	0,301%*	92	161*	164
7	80	254	5.0%	24.0%	0.083%	55	140	140
8	24	59	1.5%	28.9%	0.019%	23	105	105
9	8	8	0.5%	50.0%	0.003%	8	78	78
10	1	0	0.1%	100.0%	0.000%	1	52	52
11	0	0				0	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 2.14 SAMPLE TABLE OF TEST RESULTS

is always given against any figures which have been calculated on a reduced set. At a single term level, it can be seen that no searches were made, and therefore no figures can be estimated for precision or fallout ratios.

There were various possible procedures for estimating these figures, and these can be illustrated by reference to Fig. 2.15, which deals with the 35 questions subset searched on 1400 documents by Index Language I.5.a. Since all the questions had seven starting terms,  $z$  remains constant throughout. However, at a coordination level of 2, it is shown in column  $y$  that only 23 questions were searched. It was found that, with these 23 questions 8,565 non-relevant documents were retrieved together with 157 relevant documents. The simplest way of estimating the total non-relevant for the complete subset of 35 questions would be to scale up the above figure of 8,565 in the ratio of  $\frac{35}{23}$ , which would give a total of 13,033 non-relevant documents. On the basis of this figure the precision and fallout ratios\* could now be calculated. A second method is first to determine the precision ratio for the 23 questions searched; in this case it works out at 1.8%. It is known that the 35 questions retrieved 253 relevant documents; to maintain the precision ratio of 1.8% the total of non-relevant is scaled up by  $\frac{253}{157}$ , namely the totals of relevant documents retrieved in the full set and in the subset. This gives a figure of 13,803 and from this the fallout ratio can be calculated.

The accuracy of these scaled up results will depend on whether the sample of questions that were searched is typical of the whole set. It is unlikely that this was the case; as stated earlier, questions were not searched when they would retrieve an excessive number of non-relevant documents, so conversely the questions which were searched, and which are therefore in the sample, were those which had fewer non-relevant documents. Scaling-up from the sample could therefore be expected to give a somewhat higher precision figure than was really the case.

To check on this, we can consider the actual situation in regard to the same set of questions with Index Language I.1.a, on which, as previously mentioned, searches were made down to the single-term level.

In this language, at a coordination level of 2, the 23 questions retrieved 3871 documents. By the methods already suggested, the estimated figures would have been 6043 and 6476 respectively. In fact, the correct figure is 8086, and bears out the expectation expressed in the previous paragraph. This was also checked at the coordination level of 3, and again it was found that the remaining 12 searches retrieved

---

\*The method of calculating these ratios is discussed in Chapter 3.

Index Language I 5. a (S. T. Synonyms, Quasi-synonyms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1.400  
 Number of Questions 35 (subset 1)  
 Number of Relevant Documents 287  
 Generality Number 5.9

Coord- ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	280	(-)	97.6%	(-)	(-)	35	0	35
2	253	17,130*	88.2%	1.5%*	34.959%*	35	23*	35
3	194	7,472	67.6%	2.5%	15.339%	35	35	35
4	125	2,086	43.6%	5.6%	4.282%	34	35	35
5	65	463	22.7%	12.3%	0.950%	30	35	35
6	35	88	12.2%	28.4%	0.181%	16	35	35
7	11	18	3.9%	38.0%	0.037%	5	35	35

FIGURE 2.15 SAMPLE TABLE OF TEST RESULTS

approximately the same number of relevant documents as the original 23 searches. For this group of results, therefore, the figures at the coordination level of 2 have been estimated by doubling the total obtained for the 23 questions. Similar procedures have been used in other cases.

This can certainly be considered somewhat unsatisfactory, and it could be argued that it would have been preferable not to have attempted to obtain figures by such a dubious method. However, it is felt that they do have some value; in every case where any such action is taken, an asterisk is placed against the figure or the ratio, and, if the reader feels so inclined, these results can be ignored.

As can be seen from the example of Figs. 2.14 and 2.15, each table of results contains details of the environment in which the test was carried out. This includes the particular index language, the level of exhaustivity, the search rule, the level of relevance, the number of documents and questions, the number of relevant documents, and the generality number. The latter, and the meaning of  $x$  in the tables, is considered in Chapter 3.

### CHAPTER 3

#### Methods for Presentation of Results

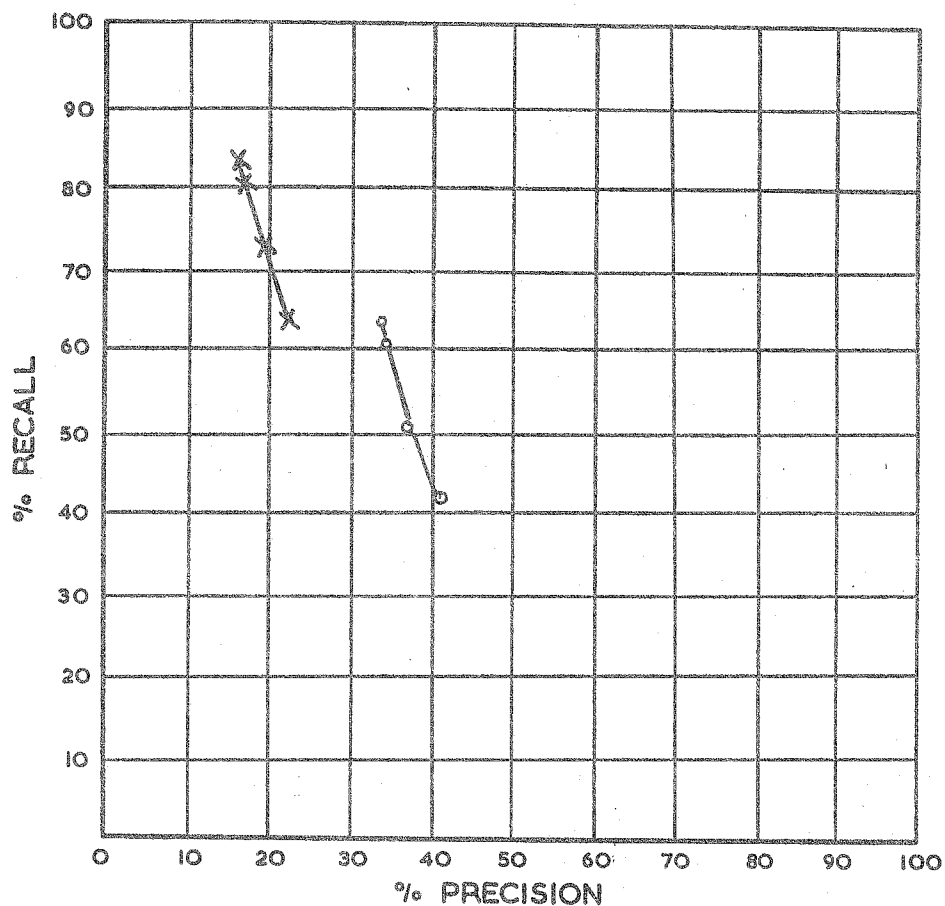
Lord Kelvin is often credited with remarking, 'When you can measure what you are speaking about and express it in numbers you know something about it, but when you cannot measure it, when you cannot express it in numbers, your knowledge is of a meagre and unsatisfactory kind'. The problem of validity - how closely do the figures relate to the 'thing we are talking about' must be separated from the problem of reliability - how accurate are the figures themselves.

L.T. Wilkins: Social Deviance, page 147

In Cranfield I, the results of the main project did no more than record what is now generally known as the Recall Ratio, which was calculated on the basis of  $\frac{100R}{C}$  where R equals the number of relevant documents retrieved and C equals the total number of documents in the collection which are relevant to the questions. In the subsequent test of the Western Reserve University Index, (Ref. 2) measurement was carried to the stage where, by making relevance assessments of all the retrieved documents, it was also possible to calculate what was originally called the Relevance Ratio, but which is now generally known as the Precision Ratio, namely  $\frac{100R}{L}$ , where L equals the total number of documents retrieved in the series of searches. In the course of this evaluation of the W.R.U. Index, the effect of varying the exhaustivity of indexing was measured, and allowed the production of the first - and, incidentally, so far the only - performance curve from the Cranfield project. It is reproduced in Fig. 3.1F and showed two interesting characteristics. The first was the inverse relationship between recall and precision, which has been considered at some length in Volume 1 of this report. The second point was that, when documents of lower relevance, were accepted, there was at any given level of indexing exhaustivity, a lower recall ratio but an improved precision ratio. It was tentatively suggested that this latter point was connected with a variation in the average number of relevant documents for each question, and that, for any given situation, it would be necessary to state also what was to be later termed the Generality Number\*, expressing it as  $\frac{1000C}{N}$ , where N equals the total number of documents in the collection.

---

\*In the earlier volume of this report, this was called Generality Ratio.



x Relevance 2

o Relevance 2 and 3

FIGURE 3.1P PERFORMANCE CURVE OBTAINED WITH FACET INDEX IN W.R.U. TEST

While the recall and precision ratios have been generally accepted as performance measures for information retrieval systems, they have also aroused some criticism. No serious attempt has been made to answer this criticism, partly because it was mostly trivial and never supported by experimental data, but mainly because an intention of Cranfield II was to investigate the performance measures which could or should be used. For this, sets of performance data were required and it was known that for every set of figures in Cranfield I, there would be hundreds of sets in Cranfield II, and it was obvious that the decisions regarding the measures to be used and the methods of presenting the test results would be of major importance. The programme of work which this aspect of the project has involved has been considerable, with many sets of results being calculated in a number of different ways. Based on this work, which has taken up a significant part of the effort during the last eighteen months of the project, the decision was finally reached that the most satisfactory method of calculating results involves three measures, namely Recall Ratio and Precision Ratio with, additionally, the new measure of Fallout Ratio. For the presentation of results on a plot, it is believed that, in the large

majority of cases, the most straightforward and most meaningful method is the Recall/Precision curve. These are, in general, the measures used in this report, although to illustrate certain points various other measures and methods of presentation are used.

A detailed account of the Cranfield work on performance measures has been presented in a thesis by M. Keen, but the following is a resumé of the more important points which led to the decisions; other matters relevant to the presentation of results in this volume are also considered.

In tests of experimental systems, it is essential that measures should be used that accurately reflect the changes in the particular component being tested, which primarily, in this particular test, was a range of index language devices. In addition, there is the strong desirability, if not the absolute necessity, that it should be possible to make direct comparison between different sets of test results.

Measures of retrieval performance may be used in experimental tests of information retrieval systems when the following requirements are met:-

1. A document collection of known size to be used in the test;
2. A set of questions, together with decisions as to exactly which documents are relevant to each question;
3. A set of results of searches made in the test; these usually give the numbers of documents retrieved in the searches, divided into the relevant and non-relevant documents.

The successive dichotomies of the total collection have been displayed by B.C. Vickery (Ref. 3, page 174) by the following table:-

TOTAL COLLECTION			
RELEVANT		NON-RELEVANT	
NOT RETRIEVED	RETRIEVED		NOT RETRIEVED
(c)	(a)	(b)	(d)

The more usual way to present the categories is in the form of a 2 x 2 contingency table as shown in Fig. 3:2. The notation given in this figure will be used throughout the remainder of this report.

	RELEVANT	NON-RELEVANT	
RETRIEVED	a	b	a + b
NOT RETRIEVED	c	d	c + d
	a + c	b + d	a + b + c + d = N (Total Collection)

FIGURE 3.2 2 x 2 CONTINGENCY TABLE

Whether it is correct to regard the values that result from retrieval tests as components of a 2 x 2 table in the statistical sense, and thus apply the principles and tests that have been developed for this situation in statistics, is an unanswered question, and at this stage, therefore, the use of this table is purely for convenience.

As mentioned earlier, there is the necessity of being able to make a comparison between several sets of results obtained in different conditions. This can only be done when it is known exactly which variables are altered in the different situations; two such situations are considered.

Assuming N (the total collection) remains constant, a, b, c and d can each vary, while a + b (total retrieved) and c + d (total not retrieved) remain constant. More common is the situation where all the above six values change, but a + c (total relevant) and b + d (total non-relevant) do not alter. This is to say that the numbers of relevant and non-relevant documents remain the same, but the numbers of retrieved and not retrieved, together with the four categories making up these groups, all vary. In such cases the change could be due to the 'cut-off' applied, that is the point in the search where the rules do not allow any further documents to be examined. At this stage the search is stopped and a record made of all the documents retrieved, both a (relevant) and b (non-relevant). A different cut-off results in a different set of values for a and b, thereby changing c and d, but without in any way affecting a + c or b + d. Alternatively, the change could be due to different indexing decisions or to different search strategies.

The second point to consider is the variables that affect a + c, b + d and N. If the decision as to what is relevant (a + c) is altered, then it must also result in a change for the total of non-relevant (b + d); if the collection size (N) is changed, other values in the table may change. Although significant changes of this nature occur rarely in operational retrieval system tests, it is necessary to consider the matter in experimental tests. Either type of change, i.e. altering the number of relevant documents or altering the collection size, can vary the number of relevant documents in relation to the collection size. Examples of the two types

of situations can be taken from these tests. Relevance decisions were based on four levels of relevance; if we consider Relevant 1 documents, there are 12 such documents relevant to the 42 questions of subset 2. Relevance 1 and 2 documents come to 57, Relevance 1, 2 and 3 documents total 154 and Relevance 1-4 documents come to 198. It can be seen that changing the decision as to the relevant documents (a + c) materially alters the proportion of relevant documents in the complete collection (N).

On the other hand, the collection size can be changed. Originally there were 1400 documents in the collection. A subset of the collection was formed which consisted of 200 documents; a characteristic of this subset was that it retained all of the 198 documents that were relevant to the 42 questions of subset 2.

While the number of relevant documents is now held constant, the proportion changes because of the reduction in the document collection from 1400 to 200 documents. It is convenient to express this variation as a parameter, and this is the aforementioned Generality number i.e.  $\frac{1000(a+c)}{N}$ , the total relevant documents divided by the collection size, with a constant. This parameter is not a measure of retrieval performance, but one which reflects the environment of the relevance decisions made; e.g. if the generality number for a set of questions is 5, this means that there are, for each question, an average of five relevant documents for every thousand documents in the collection, irrespective of what the actual size of the collection might be. For the example given above, the change from the larger to the smaller collection size (bearing in mind that there are 42 questions) changes the generality number from  $\frac{1000 \times 198}{42 \times 1400} = 3.4$  to  $\frac{1000 \times 198}{42 \times 200} = 23.6$ . Therefore, as far as retrieval performance is concerned, the significance of a change in either the relevance decisions or the collection size is that in both cases it is the generality number which alters.

The single performance measures that can be used can be listed as follows:-

- |                 |   |
|-----------------|---|
| $\frac{a}{a+c}$ | usually known as Recall Ratio; at Western Reserve University it is called 'Sensitivity', and has also been called 'Hit Rate'.                       |
| $\frac{e}{a+c}$ | complementary to recall ratio. Called by Fairthorne, 'Snobbery Ratio'.  |
| $\frac{a}{a+b}$ | now generally known as Precision Ratio, formerly called by Cranfield 'Relevance Ratio'. Also described as 'Pertinency Factor' or 'Acceptance Rate'. |

$\frac{b}{a + b}$  complementary to precision ratio. Called by Perry, 'Noise Factor'.

$\frac{b}{b + d}$  here called Fallout Ratio.

$\frac{d}{b + d}$  complementary to fallout ratio. Called by Western Reserve University, 'Specificity'.

Use of any of these single measures, either reflecting the retrieval of the relevant items or the retrieval of non-relevant items, is inadequate to reflect the performance of a system. High recall can mean very low precision, or vice versa, and the mere statement that the recall ratio is 99% means little, for it might only be achieved by retrieving more than half of the total collection.

While many different combinations of single measures have been proposed, they fall into two groups: 'twin variable measures' and 'composite measures'.

For the former, one of each of the single measures is taken and a comparison made between them by observing the relative changes in the two values, but retaining each value as a separate entity. The two major pairs of single measures are recall with precision and recall with fallout.

Examples of recall/precision ratios are given in Figs. 3.3 and 3.4. Fig. 3.3T illustrates the situation for a set of 20 searches where the variable being tested is the search coordination level, that is the number of search terms which must be matched with the index terms. At each different level, a cut-off is applied and the number of documents retrieved, relevant and non-relevant, is recorded. Since the total number of relevant documents is known, the recall and precision ratios can be calculated, as shown in the table. Alternatively these ratios can be plotted as on the graph (Fig. 3.3P) with the five performance points connected to make a recall/precision curve. In Fig. 3.4T are given the results of a series of searches with the same set of questions but with different search requirements. The particular change is incidental to the present discussion, but in fact whereas search X accepted any combination of terms, search Y would not accept certain terms unless some other given term was also present. (This matter of search strategy was discussed in Chapter 2). The result of this change was a different set of performance figures at the five coordination levels. The contrast between search X and search Y can be seen by comparing the tables or from the graph (Fig. 3.4P), which shows clearly that the maximum recall figure has fallen sharply in search Y, but on the other hand at any given recall ratio of 65% or less, search Y will give a higher precision ratio than search X.

Coordination Level	Documents Retrieved		Recall Ratio a/a + c	Precision Ratio a/a + b
	Rel (a)	Non-Rel (b)		
1	133	16,492	95.0%	0.8%
2	108	6,642	77.1%	1.6%
3	80	1,825	57.1%	4.2%
4	57	430	40.7%	11.7%
5	39	94	27.9%	29.3%

Relevant Documents (a + c) = 140

Generality number 5.

$N = 28,000$

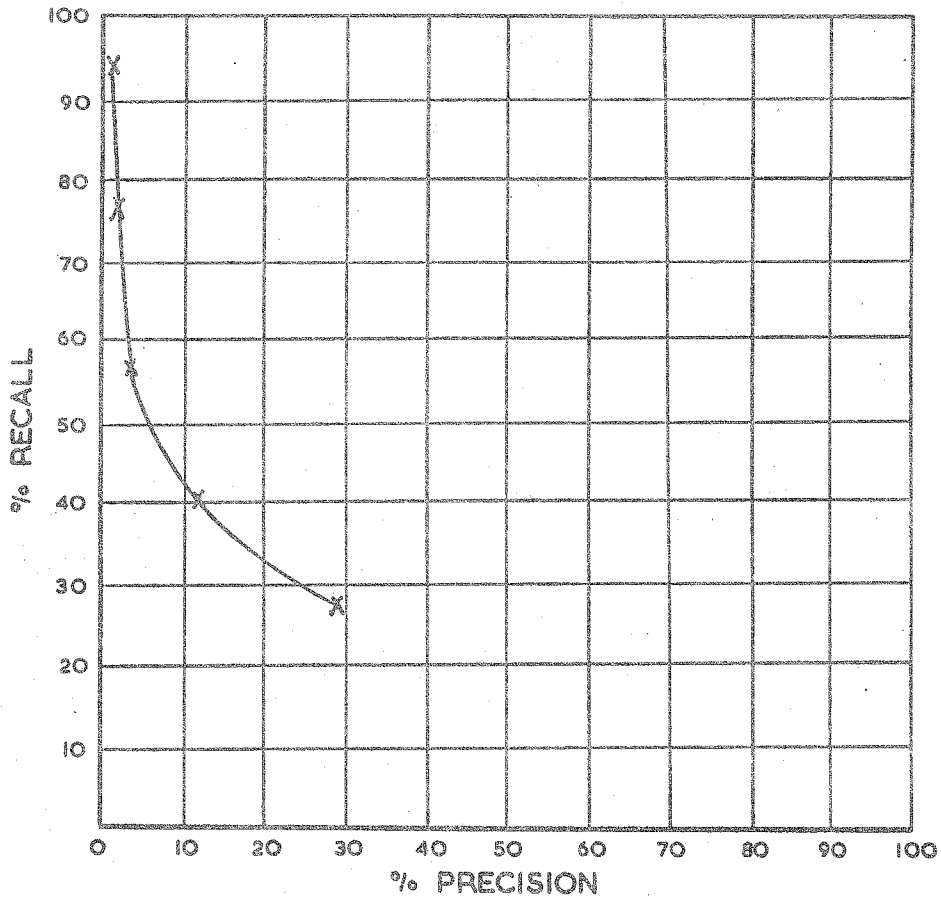


FIGURE 3.3TP

TABLE AND PLOT OF TEST RESULTS FOR 20 SEARCHES SHOWING RECALL AND PRECISION RATIOS FOR SEARCH X.

Coordination Level	Documents Retrieved		Recall Ratio a/a + c	Precision Ratio a/a + b
	Rel (a)	Non-Rel (b)		
1	97	2,674	69.3%	3.5%
2	77	788	55.0%	8.9%
3	56	220	40.0%	20.3%
4	37	30	26.4%	55.2%
5	33	17	23.6%	66.0%

Relevant documents (a + c) = 140

Generality number 5

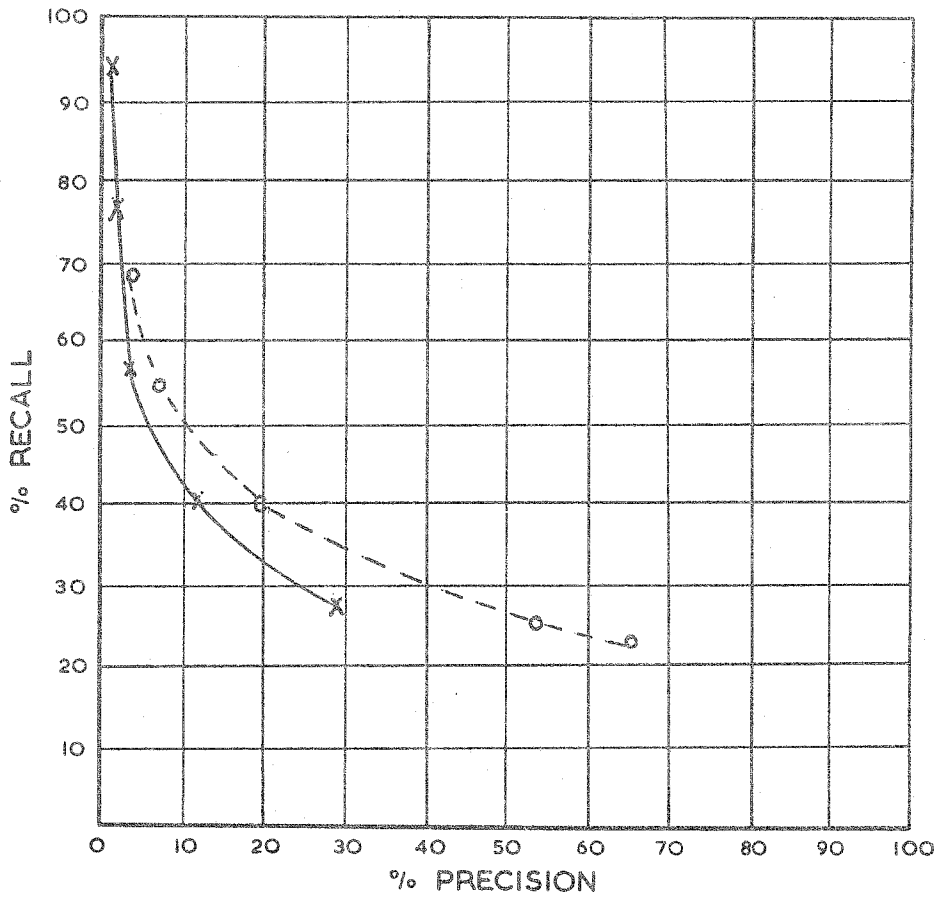


FIGURE 3.4TP

TABLE AND PLOT OF TEST RESULTS FOR 20 QUESTIONS SHOWING RECALL AND PRECISION RATIOS FOR SEARCH Y (BROKEN LINE) (SEARCH X CONTINUOUS LINE)

A comparison of the recall ratio with fallout ratio can be made in the same way. We are not aware of any previous occasions when the fallout ratio has been used for presenting test results, although Swets (Ref. 4) has discussed its possible use. In that it measures the ratio of the non-relevant retrieved to the total non-relevant in the collection

$\frac{b}{b+d}$ , it is very sensitive to N, the total number of documents in the collection. While it might not be found to be particularly satisfactory for tests on operational systems, it has an attraction in experimental testing where collections of different but known size are being tested, since it automatically compensates for the changes in size. Fig. 3.5T takes the figures of Fig. 3.3T and Fig. 3.4T and replaces the precision ratio by fallout ratio. A characteristic of fallout ratios is that they tend to be concentrated at low numbers; for this reason the figures are taken to three places of decimals and the resultant plot of recall ratio against fallout ratio is clearer if made on a semi-log scale, as in Fig. 3.5P. In this case the better performance is obtained when the curve is nearer the top left hand corner, whereas the recall precision curve is optimised towards the top right hand corner. Therefore, as in Fig. 3.4P, search Y is shown to give a generally improved performance over search X.

Either of these twin measures is satisfactory for presenting the performance of systems where the generality number is held constant, although the argument has been advanced that a plot of recall/precision is not valid since both ratios contain a (relevant retrieved). It has been incorrectly argued that in plotting  $\frac{a}{a+c}$  against  $\frac{a}{a+b}$ , all the a's cancel out, with the result that the factors being plotted are c against b. Fairthorne (Ref. 5) has said that a more reliable precision ratio is given by what he calls the 'distillation ratio' which is  $\frac{a}{a+b} - \frac{c}{d}$ .

However, he agrees that when the correction factor of  $\frac{c}{d}$  is negligible compared with the precision ratio, the latter is a valid measure. In fact, in the results presented in Fig. 3.3T, the correction factor at the coordination level of five terms is 0.0038, which can definitely be considered negligible.

Rees (Ref. 6) argues against precision ratio in favour of a measure that is complementary to fallout, namely  $\frac{d}{b+d}$ , on the grounds that it takes into account one of the vital parameters in a retrieval system - size of file. To some extent this is true, but it is a matter which has to be approached very carefully. The difficulty lies in determining exactly what is the correct value of N, that is to say how many documents can validly be considered to form the total collection in regard to any question. This matter is considered in more detail later in this chapter. It is true that the same difficulty arises in calculating the generality number, but if N is known, then it is just as easy to calculate the generality number as to

Coordination Level	Recall Ratio		Fallout Ratio	
	X	Y	X	Y
1	95.0%	69.3%	59.196%	9.598%
2	77.1%	55.0%	23.841%	2.828%
3	57.1%	40.0%	6.551%	0.790%
4	40.7%	26.4%	1.543%	0.108%
5	27.9%	23.6%	0.337%	0.061%

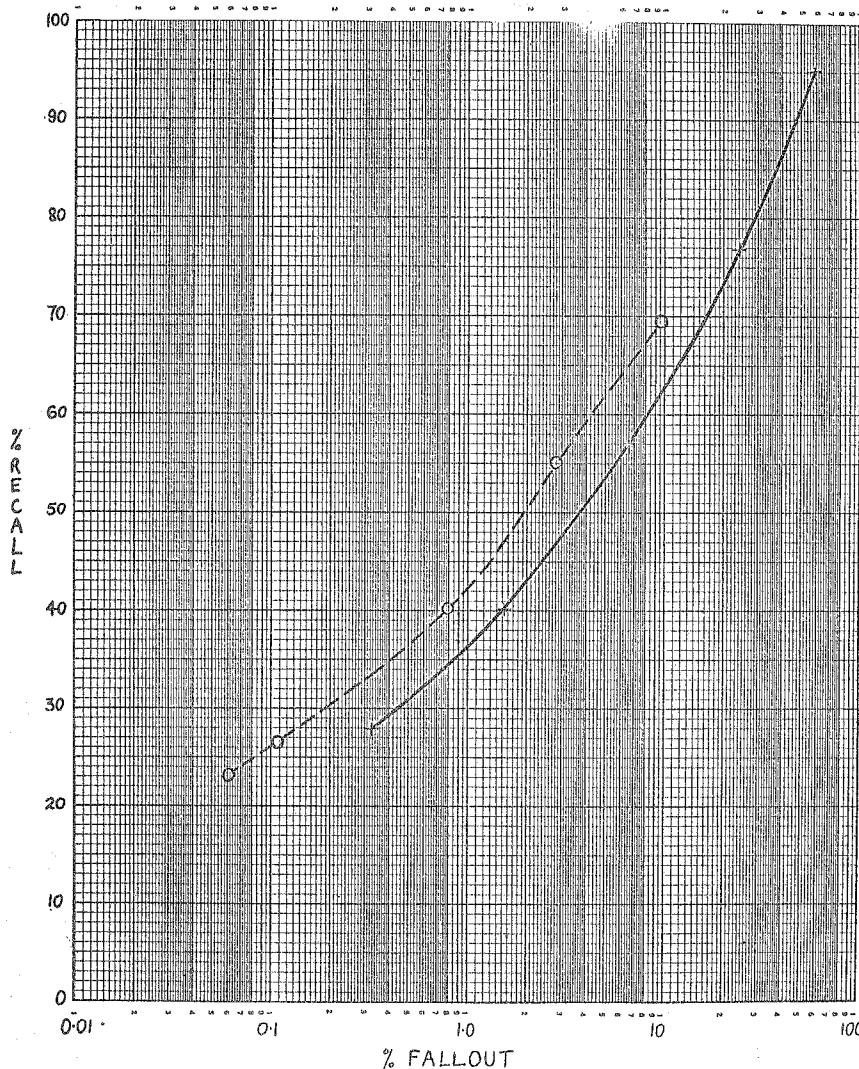


FIGURE 3. 5TP TABLE AND PLOT OF FALLOUT RATIOS DERIVED FROM FIGURES 3. 3T and 3. 4T FOR SEARCH X (CONTINUOUS LINE) AND SEARCH Y (BROKEN LINE)

calculate fallout.

A possible solution to making a full presentation of performance on a single plot is shown in Figs. 3.6P and 3.7P. Before considering these it is necessary to consider the relationship between the individual ratios of recall, fallout and precision, together with the generality number. These four ratios or parameters completely describe a given set of performance results in a retrieval table in terms of the measurements most likely to be of importance in presenting retrieval performance. However, it is only necessary to obtain any three of these in a given situation, since the fourth is then mathematically determined and can be written in terms of the other three. The four equations are:-

$$(1) \quad R \text{ (Recall Ratio)} = \frac{\left( \frac{F(1000 - G)}{1 - P} \right) \times P}{G}$$

$$(2) \quad F \text{ (Fallout Ratio)} = \frac{\left( \frac{R \times G}{P} \right) - (R \times G)}{1000 - G}$$

$$(3) \quad P \text{ (Precision Ratio)} = \frac{R \times G}{(R \times G) + F(1000 - G)}$$

$$(4) \quad G \text{ (Generality Number)} = \frac{1000}{\left( \frac{\frac{R}{P} - R}{F} \right) + 1}$$

$$\text{where } R \text{ (Recall Ratio)} = \frac{a}{a + c}$$

$$F \text{ (Fallout Ratio)} = \frac{b}{b + d}$$

$$P \text{ (Precision Ratio)} = \frac{a}{a + b}$$

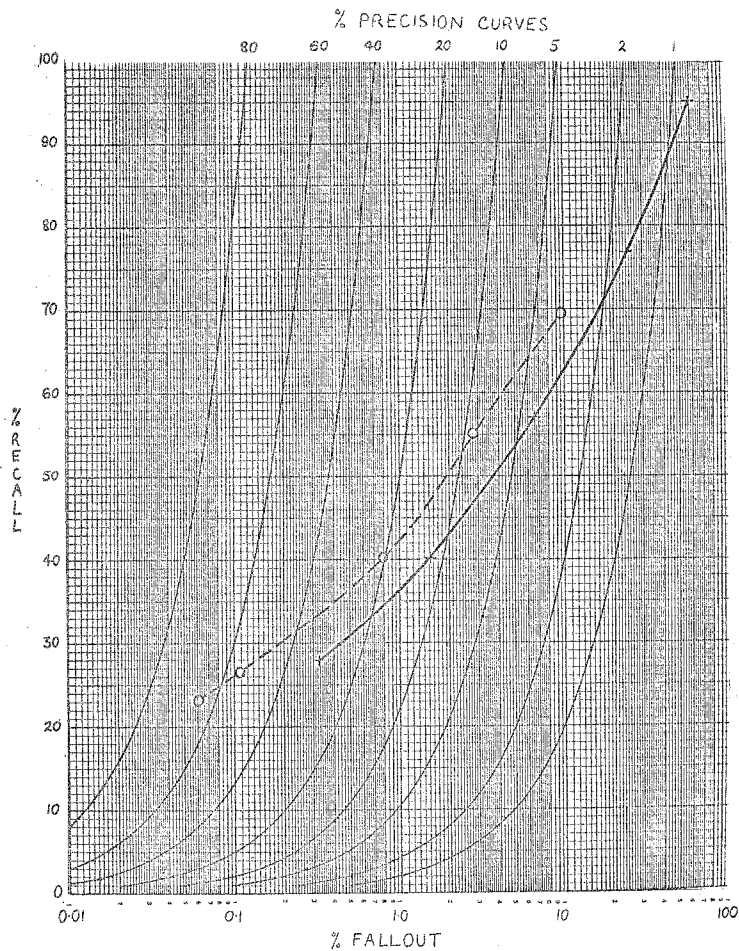


FIGURE 3.6P PLOT OF RECALL AND FALLOUT RATIOS AS FIGURE 3.5P SHOWING THE PRECISION RATIO CURVES

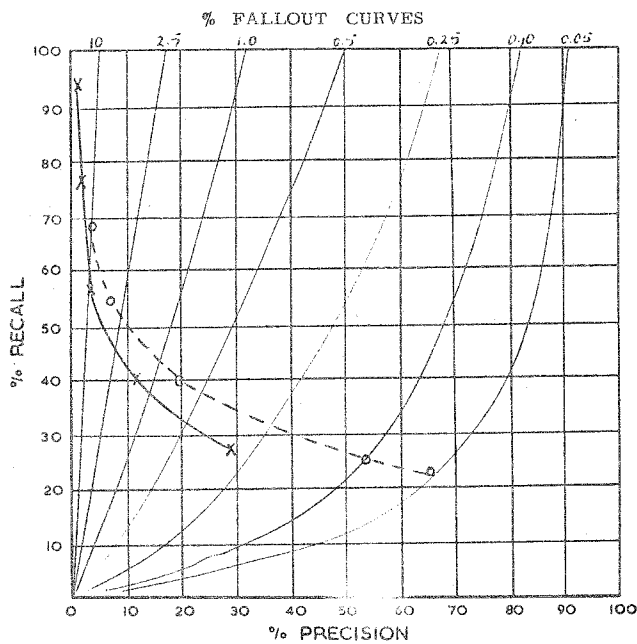


FIGURE 3.7P PLOT OF RECALL AND PRECISION RATIOS AS FIGURE 3.4P SHOWING THE FALLOUT RATIO CURVES

$$G \text{ (Generality Number)} = \frac{1000(a + c)}{N}$$

Thus equation (1) shows how, given the fallout and precision ratios together with the generality number, the recall ratio can be determined by calculation, and the other three equations show the other combinations possible. Because of this relationship, it has been possible to prepare, by computer, the figures for a series of situations where the generality number ranges from 1 - 50, recall from 5% to 100% and precision from 0.5% to 100%. In Appendix 3.3 is given this full set of tables for F (fallout) at varying generality numbers. From this set of tables, it is possible to plot on a recall/precision graph, the curves for fallout, or on a recall/fallout plot the curves for precision at all levels for any given generality number. For the example being considered, Fig. 3.6P shows the former, while Fig. 3.7P shows the precision curves on a recall/fallout graph. From either of these graphs it can be seen, for instance, that for search Y (the dotted line) at a recall ratio of 40%, precision ratio was 20% and the fallout ratio 0.8%. As the generality number for this set of searches is 5, the above figures can be confirmed from the sheet in Appendix 3.4 for generality number 5. In the column for recall of 40% and in the line for precision of 20%, fallout is 0.803%.

In a large number of situations arising in this test, comparison is made between various systems where everything is being held constant with one exception such as, for instance, the index language. In these circumstances the generality number remains constant and therefore the fallout measure does not contribute to the presentation of the results. In spite of the fact that there are some situations where comparative results are presented when the testing has been done on collections of different sizes, (with therefore, different generality numbers), the decision has been taken, as previously stated, to present the main sets of results on recall/precision graphs. The positive reason for doing this is that discussions with a number of people have led to the conclusion that such a graph can be more readily understood than a recall/fallout graph in that it more closely reflects the required performance aspects of a system. This may, of course, be due to the fact that recall/fallout graphs are unfamiliar compared with recall/precision graphs, and our decision is certainly not intended to imply that the latter are, in experimental work, basically superior to recall/fallout graphs.

In the course of this project, we have also considered a number of 'composite' measures which have been suggested. Swets (Ref. 4) argued that twin variable measures (e.g. recall/precision) were 'an unnecessarily weak procedure', but qualified this by assuming that a real retrieval system has a constant effectiveness, independent of the various forms of queries it will handle. He admitted that such an assumption is open to question, and it is clearly incorrect in an experimental situation where major variables are being changed with the result that new systems are being formed. In such tests, the twin variables are necessary to see the

changes that are taking place over the whole range of performance and even then need the additional environmental control of generality. It is difficult to understand the use of the term 'weak', since all composite measures can only present some compressed and simplified combination of the whole range of values shown by twin variable measures.

The composite measures can themselves be evaluated by recording their scale or range of values on the two twin variable plots. Any composite measure must indicate perfect retrieval in a situation of 100% recall at 100% precision at 0% fallout, and must indicate the worst retrieval in a situation of zero recall and zero precision at 100% fallout. Thus all composite measures have some scale of values between those two extremes, which can be plotted for visual examination on both recall/fallout and recall/precision plots.

Some of the measures proposed may be described as linear composite measures, when their values vary in some linear way if either the recall alters, or the precision (or fallout) alters. Perhaps the simplest composite measure suggested is the sum of the recall and precision ratios. Fig. 3.8P shows an example of this, using the simple sum of the recall and precision percentages, resulting in a range of values from 0 to 200. As can be seen, a performance of 70% recall at 10% precision would be given a value of 80, and be regarded as a better performance than 45% recall at 30% precision, or worse than a performance of 80% recall at 1% precision. The limitations of such a measure are fairly obvious, since a 70% recall at 10% precision will be rated the same as a performance as 10% recall at 70% precision or 40% recall at 40% precision, and many other different levels along the diagonal line. A simple weighting can alter the slope of the lines, e.g. if the recall ratio is weighted 1, and the precision ratio 2, the lines are more steeply positioned (Fig. 3.9P). The performance curves from Fig. 3.4P plotted on both tables are seen to have composite values which generally indicate the superior performance of search Y but, of course, the detailed differences at the cut-off points and the loss of maximum recall with search Y as against search X cannot be indicated by any composite measures.

Of measures of this type that have been suggested, J.D. Sinnott, in his thesis describing a test of role indicators, (Ref. 7), uses an effectiveness measure 'R', originally suggested by H. Borko, being  $R = 100 \left( \frac{a}{a+c} - \frac{b}{a+b} \right)$  which is the recall ratio minus the noise factor (the complement of precision). The resulting values are positioned as 45 degree diagonals on a recall/precision plot similar to Fig. 3.8P but having the range of values from -100 to +100, with the centre diagonal being 0. A second measure, put forward by Western Reserve University, is the measure 'Effectiveness', being the sum of sensitivity and specificity (Ref. 8), and appears as straight lines on a plot which reverses recall/fallout. This is shown in Fig. 3.10P, which, it should be noted, is not a semi-log plot as are the previous examples of recall/fallout.

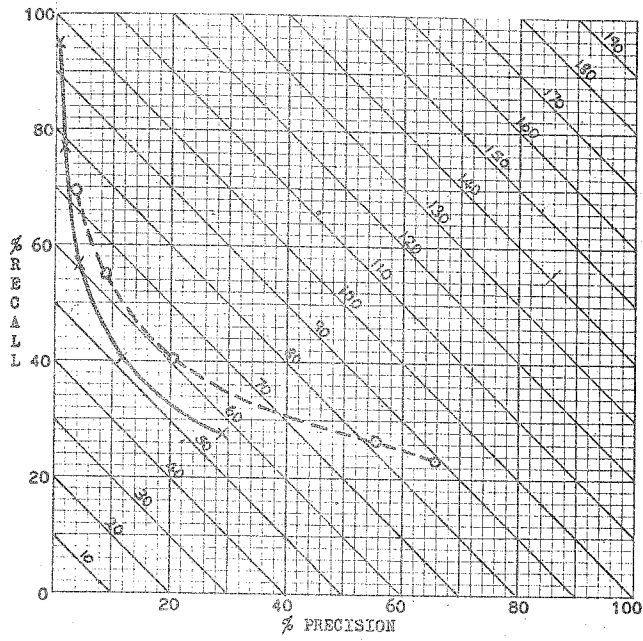


FIGURE 3.8P PLOT OF RECALL AND PRECISION AS FIGURE 3.4P SHOWING THE 'RECALL + PRECISION' LINES

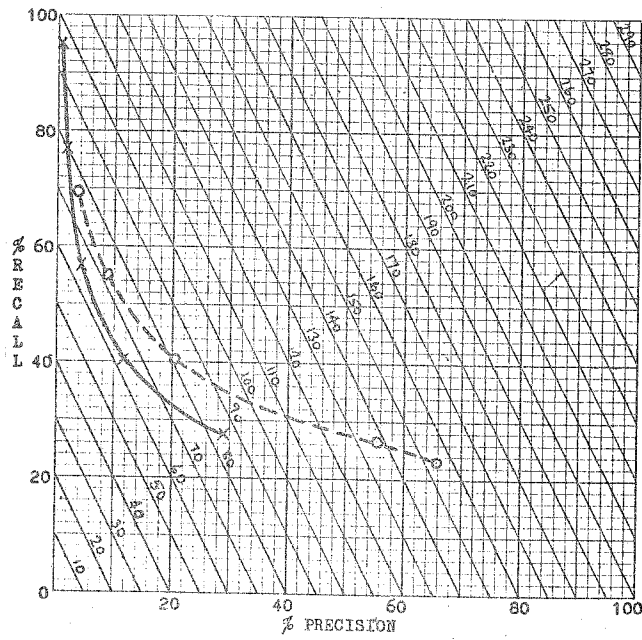


FIGURE 3.9P PLOT OF RECALL AND PRECISION AS FIGURE 3.4P SHOWING THE 'RECALL + PRECISION X2' LINES

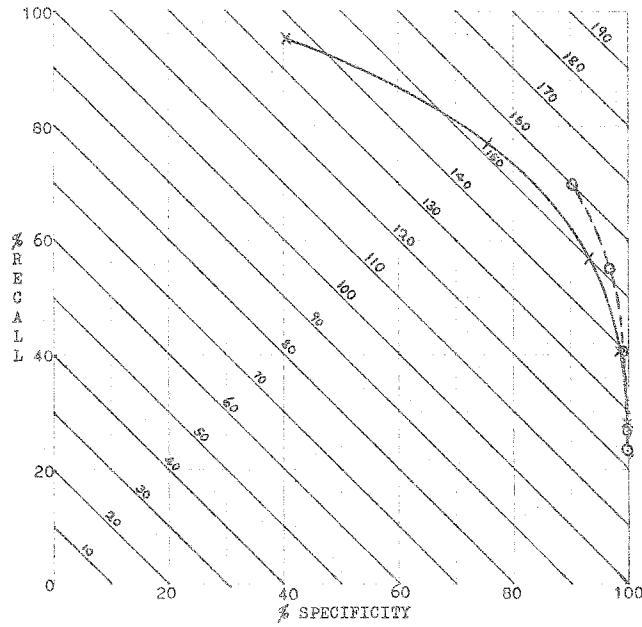


FIGURE 3.10P PLOT OF RECALL AND SPECIFICITY (ON A LINEAR SCALE) SHOWING THE 'EFFECTIVENESS' LINES

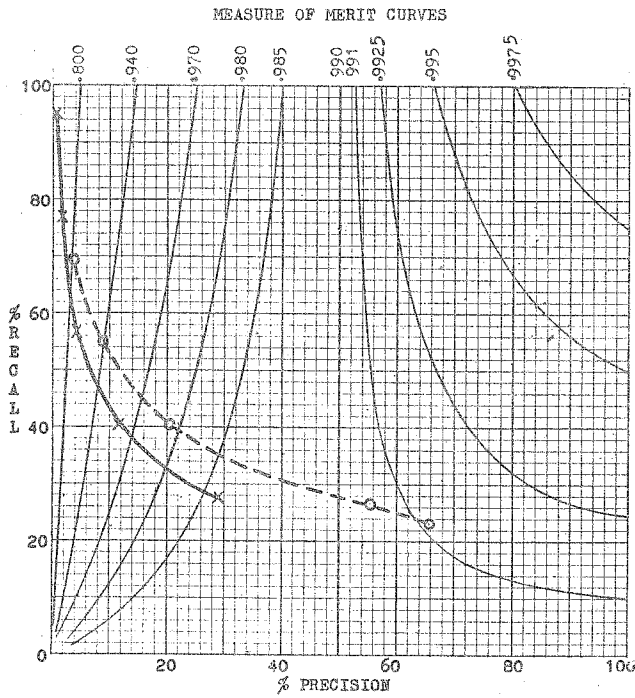


FIGURE 3.11P PLOT OF RECALL AND PRECISION AS FIGURE 3.4P SHOWING 'MEASURE OF MERIT' CURVES

MEASURE OF MERIT CURVES

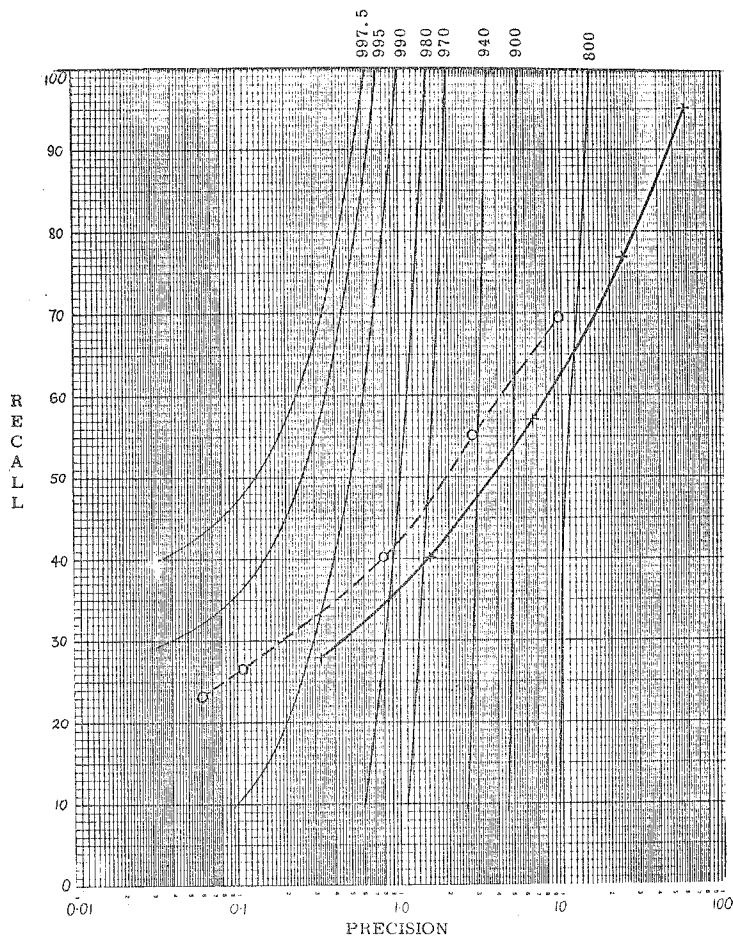


FIGURE 3.12P PLOT OF RECALL AND FALLOUT AS FIGURE 3.5P SHOWING 'MEASURE OF MERIT' CURVES

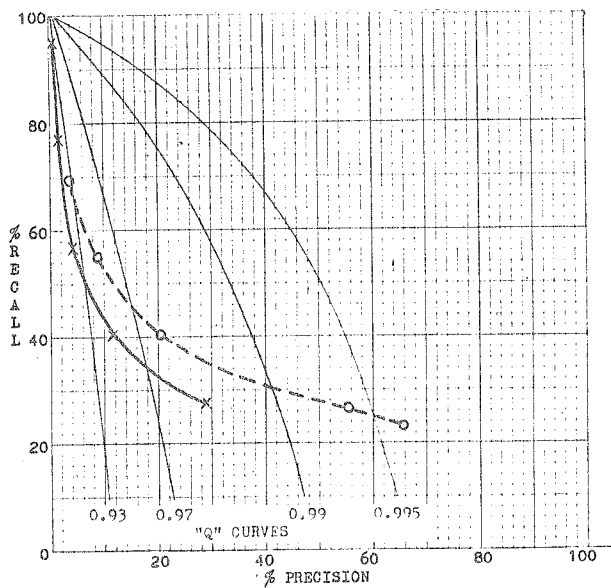


FIGURE 3.13P PLOT OF RECALL AND PRECISION AS FIGURE 3.4P SHOWING 'Q' CURVES

Other composite measures proposed can be described as non-linear composite measures, since their scale of values varies in a non-linear fashion when recall, precision, or fallout are varied, and the display of their values on the twin variable plots results in curves rather than straight lines. When a measure of this type includes d (non-relevant not retrieved) in its equation, the values and curves of the measure will be affected by the generality number. For Figs. 3.11 to 3.15 a generality of 5.0 is used in drawing the curves for the measures involved, since the performance results of searches X and Y that are plotted were obtained in a situation of that generality. The values of a composite measure of this type have been calculated in a manner similar to that adopted in making the two combined plots of recall, precision and fallout, Figs. 3.7P and 3.8P. In this case various sets of recall and fallout ratios, and also recall and precision ratios (at a generality of 5.0) were selected in advance and the resulting value of the composite measure calculated. This was done for different ratios to obtain curves of the measure that give a general indication of the range in its values.

The first of these non-linear composite measures which we consider is that proposed by J. Verhoeff and others, which is described as a 'Measure of Merit' (Ref. 9), with the basic equation:

$$M = a - b - c + d$$

This can also be written as  $M = (a + d) - (b + c)$  which is really the sum of the 'successes' minus the sum of the 'failures'. The values are shown in the two twin variable plots, Figs. 3.11P and 3.12P, with the equations divided by 'N' to obtain a range of values between 0 and 1, and it can be seen how high values of the measure occur at high recall with high precision or, to say the same thing in a different way, high recall with low fallout. The measure was intended to be used with various weights associated with the four component values, and any of the composite measures being described could incorporate this if in a given situation a meaningful set of weights can be devised. One might, for instance, hypothesise 'cost values' of failing to retrieve a relevant document or retrieving a non-relevant document. Any such weighting would alter the position of the measure's curves on the plots.

A more complex version of this is the Q factor, which has been suggested by Farradane as suitable for use in retrieval tests. This is a statistical coefficient of association proposed by Yule (Ref.10). The formula is  $Q = \frac{ad - bc}{ad + bc}$ , which can be described as the product of the successes minus the product of the failures divided by the sum of the same two products. Figs. 3.13P and 3.14P show the two graphs with Q curves plotted, with the performance curves. It has not been shown that Q curves have any significance in retrieval tests, and there does not appear to be any reason why they should.

A measure put forward in discussion by Vickery at the NATO Advanced Study Institute on Evaluation, held at The Hague, July 1965, uses the values

of a, b and c, from the retrieval table. He suggested that the measure should reflect the ability of the system to maximise a relative to b and c, described as the selectivity of the system. The proposed measure F, uses a normalisation factor S, where  $S = a + b + c$ , and

$$F = \frac{100 \frac{a}{S}}{\frac{b}{S} + \frac{c}{S} + 1}$$

F varies from 0 to 100, and is plotted on a recall/precision plot in Fig. 3.15P. The curves are symmetrical about the diagonal from the bottom left corner to the top right corner, and alter in shape as they approach the top right side.

All the composite measures described have an apparently reasonable scale of values ranging from the case of worst performance to that of best possible performance, but none of these measures can show the very large differences that occur between these two points, in the different positions at which systems actually operate. The curves in Figs. 3.4P and 3.5P are indicators of retrieval performance when a component of a system is varied to give results over the largest possible operating range, but the composite measures can only reflect one, or sometimes two, points of such curves. It is unfortunate that, in examples investigated so far, the point on the curves which determines the highest value assigned to that test by a given composite measure is usually either the point of maximum recall, or of maximum precision, neither of which may be the best points to use. It is a reasonable conclusion that for experimental tests where changes of the variables in systems are examined, the composite measures so far proposed are inadequate, although for tests where a single cut-off point is chosen, or a single cut-off is applied to two systems in a comparable manner, some of the composite measures may be useful. In experimental tests it is suggested that an 'area measure' is required; a possible solution is put forward in Chapter 5.

Having examined the main suggested performance measures, it may be asked whether any theoretical objective methods are known which could be used to evaluate the proposed measures, or whether tests and experience of actual results will be the only arbiter.

The only theoretical basis suggested so far is the use of the 2 x 2 contingency table, as already mentioned. Although the retrieval situation obviously fits the case in the sense that the resulting values of a retrieval test perfectly fit the nine categories in the table, no reasons have been advanced to show that figures from retrieval tests can benefit from the statistical tests commonly used. The retrieval situation is very different from the simple statistical one. For example, a typical 2 x 2 table taken from a popular textbook on statistics by M.J. Moroney (Ref. 11, page 264)

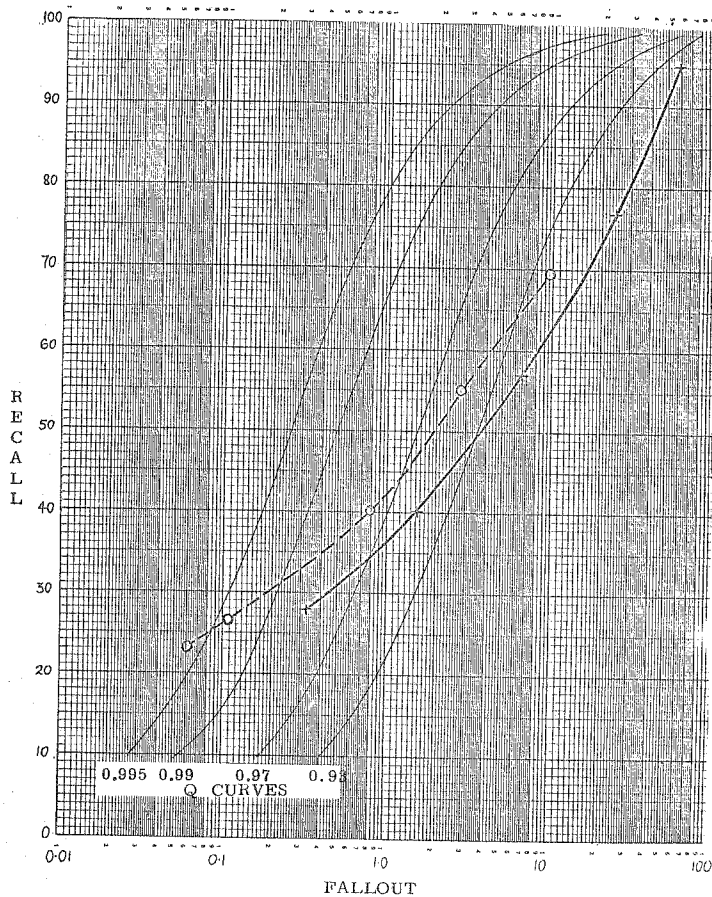


FIGURE 3.14P PLOT OF RECALL AND FALLOUT AS FIGURE 3.5P SHOWING 'Q' CURVES

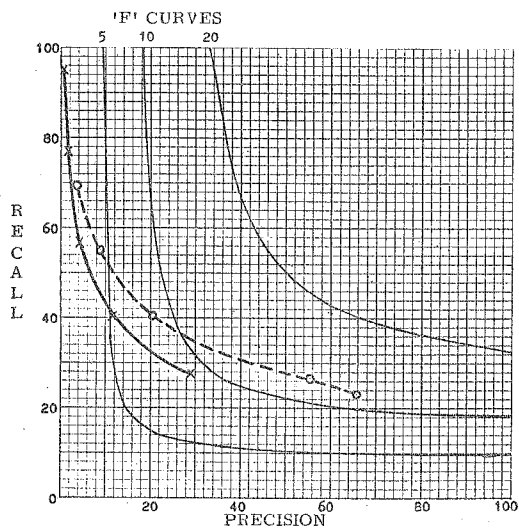


FIGURE 3.15P PLOT OF RECALL AND PRECISION AS FIGURE 3.4P SHOWING 'F' CURVES

gives data on a population of 77 people, showing the numbers that were both inoculated and not inoculated, and the numbers that were infected and not infected. The usual purpose of such a table is to ask a question of the kind, 'Is there really some degree of association between the events?', or in this particular case, 'Is the proportion of people that were not inoculated and became infected significantly different from the proportion of people that were inoculated and were infected?' In this situation, certain tests for the reality or existence of the association can be used (e.g. the chi square test), and other tests to determine the intensity of the association (e.g. the Q formula) can be applied. The form in which the question is posed, and the tests of the reality of association do not fit the retrieval case. Any question such as 'Is the proportion of relevant documents in the retrieved set significantly different from the proportion in the set not retrieved' does not make any sense in the retrieval situation. In the retrieval situation it is two sets of ratios from the table that are to be compared with one another by observing the relative changes in the ratios as conditions are changed. The actual comparative proportions do not need any test of significance. The tests of intensity of association do reflect the situation when the retrieval case is perfect, and when it is at its worst, and therefore provide one scale between the two extremes. But the deficiencies of the composite measures have been noted, and no assistance or confirmation of the twin variable measures being used seems to be given. The conclusion is that statistics does not help at all at this point.

#### Averaging sets of results

To present reliable results of performance, the figures from a set of questions must be averaged in some way. The size of the question set required in order to give reliable results will not be considered here, since there are many standard statistical tests to use in order to determine the significance level of a set of results. It is obvious that the results of individual questions will vary considerably, and some idea of the magnitude of this variation may be gained from Figs. 3.16P and 3.17P. In these plots of recall/precision, the individual results from a set of questions are plotted, where single term natural language indexing is being tested. Fig. 3.16P shows the points that result when any three out of a possible total of seven of the search terms in each of thirty-one questions are demanded in 'logical product' coordination. Fig. 3.17P shows points from thirty-five questions when the level of search terms demanded in coordination is varied from two to seven, and the scatter is quite wide, ranging from 11% recall at 1% precision in the bottom left corner, to 100% recall at 100% precision at the top right corner. However, a trend is clearly present down the left side of the plot and at the bottom right corner, with a tendency for results at a high coordination level to give high precision and low recall, and with lower coordination levels resulting in an inverse change. Two different methods of averaging these results, at each of the 'coordination levels', may be used.

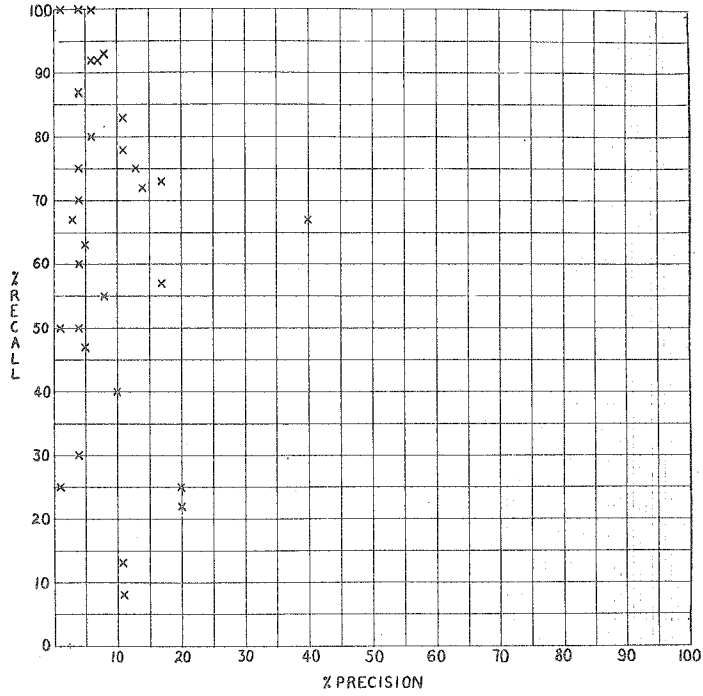
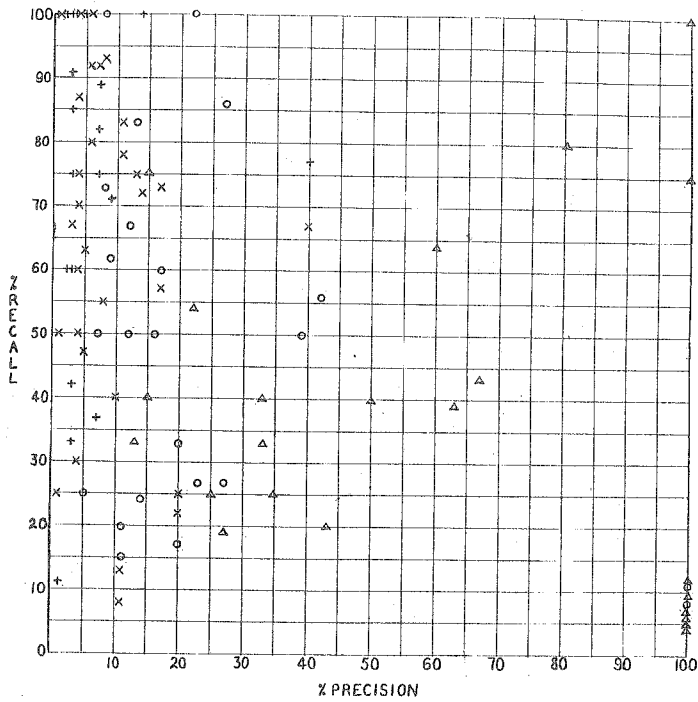


FIGURE 3.16P PLOT OF INDIVIDUAL RECALL AND PRECISION RATIOS OF 31 QUESTIONS SEARCHED AT A COORDINATION LEVEL OF 3 TERMS



+ 2+  
x 3+  
o 4+  
Δ 5+, 6+, 7.

FIGURE 3.17P PLOT OF INDIVIDUAL RECALL AND PRECISION RATIOS OF 35 QUESTIONS SEARCHED AT COORDINATION LEVELS BETWEEN 2 AND 7 TERMS

The first method, as used in Cranfield I, involves obtaining total figures of the numbers of documents involved for the whole set of questions being used in the test, and then converting the one grand total into, say, recall and precision ratios. In the case of the 35 question set, a total of 287 relevant documents is sought; at a coordination level of 3+, 157 of the relevant documents are retrieved, together with 2,865 non-relevant documents. These totals are then used to calculate the ratios of:-

$$\text{Recall} = \frac{100a}{a + c} = \frac{157}{287} \times 100 = 54.7\%$$

$$\text{Precision} = \frac{100a}{a + b} = \frac{157}{157 + 2865} \times 100 = 5.2\%$$

$$\text{Fallout} = \frac{100b}{b + d} = \frac{2865}{(35 \times 1400) - 287} \times 100 = 5.9\%$$

These ratios are obtained for all of the seven possible coordination levels, and can then be plotted as points on a graph. While this procedure of averaging the numbers was used for presenting the results of the first Aslib-Cranfield Project and the Western Reserve University test, at the time of the latter test it was realised that this method results in certain searches affecting the final figures more than others. Non-typical questions, such as those which retrieve an exceptionally large number of non-relevant documents, will exert a disproportionate influence on the final figures, and, in the W.R.U. test, separate figures were given showing the change in performance when those questions that retrieved unusually large numbers of (mainly) non-relevant documents were deleted (Ref. 2, page 13).

The second method of merging a set of results first converts the results of individual questions into recall, precision or fallout ratios and then obtains the final figures by using the average of the ratios of each question. In Fig. 3.18T are given the results of 35 questions which have been calculated in both ways, thus enabling a comparison of the 'average of numbers' and 'average of ratios' methods for these particular results. Recall, fallout and precision ratios for the two methods are compared in tabular form. It can be seen that there is no significant difference in the recall ratios between the two methods; at some coordination levels the average of ratios gives a slightly higher recall ratio, and at other levels the opposite is the case. The fallout values also show no significant difference. However, in the case of the precision ratios, it is clearly seen that the average of ratios gives a substantially higher figure for all coordination levels. Fig. 3.19P is a recall/precision plot of the two methods, where the 'better' curve results from averaging the ratios. As can be seen from the tables, a recall/fallout plot would have virtually overlapping

Coordination Level	Recall Ratio		Fallout Ratio		Precision Ratio	
	A	B	A	B	A	B
1+	93.4%	93.4%	48.329%	48.330%	1.1%	1.7%
2+	77.3%	77.3%	16.500%	16.502%	2.7%	4.2%
3+	54.7%	55.6%	6.954%	6.055%	5.2%	8.7%
4+	32.8%	31.2%	1.540%	1.538%	13.5%	24.4%
5+	16.4%	14.9%	0.547%	0.586%	25.5%	54.3%
6+	8.0%	6.9%	0.381%	0.377%	38.3%	64.2%
7	2.8%	3.3%	0.192%	0.190%	50.0%	77.8%

A Average of Numbers  
B Average of Ratios

FIGURE 3.18T COMPARISON OF RECALL, FALLOUT AND PRECISION RATIO WHEN TOTALLED BY (A) AVERAGING OF NUMBERS AND (B) AVERAGING THE RATIOS FOR 35 QUESTIONS.

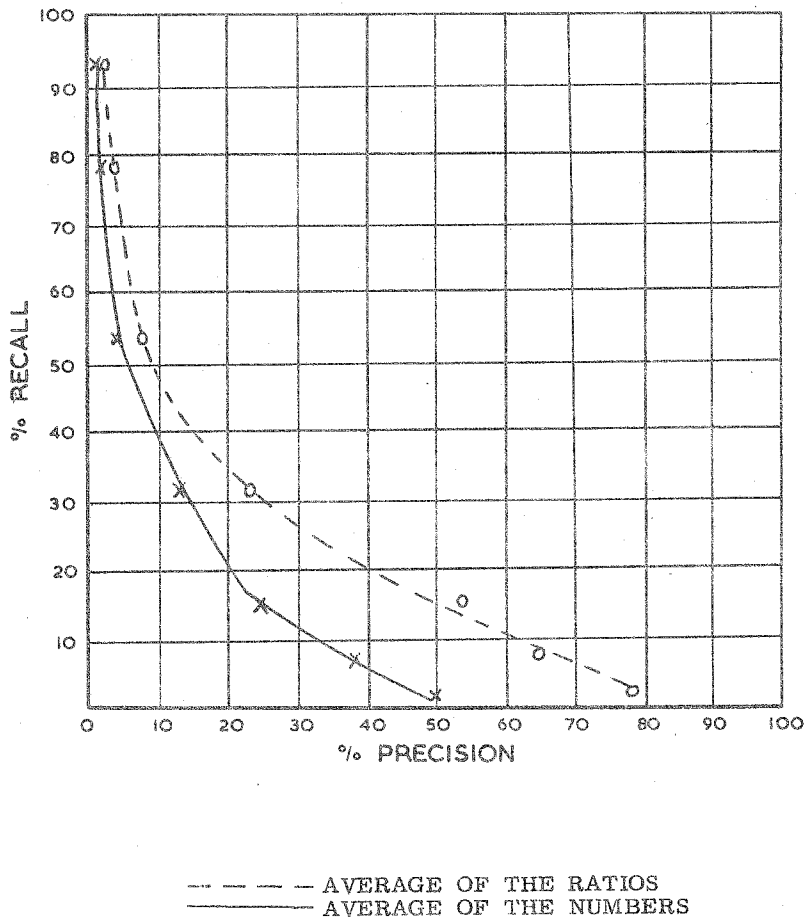


FIGURE 3.19P PLOT OF RECALL AND PRECISION RATIOS FROM FIGURE 3.18T COMPARING TOTALLED BY AVERAGING THE NUMBERS AND AVERAGING THE RATIOS

curves.

In the tests at Cranfield and in other tests where sufficient data has been available, the samples which have been processed by both methods have always shown this increase in precision with recall remaining much the same. However, we do not wish to be misquoted on this point and would emphasize that while it is probably true that the average of ratios will usually give a better performance figure, it would be wrong to assume that the proportional improvement would always be so pronounced as in the example shown.

An evaluation of the two methods which shows one method to be superior is not possible, since proponents of both methods can give good reasons for adopting one method in preference to the other. The theoretical cause of the discrepancy is the variation in the base from question to question: in the case of the recall ratio it is the number of relevant documents sought; in the precision ratio it is the total retrieved; and in the fallout ratio it is the total non-relevant. The average of numbers method weights the results of individual questions according to the base, and a larger base exerts a greater influence on the final result. The average of ratios completely ignores the base variation. In situations outside retrieval tests, where similar data has to be averaged, it is frequently advocated that the variation in base should be allowed for, and the average of numbers used (see, for instance Ref. 12, page 161). The difference in the results of the two methods is small except when the range and distribution of the variation in base becomes large, as is often the case with the precision ratio. However, both methods appear to be equally reasonable for use in retrieval situations, and the different results are really complementary viewpoints requiring careful interpretation.

A description of the different viewpoints represented by the two methods has been given by Salton (Ref. 13). He suggests that the average of ratios is 'a query-oriented viewpoint', and the average of numbers is a 'document-oriented viewpoint'; performance figures using the average of ratios indicate the performance of a single typical search question, typical that is of the set of questions used in the test. The use of average of numbers indicates the result of the whole set of questions, or indicates the success in performance of looking for a given set of relevant documents (287 in the example being used). This really ignores the actual individual questions involved, since one question with 287 relevant documents could in theory have the same result as 35 questions having in total 287 relevant documents. Thus the average of numbers gives an arithmetical mean value for a set of questions, and the average of ratios gives what approximates to a 'median' value which reflects the performance of a typical question.

Neither method appears to have any marked superiority over the other as a means of presenting results. However, the decision to use in this volume the average of numbers method was based on a most

important practical advantage, namely the comparative ease of calculation. To have used the method of the average of ratios would have increased the calculations forty-fold; work that has taken hundreds of hours would have taken hundreds of weeks. The really important matter in any test is to know which method is being used and to use it consistently in all situations.

#### Method of totalling results

Apart from deciding on whether to use the average of ratios or the average of numbers, we were faced with the additional problem which is involved in totalling results of sets of searches where the coordination level cut-off is employed; an idea of what is involved in this problem can be seen in Appendix 4A. There are given the performance results, in actual figures, for the 221 questions (subset 3), being tested on Language I.1a (single terms, natural language and coordination), Exhaustivity 3, Search rule type A, Document Relevance 1 - 4 and searched on the 1400 document collection. The questions are arranged in numerical order, and for each question is given the total number of relevant documents in the whole collection, followed by the relevant and non-relevant documents actually retrieved at each coordination level. In the final column is given the sum of the total number of postings for the search terms; the total must, of course, equal the sum of the relevant and non-relevant documents at all coordination levels. The variations between the 221 questions that affect the problem of arriving at a single result of a single performance curve for the 221 questions can be seen in tabular form in Fig. 3.20T. Here the two characteristics of the 221 questions are listed, namely the numbers of terms initially selected from the search question and used as search terms (starting terms), and the number of retrieving terms, that is the maximum number of starting terms which, used in logical product coordination, may be put to the index and will still retrieve documents (whether relevant or non-relevant).

The table shows how, for this particular test, the starting terms ranged from 2 to 15, and the retrieving terms varied from 2 to 10. Within this 14 x 9 matrix the actual number of questions involved is recorded, so it can be seen, for example, that of the 35 questions having seven starting terms (column headed 7) only three of these questions could coordinate all seven terms and still retrieve some documents. The figures in the table refer only to the particular index language in use, and a different index language such as index language I-5 which includes synonyms, word endings and quasi-synonyms, would alter the distribution of the questions in relation to the retrieving terms, while any test involving a different basic index language (such as simple concepts as compared to single terms) would alter the starting term groups also.

		Number of starting terms															
		2	3	4	5	6	7	8	9	10	11	12	13	14	15	Totals	
Number of retrieving terms	2	1	3	1	-	-	1	-	-	-	-	-	-	-	-	6	
	3		5	5	7	5	6	-	-	-	-	-	-	-	-	28	
	4			9	18	8	10	7	4	-	-	-	-	-	-	56	
	5				8	8	11	8	4	3	2	-	1	-	-	45	
	6					3	4	7	7	10	1	1	2	-	1	36	
	7						3	5	8	6	4	3	1	-	2	32	
	8							-	2	-	5	2	-	1	-	10	
	9								1	1	4	1	-	-	-	7	
	10									-	1	-	-	-	-	1	
	Totals		1	8	15	33	24	35	27	26	20	17	7	4	1	3	221

FIGURE 3.20T DISTRIBUTION OF THE 221 QUESTIONS BY STARTING TERMS AND RETRIEVING TERMS, IN ONE PARTICULAR TEST.

The table in Fig. 3.20T may be considered as showing how, in two respects, the 221 questions are a heterogeneous set of questions. Various subsets of the 221 can be picked to overcome the variations, and truly homogeneous subsets occupy each cell in the table, e.g. the five starting term group with four retrieving terms is the largest such subset, having a total of eighteen questions. A partially homogeneous subset, on the basis of one common characteristic only (either starting terms or retrieving terms), was the first to be examined in an attempt to find a method of totalling the whole set.

The subset of seven-starting-term questions was chosen and totalled by simply adding up each question at the seven possible coordination levels, resulting in seven totals. These totals are shown in Fig. 3.21T, and the recall precision percentages are recorded, these being calculated by using the average of numbers. The seven average recall and precision ratios are plotted in Fig. 3.21P, thus producing a performance curve for 35 questions, when the exhaustivity of search is altered by coordination levels. Since the characteristic of retrieving terms was ignored, not all the 35 questions provide results at all coordination levels, and, as was seen in Fig. 3.20T, one question is unable to retrieve any documents when more than two of the terms are demanded in coordination, and only three questions provide results at a coordination level of seven. The number of

Coordination Level	Documents Retrieved		Recall Ratio	Precision Ratio
	Rel.	Non-Rel.		
1	302	29,898	94.7%	1.0%
2	257	10,917	80.6%	2.3%
3	191	3,292	59.9%	5.5%
4	119	899	37.3%	11.7%
5	47	132	14.7%	26.3%

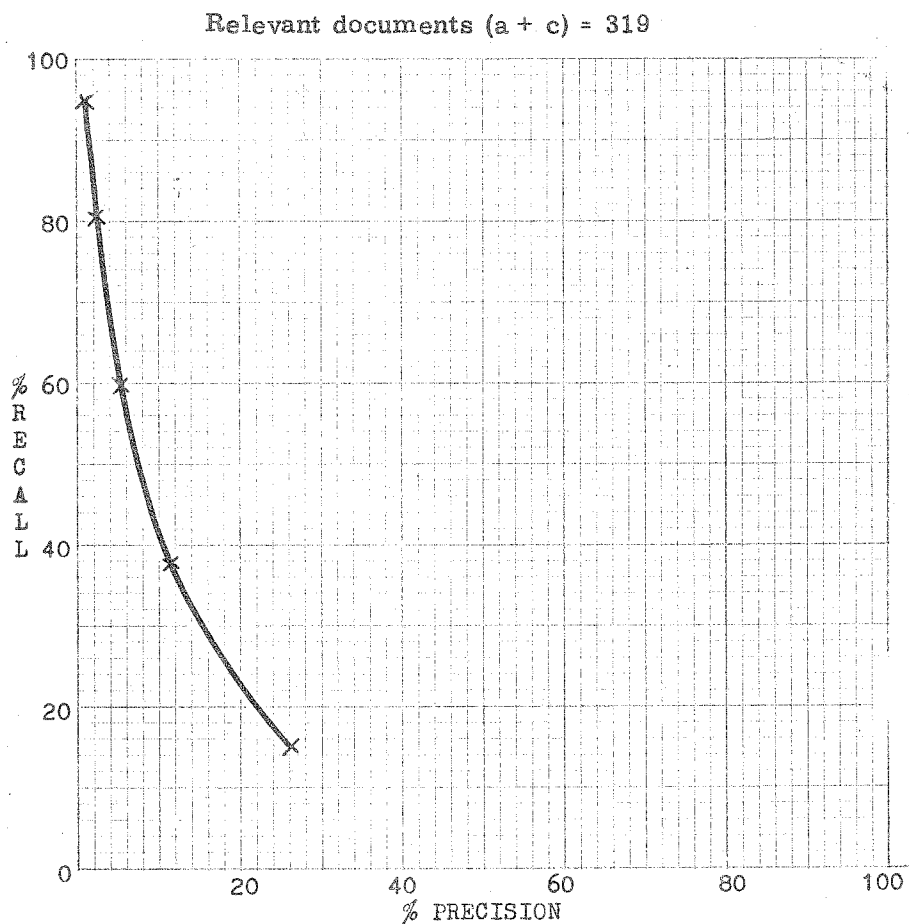


FIGURE 3.22TP TABLE AND PLOT FOR RESULTS OF 45 QUESTIONS WITH FIVE RETRIEVING TERMS TOTALLED BY COORDINATION LEVELS

Coordination Level	Recall Ratio	Precision Ratio
1	95.0%	0.9%
2	80.7%	2.2%
3	59.5%	4.1%
4	38.1%	7.6%
5	19.7%	11.6%
6	9.7%	19.0%
7	4.7%	25.5%
8	1.4%	33.8%
9	0.5%	61.5%
10	0.1%	100.0%

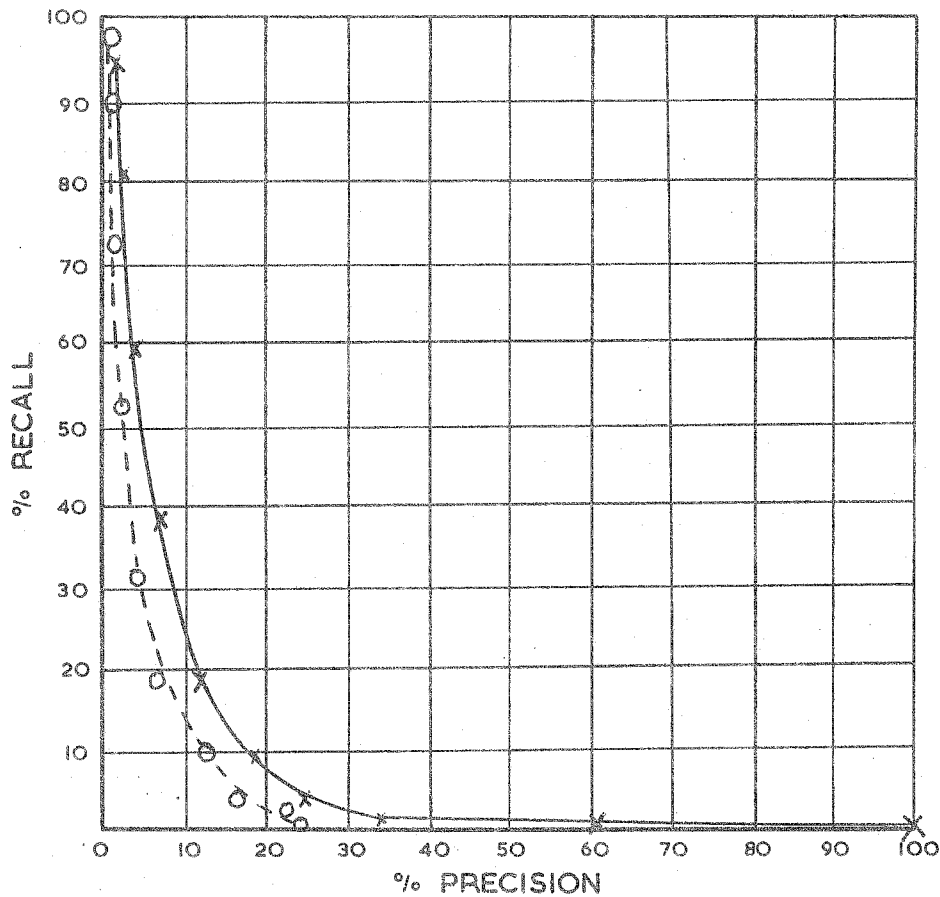


FIGURE 3.24 TP . TABLE AND PLOT FOR RESULTS OF 221 QUESTIONS TOTALLED BY METHOD 1A, COORDINATION LEVELS, FOR INDEX LANGUAGE I. 1. a. (INDEX LANGUAGE I. 6. a DOTTED LINE)

For Method 1, the questions were totalled in a similar manner to the starting term groups described above. This meant that for any given coordination level, (say, for example, four terms), the total results were obtained by adding the individual results for all the 221 questions, irrespective of the number of starting terms which each question had. Two variants of this strict coordination level totalling were considered. Method 1A involved totalling as described, and the resulting performance ratios are given in Fig. 3.24T, for Single Term Index Language I.1. The performance plot is given in Fig. 3.24P, with an additional curve of Language I.6 for comparison. In Method 1B, account is taken of the fact that at the higher coordination levels, many of the questions are not capable of contributing results, since the number of starting terms in the question is fewer than the coordination level. It is, for instance, quite impossible, at a coordination level of seven terms, to retrieve documents related to any of the questions which only have six, five, four, three or two starting terms. This effect increases, of course, with the coordination level. In this case, therefore, the recall ratio is calculated only for the questions that are capable of giving results. Fig. 3.25TP shows this, where it is seen that at a coordination level of 8+, only 704 relevant documents, i.e. less than half of the real total for this set of questions, are taken as the total of relevant documents being sought. This results in an increased recall ratio compared with Method 1A, but the precision ratio is not affected. A disadvantage of this method is that at each coordination level a change in generality occurs.

In Method 2, an attempt is made to allow for the fact that questions differ according to the number of starting terms. The strict coordination level of Method 1 can be faulted for equating, for example, the results of a five starting-term question searched at a coordination level of four terms, with the results of a ten starting-term question, also searched at four terms. The basic Method 2 can be described as 'totalling by proportional coordination levels', since it takes into account the potential range of coordination levels, which differs between questions. For example, a three starting-term question searched at a coordination level of two terms is demanding a match of two-thirds of the theoretical maximum, and in this method all questions having such a match would be included in the group. For a six starting-term question, for a nine starting-term question and for a twelve starting-term question, a two-thirds match would be four terms, six terms and eight terms respectively, although, for most other questions, no exact two-thirds match is possible. There are obviously many variations which are possible, but the example presented illustrates the use of this method when seven levels of match are chosen to obtain a total result.

There are obviously many ways in which this method could be applied; the example presented is where seven terms of match have been selected. Whatever the actual number of coordination levels in any particular question, the results are forced into the seven-term pattern. As can be seen from Figure 3.26T, this means that certain results are repeated, while for questions with more than seven starting-terms, certain results have to be omitted.

Coordination Level	Total Relevant	Relevant Retrieved	Recall Ratio	Precision Ratio
1	1590	1510	95.0%	0.9%
2	1590	1283	80.7%	2.2%
3	1583	946	59.8%	4.1%
4	1507	606	40.2%	7.6%
5	1401	314	22.4%	11.6%
6	1143	154	13.5%	19.0%
7	991	74	7.5%	25.5%
8	704	22	3.1%	33.8%
9	546	8	1.5%	61.5%
10	349	1	0.3%	100.0%

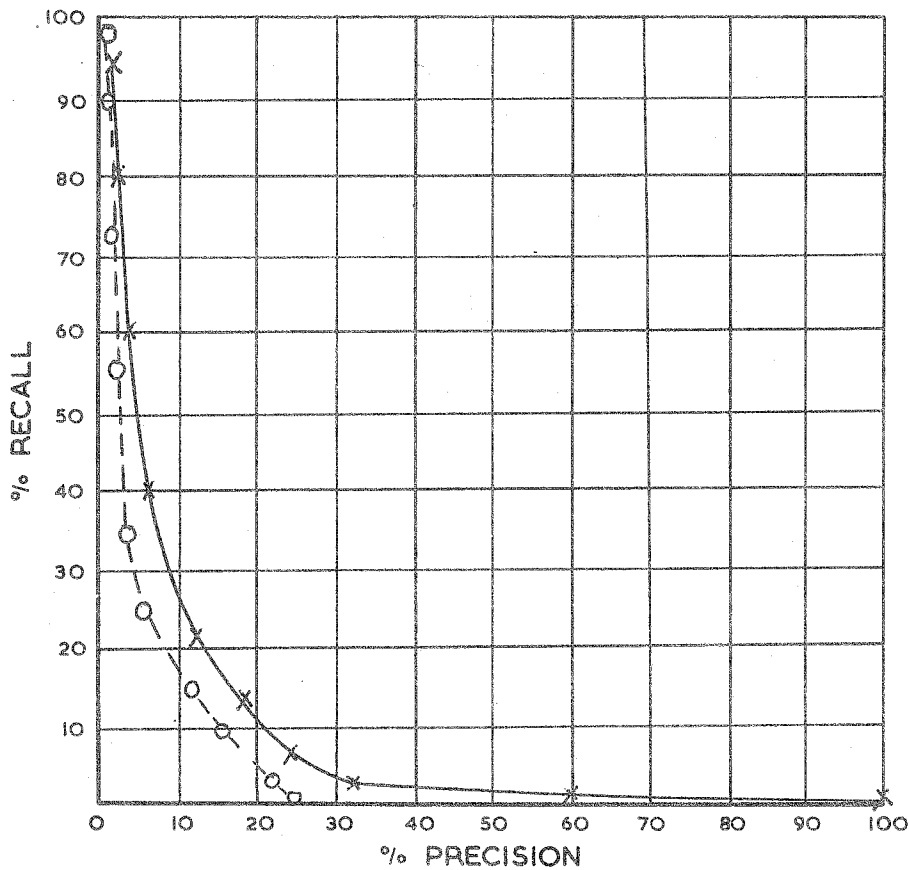


FIGURE 3.25TP

TABLE AND PLOT OF RESULTS FOR 221 QUESTIONS  
 TOTALLED BY METHOD 1B, ADJUSTED COORDINATION  
 LEVELS, FOR INDEX LANGUAGE I. 1. a.  
 (INDEX LANGUAGE I. 6. a DOTTED LINE)

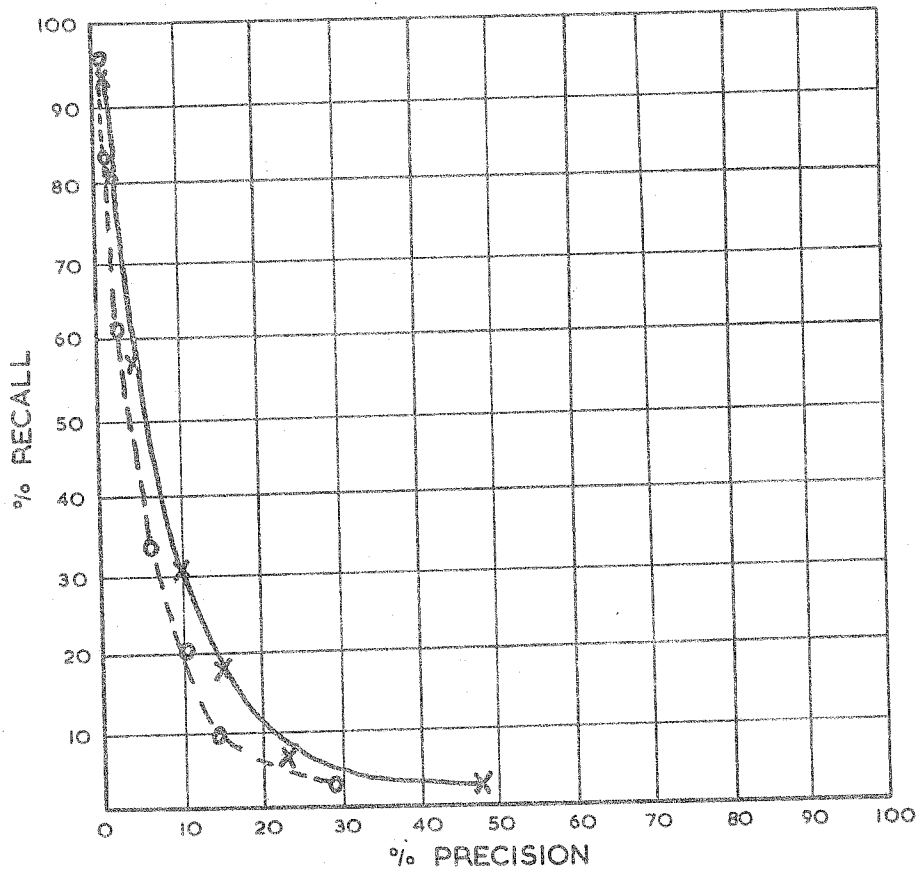
Starting Term Groups	Seven Proportional Coordination Levels						
	1/7	2/7	3/7	4/7	5/7	6/7	7/7
2/3	1	1	2	2	2	3	3
4	1	1	2	3	3	4	4
5	1	2	2	3	4	4	5
6	1	2	3	4	4	5	6
7	1	2	3	4	5	6	7
8	1	2	3	5	6	7	8
9	1	2	4	5	6	8	9
10	1	3	4	6	7	8	10
11	2	3	5	6	7	9	10
12/15	2	3	5	7	8	10	11

FIGURE 3.26T PROPORTIONAL PLACEMENT OF COORDINATION LEVELS IN STARTING TERM GROUPS FOR METHOD 2

Such a method is very arbitrary; some of the results for questions having less than seven starting terms had to be used more than once, whilst some of the results for questions having more than seven starting terms could not be used. The performance figures resulting from this method are given in Fig. 3.27TP.

For Method 3, described as (maximum starting term coordination levels; the questions were totalled by grouping according to the maximum number of starting terms. Thus the three-starting-term questions searched at a level of 3 would be totalled with the four-starting-term questions searched at 4, with the five-starting-term questions searched at 5 and so on. A single coordination level is dropped off at a time, working from right to left in the diagram given in Table 3.28T. It can be seen that questions having only a small number of starting terms are soon reduced to a single term search; therefore the results at this level are maintained together with those questions that still have terms that can be dropped off, until all questions are being searched on a single term. Results by this method are given in Fig. 3.29TP.

Coordination Level	Recall Ratio	Precision Ratio
1/7	94.4%	0.9%
2/7	81.0%	2.4%
3/7	58.7%	4.2%
4/7	31.3%	9.9%
5/7	18.7%	15.5%
6/7	8.4%	23.5%
7/7	3.2%	48.1%



3. 27TP TABLE AND PLOT FOR RESULTS BY METHOD 2, PROPORTIONAL COORDINATION LEVELS, FOR INDEX LANGUAGE I. 1. a. (INDEX LANGUAGE I. 6. a DOTTED LINE)

Starting Term Groups	Twelve Coordination Levels											
	Minus 11	Minus 10	Minus 9	Minus 8	Minus 7	Minus 6	Minus 5	Minus 4	Minus 3	Minus 2	Minus 1	Maximum
2/3	1	1	1	1	1	1	1	1	1	1	2	3
4	1	1	1	1	1	1	1	1	1	2	3	4
5	1	1	1	1	1	1	1	1	2	3	4	5
6	1	1	1	1	1	1	1	2	3	4	5	6
7	1	1	1	1	1	1	2	3	4	5	6	7
8	1	1	1	1	1	2	3	4	5	6	7	8
9	1	1	1	1	2	3	4	5	6	7	8	9
10	1	1	1	2	3	4	5	6	7	8	9	10
11	1	1	2	3	4	5	6	7	8	9	10	11
12/15	1	2	3	4	5	6	7	8	9	10	11	12

FIGURE 3.28T GROUPINGS FOR MAXIMUM STARTING TERM COORDINATION LEVELS FOR METHOD 3

All methods discussed so far have results which, at higher coordination levels, are based on increasingly smaller sets of questions. Method 4 overcomes this particular drawback, by totalling all questions at the highest coordination levels that retrieve documents in every question. Known as the method of 'maximum retrieving term coordination levels', all questions are first aligned at the highest coordination level at which, in every question, at least one document is retrieved, irrespective of whether or not it is relevant. This level will vary from question to question, and by referring back to Fig. 3.20T, it can be seen that in the particular conditions of that test, some questions only began retrieving documents when several of their starting terms had been dropped off in coordination. For example, none of the twelve starting-term questions retrieved documents at a coordination level higher than nine. When the search results have been aligned by the coordination level at which each question gives a result, the figures are totalled, the coordination level in each question being relaxed one term at a time, until every question is reduced to a single term search. The results by this method are given in Fig 3.30T, and the curve plotted in Fig. 3.30P. It will be noted that the lower end of the curve terminates at 17% recall and 22% precision; this point has been derived from the individual results

Coordination Level	Recall Ratio	Precision Ratio
Minus 11	94.5%	0.9%
Minus 10	94.2%	1.0%
Minus 9	93.2%	1.2%
Minus 8	91.2%	1.4%
Minus 7	86.3%	1.8%
Minus 6	78.8%	2.4%
Minus 5	67.7%	3.9%
Minus 4	54.2%	6.8%
Minus 3	38.7%	8.7%
Minus 2	25.3%	9.3%
Minus 1	11.0%	14.4%
Maximum	3.1%	47.6%

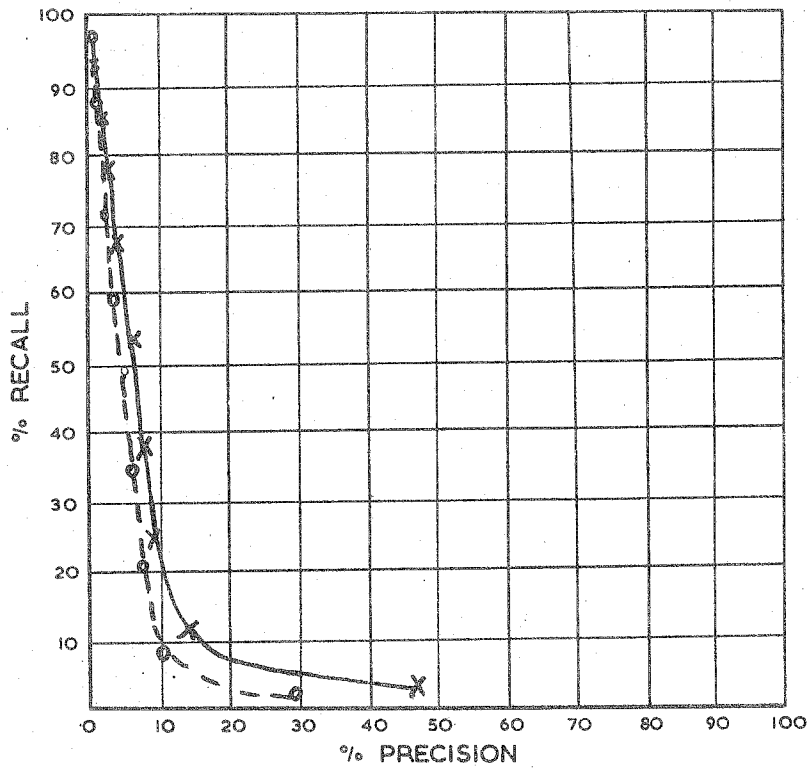


FIGURE 3. 29TP TABLE AND PLOT FOR RESULTS BY METHOD 3, MAXIMUM STARTING TERM COORDINATION LEVELS, FOR INDEX LANGUAGE I. 1, a. (INDEX LANGUAGE I. 6, a. DOTTED LINE)

Coordination Level	Recall Ratio	Precision Ratio
Minus 9	94.9%	0.9%
Minus 8	94.9%	1.0%
Minus 7	94.9%	1.1%
Minus 6	94.7%	1.3%
Minus 5	93.7%	1.6%
Minus 4	88.5%	2.7%
Minus 3	78.3%	4.6%
Minus 2	62.3%	6.0%
Minus 1	40.7%	8.9%
Maximum	16.9%	22.2%

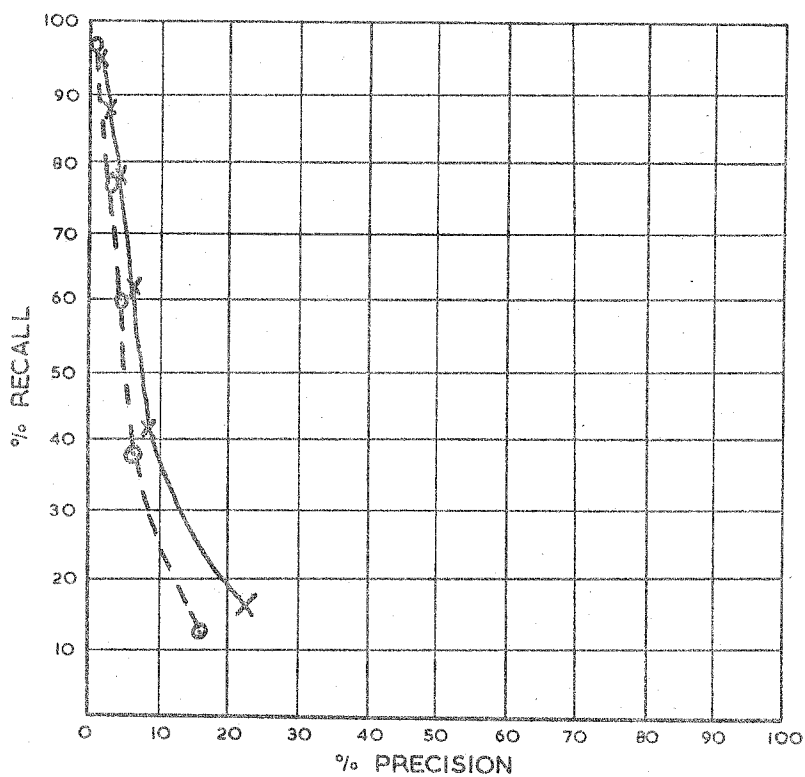


FIGURE 3.30TP TABLE AND PLOT FOR RESULTS BY METHOD 4, MAXIMUM RETRIEVING TERM COORDINATION LEVELS, FOR INDEX LANGUAGE I.1.a. (INDEX LANGUAGE I.6.a DOTTED LINE)

of every one of the 221 questions.

Method 5 differs from all other methods described so far in being based on actual retrieval results obtained in testing. The method was generally known as 'recall levels', because a series of recall ratios is chosen in advance, and the performance results closest to the chosen recall levels are used to obtain the totals, irrespective of the coordination level of the search terms. Ideally this method should be applied to each individual question in a set, with the recall and precision ratios attained by each question being recorded when closest to 5% recall, then 10% recall, and so on. The calculations by Method 5 approximated to this by using the recall levels of the nine retrieving term groups. The recall ratios of these retrieving term groups were arranged by a set of twenty-one recall levels, being 0%, 5%, 10% etc. to 100%, and then the results in figures thus arranged were used to obtain twenty-one sets of recall and precision ratios. Fig. 3.31TP gives the table and plot of results, and the large number of performance points on the plot show a slight scatter through which the performance curve is drawn.

Method 6 was known as 'Document output cutoff method', and was based on quite different principles to those already discussed. To explain this method, it is first necessary to consider the effect of the 'conventional' search cutoff method used in the test. This, as has been explained, was based on the coordination level, which is to say that with, for instance, a six-term question, the search result would be recorded for a coordination of all six terms, then it would be recorded for a coordination of five terms, then for a coordination of four terms and so on. It was this method of search cutoff, with questions having a range of different potential coordination levels, that caused the problem in totalling the results of the whole set of questions, and Method 6, involving a document output cutoff, seemed to overcome this problem.

To apply this method, it was first necessary to obtain a ranked order of documents for every question, and, in our case, this had to be based on the coordination level cutoff results. A method of doing this was developed, but it entailed a considerable amount of effort.

The decision as to which method to use for presentation of the results was not easy to make and has probably involved more discussion, both amongst ourselves and with other people, than any other single aspect of the test. The necessity for the particular series of attempts to total the results was due to the problem created by the coordination level cutoff. It seems reasonable to assume that the final method discussed, the document output cutoff method would be most satisfactory since it eliminated the basic problem of totalling different sets of results but it appeared to involve more effort than could be afforded.

Recall Levels	Recall Ratio	Precision Ratio
100%	94.4%	0.9%
95%	93.7%	1.1%
90%	93.5%	1.2%
85%	82.3%	2.5%
80%	81.4%	2.7%
75%	77.9%	3.1%
70%	67.8%	4.2%
65%	62.2%	5.2%
60%	61.7%	5.4%
55%	54.0%	7.2%
50%	52.3%	7.4%
45%	47.4%	8.2%
40%	42.6%	10.2%
35%	27.6%	15.4%
30%	27.5%	15.3%
25%	22.1%	17.7%
20%	22.1%	17.7%
15%	16.9%	22.2%
10%	16.9%	22.2%
5%	16.9%	22.2%
0%	16.9%	22.2%

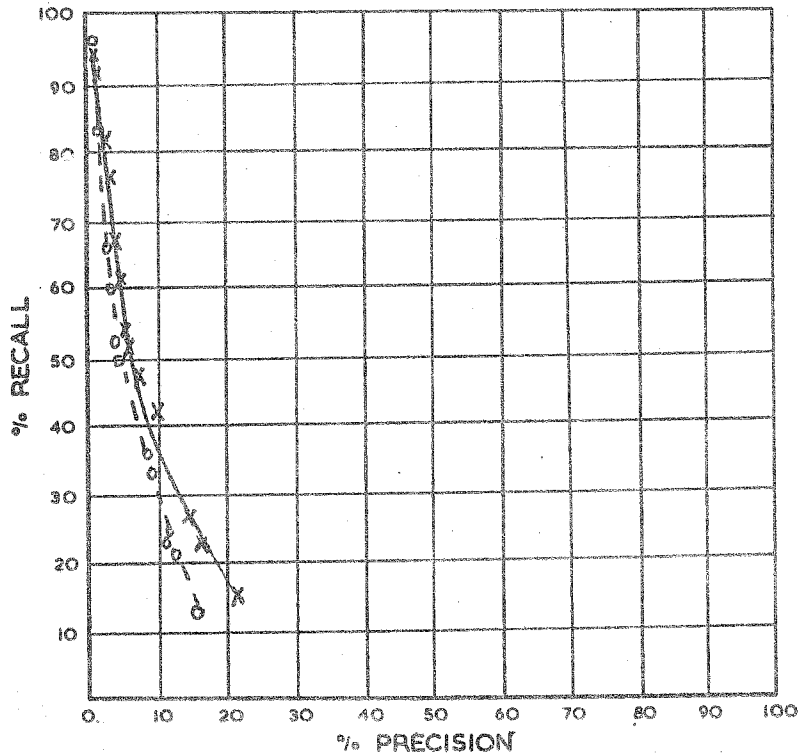


FIGURE 3.31TP

TABLE AND PLOT FOR RESULTS BY METHOD 5, RECALL LEVEL OF RETRIEVING TERM GROUPS, FOR INDEX LANGUAGE I.1. a (INDEX LANGUAGE I.6. a DOTTED LINE)

With all the methods that have been discussed and illustrated, each consistently showed Index Language I.1a to have a seemingly superior performance to Index Language I.6a. Whatever weakness there might have been in any of the methods, in no case were the results sufficiently distorted to mask this change in performance. In this situation, it again seemed most sensible to adopt a method which was relatively simple to apply, and Method 1A, using starting term coordination levels was therefore selected.

After this decision had been taken and the main sets of results had been prepared, a simpler method of obtaining a ranked output was found, and the majority of the results have been recalculated by the document output cutoff method. However, the decision to present the main results by Method 1 was not reversed, so the results obtained by this alternative Method 6 are presented separately in Chapter 5.

In view of the decision to use the starting term coordination level method, it is necessary to mention one further point. Using this method means that average results obtained at high coordination levels are based on an increasingly smaller number of questions in the set due to two reasons - firstly the variation in the number of starting terms, resulting in questions with a small number of starting terms never being capable of contributing results when the coordination level exceeds the number of starting terms. This has already been discussed in Chapter 2, where it was stated that this information was given in each table of test results in column z (see Fig. 2.15). The second reason was the variation in the number of terms that actually retrieve any documents, since the higher coordination levels in some questions demand a match that is too strong for any documents to be retrieved. Data on this point is presented in Column x which gives the total number of questions which actually retrieved any documents. As can be seen in Fig. 2.15, although z decreases at the higher coordination levels, x is smaller than z at all coordination levels of 4 or more. This was the normal experience, since there were usually some questions where the demand for a coordination of four terms would not retrieve a single document.

The generality number  $\frac{1,000(a+c)}{N}$

To return to the matter of the generality number, it is now possible to consider this in more detail. It is known that, in situations where the generality numbers are different, varying performance figures will be obtained, even though the actual operational performance may be similar. In experimental tests such situations exist when the average numbers of documents relevant to the questions differ in two cases of identical file size or, vice versa, where the file sizes are the same but the numbers of relevant documents are different. A third situation is where both the numbers of the relevant documents and the file sizes are different.

As an example, two collections are hypothesised (see Fig. 3.32T), Collection A having 1000 documents and Collection B having 10,000 documents. In both collections there are assumed to be ten relevant documents for a given question, giving a generality number of 10 for Collection A and 1 for Collection B. It is hypothesised that the recall ratio is 50% and that the proportion of non-relevant retrieved to collection size remains the same. The fact that the proportion of non-relevant retrieved remains the same means that the fallout ratio will be 1.0%\*, although the precision ratio changes from 33.3% in Collection A to 4.8% in Collection B, reflecting the decrease in the generality number. A recall/fallout plot would indicate an identical performance, concealing the information that in Collection A a fallout ratio of 1.0% means the retrieval of ten non-relevant documents and in Collection B it means the retrieval of one hundred non-relevant documents. On the other hand a plot of recall/precision would correctly indicate this change.

<u>COLLECTION A</u>	1000 DOCUMENTS			
	Relevant	Non-Relevant		Generality 10
Retrieved	5	10	15	Recall 50%
Not Retrieved	5	980	985	Fallout 1.0%
	10	990	1,000	Precision 33.3%

<u>COLLECTION B</u>	10,000 DOCUMENTS			
	Relevant	Non-Relevant		Generality 1
Retrieved	5	100	105	Recall 50%
Not Retrieved	5	9890	9895	Fallout 1.0%
	10	9990	10,000	Precision 4.8%

FIGURE 3.32T TWO SETS OF PERFORMANCE RESULTS WITH DIFFERENT GENERALITY NUMBERS AND CONSTANT RECALL AND FALLOUT RATIOS.

For a comparison of retrieval performance, it can be argued that the result revealed by the fallout ratios is more useful, since the change in precision ratio is solely due to the change in the environmental factor of the generality number. However, we have earlier stated our intention to present the main body of results with recall/precision plots, on the ground that these, in general, make a more useful and comprehensible

---

\*This is correct to one decimal place; the actual figures are, respectively, 1.0101%, recurring and 1.001001% recurring.

presentation of performance. It is therefore necessary to make adjustments to the precision ratios in certain situations (which have been considered in Chapter 2) where sets of varying generality have to be compared. This is reasonably straightforward and is obtained by the following equation:-

$$P_A \text{ (Adjusted Precision Ratio)} = \frac{R_1 \times G_2}{(R_1 \times G_2) + F_1(1000 - G_2)}$$

- where  $R_1$  = Recall ratio obtained for a given system, in a situation of a known generality number
- $F_1$  = Fallout ratio obtained for the given system, in a situation of a known generality number
- $G_2$  = Generality number to which it is desired to alter the results, to obtain the adjusted precision

Thus two sets of performance figures obtained with systems of differing generality can be compared by adjusting the precision ratio of one case, so that it is based on the generality number of the other. If the example in Fig. 3.32T were to be corrected, and if it were decided to alter the result of Collection A to fit the generality of Collection B, then, from the equation given above,

$$P_A = \frac{.50 \times 1}{(.50 \times 1) + .01(1000 - 1)} = \frac{.50}{.50 + 9.99} = .048$$

The answer, expressed as a percentage is 4.8% and this result is clearly correct, with both cases now having an identical recall ratio, fallout ratio and precision ratio,

This however, is a simplified example, and in practice the matter is complicated by what at present seems to be the most difficult problem in performance comparison, namely the determination of the correct N. (the size of the collection). To consider this, an actual result is taken from a particular set of 42 questions that were searched on collections A and B where N equals 200 and 1400 documents respectively, the documents in collection A being a subset of the documents in collection B. The details are given in Fig. 3.33T, with the two sets of performance figures obtained in exactly the same conditions. While the precision ratio for collection A has increased with the increased generality number, yet there is also a significant difference in the fallout ratio.

SYSTEMS DATA

	Collection A	Collection B
No. of documents	200	1400
No. of questions	42	42
Total No. of relevant documents	198	198
Generality Number	23.6	3.4

PERFORMANCE AT COORDINATION LEVEL OF 3

	Collection A	Collection B
Relevant retrieved	132	132
Non-relevant retrieved	761	3,984
Recall Ratio	66.7%	66.7%
Precision Ratio	14.8%	3.2%
Fallout Ratio	9.278%	6.798%

FIGURE 3.33T SYSTEMS AND PERFORMANCE DATA FOR COMPARISON OF GENERALITY NUMBERS.

If the fallout in both collections were exactly the same, this would mean that the ratio of the change of the number of non-relevant retrieved (b) would be the same as the ratio of the change of the total non-relevant (b + d) i.e.

$$\frac{b(\text{Collection B})}{b(\text{Collection A})} = \frac{(b + d)(\text{Collection B})}{(b + d)(\text{Collection A})}$$

$$\frac{b(\text{Collection B})/b(\text{Collection A})}{(b + d)(\text{Collection B})/(b + d)\text{Collection A}} = 1$$

Bearing in mind that these figures represent the sum of a series of searches for 42 questions having 198 relevant documents, the result from Fig. 3.33T is, in fact,

$$\frac{\frac{3984}{761}}{(42 \times 1400) - 198} = \frac{5.2352}{7.1448} = 0.7327$$

It is therefore shown that b(non-relevant retrieved) has increased by a factor of 5.2352 while the total number of non-relevant documents (b + d) has increased by a factor of 7.1448. Proof of the accuracy of this can be shown by assuming that collection B had retrieved 7.1448 times as many non-relevant documents as collection A in which case it would have retrieved 761 x 7.1448 = 5437 documents, as against the actual total of 3,984 documents.

The fallout ratio would now be  $\frac{5437}{58602} = 9.278\%$ .

This fallout is now identical with that of collection A in Fig. 3.33T; it should be noted however, that these figures would result in the precision ratio falling from 3.2% to 2.4%.

One has various options as to how to correct the precision ratio according to generality; it is possible to convert A to B (i.e. 23.6 to 3.4), B to A (i.e. 3.4 to 23.6) or to take a figure intermediate between A and B, such as 11. The effect of these three possible changes would result in the following figures:-

	Uncorrected Precision Ratio	Adjusted Precision Ratio			Fallout Ratio
		G = 3.4	G = 23.6	G = 11	
Collection A	14.8%	2.4%	14.8%	7.3%	9.278%
Collection B	3.2%	3.2%	19.0%	9.7%	6.798%

Whereas uncorrected precision ratio shows A to be superior, all adjusted precision ratios show B to be superior. To discover what is the factor which, in terms of the two collections, causes the difference in performance, Collection A will be taken as giving the expected result, and we will investigate the reasons why B should show the improved performance after precision ratio has been adjusted.

The problem is why, with collection B, fewer non-relevant documents are retrieved than expected. This can be explained by saying that there is more diversification in the indexing terms (and, therefore, presumably of the subject) of some of the documents in the larger file in relation to the search terms of the questions. The 42 questions in the test were all specifically on aerodynamics, as were all the 200 documents in collection A. However, it is known that 257 of the documents in collection B were included in relation to questions dealing with the theory of aircraft structures; if it is assumed that these were never retrieved by any of the 42 questions on aerodynamics, then this would reduce N for collection B from 1400 to 1143, which is shown as B<sub>1</sub> in Fig. 3.34T, where the new generality number and fallout ratio are given. The fallout, at 8.333%, is now closer to, but still does not reach, the level for collection A.

It is therefore clear that if the performances are to be equated, it is necessary to hypothesise that in collection B there is a further subset of documents which are not retrieved by the questions. This number can be found by calculating the size of a hypothetical collection, B<sub>2</sub>, which would result in an identical performance as collection A; the size of this

collection, B<sub>2</sub>, is calculated to be 1027 documents, which means that a further 116 documents must be deleted from collection B<sub>1</sub>. As will be seen in Fig. 3.34T, the collections are now equated with the fallout ratio being, in both cases, 9.278%.

Collection	No. of Documents	Generality Number	Fallout Ratio	Collection	No. of Documents	Generality Number	Fallout Ratio
A	200	23.6	9.278%	B	1400	3.4	6.798%
				B <sub>1</sub>	1143	4.1	8.333%
A <sub>1</sub>	271	17.4	6.798%	B <sub>2</sub>	1027	4.6	9.278%

FIGURE 3.34T CORRECTED COLLECTION SIZES TO FIT GENERALITY NUMBERS.

If instead of correcting collection B to collection A, the reverse step had been taken, then it can be seen that it would have meant adding 71 documents to collection A making A<sub>1</sub>, which would then have a fallout of 6.798%, the same as the original collection B.

As a result of doing this, the precision ratio of collection A, can now be converted by the equation given earlier, and, since recall, fallout and generality are equal, the adjusted precision ratio must be 3.2% as for collection B.

While the above may seem to be somewhat involved, it is, in fact, a simplification of the real situation in that 42 questions have been taken as a block. A more detailed analysis would require that each question should be treated separately. Then, again, the analysis has been done in a single fixed situation, namely a certain index language at a certain level of coordination, and clearly it could be repeated over many hundreds of situations of a similar type. However the implications of such analysis are far-reaching, going beyond the scope of this chapter, so they will be considered later in this report.

In addition to explaining the performance measures adopted in this report, this chapter has also attempted to cover, albeit in a non-exhaustive manner, the main considerations regarding their use and effect. For ourselves, we feel that it is foolish, at the present stage of development, to be dogmatic on this subject. Wherever it has been necessary to make a choice between different methods, in most cases the decision has been taken for reasons which could be considered peculiar to this project. Other

experimenters may well find that different measures better suit their purpose; hopefully, in this survey the relationship between different measures has now been established, and so long as complete sets of figures are given in reporting test results, there should be no serious difficulty in converting from one set of measures to another. Ultimately, one assumes that something approaching general agreement will be reached on the measures to be used. All that we would claim is that the measures used in this report appear to be as good as any others so far proposed.

## CHAPTER 4

### Main Test Results

Efficiency is something which cannot be achieved without effort expended in the appropriate ways, and measurement is one of the ways towards improved efficiency. It is difficult to see how efficiency can be improved without some basis in measurement.

L.T. Wilkins: Social Deviance, page 8.

In the two previous chapters have been outlined the environment in which the tests were carried out and the measures which are used for presenting the results. For those who did not wish to work through these two chapters, a brief summary is now given to assist in the interpretation of the test results presented in this chapter. The simplest method of doing this is to illustrate the various points with an example of the tables which present the test results, as in Figure 4.001, to which reference should be made.

(1) There are four main groups of index languages and these are: identified by roman numerals:

- I Single term index languages (eight languages tested)
- II Simple concept index languages (fifteen languages tested)
- III Controlled term index languages (six languages tested)
- IV Abstract and title searches (four languages tested)

With each index language there are a number of different recall devices. These are identified by arabic numerals, and the full range of these index languages is set out in Figures 2.5, 2.6 and 2.7.

The lower case letter identifies the precision device which is being used. The basic device is coordination, shown by a. The other precision devices are fully explained in Figure 2.8.

(2) At the indexing stage, each term was given a rating in relation to its importance in each particular document, and this permits tests to be made with different levels of exhaustivity of indexing. This is to say that searches can be made on the full indexing, which averaged 31 terms for each document, and this is shown as Exhaustivity 3. Alternatively it can be done on a restricted set of terms, where the average was 25 terms per document; this is shown as Exhaustivity 2. Finally searches can be made on the least exhaustive indexing, where the average was 13 terms per document, and this is shown as Exhaustivity 1.

(3) Various search rules that were used are explained in Figure 2.9. The basic search rule permitted the combination of any terms and is shown by A. Other search rules are explained in the appropriate section of the tables.

- 1 Index Language L 5.a (S. T. Synonyms, Quasi-synonyms. Coordination)
- 2 Exhaustivity of Indexing 3
- 3 Search Rule A
- 4 Document Relevance 1 - 4
- 5 Number of Documents in Collection 1,400
- 6 Number of Questions 221 (Subset 3)
- 7 Number of Relevant Documents 1,590
- 8 Generality Number 5.1
- 9
- 10
- 11
- 12

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x y z		
	Rel.	Non-rel.						
1	1,548	(-)	97.4%	(-)	(-)	221	0	221
2	1,406	114,265*	88.4%	1.2%*	37.099%*	221	44*	221
3	1,121	42,364*	70.5%	2.6%*	13.755%*	218	109*	220
4	802	16,191*	50.4%	4.7%*	5.257%*	204	142*	212
5	475	8,164*	29.9%	5.5%*	2.651%*	164	177*	197
6	265	3,013*	16.7%	8.1%*	0.378%*	114	161*	164
7	131	910	8.2%	12.6%	0.296%	79	140	140
8	50	266	3.1%	15.8%	0.086%	44	105	105
9	19	56	1.2%	25.3%	0.018%	20	78	78
10	2	12	0.1%	14.3%	0.004%	6	52	52
11	0	0				0	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 4.001 SAMPLE TEST RESULTS SHEET

- (4) The relevance of the documents to the questions was assessed by the questioners in four grades of relevance (see Vol. I, p.21). Most of the test results are given for documents of all grades of relevance (shown as 1-4), except in Section 4.6 which specifically deals with the effect of relevance.
- (5) The main test collection had 1,400 documents, but two other document sets were used. These were of 200 and 350 documents, a characteristic of these smaller sets being that all the documents were concerned with aerodynamics, whereas the main set contained some 300 documents on theory of aircraft structures.
- (6) The largest set of questions in the test had 221 questions. Most of the results are based on a subset of 42 questions, all of which were concerned with aerodynamics. Another subset had 35 questions, the characteristic of which was that each question had seven starting terms. Other sets were used in special cases; full details are given in Figure 2.12, and also in the appropriate section of the tables.
- (7) The number of relevant documents will vary with the document set, the question set and the relevance grade. This number must be known for calculating the recall ratio.
- (8) The generality number is a function of the number of relevant documents and the number of documents in the test collection. With Question Subset 2 (for which there are 198 relevant documents), when the search is made on the 1400 document collection, the generality number is 3.4. When searched on the 200 document collection, the generality number is 23.6. The effect of the increase in generality number is to bring about an apparent improvement in the performance figures. The matter of generality is fully discussed in Chapter 3.
- (9) All the test results given in this chapter were based on searches where the coordination level was progressively decreased from the maximum down to a single term. The maximum number of single terms in any question was 15, while the lowest number of terms was 2.
- (10) In most situations, for the reason stated in the previous paragraph, the number of questions which can be searched at a given coordination level will be limited by the number of questions having that number of starting terms. This information is given in column z, which shows, for example, that at a coordination level of 6, there are 164 questions which, having six or more starting terms, can be searched at this level.

In certain searches, the number of questions actually searched at a given coordination level was less than the theoretical maximum possible. This was because of the large clerical effort required and the number of questions actually searched is given in column y.

It will be seen from Figure 4.001 that, at a coordination level of 6, only 161 of the possible 164 questions were searched, and from this stage there is an increasing disparity between the figures in column y and z.

Column x shows the number of questions that were able, at any given coordination level, to retrieve any documents, whether relevant or non-relevant. It will be noted, for instance, that it was not until the coordination level had dropped to ten terms was it possible to retrieve a single document; at this level, as can be seen from the figures in columns x and z, of the 52 questions having ten or more starting terms, 6 questions retrieved documents.

(11) In the columns of documents retrieved are shown the total numbers of relevant and non-relevant documents retrieved at the various coordination levels. These figures have been obtained by summing the results for each individual question in the question set.

As mentioned in the previous section, in some cases the searches were not completed at the lower coordination levels. The result is that the figures for non-relevant documents retrieved are estimated by a method described on page 28. All such estimated figures are indicated by an asterisk. However, it should be noted that the figures for the retrieval of relevant documents are always correct.

(12) The actual performance measures presented in the tables are recall ratio, precision ratio and fallout ratio. These are derived from the following:-

- a. Relevant documents retrieved
- b. Non-relevant documents retrieved
- c. Relevant documents not retrieved
- d. Non-relevant documents not retrieved

Recall ratio is  $\frac{100a}{a + c}$ , that is relevant documents retrieved over the total relevant. All such figures in the tables are correct, but for precision ratio ( $\frac{100a}{a + b}$ ) and fallout ratio ( $\frac{100b}{b + d}$ ), it was in some cases necessary to use the estimated figures discussed in the previous section. Where this has been done, an asterisk is placed against the figure in the table.

#### TABLE OF RESULTS

The tables of results are presented in nine main sections. Details are given before each section of tables, but the following is a brief resumé.

In the first section the results are given for the single term index languages with the largest set of documents and questions. Because of the doubts regarding the totalling method used, the 35 seven-term questions set is given for comparison. This is followed by the 42 aerodynamic questions on the 1400 document collection, with a final table for this set of questions with the 200 document collection. The purpose of this group of tables is first to justify the totalling procedure, then the reduced set of questions and finally the reduced set of documents.

Section 2 gives the results for the 42-question and 200-document sets for the eight single term languages. The tests using precision devices of interfixing and partitioning are given in Section 3, and are followed in Section 4 by comparative results at different levels of exhaustivity. The effect of search rules is shown in Section 5 and the different levels of relevance are tested in Section 6. Results for fifteen simple concepts index languages are presented in Section 7, while Section 8, which deals with the six controlled term index languages, includes results for recall devices, precision devices and search rules. In Section 9 are presented the results of searches on the titles and abstracts,

The tables and plots are numbered according to the above nine sections, 'T' indicating a table and 'P' indicating a plot. In all cases unless otherwise stated, the plots are of recall and precision ratios.

## Section 1 Introductory Tables

The first set of results is based on 221 questions (question subset 3) searched on the full collection of 1400 documents. In Tables 4.100T to 4.104T, five of the languages investigated on the single term indexing are presented, namely languages I.1.a, I.2.a, I.3.a, I.5.a, and I.6.a, showing the effect of recall devices. The results for these five languages are presented on a single plot in Figure 4.105P, with a performance curve for each of the languages plotted.

It was very difficult to find a satisfactory method of totalling the results with these 221 questions, because of the large variation in the number of starting terms. This problem was fully discussed in Chapter 3 (pages 51 - 71), but in order to validate the selected method, a subset of 35 questions was selected, the characteristic of which was that each question had seven starting terms, and in this respect was an average set. The results for searches for the same five index languages are presented in Figures 4.110T - 4.114T. The questions are again searched on the 1400 document collection, and a single plot of the five curves is given in Fig. 4.115P.

A further subset of questions, 42 in number (subset 2), is used for the results given in Figs. 4.120T to 4.124T. The same five languages are used, and the 1400 document collection is searched. These questions pose the same problems in totalling as did the 221 questions, since the 42 questions have varying numbers of starting terms, ranging from three to twelve. A single plot of the five languages is given in Fig. 4.125P.

In these three sets of questions, a progression may be observed from the largest set of questions (221) to a smaller set of 35 questions specially selected to minimise the totalling problem and to another small set of 42 questions that has the same problem of totalling, due to the variation in number of starting terms, as the collection of 221 questions. The differences in question sample size may be expected to affect any direct comparison of the three sets of questions, in addition to the totalling method problem. The effect of this can be seen in Fig. 4.130P where the natural language results (Language I.1.a) are compared for the three subsets of questions, the three curves being based on the results in Figs. 4.100T, 4.110T, and 4.120T. Although the comparison is accurate in terms of recall and precision as calculated, the comparison of three different sets of questions brings in a new variation, namely the generality number (G). For the 221 questions G is 5.1, for the 35 questions G is 5.9, and for the 42 questions G is 3.4. The need to allow for this difference in generality has been discussed, and Fig. 4.131P is a graph that allows for this by use of a plot of recall and fallout.

It can now be argued that the performance results based on the smaller sets of questions give valid results. First, there is a close similarity in the relative differences obtained when five recall languages are compared, whichever set of questions is used. This may be seen by comparing Figs. 4.105P, 4.115P and 4.125P where the relative differences between the five languages are very similar. A further comparison is made on a recall fallout plot, in Fig. 4.131P, where generality is allowed for, and the language is held constant as type I.1.a. Some small variation is to be expected when the question sets are altered in size, and when a universal totalling method is to be used, but the subset of 42 questions, on which many of the later results are based, is seen to be representative of both the larger set of 221 questions and of the set with the chosen characteristic of each question having the same number of starting terms.

In the question sets shown so far, the collection size has remained constant at 1400 documents. Most of the later results have been obtained on the 200 document collection (collection subset 1), and the validity of results based on the smaller collection size will be considered next.

Table 4.140T gives the results for a search on language I.1.a made with the 42 questions searched on the 200 document collection. The results from the table are plotted as a performance curve in Fig. 4.140P. Also shown on this plot are the results from Fig. 4.120T, which are based on the same question set and language but searched on the 1400 collection. It would be expected that with a recall/precision plot, the increased generality number would result in a better performance for the searches on the 200 collection as compared to those on the 1400 collection. This is, in fact, the case, for while the recall at each coordination level is seen to be identical, the expected large increase in precision is seen when the 200 collection is searched.

The effect of the change in generality on the precision ratio, as discussed in Chapter 3, is allowed for in Fig. 4.141P, which is a plot of the two curves using recall and fallout ratios, and Fig. 4.142P which plots the two curves on a recall/precision plot with generality adjusted to a constant of 23.6, this number being the generality of the situation in the 200 collection. The result for the 1400 collection is now no longer inferior to the 200 collection, and in fact the situation is reversed. The reason why the 200 collection now has a somewhat worse performance has been investigated in Chapter 3, where it has been shown that the cause of the difference can be adequately explained and allowed for; this was, in fact, done and the result was shown in Figure 3.34T.

The purpose of the introductory results presented is to demonstrate that test results based on a relatively small collection and set of questions do give valid results. The variations observed between the three different

sets of questions and the two collections can all be adequately explained, and can be allowed for by methods of adjustment that have been developed.

LIST OF FIGURES

	Index Language	No. of Questions	Question Subset	Document Collection	Plots
4.100T	I.1.a	221	3	1400	
4.101T	I.2.a	221	3	1400	
4.102T	I.3.a	221	3	1400	
4.103T	I.5.a	221	3	1400	
4.104T	I.6.a	221	3	1400	
4.105P	I.1.a	221	3	1400	Plot 4.100T, 4.101T, 4.102T, 4.103T, 4.104T.
	I.2.a				
	I.3.a				
	I.5.a				
	I.6.a				
4.110T	I.1.a	35	1	1400	
4.111T	I.2.a	35	1	1400	
4.112T	I.3.a	35	1	1400	
4.113T	I.5.a	35	1	1400	
4.114T	I.6.a	35	1	1400	
4.115P	I.1.a	35	1	1400	Plot 4.110T, 4.111T, 4.112T, 4.113T, 4.114T.
	I.2.a				
	I.3.a				
	I.5.a				
	I.6.a				
4.120T	I.1.a	42	2	1400	
4.121T	I.2.a	42	2	1400	
4.122T	I.3.a	42	2	1400	
4.123T	I.5.a	42	2	1400	
4.124T	I.6.a	42	2	1400	
4.125P	I.1.a	42	2	1400	Plot 4.120T, 4.121T, 4.122T, 4.123T, 4.124T.
	I.2.a				
	I.3.a				
	I.5.a				
	I.6.a				
4.130P	I.1.a	221	3	1400	Plot 4.100T, 4.110T, 4.120T,
		35	1		
		42	2		
4.131P	(same as 4.130P, but recall/fallout)				

	Index Language	No. of Questions	Question Subset	Document Collection	Plots
4.140T	I.1.a	42	2	200	
4.141P	I.1.a	42	2	1400 200	Plot 4.120T, 4.140T
4.142P	(same as 4.141P, but recall/adjusted precision)				

FIGURE 4.100T

Index Language I.1.a. (S. T. Natural language. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1-4  
 Number of Documents in Collection 1,400  
 Number of Questions 221 (Subset 3)  
 Number of Relevant Documents 1,590  
 Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,510	159,122	95.0%	0.9%	51.696%	221	221	221
2	1,283	58,122	80.7%	2.2%	18.883%	221	221	221
3	946	21,933	59.5%	4.1%	7.125%	215	220	220
4	606	7,359	38.1%	7.6%	2.390%	187	212	212
5	314	2,380	19.7%	11.6%	0.773%	131	197	197
6	154	699	9.7%	18.0%	0.227%	86	164	164
7	74	216	4.7%	25.5%	0.070%	50	140	140
8	22	43	1.4%	33.8%	0.014%	18	105	105
9	8	5	0.5%	61.5%	0.002%	8	78	78
10	1	0	0.1%	100.0%	0.000%	1	52	52
11	0	0				0	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 4.101T

Index Language I.2.a (S. T. Synonyms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1-4  
 Number of Documents in Collection 1,400  
 Number of Questions 221 (Subset 3)  
 Number of Relevant Documents 1,590  
 Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,514	(-)	95.2%	(-)	(-)	221	0	221
2	1,313	59,734*	82.6%	2.2%*	19.406%*	221	44*	221
3	981	23,654*	61.7%	4.0%*	7.680%*	216	109*	220
4	644	8,850*	40.5%	6.8%*	2.873%*	192	142*	212
5	355	2,946*	22.3%	10.4%*	0.957%*	139	177*	197
6	169	928*	10.6%	15.4%*	0.301%*	92	161*	164
7	80	254	5.0%	24.0%	0.083%	55	140	140
8	24	59	1.5%	28.9%	0.019%	23	105	105
9	8	8	0.5%	50.0%	0.003%	8	78	78
10	1	0	0.1%	100.0%	0.000%	1	52	52
11	0	0				0	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 4.102T

Index Language I.3.a (S.T. Word forms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 221 (Subset 3)  
 Number of Relevant Documents 1,590  
 Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,533	(-)	96.4%	(-)	(-)	221	0	221
2	1,338	62,765*	84.2%	2.1%*	20.378%*	221	44*	221
3	1,017	24,726*	64.0%	3.9%*	8.028%*	217	109*	220
4	677	9,565*	42.6%	6.6%*	2.530%*	192	142*	212
5	374	3,084*	23.5%	10.8%*	1.001%*	139	177*	197
6	192	1,112*	12.1%	14.8%*	0.361%*	99	164*	161
7	96	333	6.0%	22.4%	0.108%	64	140	140
8	34	87	2.1%	28.1%	0.028%	28	105	105
9	13	15	0.8%	46.4%	0.005%	17	78	78
10	2	0	0.1%	100.0%	0.000%	2	52	52
11	0	0				0	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 4.103T

Index Language I.5.a (S.T. Synonyms, Quasi-synonyms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 221 (Subset 3)  
 Number of Relevant Documents 1,590  
 Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,548	(-)	97.4%	(-)	(-)	221	0	221
2	1,406	114,265*	88.4%	1.2%*	37.099%*	221	44*	221
3	1,121	42,364*	70.5%	2.6%*	13.755%*	218	109*	220
4	802	16,191*	50.4%	4.7%*	5.257%*	204	142*	212
5	475	8,164*	29.9%	5.5%*	2.651%*	164	177*	197
6	265	3,013*	16.7%	8.1%*	0.378%*	114	161*	164
7	131	910	8.2%	12.6%	0.296%	79	140	140
8	50	266	3.1%	15.8%	0.086%	44	105	105
9	19	56	1.2%	25.3%	0.018%	20	78	78
10	2	12	0.1%	14.3%	0.004%	6	52	52
11	0	0				0	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 4.104T

Index Language I.6.a (S. T. Synonyms, Quasi-synonyms, Word forms, Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 221 (Subset 3)  
 Number of Relevant Documents 1,590  
 Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,557	(-)	97.9%	(-)	(-)	221	0	221
2	1,430	116,374*	89.9%	1.2%*	37.783%*	231	44*	221
3	1,165	45,101*	73.3%	2.5%*	14.643%*	218	109*	220
4	848	18,373*	53.3%	4.4%*	5.965%*	206	142*	212
5	503	8,895*	31.6%	5.4%*	2.883%*	169	177*	197
6	295	3,874*	18.6%	7.1%*	1.257%*	119	161*	164
7	161	1,136	10.1%	12.4%	0.369%	83	140	140
8	72	344	4.5%	17.3%	0.112%	54	105	105
9	24	82	1.5%	22.6%	0.027%	25	78	78
10	6	18	0.4%	25.0%	0.006%	12	52	52
11	0	1	0.0%	0.0%	0.0003%	1	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

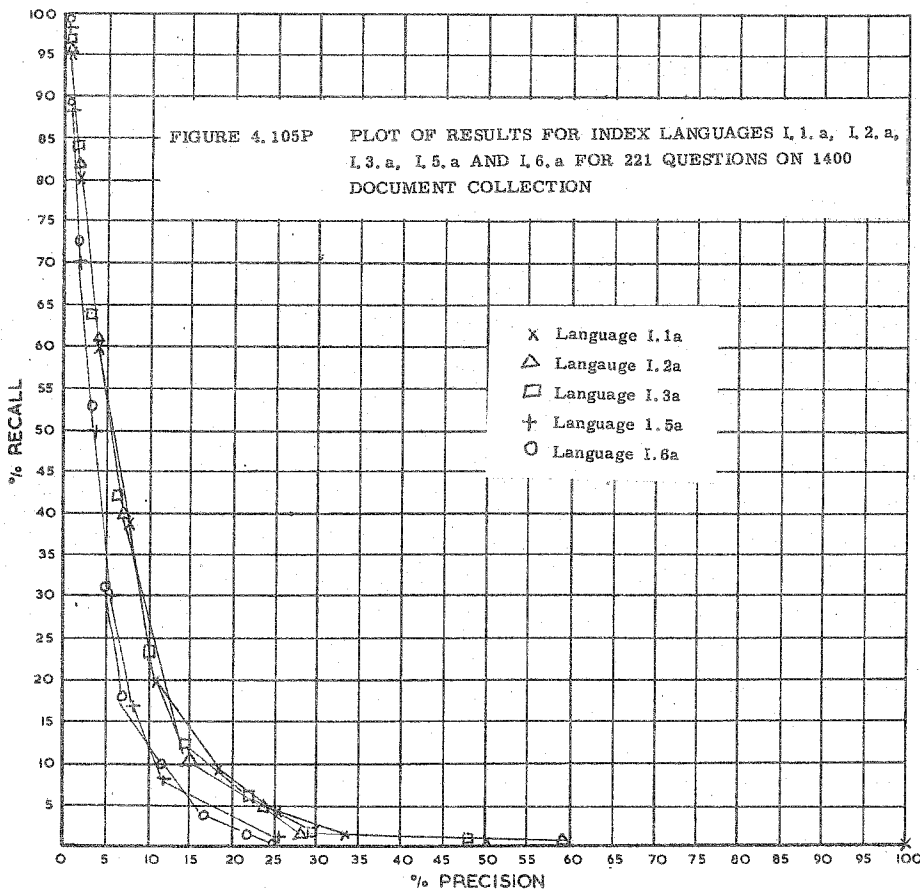


FIGURE 4.110T

Index Language I.1.a (S. T. Natural language. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 35 (Subset 1)  
 Number of Relevant Documents 287  
 Generality Number 5.9

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	258	23,681	93.4%	1.1%	48.329%	35	35	35
2	221	8,086	77.0%	2.7%	16.500%	35	35	35
3	157	2,865	54.7%	5.2%	5.821%	34	35	35
4	94	600	32.8%	13.5%	1.232%	28	35	35
5	46	147	16.4%	23.8%	0.281%	18	35	35
6	22	37	7.7%	37.2%	0.076%	7	35	35
7	8	8	2.8%	50.0%	0.016%	3	35	35

FIGURE 4.111T

Index Language I.2.a (S. T. Synonyms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 35 (Subset 1)  
 Number of Relevant Documents 287  
 Generality Number 5.9

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	269	(-)	93.8%	(-)	(-)	35	0	35
2	226	8,514*	76.8%	2.6%*	17.355%*	35	23*	35
3	164	3,026	57.2%	5.1%	6.212%	34	35	35
4	99	650	34.5%	13.2%	1.334%	28	35	35
5	51	151	17.8%	25.2%	0.310%	19	35	35
6	23	39	8.0%	37.1%	0.080%	6	35	35
7	8	8	2.8%	50.0%	0.016%	3	35	35

FIGURE 4.112T

Index Language I.3.a (S. T. Word forms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 35 (Subset 1)  
 Number of Relevant Documents 287  
 Generality Number 5.9

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	270	(-)	94.1%	(-)	(-)	35	0	35
2	227	9,174*	79.1%	2.4%*	18,518%*	35	23*	35
3	170	3,408	59.3%	4.7%	6.696%	35	35	35
4	103	768	35.9%	11.8%	1.577%	30	35	35
5	55	184	19.2%	23.0%	0.378%	21	35	35
6	25	44	8.8%	36.2%	0.090%	12	35	35
7	8	8	2.8%	50.0%	0.016%	3	35	35

FIGURE 4.113T

Index Language I 5.a (S. T. Synonyms, Quasi-synonyms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 35 (subset 1)  
 Number of Relevant Documents 287  
 Generality Number 5.9

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	280	(-)	97.6%	(-)	(-)	35	0	35
2	253	17,130*	88.2%	1.5%*	34.959%*	35	23*	35
3	194	7,472	67.6%	2.5%	15.339%	35	35	35
4	125	2,086	43.6%	5.6%	4.282%	34	35	35
5	65	463	22.7%	12.3%	0.950%	30	35	35
6	35	88	12.2%	28.4%	0.181%	18	35	35
7	11	18	3.9%	38.0%	0.037%	5	35	35

FIGURE 4.114T

Index Language I 6.a (S. T. Synonyms, Quasi-synonyms, Word forms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 35 (Subset 1)  
 Number of Relevant Documents 287  
 Generality Number 5.9

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	280	(-)	97.6%	(-)	(-)	35	0	35
2	255	19,264*	88.9%	1.3%*	39.314%*	35	23*	35
3	202	6,070	70.4%	2.4%	16.568%	35	35	35
4	130	2,426	45.3%	5.1%	4.980%	35	35	35
5	72	571	25.1%	11.2%	1.172%	30	35	35
6	38	109	13.3%	25.8%	0.224%	18	35	35
7	13	19	4.6%	40.6%	0.039%	6	35	35

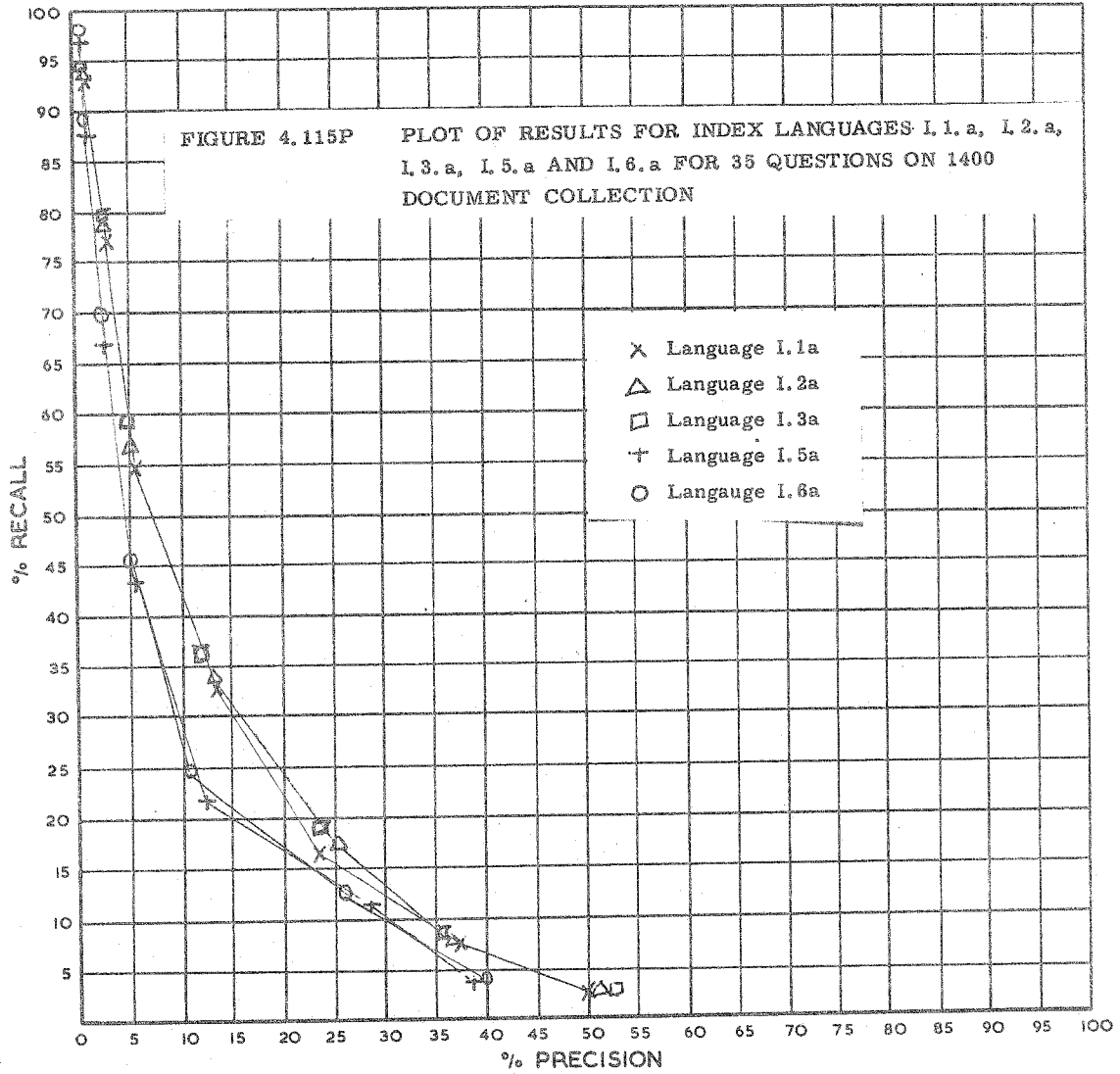


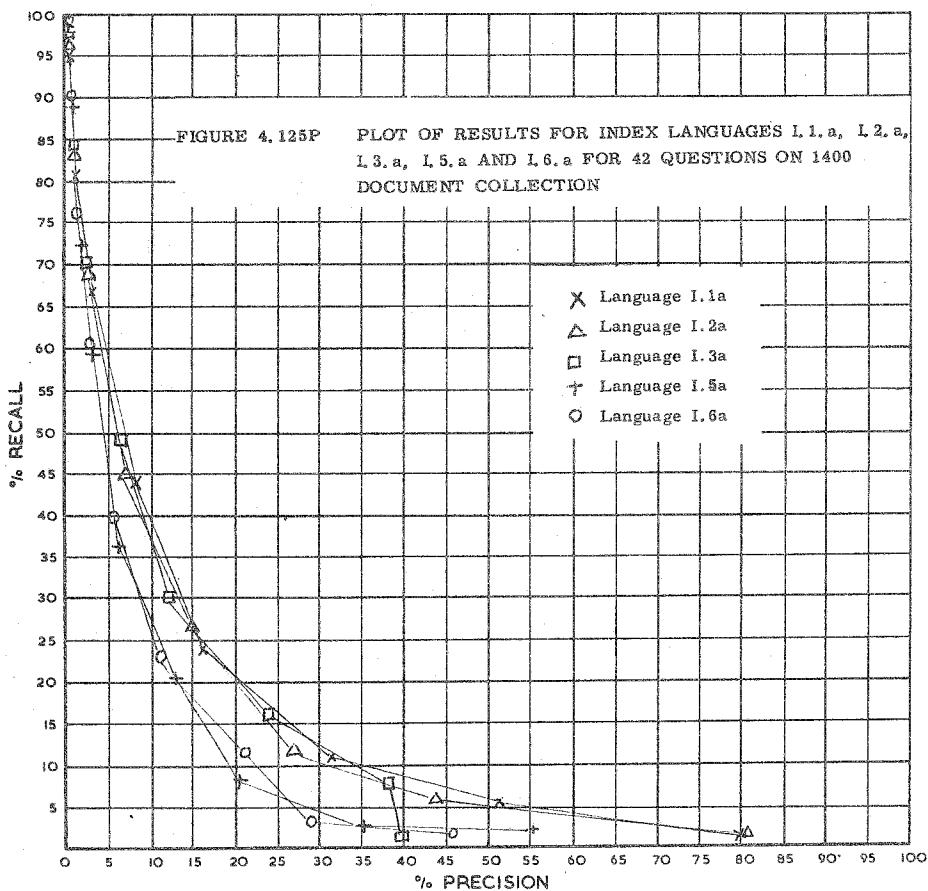




FIGURE 4.124T

Index Language I.6.a (S. T. Synonyms, Quasi-synonyms, Word forms, Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 3.4

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	195	(-)	98.5%	(-)	(-)	42	0	42
2	179	13,826*	90.4%	1.3%*	23.120%*	42	6*	42
3	151	10,506	76.3%	1.4%	17.928%	42	42	42
4	120	3,875	60.6%	3.0%	6.612%	40	41	41
5	79	1,299	39.9%	5.7%	2.217%	34	39	39
6	47	362	23.7%	11.5%	0.618%	24	33	33
7	27	91	13.6%	22.9%	0.155%	16	27	27
8	7	17	3.5%	29.2%	0.029%	7	18	18
9	5	6	2.5%	45.5%	0.010%	5	11	11
10	0	3	0.0%	0.0%	0.005%	2	7	7
11	0	0				0	3	3
12	0	0				0	1	1
13								
14								
15								



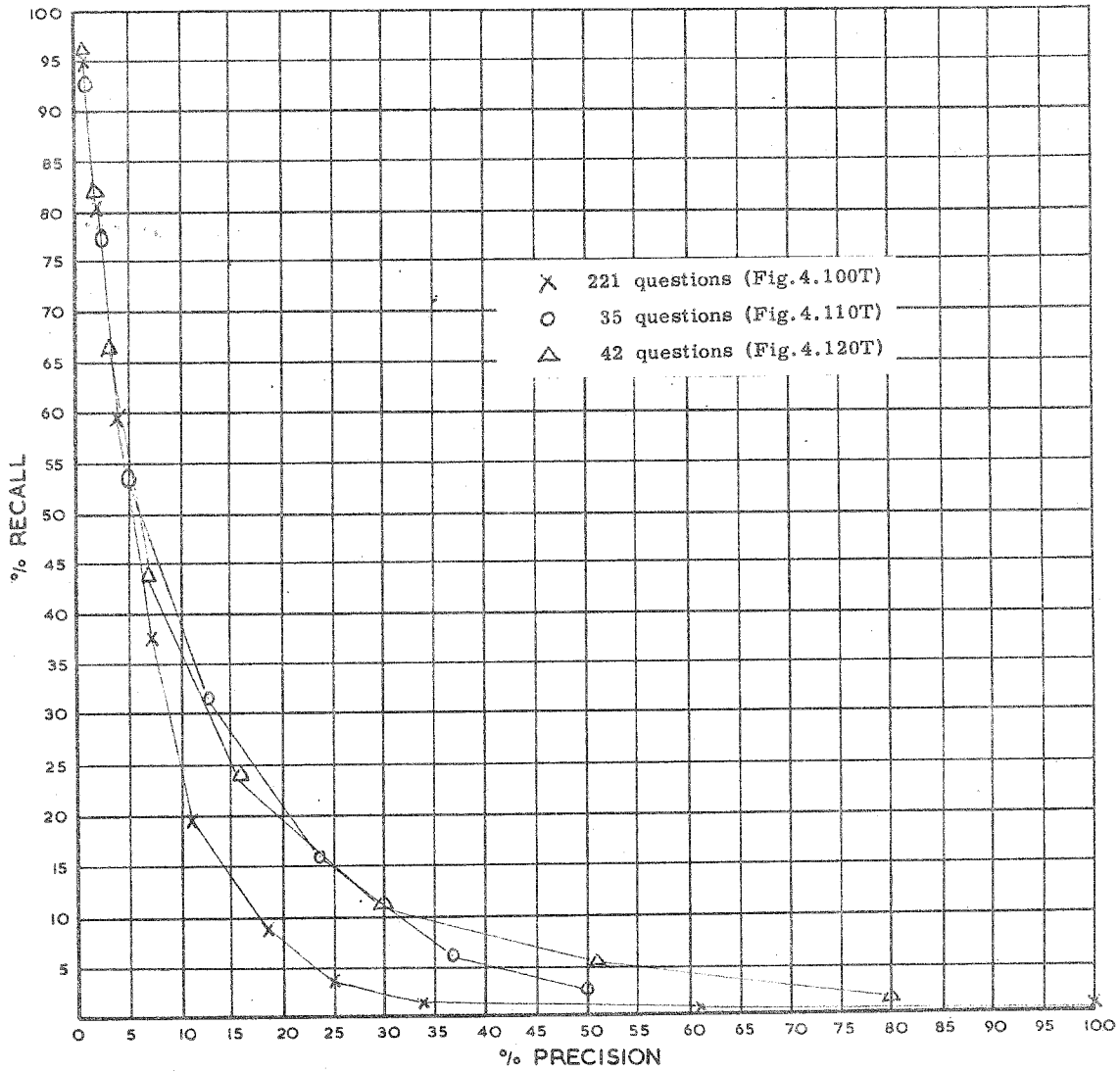


FIGURE 4.130P PLOT OF RESULTS FOR INDEX LANGUAGE I.1. a FOR 221 QUESTIONS, 35 QUESTIONS AND 42 QUESTIONS ON 1400 DOCUMENT COLLECTION

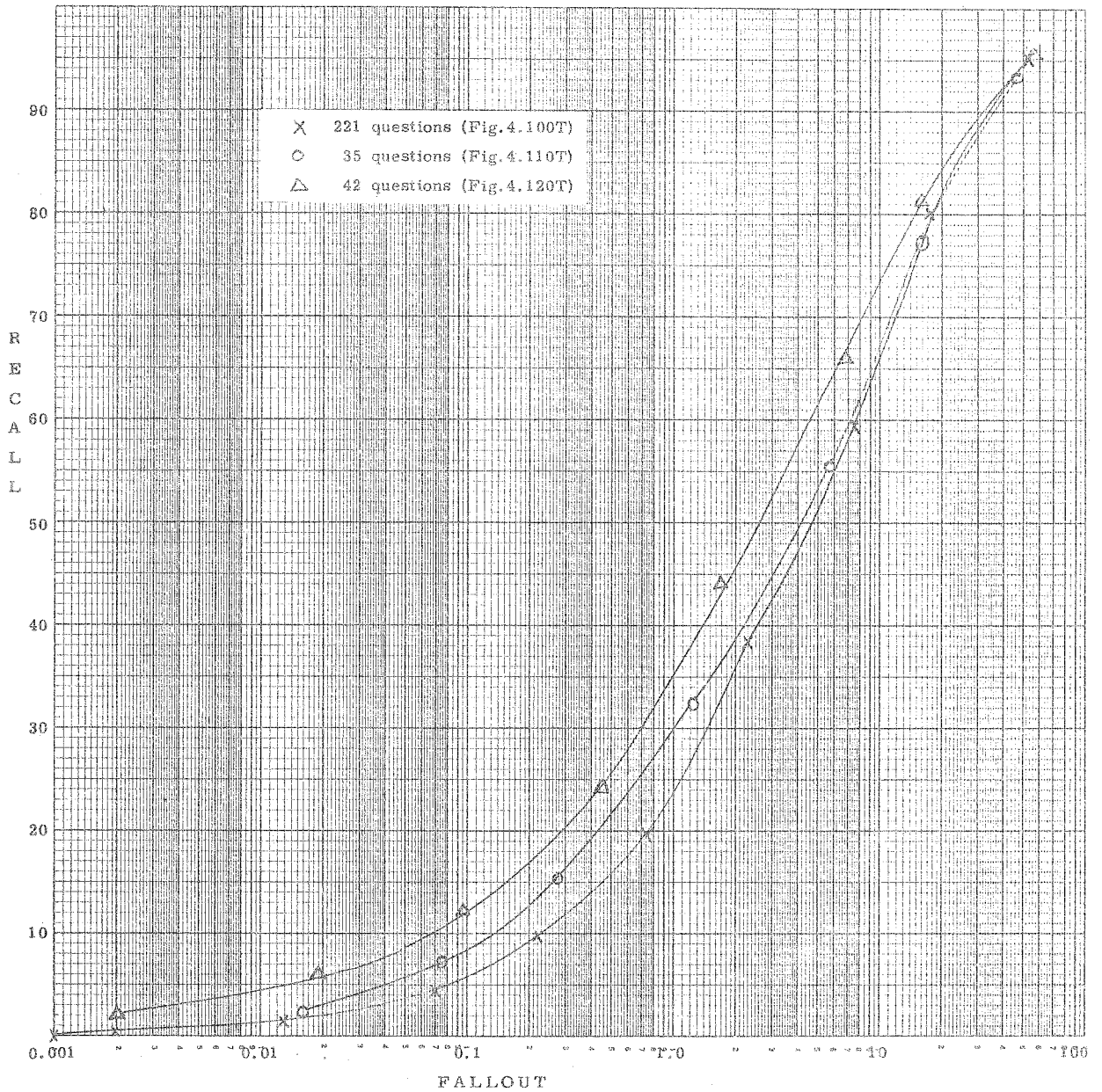


FIGURE 4.131P RECALL/FALLOUT PLOT OF RESULTS FOR INDEX LANGUAGE I.1.a FOR 221 QUESTIONS, 35 QUESTIONS AND 42 QUESTIONS ON 1400 DOCUMENT COLLECTION

FIGURE 4.140T

Index Language I.1.a (S. T. Natural language. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

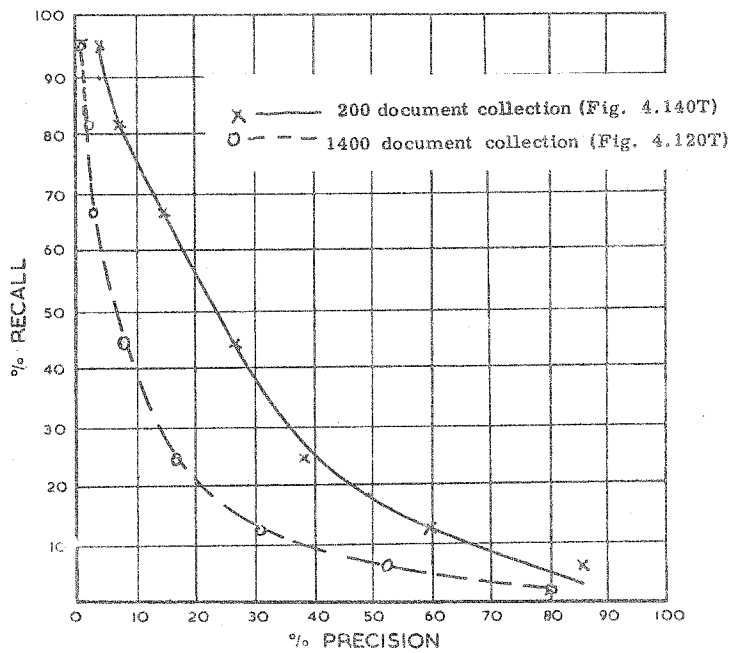
Generality Number 23.6

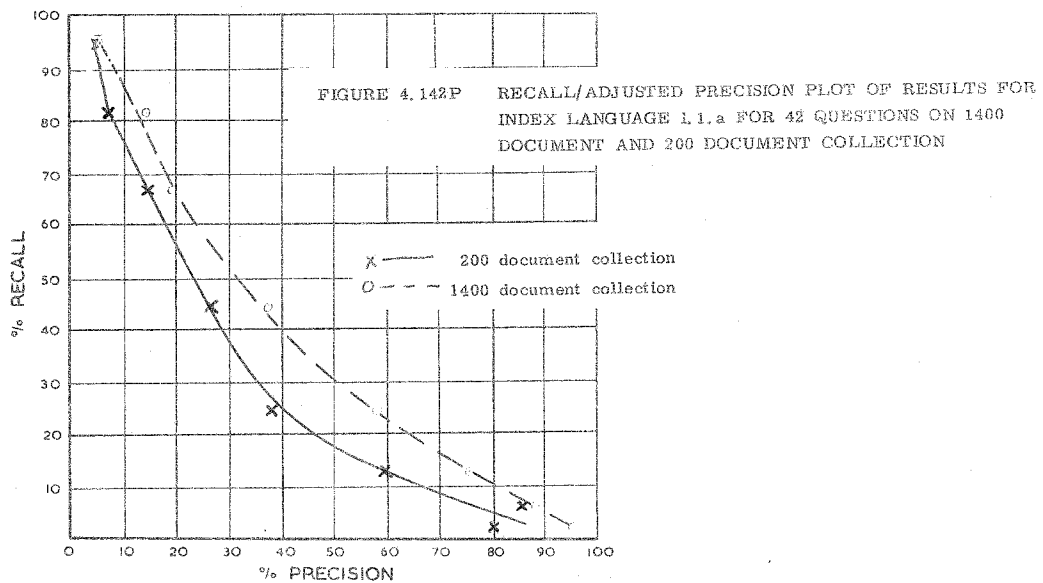
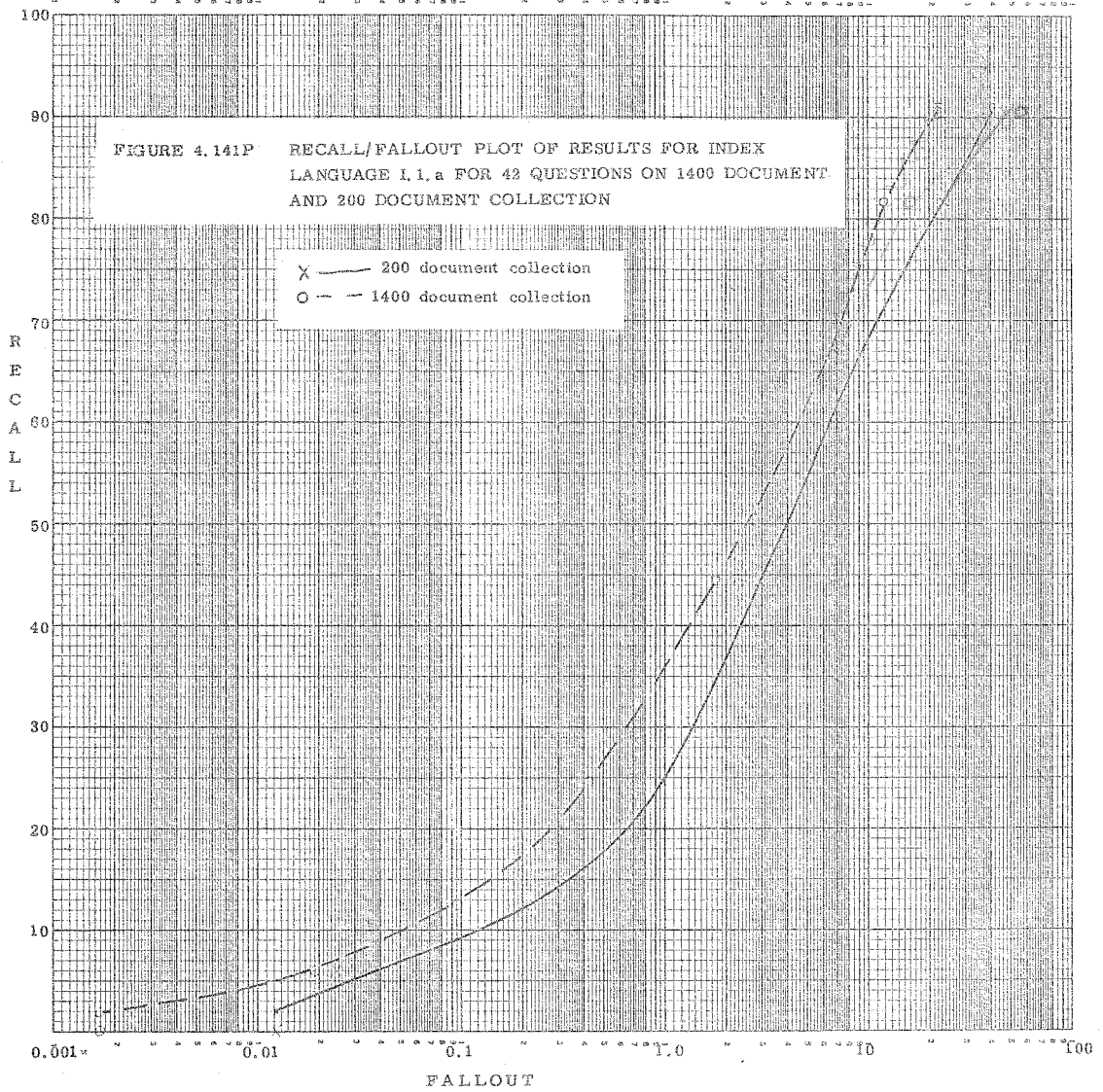
Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	189	4,072	95.5%	4.4%	49.632%	42	42	42
2	162	1,929	81.8%	7.7%	23.519%	42	42	42
3	132	781	66.7%	14.8%	9.278%	42	42	42
4	88	241	44.4%	26.7%	2.938%	34	41	41
5	49	80	24.7%	38.0%	0.975%	23	39	39
6	25	17	12.6%	59.5%	0.207%	15	33	33
7	12	2	6.1%	85.7%	0.024%	6	27	27
8	4	1	2.0%	80.0%	0.012%	2	18	18
9	0	1	0.0%	0.0%	0.012%	1	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1

7765

207

306





Section 2 Recall Devices

The tests made on different recall devices on the single term index languages are based on the 42 questions, (subset 2), searched on the 200 document collection. Tables 4.200T to 4.206T give the performance results, and the recall/precision plots given in Figs. 4.200P to 4.206P show the curve for the particular language being tested, together with the curve, (shown by a broken line) obtained by the device-less search in natural language, which was presented in the previous section as Fig. 4.140P.

A combined plot of all the seven devices or their aggregates, together with natural language, is given in Fig. 4.207P.

In presenting these results, all other possible variables of search rules, indexing exhaustivity and document relevance are held constant, and are fixed to the one particular variation indicated at the head of the tables.

LIST OF FIGURES

	Index Language	No. of Questions	Question Subset	Document Collection	Plots
4.200TP	I.2.a	42	2	200	Plot + 4.140T
4.201TP	I.3.a	42	2	200	" " "
4.202TP	I.5.a	42	2	200	" " "
4.203TP	I.6.a	42	2	200	" " "
4.204TP	I.7.a	42	2	200	" " "
4.205TP	I.8.a	42	2	200	" " "
4.206TP	I.9.a	42	2	200	Plot + 4.140T
4.207P	I.1.a	42	2	200	Plot 4.140T,
	I.2.a				4.200T,
	I.3.a				4.201T,
	I.5.a				4.202T,
	I.6.a				4.203T,
	I.7.a				4.204T,
	I.8.a				4.205T,
	I.9.a				4.206T.

FIGURE 4.200T

Index Language I.2a (S.T. Synonyms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	189	4,193	95.5%	4.3%	51.147%	42	42	42
2	165	2,092	83.3%	7.3%	25.506%	42	42	42
3	134	837	67.7%	13.8%	10.205%	42	42	42
4	90	282	45.5%	24.2%	3.438%	34	41	41
5	55	94	27.8%	36.9%	1.146%	24	39	39
6	25	23	12.6%	52.1%	0.280%	15	33	33
7	12	2	6.1%	85.7%	0.024%	6	27	27
8	4	1	2.0%	80.0%	0.012%	2	18	18
9	0	1	0.0%	0.0%	0.012%	1	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1
13								
14								
15								

8199

304 206

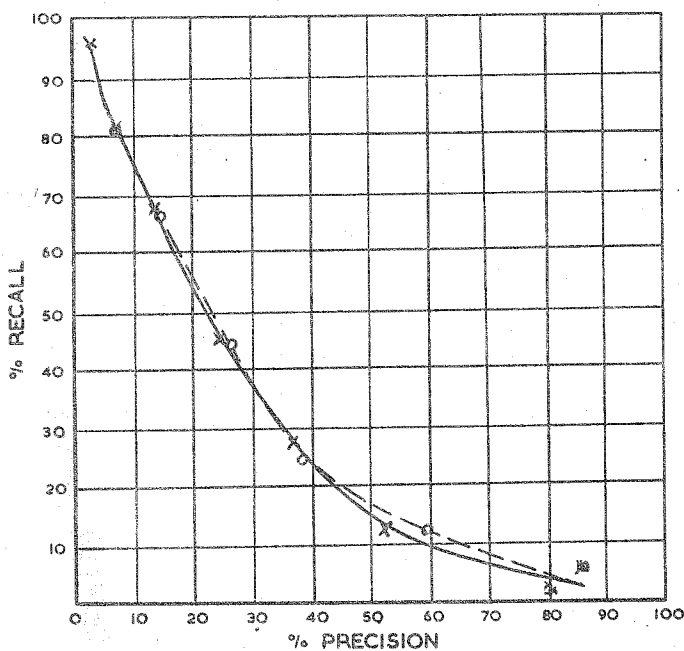


FIGURE 4.201T

Index Language I.3a (S.T. Word forms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	191	4,416	96.5%	4.1%	53.826%	42	42	42
2	166	2,221	83.8%	7.0%	27.079%	42	42	42
3	139	913	70.2%	13.2%	11.131%	42	42	42
4	97	315	49.0%	23.5%	3.841%	35	41	41
5	59	102	29.8%	36.6%	1.244%	25	39	39
6	32	28	16.2%	53.3%	0.341%	17	33	33
7	14	3	7.1%	82.4%	0.037%	8	27	27
8	4	1	2.0%	80.0%	0.012%	2	18	18
9	0	1	0.0%	0.0%	0.012%	1	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1
13								
14								
15								

8702

214 306

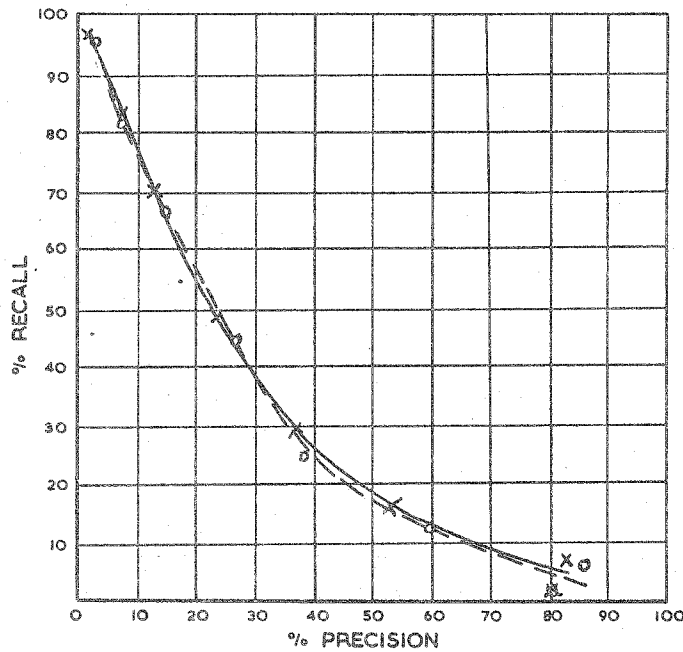


FIGURE 4.202T

Index Language I.5a (S.T. Synonyms, Quasi-synonyms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	195	6,126	98.5%	3.0%	74.715%	42	42	42
2	177	3,565	89.4%	4.7%	43.465%	42	42	42
3	144	1,613	72.7%	8.2%	19.666%	42	42	42
4	117	660	59.1%	15.1%	8.047%	37	41	41
5	73	234	36.9%	23.8%	2.853%	30	39	39
6	41	69	20.7%	37.3%	0.841%	23	33	33
7	17	14	8.6%	54.8%	0.171%	13	27	27
8	6	1	3.0%	85.7%	0.012%	3	18	18
9	5	1	2.5%	83.3%	0.012%	2	11	11
10	0	1	0.0%	0.0%	0.012%	1	7	7
11	0	0				0	3	3
12	0	0				0	1	1
13								
14								
15								

13059

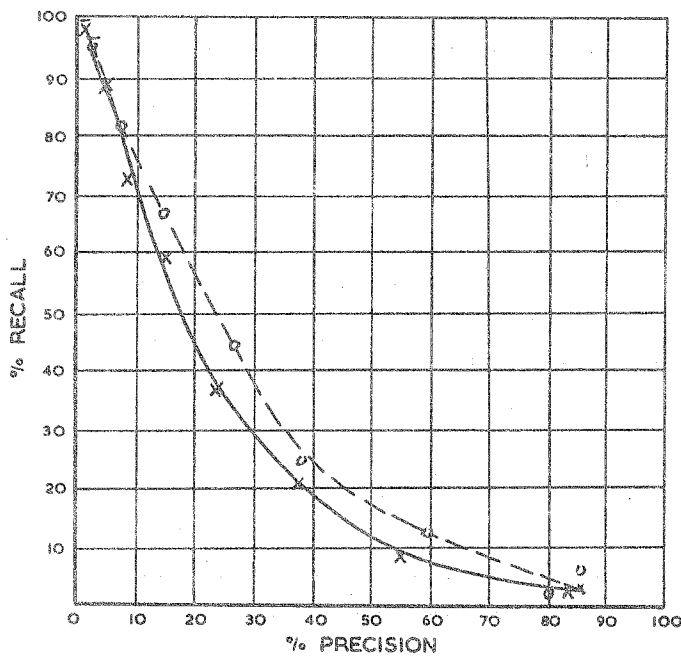


FIGURE 4.203T

Index Language I.6a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	195	6,347	98.5%	2.9%	77.428%	42	42	42
2	179	3,763	90.4%	4.5%	45.879%	42	42	42
3	151	1,785	76.3%	7.8%	21.763%	42	42	42
4	120	768	60.6%	13.5%	9.364%	37	41	41
5	79	280	39.8%	22.0%	3.414%	31	39	39
6	47	87	23.7%	35.1%	1.061%	23	33	33
7	27	20	13.6%	57.4%	0.244%	15	27	27
8	7	1	3.5%	87.5%	0.012%	3	18	18
9	5	1	2.5%	83.3%	0.012%	2	11	11
10	0	1	0.0%	0.0%	0.012%	1	7	7
11	0	0				0	3	3
12	0	0				0	1	1
13								
14								
15								

13863

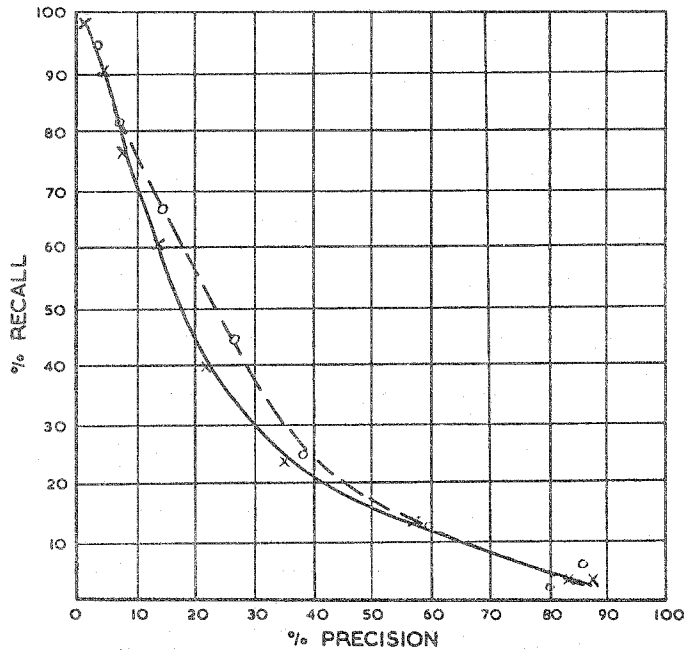


FIGURE 4.204T

Index Language I.7a (S.T. Synonyms, First hierarchical reduction. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	193	4,693	97.5%	3.9%	57.235%	42	42	42
2	172	2,366	86.9%	6.8%	28.847%	42	42	42
3	143	989	72.2%	12.6%	12.058%	42	42	42
4	101	353	51.0%	22.2%	4.304%	36	41	41
5	62	118	31.3%	34.4%	1.439%	27	39	39
6	31	31	15.7%	50.0%	0.378%	20	33	33
7	11	4	5.6%	73.3%	0.049%	6	27	27
8	4	1	2.0%	80.0%	0.012%	2	18	18
9	1	0	0.5%	100.0%	0.000%	1	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1
13								
14								
15								

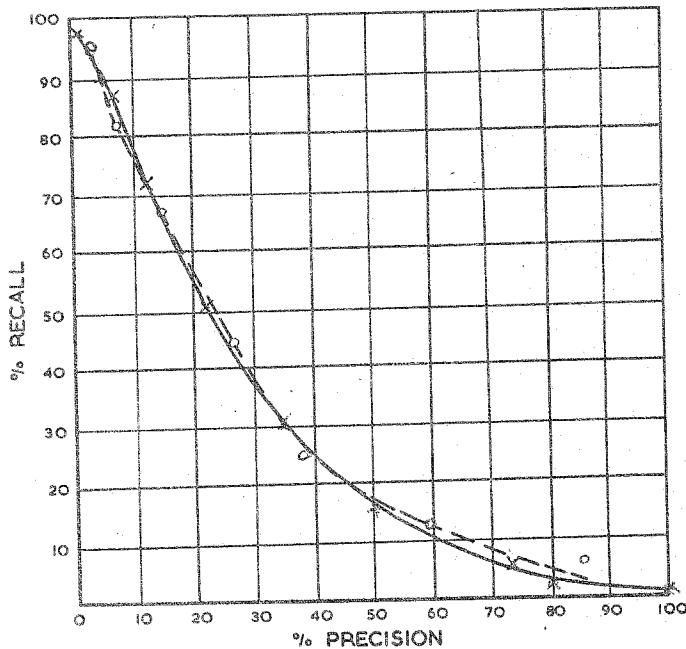


FIGURE 4.205T

Index Language I.8a (S.T. Synonyms, Second hierarchical reduction. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	193	4,954	97.5%	3.7%	60.409%	42	42	42
2	173	2,551	87.4%	6.4%	31.102%	42	42	42
3	146	1,081	73.7%	11.9%	13.180%	42	42	42
4	104	418	52.5%	19.9%	5.096%	37	41	41
5	67	146	33.8%	31.5%	1.780%	27	39	39
6	33	37	16.7%	47.1%	0.451%	21	33	33
7	12	7	6.1%	63.2%	0.085%	8	27	27
8	5	1	2.5%	83.3%	0.012%	2	18	18
9	2	0	1.0%	100.0%	0.000%	1	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1
13								
14								
15								

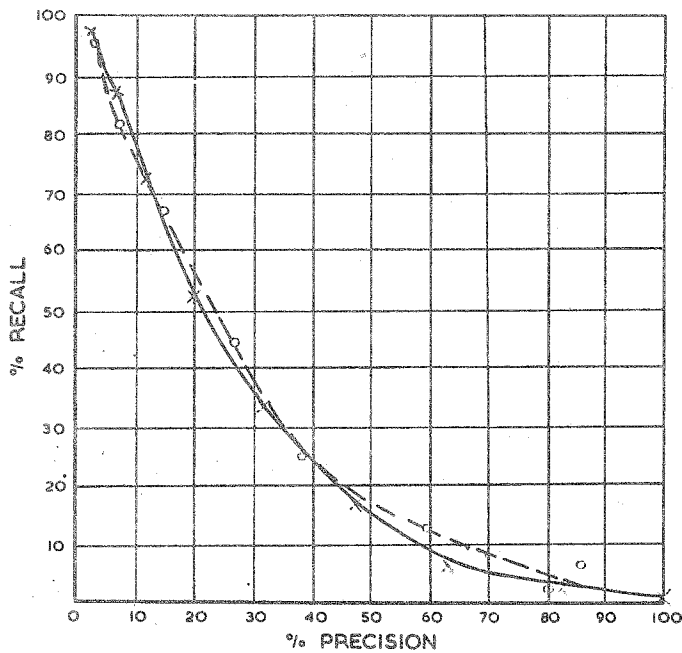


FIGURE 4.206T

Index Language I.9a (S.T. Synonyms, Third hierarchical reduction. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

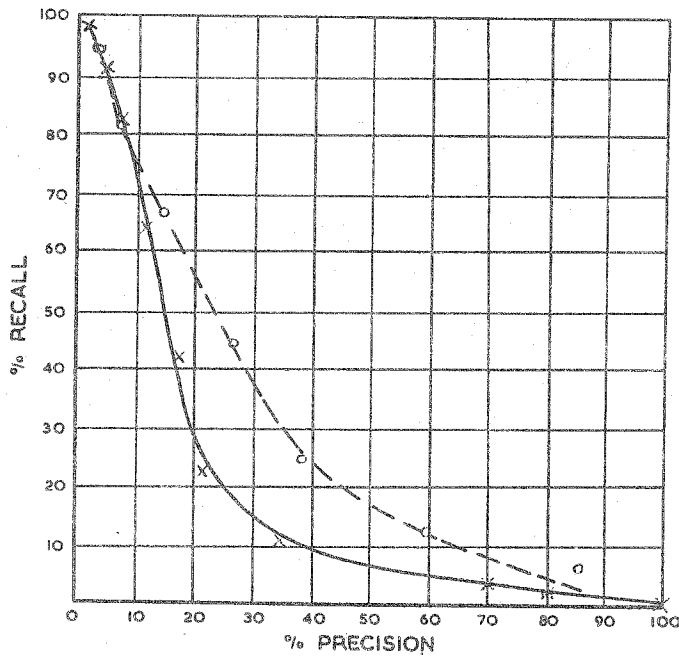
Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	195	6,576	98.5%	2.8%	78.936%	42	42	42
2	181	3,990	91.4%	4.3%	48.647%	42	42	42
3	163	2,034	82.3%	7.4%	24.799%	42	42	42
4	127	959	64.1%	11.7%	11.692%	39	41	41
5	83	403	41.9%	17.1%	4.813%	32	39	39
6	46	171	23.2%	21.2%	2.085%	25	33	33
7	22	42	11.1%	34.4%	0.512%	14	27	27
8	7	3	3.5%	70.0%	0.037%	3	18	18
9	4	1	2.0%	80.0%	0.012%	2	11	11
10	1	0	0.5%	100.0%	0.000%	1	7	7
11	0	0				0	3	3
12	0	0				0	1	1
13								
14								
15								

15008

5.76 7.29



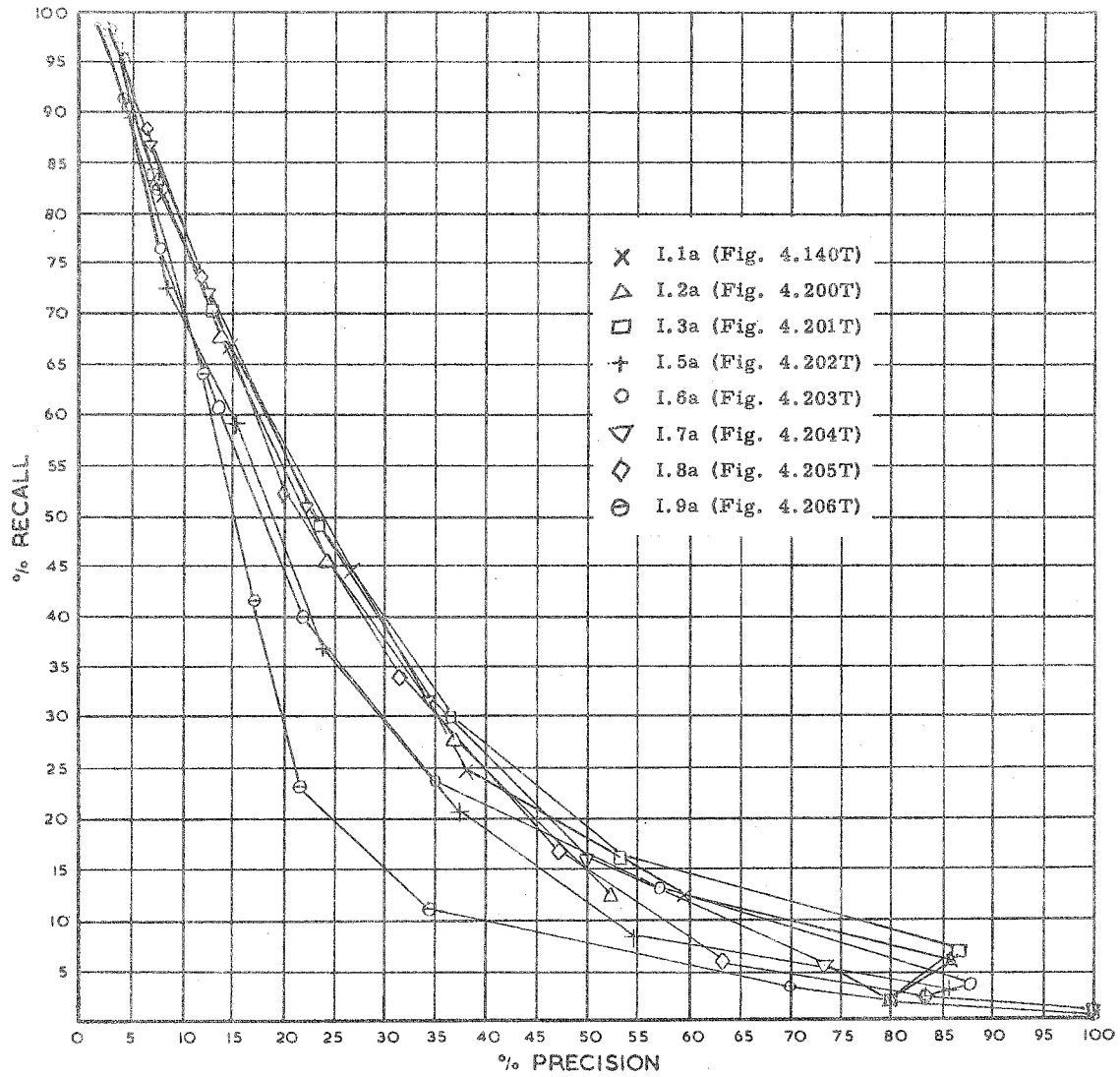


FIGURE 4.207P PLOT OF RESULTS FOR INDEX LANGUAGES I.1a, I.2a, I.3a, I.5a, I.6a, I.7a, I.8a and I.9a for 42 QUESTIONS ON 200 DOCUMENT COLLECTION.

Section 3 Precision Devices

Two subsets of questions have been used to investigate precision devices. Subset 4 consisted of 19 questions all having seven starting terms, while Subset 5 had 17 questions all having eleven starting terms. For reasons explained in Chapter 2, these tests made on precision devices required a different search rule (B) to what has been used so far. This rule required that the main term in a concept must be present if the qualifying term was to be accepted, e.g. with the concept 'Hinged flap', 'Hinged' would not be accepted unless 'Flap' was also present. Therefore it is necessary first to present the results of searches made on these two subsets with search rule B, so that later comparison can be made. Each subset is tested with Index Language I.1.a and I.6.a.

There now follow the results of tests with (b) partitioning, (Figs. 4.310T - 4.313T), (c) interfixing, (Figs. 4.320T - 4.323T) and (a) partitioning plus interfixing (Figs. 4.330T - 4.333T), with both subsets of questions presented separately. For each precision device results are given when two recall languages are used, I.1 Natural Language and I.6. the aggregate of synonyms, quasi-synonyms and word forms.

Because of the large clerical effort required to obtain these results, the precision ratios were not obtained at the lower levels of coordination. This, combined with the relatively small sets of questions used, limits the usefulness of any graphical presentation, for the results show little consistency. Therefore in the plots 4.340P, 4.341P, 4.342P and 4.343P, only a single generalised curve has been drawn for the four sets of data that are presented in each curve. The relative positions of the various symbols give an indication of the performance of the particular system.

LIST OF FIGURES

	Precision Device	Index Language	No. of Questions	Question Subset	Document Collection
4.300T		I.1.a	19	4	1400 (Search
4.301T		I.1.a	17	5	1400 Rule B)
4.302T		I.6.a	19	4	1400
4.303T		I.6.a	17	5	1400
4.310T	PARTITIONING	I.1.b	19	4	1400
4.311T	"	I.1.b	17	5	1400

List of Figures (Cont)

	Precision Device	Index Language	No. of Questions	Question Subset	Document Collection	Plots
4.312T	PARTITIONING	I.6.b	19	4	1400	
4.313T	"	I.6.b	17	5	1400	
4.320T	INTERFIXING	I.1.c	19	4	1400	
4.321T	"	I.1.c	17	5	1400	
4.322T	"	I.6.c	19	4	1400	
4.323T	"	I.6.c	17	5	1400	
4.330T	PART. + INT.	I.1.d	19	4	1400	
4.331T	"	I.1.d	17	5	1400	
4.332T	"	I.6.d	19	4	1400	
4.333T	"	I.6.d	17	5	1400	
4.340P	PART., INT., & PART. + INT.	I.1.a,b, c,d	19	4	1400	Plot 4.300T, 4.310T, 4.320T, 4.330T.
4.341P	"	I.1.a,b, c,d	17	5	1400	Plot 4.301T, 4.311T, 4.321T, 4.331T.
4.342P	"	I.6.a,b, c,d	19	4	1400	Plot 4.302T, 4.312T, 4.322T, 4.332T.
4.343P	"	I.6.a,b, c,d	17	5	1400	Plot 4.303T, 4.313T, 4.323T, 4.333T.

















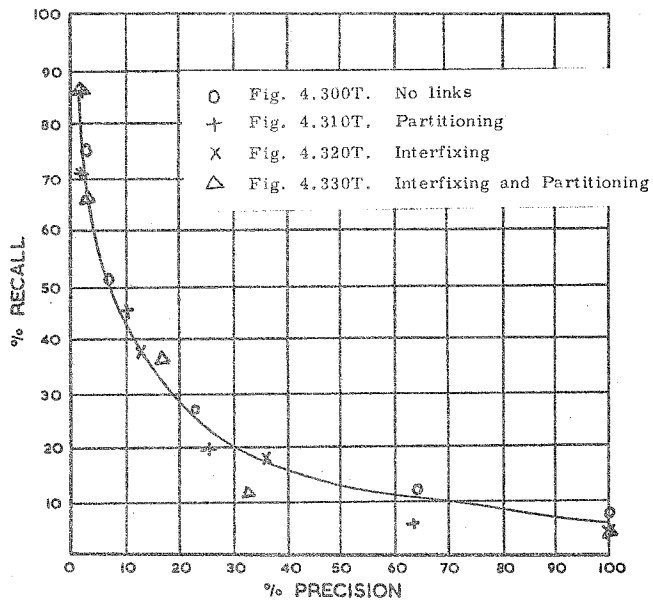


FIGURE 4.340P EFFECT OF LINKING DEVICES FOR INDEX LANGUAGE 1.1 WITH 19 QUESTIONS

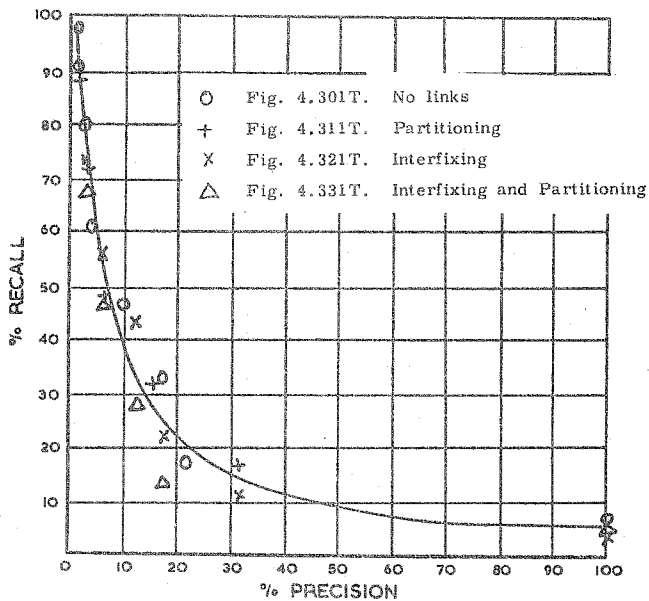


FIGURE 4.341P EFFECT OF LINKING DEVICES FOR INDEX LANGUAGE 1.1 WITH 17 QUESTIONS

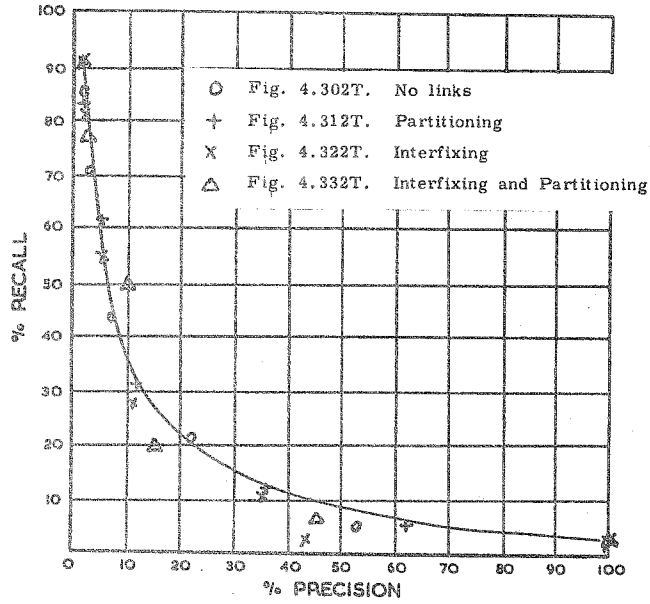


FIGURE 4.342P EFFECT OF LINKING DEVICES FOR INDEX LANGUAGE I.6 WITH 19 QUESTIONS

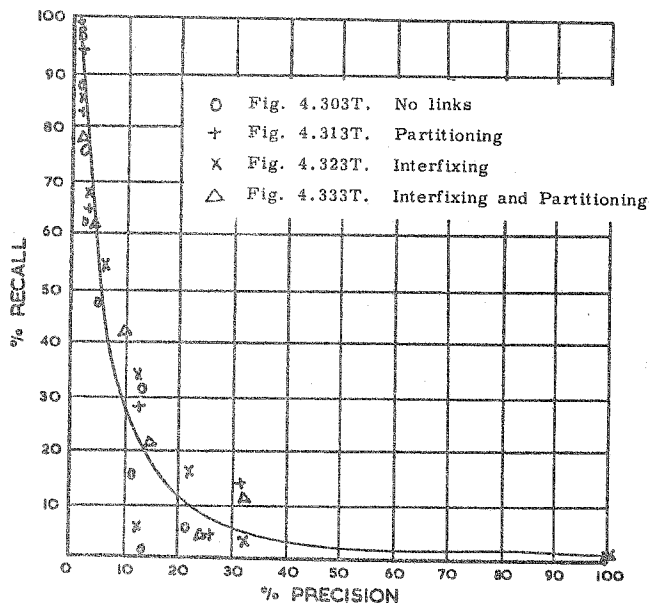


FIGURE 4.343P EFFECT OF LINKING DEVICES FOR INDEX LANGUAGE I.6 WITH 17 QUESTIONS

Section 4 Indexing Exhaustivity

Results when the three different levels of indexing exhaustivity are compared are given in Figs. 4.400T - 4.405T and 4.410T - 4.415T. Three sets of questions are presented separately, being the 221 questions (subset 3) searched on the 1400 collection, the 35 seven starting term questions (subset 1) searched on the 1400 collection, and the 42 questions (subset 2) searched on the 200 collection. Results for index languages I.1.a and I.6.a are given, and the three different performance curves for each question subset and language are given in Figs. 4.420P - 4.425P. No adjustment has been made in these graphs for the differing generality number amongst the different plots, since only the relative position of the curves on each plot is important.

LIST OF FIGURES

	Exhaustivity	Index Language	No. of Questions	Question Subset	Document Collection	Plots
4.400T	2	I.1.a	221	3	1400	
4.401T	1	I.1.a	221	3	1400	
4.402T	2	I.1.a	35	1	1400	
4.403T	1	I.1.a	35	1	1400	
4.404T	2	I.1.a	42	2	200	
4.405T	1	I.1.a	42	2	200	
4.410T	2	I.6.a	221	3	1400	
4.411T	1	I.6.a	221	3	1400	
4.412T	2	I.6.a	35	1	1400	
4.413T	1	I.6.a	35	1	1400	
4.414T	2	I.6.a	42	2	200	
4.415T	1	I.6.a	42	2	200	
4.420P	3, 2, 1	I.1.a	221	3	1400	Plot 4.100T 4.400T 4.401T
4.421P	3, 2, 1	I.1.a	35	1	1400	Plot 4.110T 4.402T 4.403T
4.422P	3, 2, 1	I.1.a	42	2	200	Plot 4.140T 4.404T 4.405T
4.423P	3, 2, 1	I.6.a	221	3	1400	Plot 4.104T 4.410T 4.411T
4.424P	3, 2, 1	I.6.a	35	1	1400	Plot 4.114T 4.412T 4.414T
4.425P	3, 2, 1	I.6.a	42	2	200	Plot 4.203T 4.414T 4.415T

FIGURE 4.400T

Index Language I.1.a (S.T. Natural language. Coordination)

Exhaustivity of Indexing 2

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 221 (Subset 3)

Number of Relevant Documents 1,590

Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,480	(-)	93.1%	(-)	(-)	221	0	221
2	1,227	49,372*	77.2%	2.4%*	16.040%*	221	44*	221
3	886	18,764*	55.7%	4.5%*	6.092%*	210	109*	220
4	552	6,813*	34.7%	7.5%*	2.212%*	177	142*	212
5	279	2,301*	17.5%	10.8%*	0.747%*	121	177*	197
6	131	514*	8.2%	20.3%*	0.167%*	76	161*	164
7	58	133	3.6%	30.4%	0.043%	41	140	140
8	18	28	1.1%	39.1%	0.008%	14	105	105
9	8	0	0.5%	100.0%	0.000%	6	78	78
10	1	0	0.1%	100.0%	0.000%	1	52	52
11	0	0				0	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 4.401T

Index Language I.1.a (S.T. Natural language. Coordination)

Exhaustivity of Indexing 1

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 221 (Subset 3)

Number of Relevant Documents 1,590

Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,373	(-)	86.4%	(-)	(-)	221	0	221
2	1,055	23,164*	66.4%	4.4%*	7.525%*	221	44*	221
3	655	7,823*	41.2%	7.7%*	2.540%*	199	109*	220
4	331	2,169*	20.8%	13.2%*	0.704%*	154	142*	212
5	148	477*	9.3%	23.7%*	0.155%*	88	177*	197
6	67	134*	4.2%	33.2%*	0.043%*	49	161*	164
7	29	38	1.8%	43.3%	0.012%	25	140	140
8	7	9	0.4%	43.8%	0.003%	9	105	105
9	2	0	0.1%	100.0%	0.000%	2	78	78
10	0	0				0	52	52
11	0	0				0	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 4.402T

Index Language I.1.a (S.T. Natural language. Coordination)

Exhaustivity of Indexing 2

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 35 (Subset 1)

Number of Relevant Documents 287

Generality Number 5.9

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	263	(-)	91.7%	(-)	(-)	35	0	35
2	211	6,823*	73.6%	3.0%*	13.924%*	35	23*	35
3	141	2,442	49.2%	5.4%	5.013%	34	35	35
4	85	488	29.7%	15.0%	1.002%	25	35	35
5	43	94	15.0%	31.4%	0.193%	17	35	35
6	23	20	8.0%	53.5%	0.041%	6	35	35
7	8	8	2.8%	50.0%	0.016%	3	35	35

FIGURE 4.403T

Index Language I.1.a (S.T. Natural language. Coordination.)

Exhaustivity of Indexing 1

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 35 (Subset 1)

Number of Relevant Documents 287

Generality Number 5.9

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	237	(-)	82.6%	(-)	(-)	35	0	35
2	184	3,532*	64.1%	4.9%*	7.208%*	35	23*	35
3	111	1,045	38.7%	9.6%	2.145%	34	35	35
4	54	168	18.9%	24.3%	0.345%	22	35	35
5	22	19	7.7%	53.7%	0.039%	6	35	35
6	10	2	3.5%	83.3%	0.004%	5	35	35
7	5	2	1.8%	71.4%	0.004%	3	35	35



FIGURE 4.410T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)

Exhaustivity of Indexing 2

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 221 (Subset 3)

Number of Relevant Documents 1,590

Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,530	(-)	96.2%	(-)	(-)	221	0	221
2	1,389	98,743*	87.4%	1.4%*	32.089%*	221	44*	221
3	1,104	36,132*	69.4%	3.0%*	11.739%*	215	109*	220
4	782	13,974*	49.2%	5.3%*	4.540%*	203	142*	212
5	467	6,701*	29.3%	6.5%*	2.178%*	158	177*	197
6	261	2,549*	16.4%	9.3%*	0.828%*	111	161*	164
7	128	663	8.1%	16.2%	0.215%	74	140	140
8	53	197	3.3%	21.2%	0.064%	40	105	105
9	19	48	1.2%	28.4%	0.016%	18	78	78
10	5	7	0.3%	41.7%	0.002%	6	52	52
11	0	1	0.0%		0.0003%	1	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 4.411T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)

Exhaustivity of Indexing 1

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 221 (Subset 3)

Number of Relevant Documents 1,590

Generality Number 5.1

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	1,455	(-)	91.5%	(-)	(-)	221	0	221
2	1,226	54,263*	77.1%	2.3%*	17.618%*	221	44*	221
3	861	15,346*	54.2%	5.3%*	4.982%*	212	109*	220
4	520	4,903*	32.7%	9.6%*	1.462%*	186	142*	212
5	261	2,214*	16.4%	10.6%*	0.719%*	124	177*	197
6	122	586*	7.7%	17.2%*	0.190%*	77	161*	164
7	54	141	3.4%	27.7%	0.046%	43	140	140
8	20	38	1.3%	34.5%	0.012%	22	105	105
9	6	7	0.4%	46.2%	0.002%	7	78	78
10	1	1	0.1%	50.0%	0.0003%	2	52	52
11	0	0				0	32	32
12	0	0				0	15	15
13	0	0				0	8	8
14	0	0				0	4	4
15	0	0				0	3	3

FIGURE 4.412T

Index Language I.G.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)  
 Exhaustivity of Indexing 2  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 35 (Subset 1)  
 Number of Relevant Documents 287  
 Generality Number 5.0

Coord- ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	276	(-)	96.2%	(-)	(-)	35	0	35
2	247	15,946*	86.1%	1.5%*	32.542%*	35	23*	35
3	187	6,567	65.2%	2.8%	13.461%	35	35	35
4	114	1,742	39.8%	6.1%	3.675%	32	35	35
5	60	397	21.0%	13.1%	0.815%	25	35	35
6	37	62	12.9%	37.4%	0.127%	16	35	35
7	12	13	4.2%	48.0%	0.027%	5	35	35

FIGURE 4.413T

Index Language I.G.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)  
 Exhaustivity of Indexing 1  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 35 (Subset 1)  
 Number of Relevant Documents 287  
 Generality Number 5.0

Coord- ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	257	(-)	89.6%	(-)	(-)	35	0	35
2	215	8,416*	75.0%	2.5%*	17.176%*	35	23*	35
3	142	2,901	49.5%	4.7%	5.955%	34	35	35
4	74	539	25.8%	12.1%	1.106%	30	35	35
5	32	85	11.2%	27.4%	0.174%	14	35	35
6	12	5	4.2%	70.6%	0.010%	6	35	35
7	6	2	2.1%	75.0%	0.004%	3	35	35

FIGURE 4.414T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)

Exhaustivity of Indexing 2

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	193	5,693	97.5%	3.4%	67.773%	42	42	42
2	175	3,253	88.4%	5.1%	39.661%	42	42	42
3	145	1,431	73.2%	9.2%	17.447%	41	42	42
4	113	544	57.1%	17.2%	6.933%	36	41	41
5	64	178	32.3%	28.7%	2.146%	30	39	39
6	42	40	21.2%	51.2%	0.488%	20	33	33
7	21	13	10.6%	61.8%	0.158%	13	27	27
8	5	1	2.5%	83.3%	0.012%	2	18	18
9	4	0	2.0%	100.0%	0.000%	2	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1

FIGURE 4.415T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)

Exhaustivity of Indexing 1

Search Rule A

Document Relevance 1 - 4

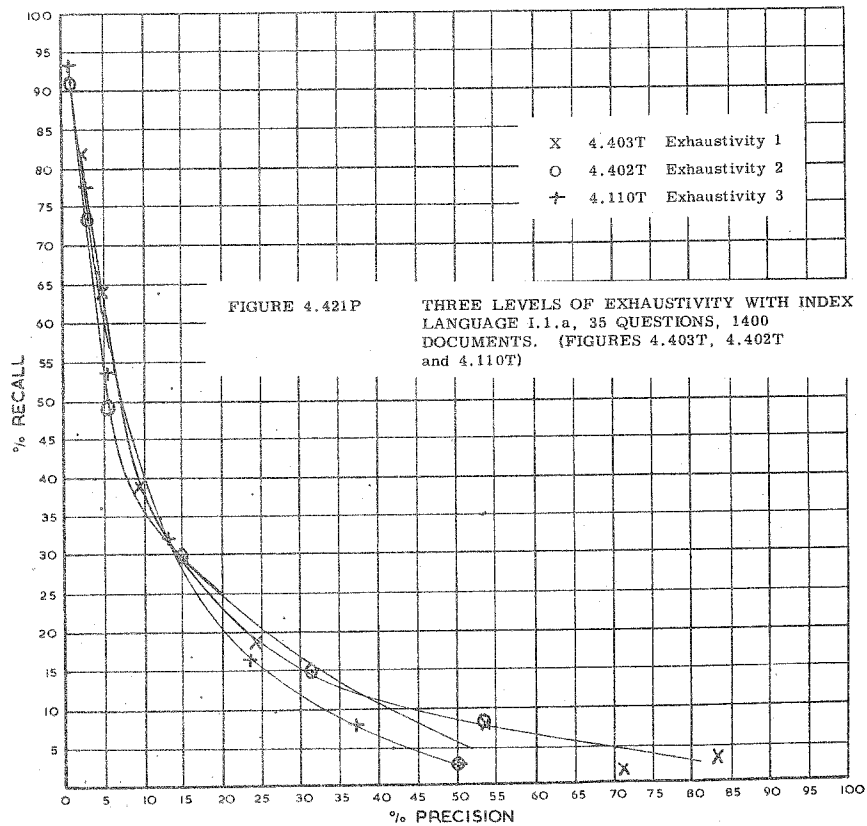
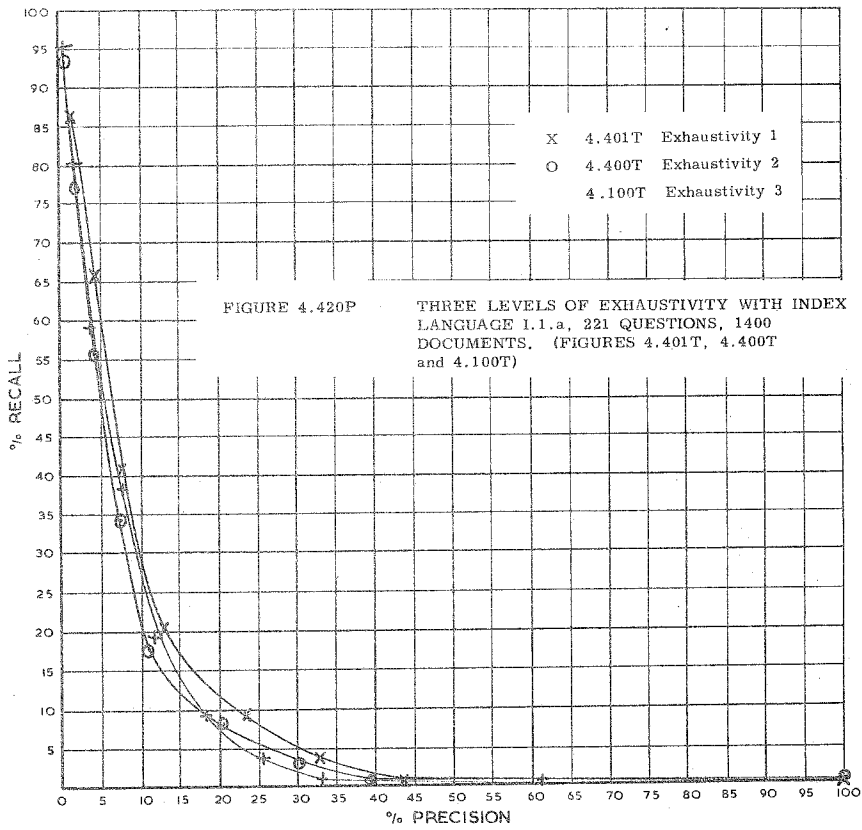
Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	163	4,127	92.4%	4.2%	49.130%	42	42	42
2	155	2,118	78.3%	6.8%	25.823%	42	42	42
3	117	685	59.1%	14.6%	8.352%	39	42	42
4	81	190	40.9%	29.9%	2.317%	32	41	41
5	47	51	23.7%	48.0%	0.622%	23	39	39
6	24	11	12.1%	88.8%	0.134%	14	33	33
7	10	3	5.1%	76.9%	0.037%	6	27	27
8	2	0	1.0%	100.0%	0.000%	1	18	18
9	2	0	1.0%	100.0%	0.000%	1	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1



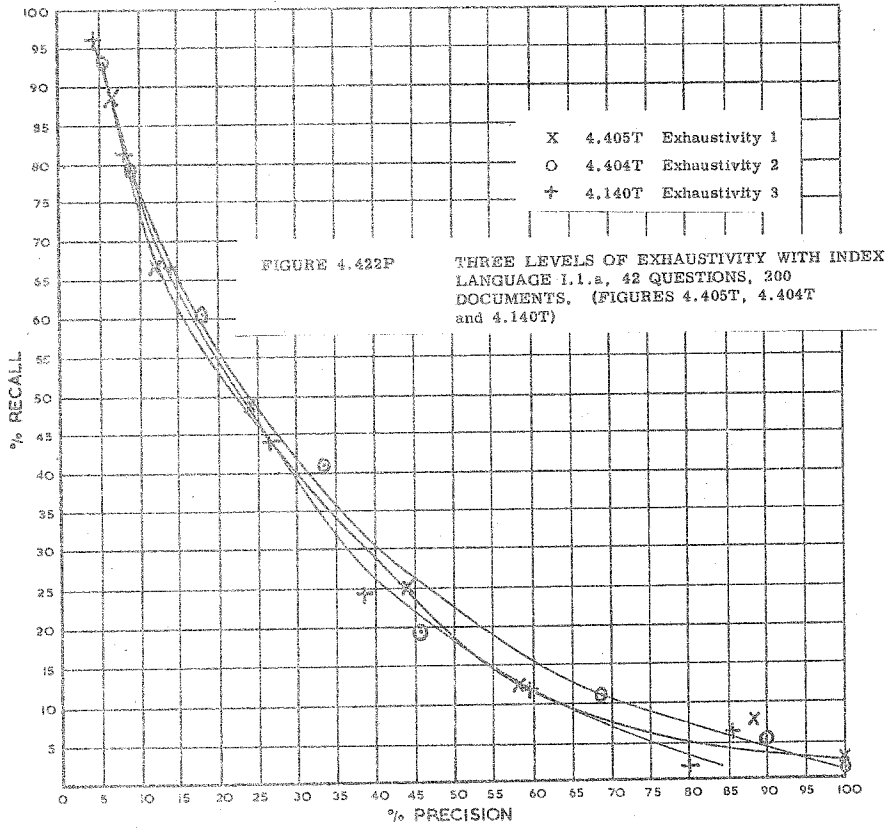
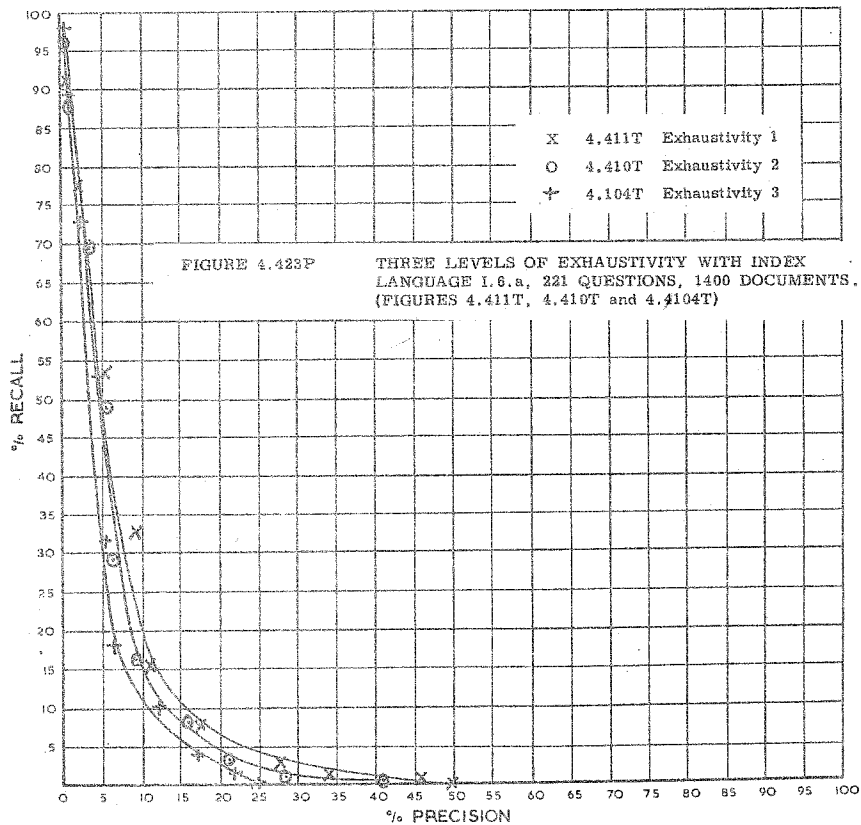
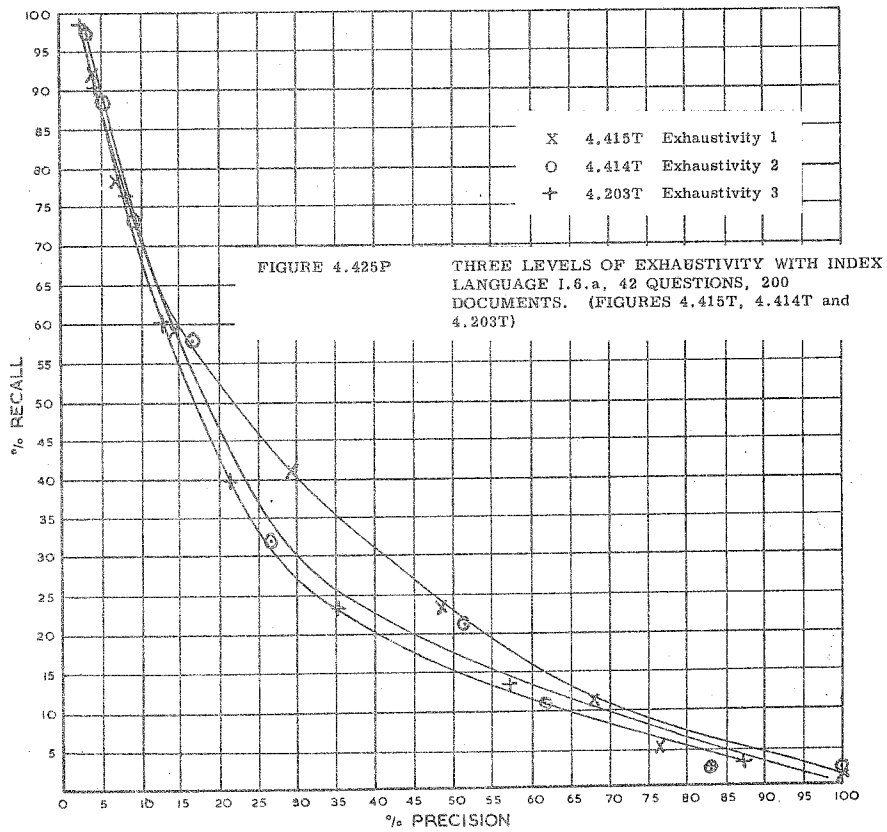
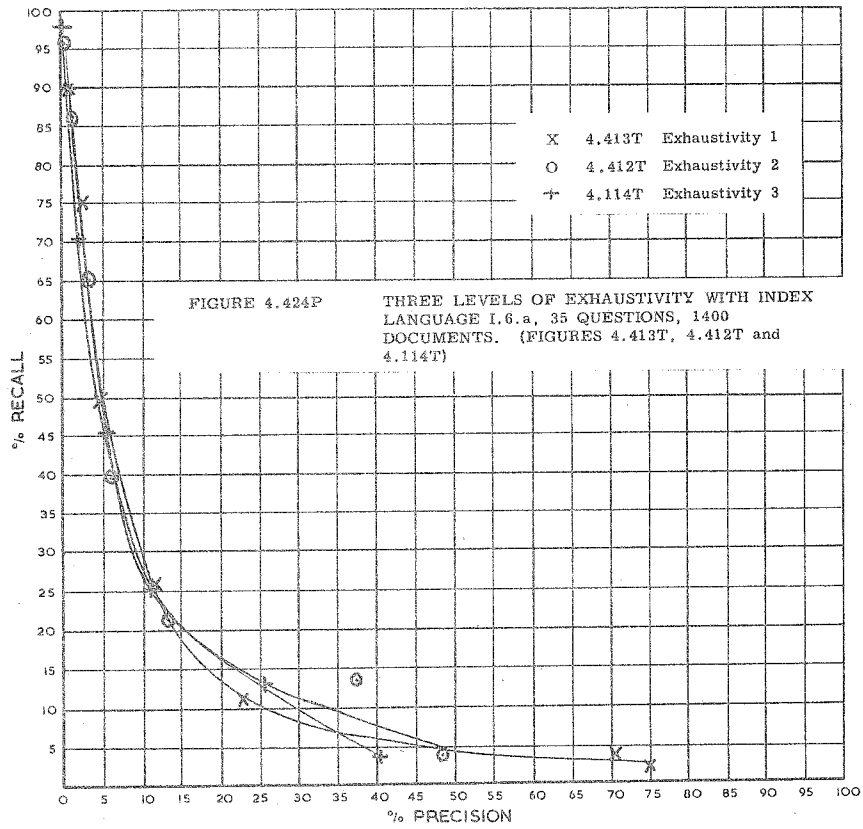


Fig 503 P





## Section 5 Search Rules

The first search rule investigated is type B, which, as explained in Chapter 2, was used to provide a minimal amount of intellect prior to making the test of precision devices. The sets of questions (subsets 4 and 5) and languages (I.1.a and I.6.a) that were used to test precision devices are investigated. Figs. 4.500T - 4.503T give tables of results for search A, this being the basic type of search used in presenting the results on, for instance, the recall devices. This particular series of results has to be given because of the different set of questions which was used for testing Search B; the results of these latter searches have already been presented in Figs. 4.300T to 4.303T in Section 3. Graphs comparing the curves of search A with search B are given in Figs. 4.504P to 4.507P. Search Rule C demanded a selection of the original starting terms taken from the question. This eliminated such terms as Problem, Influence, Comparison, etc., and the result was to reduce the number of starting terms in every question. Whereas originally, in this subset, each question had seven starting terms, when search rule C was in force, the number varied from two to six.

Search Rule D retained the terms selected for search rule C but imposed restrictions concerning the combinations of terms that would be accepted at every coordination level, thereby eliminating non-sensical combinations.

For search rules C and D a set of 20 questions (subset 6), all taken from the 35 questions having seven starting terms (subset 1) is used, searched on the 1400 collection. Languages I.1.a and I.6.a are tested and Figs. 4.510T and 4.511T give the results with search rule A.

Results for searches C and D were totalled by method 1B, (see chapter 3, page 62), since the ordinary strict coordination level method gave unsatisfactory results due to the small number of questions and variation in starting terms that resulted from the rules used in the searches. Figs. 4.512T and 4.513T give the results of search C, and Figs. 4.514T and 4.515T the results of search D; because of the different totalling methods the tables of results have to include the additional data which shows the reduced number of documents regarded as relevant at the higher coordination levels, and also the resultant change in generality number. Two plots Figs. 4.516P and 4.517P present the result for searches A, C and D with the two index languages.

Although the comparison between search A on the one hand and searches C and D on the other hand might seem to be influenced by the different totalling method used, all possible methods were tried and that used was found to be the most satisfactory.

	Search Rule	Index Language	No. of Questions	Question Subset	Document Collection	Plot
4.500T	A	I.1.a	19	4	1400	
4.501T	A	I.1.a	17	5	1400	
4.502T	A	I.6.a	19	4	1400	
4.503T	A	I.6.a	17	5	1400	
4.504P	A & B	I.1.a	19	4	1400	Plot 4.300T 4.500T
4.505P	A & B	I.1.a	17	5	1400	Plot 4.301T 4.501T
4.506P	A & B	I.6.a	19	4	1400	Plot 4.302T 4.502T
4.507P	A & B	I.6.a	17	5	1400	Plot 4.303T 4.503T
4.510T	A	I.1.a	20	6	1400	
4.511T	A	I.6.a	20	6	1400	
4.512T	C	I.1.a	20	6	1400	
4.513T	C	I.6.a	20	6	1400	
4.514T	D	I.1.a	20	6	1400	
4.515T	D	I.6.a	20	6	1400	
4.516P	A, C & D	I.1.a	20	6	1400	Plot 4.510T 4.512T 4.514T
4.517P	A, C & D	I.6.a	20	6	1400	Plot 4.511T 4.513T 4.515T

FIGURE 4.500T

Index Language I.1.a (S.T. Natural language. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 19 (Subset 4)

Number of Relevant Documents 131

Generality Number 4.9

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	118	12,700	90.1%	0.9%	47.744%	19	19	19
2	100	4,294	76.3%	2.4%	16.143%	19	19	19
3	72	1,426	55.0%	5.0%	5.360%	19	19	19
4	39	188	29.8%	17.2%	0.710%	15	19	19
5	18	16	13.7%	52.9%	0.060%	8	19	19
6	10	3	7.6%	76.9%	0.011%	2	19	19
7	4	0	3.1%	100.0%	0.000%	2	19	19

FIGURE 4.501T

Index Language I.1.a (S.T. Natural language. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 17 (Subset 5)

Number of Relevant Documents 109

Generality Number 4.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	108	14,710	99.1%	0.7%	61.807%	17	17	17
2	105	8,030	96.3%	1.3%	33.739%	17	17	17
3	91	3,652	83.5%	2.4%	15.344%	17	17	17
4	74	1,559	67.9%	4.5%	6.550%	17	17	17
5	52	618	47.7%	7.8%	2.609%	17	17	17
6	39	237	35.8%	14.1%	1.000%	15	17	17
7	22	100	20.2%	18.0%	0.422%	14	17	17
8	13	22	11.9%	37.1%	0.093%	10	17	17
9	6	4	5.5%	60.0%	0.017%	5	17	17
10	1	0	0.9%	100.0%	0.000%	1	17	17
11	0	0				0	17	17

FIGURE 4.502T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word Forms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 19 (Subset 4)

Number of Relevant Documents 131

Generality Number 4.9

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	126	(-)	96.2%	(-)	(-)	19	0	19
2	114	9,568	87.0%	1.2%	35.969%	19	19	19
3	98	3,962	74.8%	2.5%	14.895%	19	19	19
4	61	1,142	46.6%	5.1%	4.314%	19	19	19
5	33	192	25.2%	14.7%	0.725%	16	19	19
6	18	30	13.7%	37.5%	0.113%	8	19	19
7	8	1	6.1%	88.9%	0.004%	3	19	19

FIGURE 4.503T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms, Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 17 (Subset 5)

Number of Relevant Documents 109

Generality Number 4.6

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	109	(-)	100.0%	(-)	(-)	17	0	17
2	108	13,865	99.1%	0.8%	58,256%	17	17	17
3	101	7,264	92.7%	1.4%	30,521%	17	17	17
4	91	4,124	83.5%	2.2%	17,328%	17	17	17
5	73	2,353	67.0%	3.0%	9.932%	17	17	17
6	54	1,026	49.5%	5.0%	4.331%	17	17	17
7	37	361	34.0%	9.3%	1.524%	17	17	17
8	23	136	21.1%	14.5%	0.574%	16	17	17
9	9	41	8.3%	18.0%	0.173%	10	17	17
10	2	12	1.8%	14.3%	0.051%	4	17	17
11	0	0				0	17	17

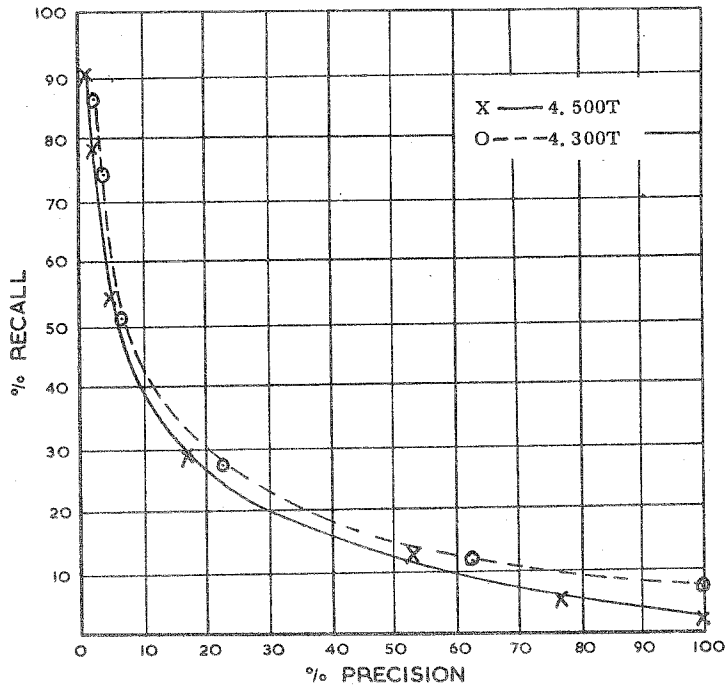


FIGURE 4.504P SEARCH RULES A and B, 19 QUESTIONS, INDEX LANGUAGE I.1.a (FIGURES 4.500T and 4.300T)

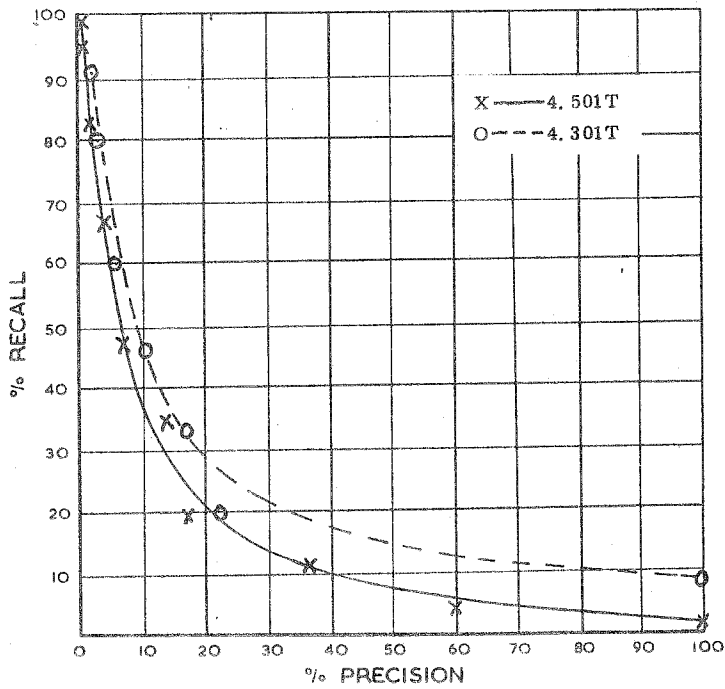


FIGURE 4.505P SEARCH RULES A and B, 17 QUESTIONS, INDEX LANGUAGE I.6.a (FIGURES 4.501T and 4.301T)

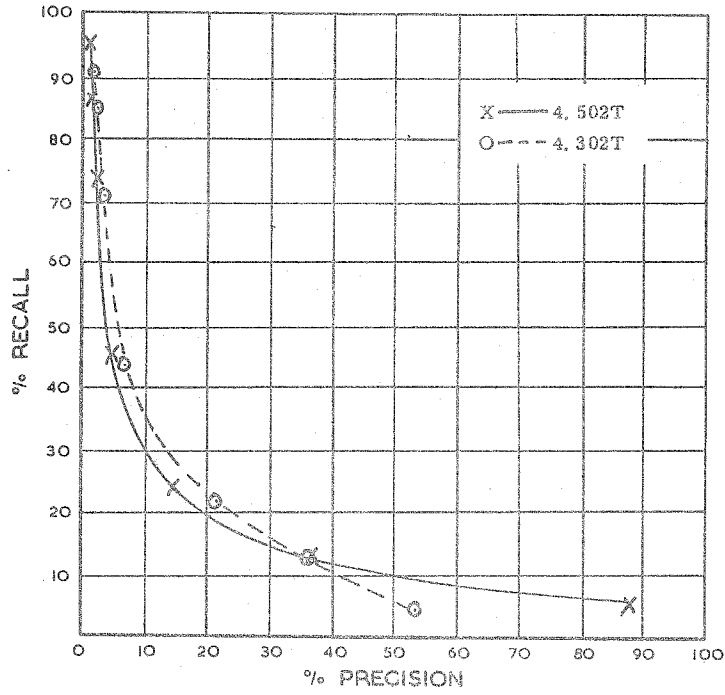


FIGURE 4.506P SEARCH RULES A and B, 19 QUESTIONS, INDEX LANGUAGE I.1.a (FIGURES 4.502T and 4.302T)

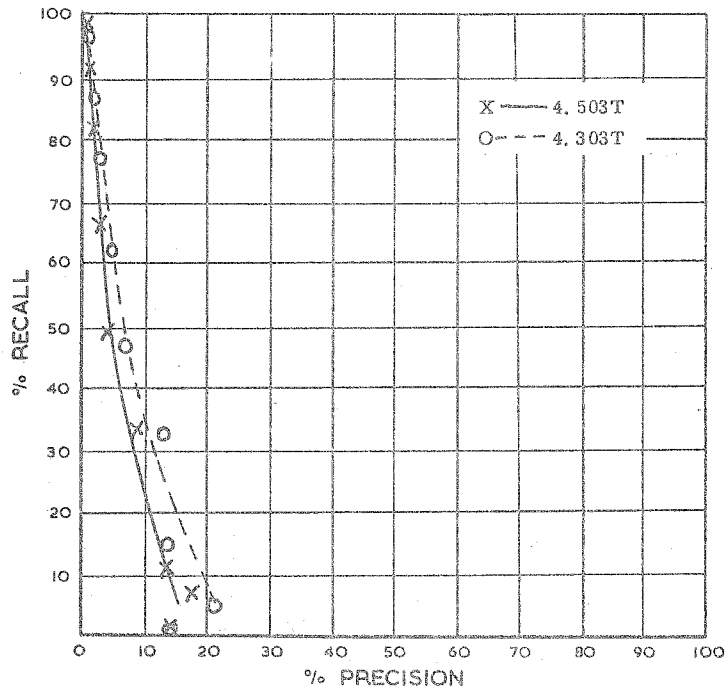


FIGURE 4.507P SEARCH RULES, A and B, 17 QUESTIONS, INDEX LANGUAGE I.6.a (FIGURES 4.503T and 4.303T)

FIGURE 4.510T

Index Language I.i.a (S.T. Natural language. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 20 (Subset 6)  
 Number of Relevant Documents 147  
 Generality Number 5.3

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	133	14,225	90.5%	0.9%	50.804%	20	20	20
2	114	4,704	77.6%	2.4%	16.808%	20	20	20
3	82	1,630	55.9%	4.9%	5.852%	20	20	20
4	47	230	32.0%	17.0%	0.826%	16	20	20
5	21	24	14.3%	46.7%	0.086%	9	20	20
6	10	3	6.8%	76.9%	0.011%	2	20	20
7	4	0	2.7%	100.0%	0.000%	2	20	20

FIGURE 4.511T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 20 (Subset 6)  
 Number of Relevant Documents 147  
 Generality Number 5.3

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	142	(-)	96.6%	(-)	(-)	20	0	20
2	129	8,088*	87.8%	1.6%*	29.038%*	20	13*	20
3	109	4,657	74.1%	2.3%	16.720%	20	20	20
4	69	1,276	46.9%	5.1%	4.581%	20	20	20
5	36	226	24.5%	13.7%	0.811%	16	20	20
6	19	32	12.9%	37.3%	0.115%	8	20	20
7	8	1	5.4%	88.9%	0.004%	3	20	20

FIGURE 4.512T

Index Language I.i.a (S.T. Natural language. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule C  
 Document Relevance 1 - 4  
 Number of Documents in Collection 1,400  
 Number of Questions 20 (Subset 6)  
 Number of Relevant Documents see table  
 Generality Number see table

Coordination Level	Relevant Documents	Generality Number	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio	x	y	z
			Rel.	Non-rel.						
1	147	5.3	128	(-)	87.1%	(-)	(-)	20	0	20
2	147	5.3	96	1,809*	65.3%	5.0%*	6.495%*	20	13*	20
3	119	5.0	57	533	47.9%	9.7%	2.251%	17	17	17
4	94	4.8	26	63	27.7%	29.2%	0.323%	9	14	14
5	43	4.4	7	8	16.3%	46.7%	0.082%	2	7	7
6	21	5.0	1	2	4.8%	33.3%	0.048%	1	3	3

FIGURE 4.513T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)

Exhaustivity of Indexing 3

Search Rule C

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 20 (Subset 6)

Number of Relevant Documents see table

Generality Number see table

Coordination Level	Relevant Documents	Generality Number	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio	x	y	z
			Rel.	Non-rel.						
1	147	5.3	140	(-)	95.2%	(-)	(-)	20	0	20
2	147	5.3	113	3,323*	76.9%	3.3%*	11.930%*	20	13*	20
3	119	5.0	68	1,544	57.1%	4.2%	6.519%	17	17	17
4	94	4.8	38	299	40.4%	11.3%	1.533%	11	14	14
5	43	4.4	12	33	27.9%	26.7%	0.338%	6	7	7
6	21	5.0	2	2	9.5%	50.0%	0.048%	1	3	3

FIGURE 4.514T

Index Language I.1.a (S.T. Natural language. Coordination)

Exhaustivity of Indexing 3

Search Rule D

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 20 (Subset 6)

Number of Relevant Documents see table

Generality Number see table

Coordination Level	Relevant Documents	Generality Number	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio	x	y	z
			Rel.	Non-rel.						
1	147	5.3	97	2,213	66.0%	4.5%	7.945%	20	20	20
2	147	5.3	73	653	49.7%	10.1%	2.344%	20	20	20
3	119	5.0	42	163	35.3%	20.5%	0.688%	15	17	17
4	94	4.8	20	14	21.3%	58.8%	0.072%	8	14	14
5	43	4.4	7	4	16.3%	63.6%	0.041%	2	7	7
6	21	5.0	1	2	4.8%	33.3%	0.048%	1	3	3

FIGURE 4.515T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)

Exhaustivity of Indexing 3

Search Rule D

Document Relevance 1 - 4

Number of Documents in Collection 1,400

Number of Questions 20 (Subset 6)

Number of Relevant Documents see table

Coordination Level	Relevant Documents	Generality Number	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio	x	y	z
			Rel.	Non-rel.						
1	147	5.3	102	2,668	69.4%	3.7%	9.579%	20	20	20
2	147	5.3	81	791	55.1%	9.3%	2.840%	20	20	20
3	119	5.0	48	189	40.3%	20.3%	0.798%	15	17	17
4	94	4.8	25	21	26.6%	54.3%	0.108%	11	14	14
5	43	4.4	10	6	23.3%	62.5%	0.061%	4	7	7
6	21	5.0	1	2	4.8%	33.3%	0.048%	1	3	3

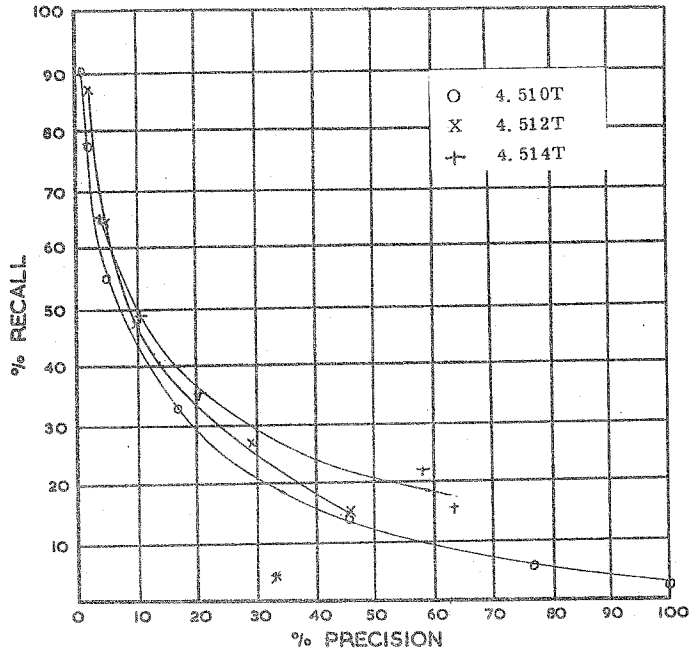


FIGURE 4.516P SEARCH RULES A, C, and D, 20 QUESTIONS, INDEX LANGUAGE I.1.a (FIGURES 4.510T, 4.512T and 4.514T)

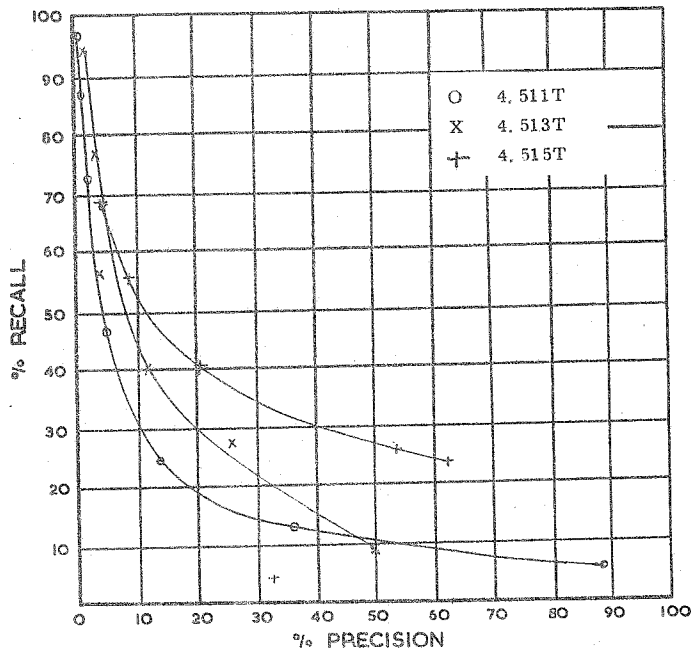


FIGURE 4.517P SEARCH RULES, A, C and D, 20 QUESTIONS, INDEX LANGUAGE I.6.a (FIGURES 4.511T, 4.513T and 4.515T)

Section 6 Document Relevance

The four grades of document relevance were first tested on the 35 questions (subset 1), and results are given for languages I.1.a and I.6.a. Figs. 4.600T - 4.605T give the tables of results, and Figs. 4.606P - 4.607P the plots of performance curves. The change in the different relevance decisions causes a change in the generality number, which is shown in each table, and therefore recall/fallout plots are presented. Although the set of questions used was typical, the test results of the two higher grades of relevance (1-2 and 1) suffer from the fact that not all questions had any relevant documents of these grades. 27 questions had relevant documents when grade 1-2 was tested, but only 11 questions when grade 1 was tested. All 35 questions, however, were kept in the set and the non-relevant documents retrieved by these questions were still counted. It can be noted that the change in document relevance from low to high merely transfers some of the relevant documents from the relevant retrieved to the non-relevant retrieved category, while the total documents retrieved always stays constant.

Because of the small total of documents having relevance 1 documents, a further test was made on a set of fifty questions, for which the criterion of selection was that each question must have a relevance 1 document. These were tested on the 1400 document collection with index language I.1.a and the results are shown in Figs. 4.610T - 4.613T, with a recall/fallout plot as Figure 4.614P.

LIST OF FIGURES

	Relevance	Index Language	No. of Questions	Question Subset	Document Collection	Plots
4.600T	1-3	I.1.a	35	1	1400	
4.601T	1-3	I.6.a	35	1	1400	
4.602T	1-2	I.1.a	35	1	1400	
4.603T	1-2	I.6.a	35	1	1400	
4.604T	1	I.1.a	35	1	1400	
4.605T	1	I.6.a	35	1	1400	
4.606P	1-4	I.1.a	35	1	1400	Plot 4.110T
	1-3					4.600T
	1-2					4.602T
	and 1					4.604T
4.607P	1-4	I.6.a	35	1	1400	Plot 4.114T
	1-3					4.601T
	1-2					4.603T
	and 1					4.605T
4.610T	1-4	I.1.a	50	9	1400	
4.611T	1-3	I.1.a	50	9	1400	
4.612T	1-2	I.1.a	50	9	1400	
4.613T	1	I.1.a	50	9	1400	
4.614P	1-4	I.1.a	50	9	1400	Plot 4.610T
	1-3					4.611T
	1-2					4.612T
	1					4.613T

FIGURE 4.600T

Index Language I.1.a (S.T. Natural language. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 3

Number of Documents in Collection 1,400

Number of Questions 35 (subset 1)

Number of Relevant Documents 212

Generality Number 4.3

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	194	23,755	91.5%	0.6%	48.520%	35	35	35
2	158	8,151	73.8%	1.9%	19.259%	35	35	35
3	112	2,914	52.8%	3.7%	5.965%	34	35	35
4	69	625	32.5%	8.9%	1.281%	29	35	35
5	34	150	18.0%	18.5%	0.307%	13	35	35
6	17	43	8.0%	28.3%	0.088%	7	35	35
7	6	10	2.8%	37.5%	0.020%	3	35	35

FIGURE 4.601T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 3

Number of Documents in Collection 1,400

Number of Questions 35 (Subset 1)

Number of Relevant Documents 212

Generality Number 4.3

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	205	(-)	96.7%	(-)	(-)	35	0	35
2	186	19,333*	87.7%	1.0%*	39.455%*	35	23*	35
3	146	8,126	68.9%	1.8%	16.656%	35	35	35
4	93	2,463	43.8%	3.6%	5.040%	35	35	35
5	50	593	23.6%	7.8%	1.215%	30	35	35
6	28	119	13.2%	19.0%	0.244%	18	35	35
7	10	22	4.7%	31.3%	0.045%	8	35	35

FIGURE 4.602T

Index Language I.1.a (S.T. Natural language. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 2

Number of Documents in Collection 1,400

Number of Questions 35 (8 questions had no relevant documents)

Number of Relevant Documents 79

Generality Number 1.8

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	73	23,876	92.4%	0.3%	48.726%	35	35	35
2	55	8,252	89.6%	0.7%	20.176%	35	35	35
3	43	2,983	54.4%	1.4%	6.088%	34	35	35
4	30	664	38.0%	4.3%	1.357%	28	35	35
5	16	177	20.3%	8.3%	0.361%	18	35	35
6	9	51	11.4%	15.0%	0.104%	7	35	35
7	3	13	3.8%	18.8%	0.027%	3	35	35

FIGURE 4.603T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 2

Number of Documents in Collection 1,400

Number of Questions 35 (8 questions had no relevant documents)

Number of Relevant Documents 79

Generality Number 1.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	78	(-)	98.7%	(-)	(-)	35	0	35
2	70	19,449*	88.6%	0.4%*	39.691%*	35	23*	35
3	53	8,218	67.1%	0.6%	16.801%	35	35	35
4	37	2,518	46.8%	1.4%	5.149%	35	35	35
5	24	810	30.4%	3.7%	1.265%	30	35	35
6	17	130	21.5%	11.6%	0.266%	18	35	35
7	8	26	7.6%	18.8%	0.053%	6	35	35

FIGURE 4.604T

Index Language I.1.a (S.T. Natural language. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1

Number of Documents in Collection 1,400

Number of Questions 35 (24 questions had no relevant documents)

Number of Relevant Documents 18

Generality Number 0.4

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	17	23,932	94.4%	0.1%	48.840%	35	35	35
2	16	8,291	88.9%	0.2%	16.920%	35	35	35
3	12	3,014	66.7%	0.4%	6.145%	34	35	35
4	8	686	44.4%	1.2%	1.401%	28	35	35
5	5	188	27.8%	2.6%	0.384%	18	35	35
6	2	58	11.1%	3.3%	0.118%	7	35	35
7	2	14	11.1%	12.5%	0.022%	3	35	35

FIGURE 4.605T

Index Language I.6.a (S.T. Synonyms, Quasi-synonyms, Word forms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1

Number of Documents in Collection 1,400

Number of Questions 35 (24 questions had no relevant documents)

Number of Relevant Documents 18

Generality Number 0.4

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	18	(-)	100.0%	(-)	(-)	35	0	35
2	16	19,503*	88.9%	0.1%*	39.802%	35	23*	35
3	13	8,259	72.2%	0.2%	16.861%	35	35	35
4	10	2,548	55.6%	0.3%	5.198%	35	35	35
5	6	636	50.0%	1.4%	1.298%	30	35	35
6	6	141	33.3%	4.1%	0.288%	18	35	35
7	3	29	16.7%	9.4%	0.058%	6	35	35

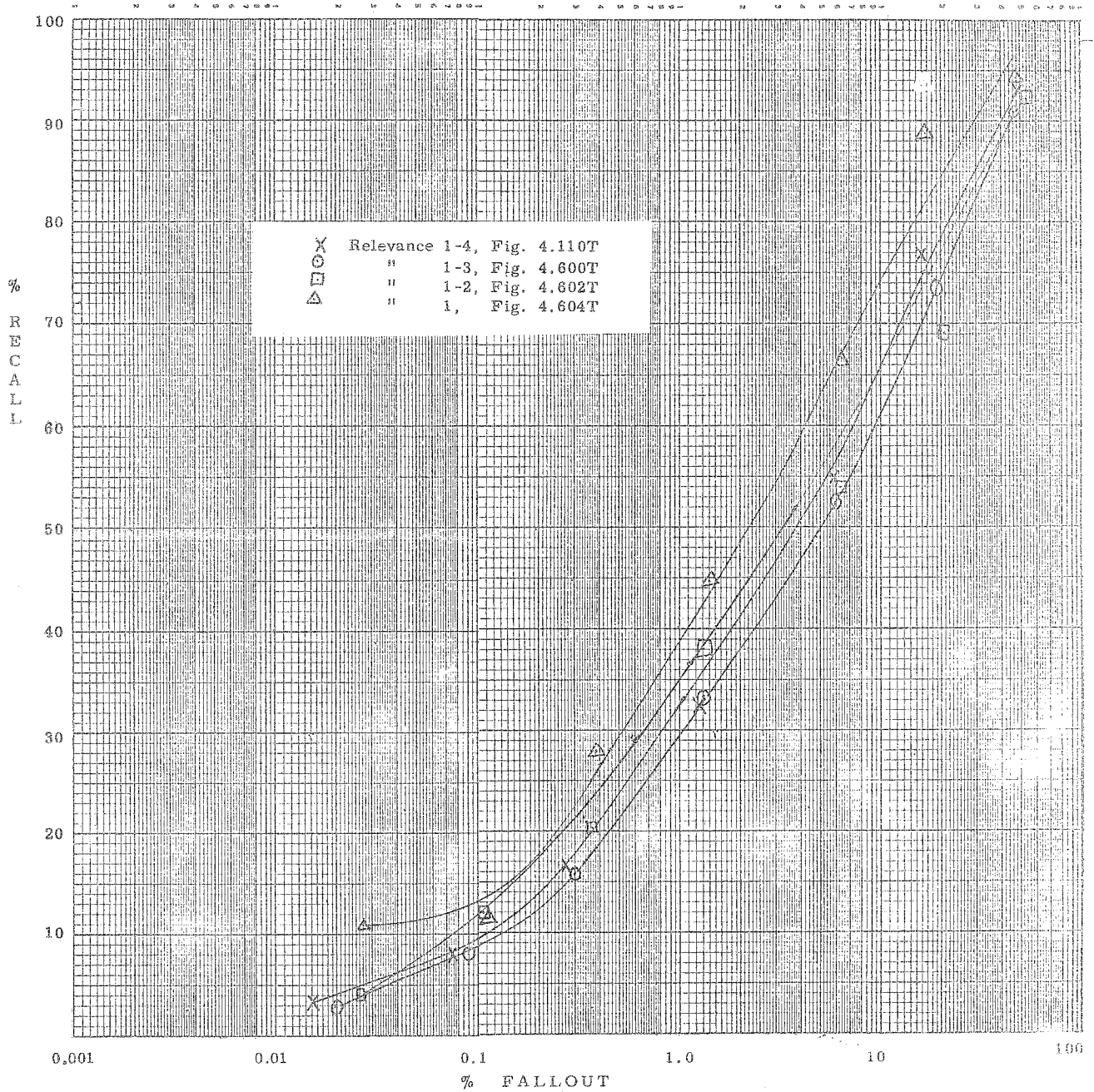


FIGURE 4.606P RECALL-FALLOUT PLOT SHOWING EFFECT OF CHANGES IN RELEVANCE ON 42 QUESTIONS WITH 1400 DOCUMENT COLLECTION, INDEX LANGUAGE I.1.a

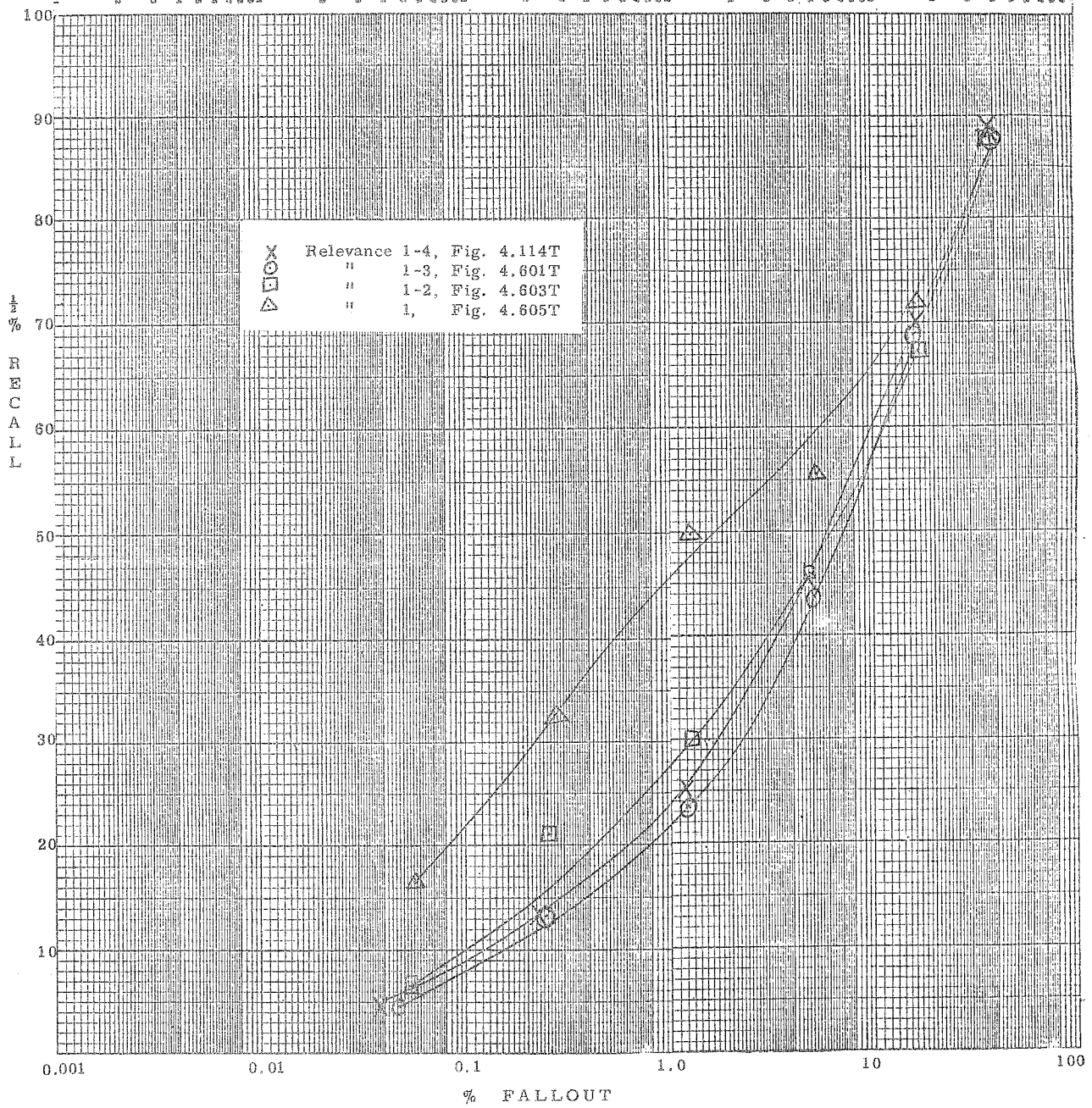


FIGURE 4.607P RECALL-FALLOUT PLOT SHOWING EFFECT OF CHANGES IN RELEVANCE ON 42 QUESTIONS WITH 1400 DOCUMENT COLLECTION, INDEX LANGUAGE I.S.a.

FIGURE 4.610T

Index Language 1.1.a

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1-4

Number of Documents in Collection 1400

Number of Questions 50

Number of Relevant Documents 361

Generality Number 5.2

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	340	30,778	94.2%	1.1%	44.196%	50	50	50
2	281	11,572	77.8%	2.4%	16.617%	50	50	50
3	176	4,717	48.8%	3.6%	6.759%	47	50	50
4	107	1,727	29.6%	5.8%	2.480%	39	46	46
5	59	717	16.3%	7.6%	1.015%	21	40	40
6	35	273	9.7%	11.4%	0.392%	12	30	30
7	19	80	5.3%	19.2%	0.115%	8	23	23
8	7	17	1.9%	29.2%	0.024%	4	12	12
9	5	2	1.4%	71.4%	0.003%	3	10	10
10	1	0	0.3%	100.0%	0.000%	1	8	8
11	0	0				0	6	6
12	0	0				0	1	1

FIGURE 4.611T

Index Language 1.1.a

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1-3

Number of Documents in Collection 1400

Number of Questions 50

Number of Relevant Documents 297

Generality Number 4.2

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	277	30,841	93.2%	0.9%	44.246%	50	50	50
2	235	11,618	79.1%	2.0%	16.651%	50	50	50
3	162	4,731	54.5%	3.4%	6.787%	47	50	50
4	97	1,737	32.7%	5.3%	2.492%	39	46	46
5	49	724	16.5%	6.4%	1.039%	21	40	40
6	29	279	9.8%	9.2%	0.400%	12	30	30
7	16	83	5.3%	16.2%	0.119%	8	23	23
8	6	18	2.0%	25.0%	0.026%	4	12	12
9	4	3	1.3%	57.1%	0.004%	3	10	10
10	1	0	0.3%	100.0%	0.000%	1	8	8
11	0	0				0	6	6
12	0	0				0	1	1

FIGURE 4.612T

Index Language I.1.a  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1-2  
 Number of Documents in Collection 1400  
 Number of Questions 50  
 Number of Relevant Documents 155  
 Generality Number 2.2

Coord- ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	146	30,918	94.2%	0.5%	44.267%	50	50	50
2	125	11,726	80.6%	1.0%	16.789%	50	50	50
3	87	4,806	56.1%	1.8%	6.958%	47	50	50
4	58	1,776	37.4%	3.2%	2.543%	39	46	46
5	33	740	21.3%	4.3%	1.059%	21	40	40
6	21	287	13.5%	6.9%	0.411%	12	30	30
7	10	89	6.5%	10.2%	0.127%	8	23	23
8	6	18	3.9%	25.0%	0.026%	4	12	12
9	4	3	2.6%	57.1%	0.004%	3	10	10
10	1	0	0.6%	100.0%	0.000%	1	8	8
11	0	0				0	6	6
12	0	0				0	1	1

FIGURE 4.613 T

Index Language I.1.a  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1  
 Number of Documents in Collection 1400  
 Number of Questions 50  
 Number of Relevant Documents 95  
 Generality Number 1.4

Coord- ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	90	30,974	94.7%	0.3%	44.309%	50	50	50
2	81	11,770	85.3%	0.7%	16.837%	50	50	50
3	57	4,836	60.0%	1.2%	6.918%	47	50	50
4	40	1,794	42.1%	2.2%	2.566%	39	46	46
5	24	749	25.3%	3.1%	1.071%	21	40	40
6	14	294	14.7%	4.5%	0.420%	12	30	30
7	7	92	7.4%	7.1%	0.153%	8	23	23
8	3	21	3.2%	14.3%	0.030%	4	12	12
9	2	5	2.1%	28.6%	0.007%	3	10	10
10	0	1	0.0%	0.0%	0.001%	1	8	8
11	0	0				0	6	6
12	0	0				0	1	1

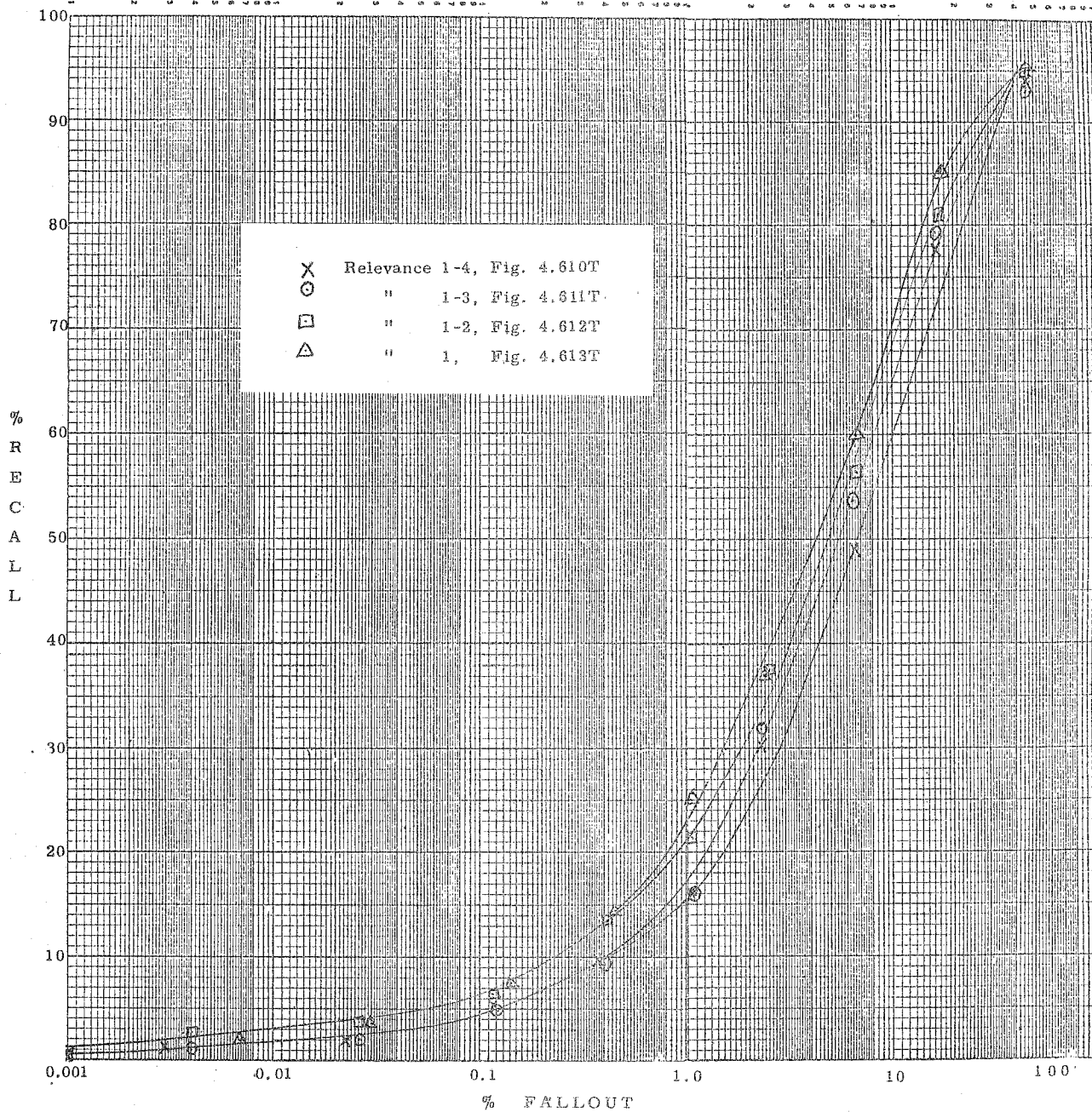


FIGURE 4.614P RECALL-FALLOUT PLOT SHOWING EFFECT OF CHANGES IN RELEVANCE ON 50 QUESTIONS WITH 1400 DOCUMENT COLLECTION, INDEX LANGUAGE 1.1.a

Section 7 Simple Concepts:- Recall Devices

The simple concept languages are tested by the 42 questions (subset 2) searched on the 200 collection (subset 1). Fifteen separate languages are presented in the tables, 4.700T - 4.714T. Index language I.1.a uses the basic natural terms, and index language 1.2.a incorporates the recall device of synonym control. All subsequent index languages are based on index language I.2.a,

Index languages I.3.a, I.4.a, I.5.a, I.6.a, I.7.a and I.8.a investigate different selections of classes based on the schedules as given in Appendix 5.4 of Volume I. Index languages I.9.a, I.10.a and I.11.a represent selections from the alphabetical index as given in Appendix 5.5 of Volume I. Index languages I.12.a, I.13.a and I.14.a investigate the complete classes of the schedules of Appendix 5.4, while index language I.15.a brings in both the alphabetical and the hierarchical languages.

A table presenting the fifteen languages and their relationships appears as Figure 2.6. (page 10).

Seven graphs are included in this section, each plotting two or more of the sets of test results; these graphs fall into three main groups and provide comparison of

- (1) the different hierarchical devices, including the selective and obligatory types, (Figs. 4.715P, 4.716P, 4.718P, 4.719P and 4.721P).
- (2) the alphabetical device, (Figs. 4.717P and 4.721P).
- (3) hierarchical and alphabetical devices (Figs. 4.719P and 4.721P).

LIST OF FIGURES

	Index Language	No. of Questions	Document Subset	Document Collection
4.700T	II.1.a	42	2	200
4.701T	II.2.a	42	2	200
4.702T	II.3.a	42	2	200
4.703T	II.4.a	42	2	200
4.704T	II.5.a	42	2	200
4.705T	II.6.a	42	2	200
4.706T	II.7.a	42	2	200
4.707T	II.8.a	42	2	200
4.708T	II.9.a	42	2	200
4.709T	II.10.a	42	2	200
4.710T	II.11.a	42	2	200
4.711T	II.12.a	42	2	200

	Index Language	No. of Questions	Document Subset	Document Collection	Plot
4.712T	II.13.a	42	2	200	
4.713T	II.14.a	42	2	200	
4.714T	II.15.a	42	2	200	
4.715P	II.1.a, II.2.a, II.3.a, II.4.a,	42	2	200	Plots 4.700T, 4.701T, 4.702T, 4.703T.
4.716P	II.2.a, II.6.a, II.7.a,	42	2	200	Plots 4.701T, 4.705T, 4.706T.
4.717P	II.2.a, II.9.a, II.10.a,	42	2	200	Plots 4.701T, 4.708T, 4.709T.
4.718P	II.2.a, II.12.a, II.13.a, II.14.a,	42	2	200	Plots 4.701T, 4.711T, 4.712T, 4.713T.
4.719P	II.5.a, II.6.a, II.7.a, II.9.a,	42	2	200	Plots 4.704T, 4.705T, 4.706T, 4.708T.
4.720P	II.8.a, II.10.a,	42	2	200	Plots 4.707T, 4.709T.
4.721P	II.8.a, II.11.a, II.14.a, II.15.a,	42	2	200	Plots 4.707T, 4.710T, 4.713T, 4.714T.

FIGURE 4.700T

Index Language II.1.a (Simple concepts. Natural language. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	73	459	36.9%	13.7%	5.597%	42	42	42
2	28	27	14.1%	50.9%	0.329%	18	41	41
3	10	0	5.1%	100.0%	0.000%	5	38	38
4	0	0				0	32	32
5	0	0				0	20	20
6	0	0				0	9	9
7	0	0				0	3	3

FIGURE 4.701T

Index Language II.2.a (Simple Concepts. Synonyms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	101	635	51.0%	13.7%	7.743%	42	42	42
2	44	65	22.2%	40.4%	0.792%	25	41	41
3	10	2	5.1%	83.3%	0.024%	6	38	38
4	2	0	1.0%	100.0%	0.000%	2	32	32
5	0	0				0	20	20
6	0	0				0	9	9
7	0	0				0	3	3

FIGURE 4.702T

Index Language II.3.a (Simple Concepts. Species (selected). Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	115	1,010	58.1%	10.2%	12.317%	42	42	42
2	60	115	30.3%	34.3%	1.890%	32	41	41
3	24	8	12.1%	75.0%	0.098%	12	38	38
4	4	0	2.0%	100.0%	0.000%	3	32	32
5	0	0				0	20	20
6	0	0				0	9	9
7	0	0				0	3	3

FIGURE 4.703T

Index Language II.4.a (Simple concepts. Superordinate (Selected). Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	129	1,336	65.2%	8.8%	18.292%	42	42	42
2	67	255	33.8%	20.8%	3.109%	36	41	41
3	19	40	9.5%	32.2%	0.488%	16	38	38
4	2	8	1.0%	20.0%	0.098%	4	32	32
5	0	2	0.0%	0.0%	0.024%	1	20	20
6	0	1	0.0%	0.0%	0.012%	1	9	9
7	0	0				0	3	3

FIGURE 4.704T

Index Language II.5.a (Simple concepts. Species and Superordinate (Selected). Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 300 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	138	1,614	69.7%	7.8%	19.682%	42	42	42
2	81	359	40.9%	18.4%	4.377%	39	41	41
3	32	60	16.2%	34.8%	0.732%	22	38	38
4	6	10	3.0%	37.5%	0.122%	5	32	32
5	0	4	0.0%	0.0%	0.049%	1	20	20
6	0	1	0.0%	0.0%	0.012%	1	9	9
7	0	0				0	3	3

FIGURE 4.705T

Index Language II.6.a (Simple Concepts. Co-ordinate (Selected). Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	118	1,124	59.6%	9.5%	13.707%	42	42	42
2	57	155	28.8%	26.8%	1.890%	33	41	41
3	18	26	9.1%	40.9%	0.219%	11	38	38
4	4	5	2.0%	44.4%	0.061%	4	32	32
5	0	0				0	20	20
6	0	0				0	9	9
7	0	0				0	3	3

FIGURE 4.706T

Index Language II.7.a (Simple Concepts. Co-ordinate and Collateral (Selected). Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	137	1,532	69.2%	8.2%	18.682%	42	42	42
2	75	318	37.9%	19.1%	3.877%	33	41	41
3	23	50	11.6%	31.5%	0.810%	14	38	38
4	6	9	3.0%	40.0%	0.110%	4	32	32
5	1	0	0.5%	100.0%	0.000%	1	20	20
6	0	0				0	9	9
7	0	0				0	3	3

FIGURE 4.707T

Index Language II.8.a (Simple Concepts. Species, Superordinate, Coordinate and Collateral (Selected). Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	159	2,133	80.3%	6.9%	28.012%	42	42	42
2	89	605	50.0%	14.1%	7.378%	40	41	41
3	38	123	19.2%	23.8%	1.500%	23	38	38
4	14	33	7.1%	29.8%	0.402%	10	32	32
5	5	8	2.5%	38.5%	0.088%	2	20	20
6	0	1	0.0%	0.0%	0.012%	1	9	9
7	0	0				0	3	3

FIGURE 4.708T

Index Language II.9.a (Simple Concepts. First alphabetical collateral (Selected). Coordination).  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	135	1,524	68.2%	8.1%	18.585%	42	42	42
2	75	396	37.9%	20.2%	3.609%	33	41	41
3	37	45	18.7%	45.1%	0.549%	19	38	38
4	13	6	6.6%	88.4%	0.073%	6	32	32
5	1	2	0.5%	33.6%	0.024%	3	20	20
6	0	0				0	9	9
7	0	0				0	3	3

ABC  
DEF  
SM  
A...c  
A...c  
A...c

FIGURE 4.709T

Index Language II.10.a (Simple Concepts. Second alphabetical collateral (selected), Coordination.)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	182	2,892	91.9%	5.9%	35.268%	42	42	42
2	129	942	65.2%	12.0%	11.485%	32	41	41
3	79	253	40.0%	23.8%	3.085%	34	38	38
4	27	44	13.6%	38.0%	0.536%	18	32	32
5	7	6	3.5%	53.8%	0.073%	6	20	20
6	1	0	0.5%	100.0%	0.000%	1	9	9
7	0	0				0	3	3

FIGURE 4.710T

Index Language II.11.a (Simple Concepts. Species, superordinate, co-ordinate, collateral, and second alphabetical collateral (selected). Coordination.)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	183	3,317	92.4%	5.2%	40.451%	42	42	42
2	146	1,259	73.7%	10.4%	15.350%	41	41	41
3	94	408	47.5%	18.72%	4.974%	37	38	38
4	34	99	17.2%	25.6%	1.207%	19	32	32
5	14	25	7.1%	35.9%	0.305%	9	20	20
6	3	1	1.5%	75.0%	0.012%	3	9	9
7	0	0				0	3	3

FIGURE 4.711T

Index Language II.12.a (Simple Concepts. Species (complete). Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	128	1,432	94.6%	8.2%	17.463%	42	42	42
2	72	287	36.4%	20.1%	3.500%	34	41	41
3	27	32	13.6%	45.8%	0.390%	13	38	38
4	9	3	4.5%	75.0%	0.037%	5	32	32
5	0	0				0	20	20
6	0	0				0	9	9
7	0	0				0	3	3

FIGURE 4.712T

Index Language II.13.a (Simple Concepts. Species, Superordinate and Co-ordinate (complete).  
 Exhaustivity of Indexing 3 Coordination).  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	168	2,490	84.6%	6.3%	30.365%	42	42	42
2	117	937	59.1%	11.1%	11.424%	40	41	41
3	49	244	24.7%	16.7%	2.975%	29	38	38
4	20	54	10.1%	27.0%	0.658%	11	32	32
5	8	21	4.0%	27.6%	0.256%	3	20	20
6	0	1	0.0%	0.0%	0.012%	1	9	9
7	0	0				0	3	3

FIGURE 4.713T

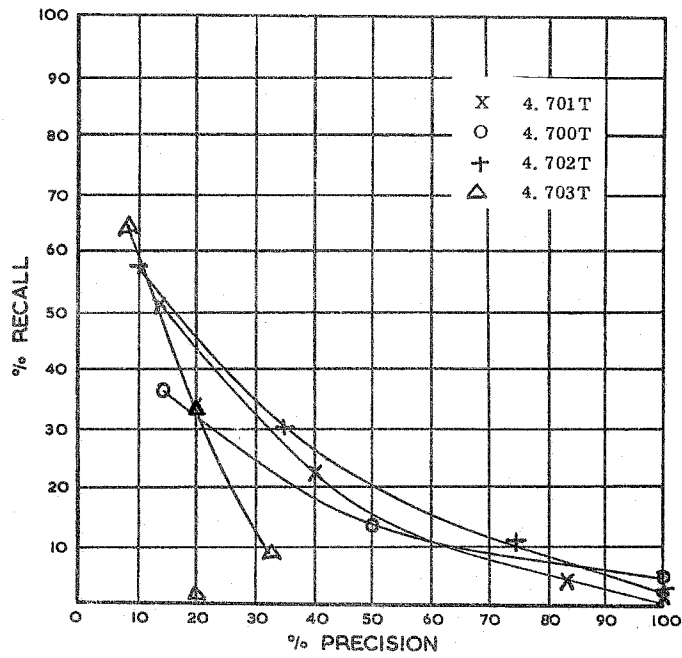
Index Language II.14.a (Simple Concepts. Species, Superordinate, Coordinate and Collateral (complete). Coordination).  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	180	3,684	90.9%	4.6%	44.926%	42	42	42
2	143	2,047	72.2%	6.5%	24.957%	40	41	41
3	66	600	33.3%	9.9%	7.315%	33	38	38
4	26	179	13.1%	12.7%	2.182%	19	32	32
5	13	59	6.6%	18.1%	0.719%	7	20	20
6	0	11	0.0%	0.0%	0.134%	2	9	9
7	0	0				0	3	3

FIGURE 4.714T

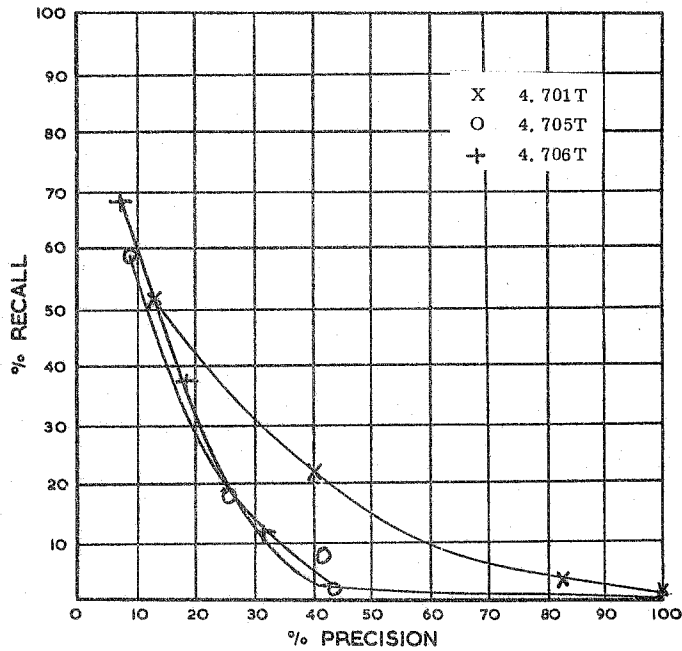
Index Language II.15.a (Simple Concepts. Species, Superordinate, Coordinate, and Collateral (complete), Second alphabetical Collateral (selected). Coordination).  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	190	4,542	96.0%	4.0%	55.390%	42	42	42
2	168	2,590	84.8%	6.1%	31.578%	41	41	41
3	117	1,020	59.1%	10.3%	12.436%	37	38	38
4	51	325	25.8%	13.6%	3.962%	26	32	32
5	21	96	10.6%	17.9%	1.170%	12	20	20
6	4	15	2.0%	21.1%	0.183%	5	9	9
7	0	0				0	3	3



	Rec	Prec
II2	31.8	27.5
II1	22.2	25.7
II3	30.3	31.1
II4	29.3	23.4

FIGURE 4.715P SIMPLE CONCEPT INDEX LANGUAGES II.1.a, II.2.a, II.3.a, and II.4.a (FIGURES 4.700T - 4.703T)



	Rec	Prec
II2	31.8	27.5
II6	21.2	22.7
II7	32.3	20.3

FIGURE 4.716P SIMPLE CONCEPT INDEX LANGUAGES II.2.a, II.6.a and II.7.a (FIGURES 4.701T, 4.705T and 4.706T)

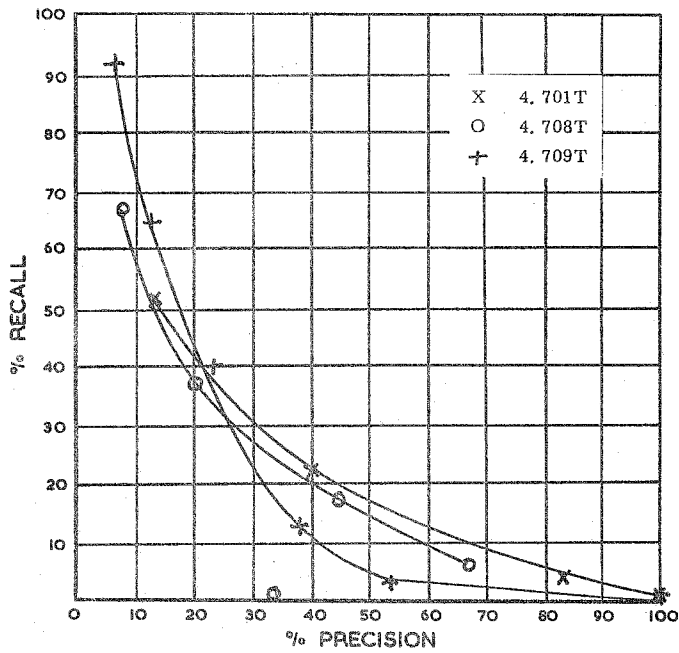


FIGURE 4.717P SIMPLE CONCEPT INDEX LANGUAGES II.2.a, II.9.a, and II.10.a (FIGURES 4.701T, 4.708T and 4.709T)

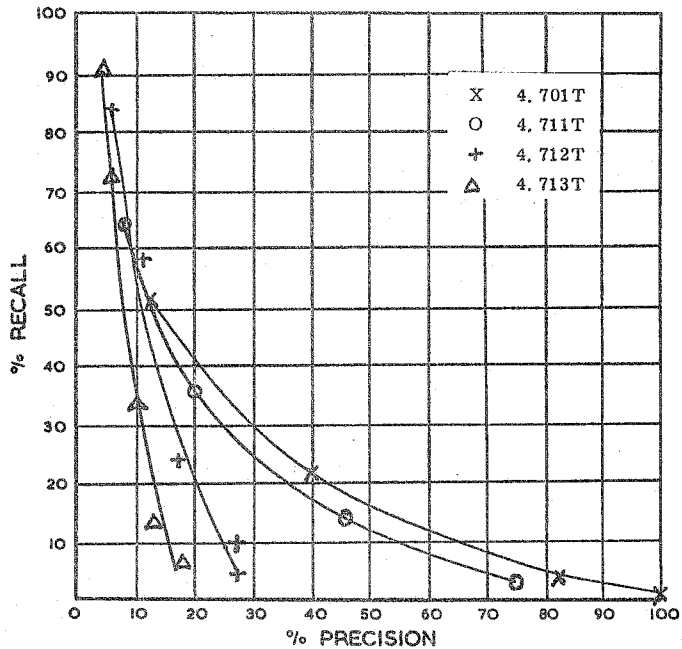
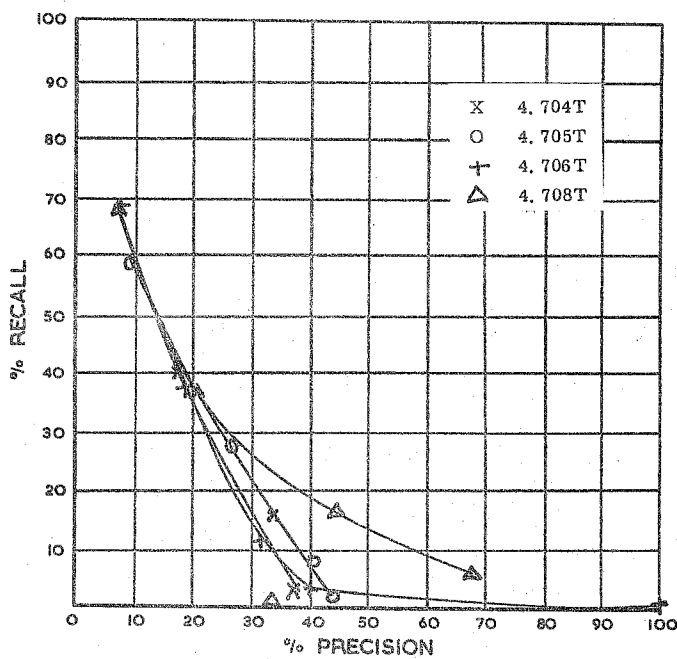


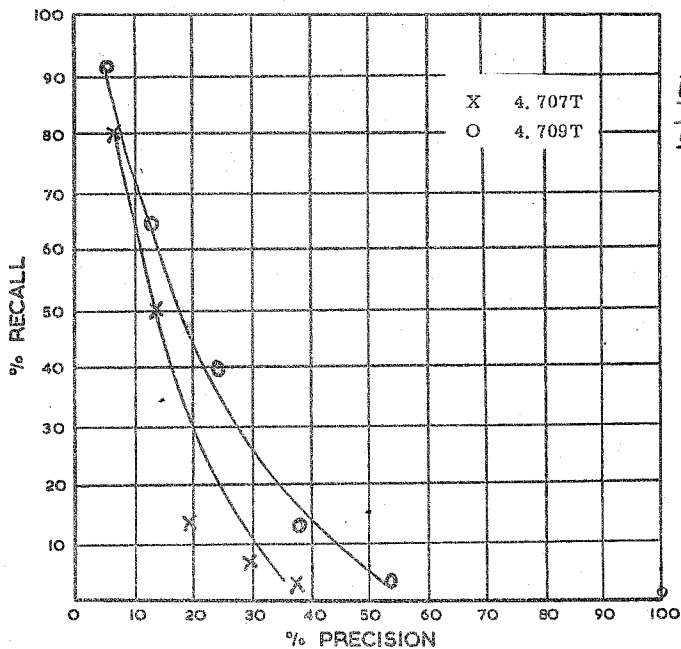
FIGURE 4.718P SIMPLE CONCEPT INDEX LANGUAGES II.2.a, II.12.a, II.13.a and II.14.a (FIGURES 4.701T, 4.711T, 4.712T and 4.713T)



	Rec	Prec
45	30.8	29.8
46	21.2	22.7
47	32.3	20.3
49	31.3	29.7

FIGURE 4.719P

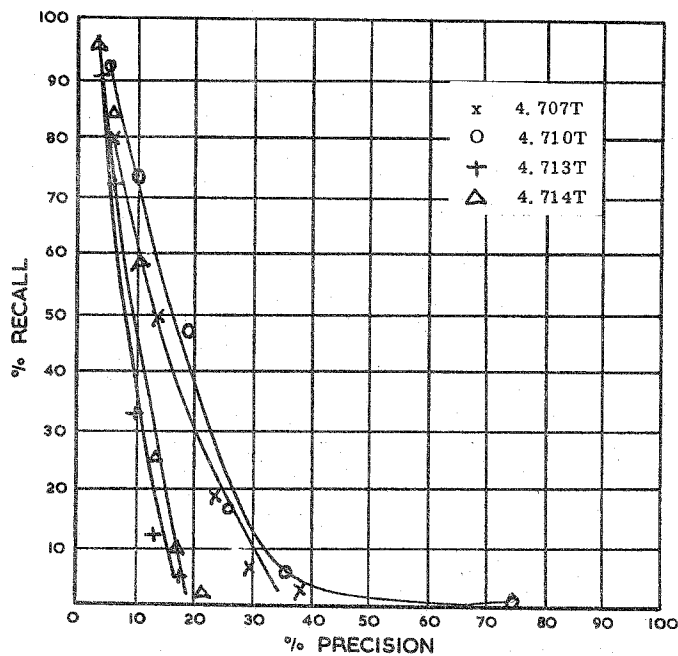
SIMPLE CONCEPT INDEX LANGUAGES II.5.a, II.6.a, II.7.a and II.9.a (FIGURES 4.704T, 4.705T, 4.706T and 4.708T)



	Rec	Prec
48a	27.8	21.8
10a	31.3	22.9

FIGURE 4.720P

SIMPLE CONCEPT INDEX LANGUAGES II.8.a, II.10.a (FIGURES 4.707T and 4.709T)



	Rec	Prec
□ 8	27.8	21.8
□ 11	32.3	22.4
□ 14	27.3	18.9
□ 15	34.8	17.3

FIGURE 4.721P

SIMPLE CONCEPT INDEX LANGUAGES II.8.a, II.11.a, II.14.a and II.15.a (FIGURES 4.707T, 4.710T, 4.713T and 4.714T)

## Section 8 Controlled terms

### Recall devices

The investigations on the controlled term languages were carried out on two sets of questions and two collection sizes: the 42 questions (subset 2) searched on the 200 document collection (subset 1), and the 77 questions (subset 7) searched on the 350 document collection (subset 2). Six recall languages were tested, and the effect of each recall device is seen in the tables and plots given in Figs. 4.800TP - 4.805TP, tested with a search rule A. A plot of the six curves is given in Figure 4.806P. The recall devices are shown in Fig. 2.7 (p.11)

### Search Rules

Further comparison of the six index languages is given with the tests made using Search Rule E on both the 42 questions with the 200 document collection and the 77 questions with the 350 document collection. The results are presented in Figs. 4.810TP - 4.815TP and 4.820T - 4.825T, with plots covering the various index languages as Figs. 4.816P and 4.826P. With Search Rule E all the combinations of acceptable terms were selected for each coordination level. Examples of search formulations for a number of questions are given in Appendix 8.1.

Tests with Search Rule F were done on index languages III.5.a and III.6.a with 42 questions on the 200 document collection. This search was superimposed on Search E, and the results are presented at the various coordination levels according to the number of basic terms as apart from related terms. Thus, at a coordination level of 4, the tables (Figs. 4.850T and 4.851T) show the results when all terms were basic terms, where one term was a related term, where two terms were related terms and so on. It is obvious that when all the terms are basic terms, then the results must be the same as for index language III.1.a with Search E; when all the terms are related terms, then the results must be the same as for the corresponding index language with the basic Search E. Therefore in the plot 4.850P, the two main curves represent index languages III.1.a and III.5.a with Search E and in Figure 4.851P the main curves represent index languages III.1.a and III.6.a, again with Search E. The additional results obtained with Search F now produce a series of secondary curves at each coordination level which span the main curves.

### Precision Device

The precision device of weighting was tested, in which search questions were weighted and tested on the most exhaustive index language, using the weights assigned in indexing. Figs. 4.830TP and 4.831TP give results for languages III.1.a and III.6.a respectively, and in the plots a comparison is made of the weighted and unweighted searches, carried out with search rule E.

SECTION 8 LIST OF FIGURES

	Index Language	Search Rule	No. of Questions	Document Collection	Plot
4.800TP	III.1.a	A	42 (2)	200	
4.801TP	III.2.a	A	42 (2)	200	
4.802TP	III.3.a	A	42 (2)	200	
4.803TP	III.4.a	A	42 (2)	200	
4.804TP	III.5.a	A	42 (2)	200	
4.805TP	III.6.a	A	42 (2)	200	
4.806P	III.1.a, III.2.a, III.3.a, III.4.a, III.5.a, III.6.a	A	42 (2)	200	4.800T, 4.801T, 4.802T, 4.803T, 4.804T, 4.805T
4.810TP	III.1.a	E	42 (2)	200	
4.811TP	III.2.a	E	42 (2)	200	
4.812TP	III.3.a	E	42 (2)	200	
4.813TP	III.4.a	E	42 (2)	200	
4.814TP	III.5.a	E	42 (2)	200	
4.815TP	III.6.a	E	42 (2)	200	
4.816P	III.1.a, III.2.a, III.3.a, III.4.a, III.5.a, III.6.a	E	42 (2)	200	4.810T, 4.811T, 4.812T, 4.813T, 4.814T, 4.815T
4.820T	III.1.a	E	77 (7)	350	
4.821T	III.2.a	E	77 (7)	350	
4.822T	III.3.a	E	77 (7)	350	
4.823T	III.4.a	E	77 (7)	350	
4.824T	III.5.a	E	77 (7)	350	
4.825T	III.6.a	E	77 (7)	350	
4.826P	III.1.a, III.2.a, III.3.a, III.4.a, III.5.a, III.6.a	E	77 (7)	350	4.820T, 4.821T, 4.822T, 4.823T, 4.824T, 4.825T
4.830TP	III.1.e (weighting)	E	42 (2)	200	+ 4.810T
4.831TP	III.6.e (weighting)	E	42 (2)	200	+ 4.815T
4.840P	III.1.a	A & E	42 (2)	200	4.800T, 4.810T
4.841P	III.2.a	A & E	42 (2)	200	4.801T, 4.811T
4.842P	III.3.a	A & E	42 (2)	200	4.802T, 4.812T
4.843P	III.4.a	A & E	42 (2)	200	4.803T, 4.813T
4.844P	III.5.a	A & E	42 (2)	200	4.804T, 4.814T
4.845P	III.6.a	A & E	42 (2)	200	4.805T, 4.815T
4.850TP	III.5.a	F	42 (2)	200	+ 4.814T
4.851TP	III.6.a	F	42 (2)	200	+ 4.815T

FIGURE 4.800T

Index Language III.1.a (Controlled terms. Basic terms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	174	2,884	87.9%	5.6%	35.170%	42	42	42
2	136	946	68.7%	12.6%	11.534%	42	42	42
3	79	213	39.9%	27.1%	2.597%	31	41	41
4	36	49	18.2%	42.4%	0.597%	20	34	34
5	17	7	8.6%	70.8%	0.085%	6	24	24
6	5	3	2.5%	62.5%	0.037%	3	13	13
7	2	2	1.0%	50.0%	0.024%	2	8	8
8	1	1	0.5%	50.0%	0.012%	2	4	4
9	0	0				0	3	3

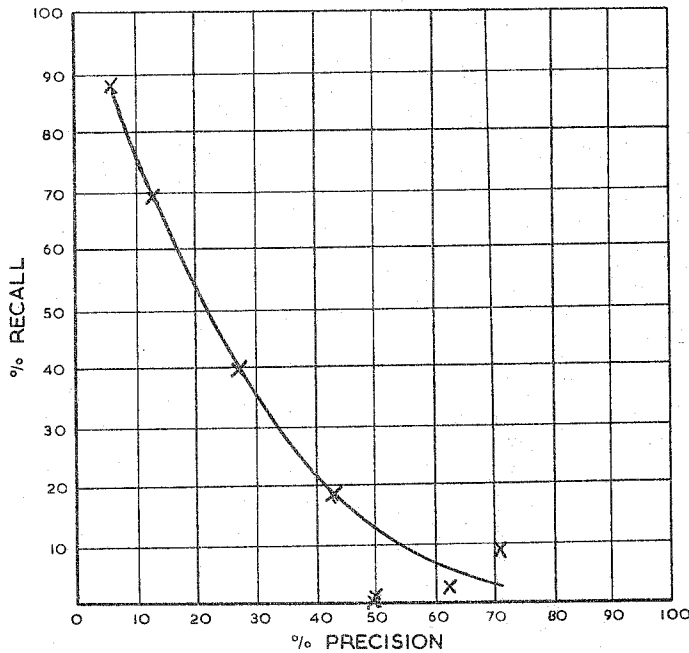


FIGURE 4.800P INDEX LANGUAGE III.1.a

FIGURE 4.801T

Index Language III.2.a (Controlled terms. Narrower terms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	176	3,062	88.9%	5.4%	37.34%	42	42	42
2	137	1,024	69.2%	11.8%	12.485%	42	42	42
3	83	232	41.9%	26.3%	2.829%	31	41	41
4	38	53	19.2%	41.8%	0.646%	20	34	34
5	17	9	8.6%	65.4%	0.110%	6	24	24
6	6	3	3.0%	66.7%	0.037%	3	13	13
7	2	2	1.0%	50.0%	0.012%	2	8	8
8	1	1	0.5%	50.0%	0.012%	2	4	4
9	0	0				0	3	3

4846

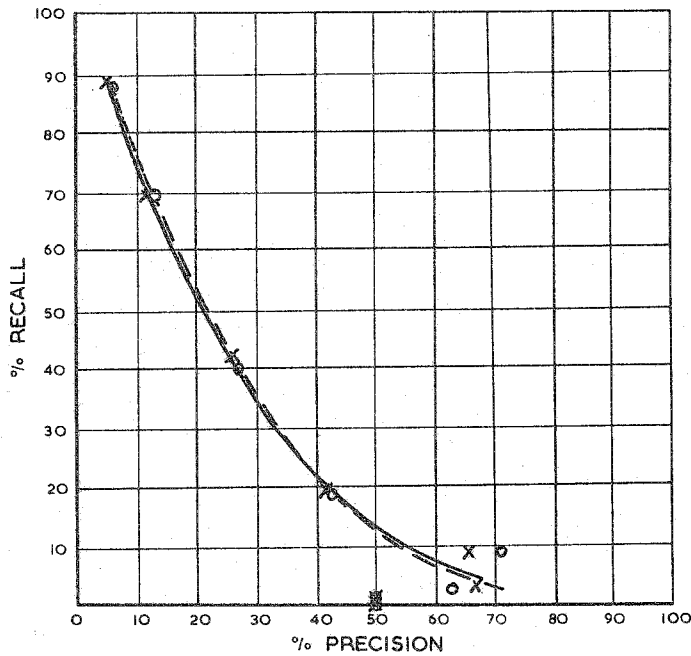


FIGURE 4.801P INDEX LANGUAGE III.2.a  
 (Index Language III.1.a Broken line)

FIGURE 4.802T

Index Language III.3.a (Controlled terms. Broader Terms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	178	3,565	89.9%	4.7%	43.475%	42	42	42
2	147	1,575	74.2%	8.5%	19.203%	42	42	42
3	90	364	45.5%	19.8%	4.438%	37	41	41
4	43	78	21.7%	35.5%	0.951%	22	34	34
5	18	8	9.1%	69.2%	0.098%	6	24	24
6	7	3	3.7%	70.0%	0.037%	4	13	13
7	2	2	1.0%	50.0%	0.024%	4	8	8
8	1	1	0.5%	50.0%	0.012%	2	4	4
9	0	0				0	3	3

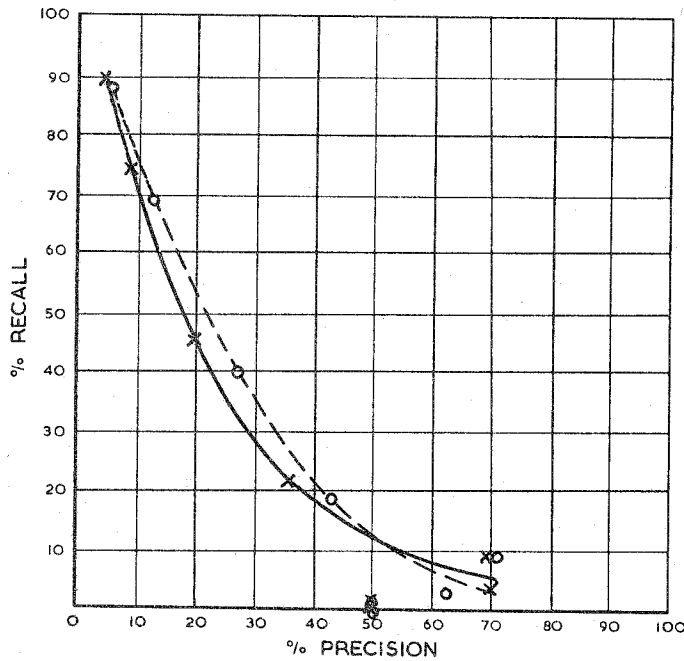


FIGURE 4.802P INDEX LANGUAGE III.3.a  
(Index Language III.1.a Broken line)

FIGURE 4.803T

Index Language III.4.a (Controlled terms. Narrower and Broader terms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	180	3,614	90.9%	4.7%	44.073	42	42	42
2	148	1,661	74.7%	8.2%	20.251%	42	42	42
3	94	386	47.5%	19.6%	4.706%	37	41	41
4	44	32	22.2%	34.9%	1.000%	23	34	34
5	18	12	9.1%	60.0%	0.146%	8	24	24
6	8	3	4.0%	72.7%	0.037%	4	13	13
7	2	2	1.0%	50.0%	0.024%	2	8	8
8	1	1	0.5%	50.0%	0.012%	2	4	4
9	0	0				0	3	3

6256

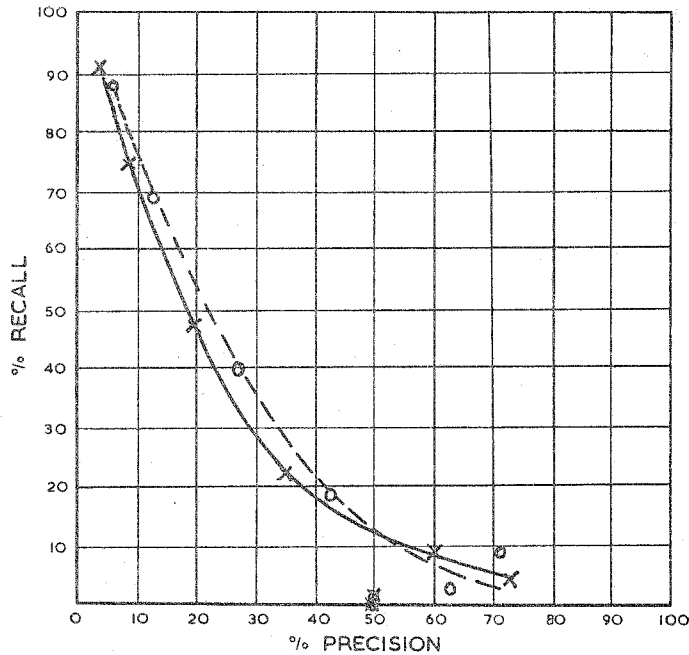


FIGURE 4.803P INDEX LANGUAGE III.4.a  
(Index Language III.1.a Broken line)

FIGURE 4.804T

Index Language III.5.a (Controlled terms. Related terms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents ~~42 (Subset 2)~~ 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	195	5,424	98.5%	3.4%	66.146%	42	42	42
2	184	3,044	92.9%	5.7%	37.113%	42	42	42
3	133	1,252	67.2%	9.6%	15.265%	41	41	41
4	62	420	31.3%	12.9%	5.121%	29	34	34
5	35	122	17.7%	22.3%	1.487%	19	24	24
6	12	26	6.1%	31.6%	0.317%	7	13	13
7	4	7	2.0%	36.4%	0.085%	3	8	8
8	2	3	1.0%	40.0%	0.037%	2	4	4
9	2	1	1.0%	66.7%	0.012%	2	3	3

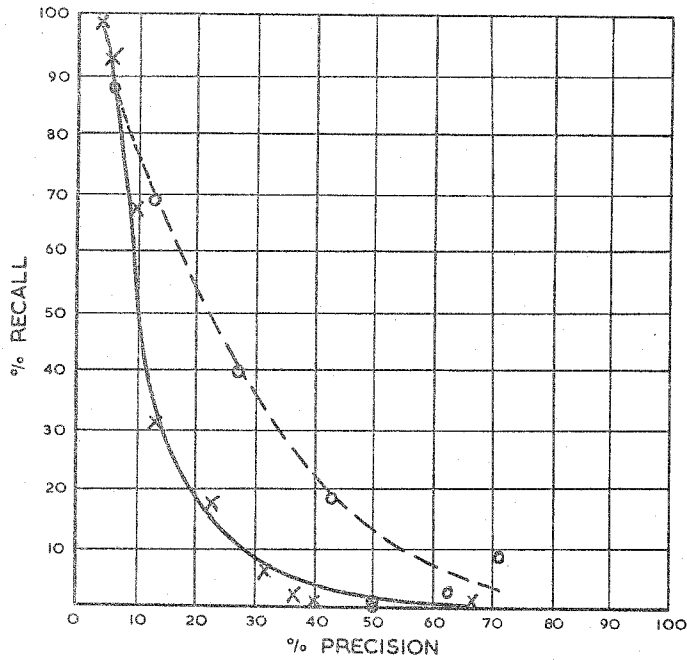


FIGURE 4.804P INDEX LANGUAGE III.5.a  
(Index Language III.1.a Broken line)

FIGURE 4.805T

Index Language III.6.a (Controlled terms. Narrower, broader and related terms. Coordination)

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	196	5,665	99.0%	3.3%	69.085%	42	42	42
2	187	3,483	94.4%	5.1%	42.465%	42	42	42
3	149	1,609	75.3%	8.5%	19.617%	41	41	41
4	72	527	36.4%	12.0%	6.425%	31	34	34
5	38	140	19.2%	21.3%	1.707%	20	24	24
6	13	31	6.6%	29.5%	0.378%	7	13	13
7	4	7	2.0%	36.4%	0.085%	3	8	8
8	2	3	1.0%	40.0%	0.037%	2	4	4
9	2	1	1.0%	66.7%	0.012%	2	3	3

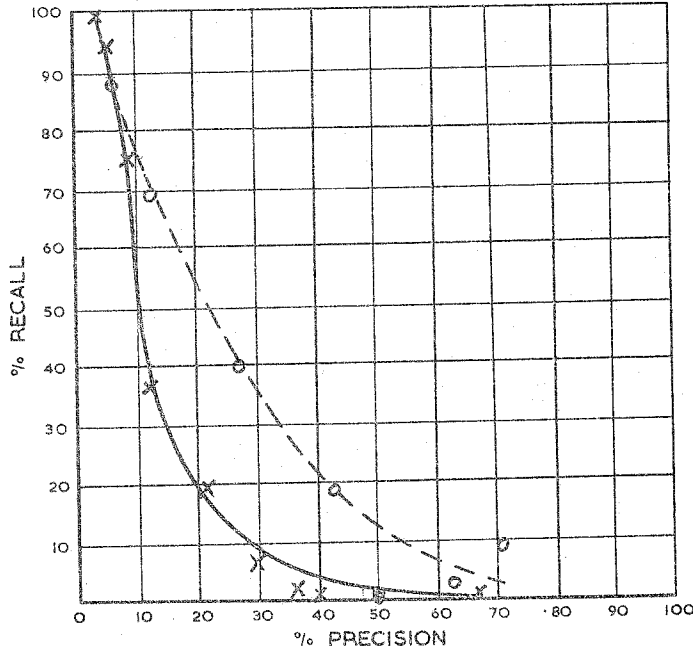


FIGURE 4.805P INDEX LANGUAGE III.6.a  
(Index Language III.1.a Broken line)

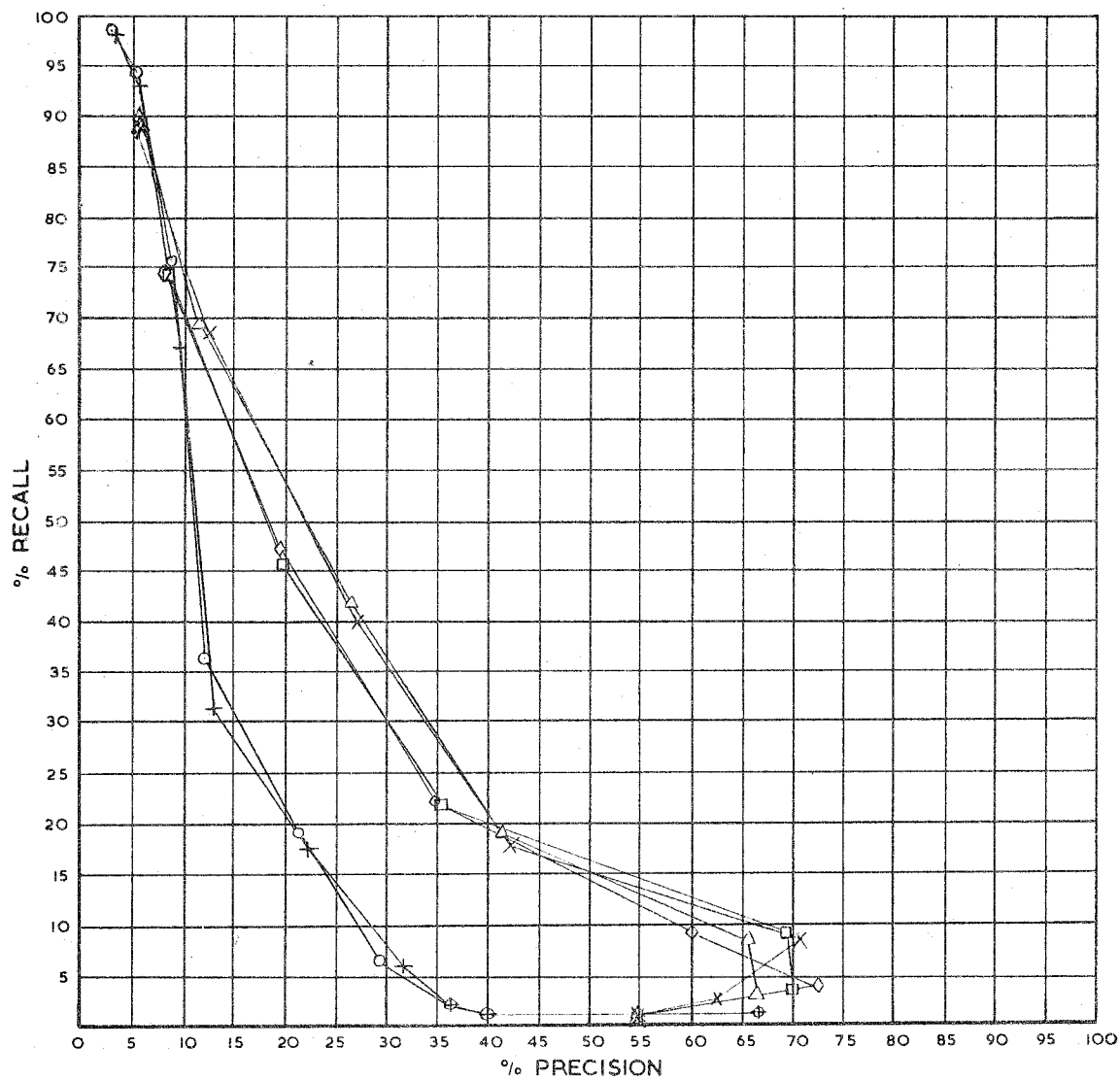


FIGURE 4.806P INDEX LANGUAGES III. 1. a (X), III. 2. a (Δ),  
III. 3. a (□), III. 4. a (◇), III. 5. a (+) and  
III. 6. a (O). SEARCH A (Figures 4.800T-4.805T)

FIGURE 4.810T

Index Language III.1.a (Controlled terms. Basic terms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule E  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	169	1,933	85.4%	8.0%	23,573%	42	42	42
2	114	428	57.6%	21.0%	5.218%	39	42	42
3	62	92	31.3%	40.3%	1.122%	29	41	41
4	25	13	12.6%	65.8%	0.158%	15	34	34
5	10	4	5.1%	71.4%	0.049%	6	24	24
6	0	2	0.0%	0.0%	0.024%	1	13	13
7	0	1	0.0%	0.0%	0.012%	1	8	8
8	0	1	0.0%	0.0%	0.012%	1	4	4
9	0	0	0.0%	0.0%	0.012%	0	3	3

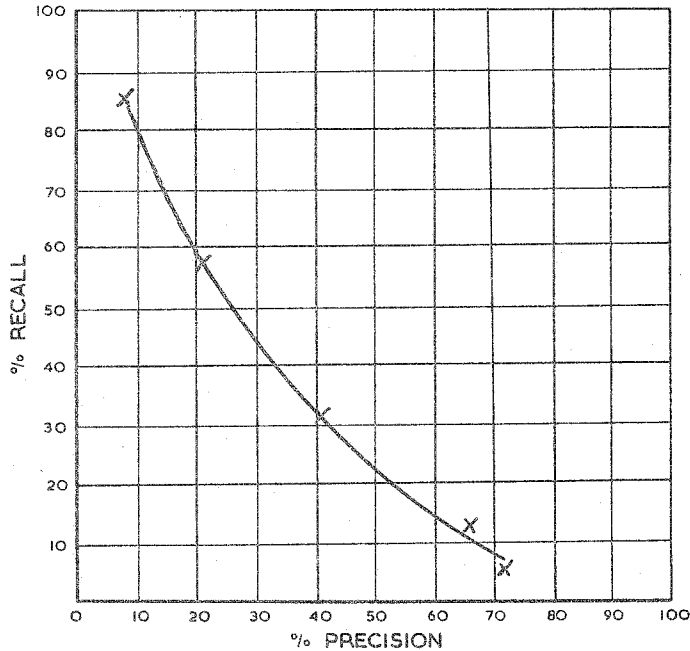


FIGURE 4.810P INDEX LANGUAGE III.1.a SEARCH E  
 200 DOCUMENTS

FIGURE 4.811T

Index Language III.2.a (Controlled terms. Narrower terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	171	2,064	86.4%	7.6%	25.414%	42	42	42
2	116	455	58.6%	20.3%	5.547%	39	42	42
3	64	100	32.3%	39.0%	1.219%	29	41	41
4	28	14	14.1%	66.7%	0.171%	15	34	34
5	11	4	5.6%	73.3%	0.049%	6	24	24
6	0	2	0.0%	0.0%	0.024%	1	13	13
7	0	1	0.0%	0.0%	0.012%	1	8	8
8	0	1	0.0%	0.0%	0.012%	1	4	4
9	0	0	0.0%	0.0%	0.012%	0	3	3

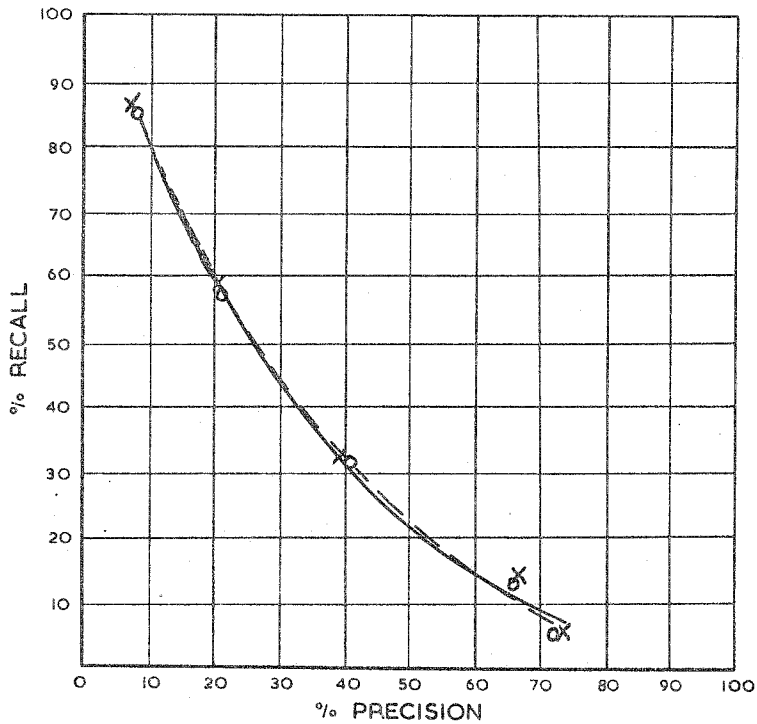


FIGURE 4.811P INDEX LANGUAGE III.2.a SEARCH E  
200 DOCUMENTS  
(Index Language III.1.a Broken line)

FIGURE 4.812T

Index Language III.3.a (Controlled terms. Broader terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	173	2,758	87.4%	5.9%	33.634%	42	42	42
2	117	762	59.1%	13.3%	9.290%	39	42	42
3	79	172	39.9%	31.5%	2.097%	35	41	41
4	31	24	15.7%	56.4%	0.2926%	19	34	34
5	12	5	6.1%	70.6%	0.0609%	7	24	24
6	2	2	1.0%	50.0%	0.024%	2	13	13
7	0	1	0.0%	0.0%	0.012%	1	8	8
8	0	1	0.0%	0.0%	0.012%	1	4	4
9	0	0				0	3	3

4137

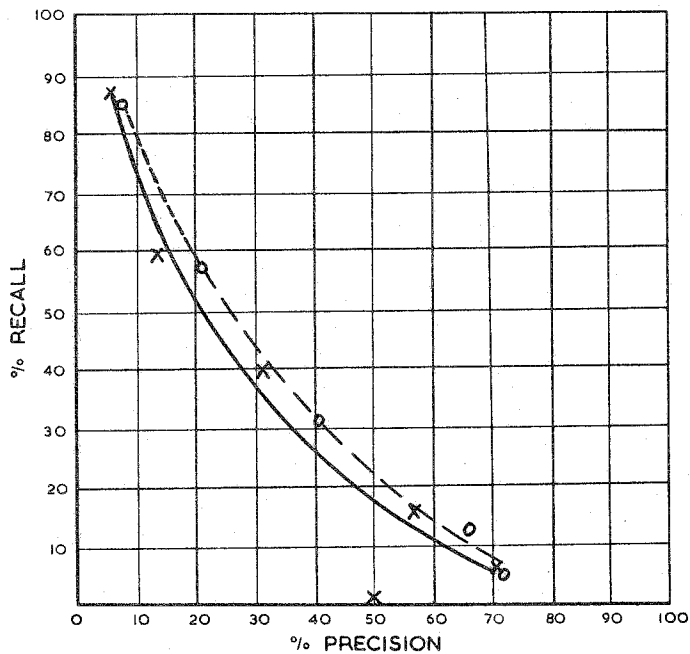


FIGURE 4.812P INDEX LANGUAGE III.3.a SEARCH E  
200 DOCUMENTS  
(Index Language III.1.a Broken line)

FIGURE 4.813T

Index Language III.4.a (Controlled terms. Narrower and broader terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	175	2,854	88.4%	5.8%	34.804%	42	42	42
2	129	797	65.2%	13.9%	9.717%	39	42	42
3	81	180	40.9%	31.0%	2.195%	35	41	41
4	34	24	17.2%	58.6%	0.2926%	19	34	34
5	13	6	6.6%	68.4%	0.073%	7	24	24
6	2	2	1.0%	50.0%	0.024%	2	13	13
7	0	1	0.0%	0.0%	0.012%	1	8	8
8	0	1	0.0%	0.0%	0.012%	1	4	4
9	0	0				0	3	3

4299

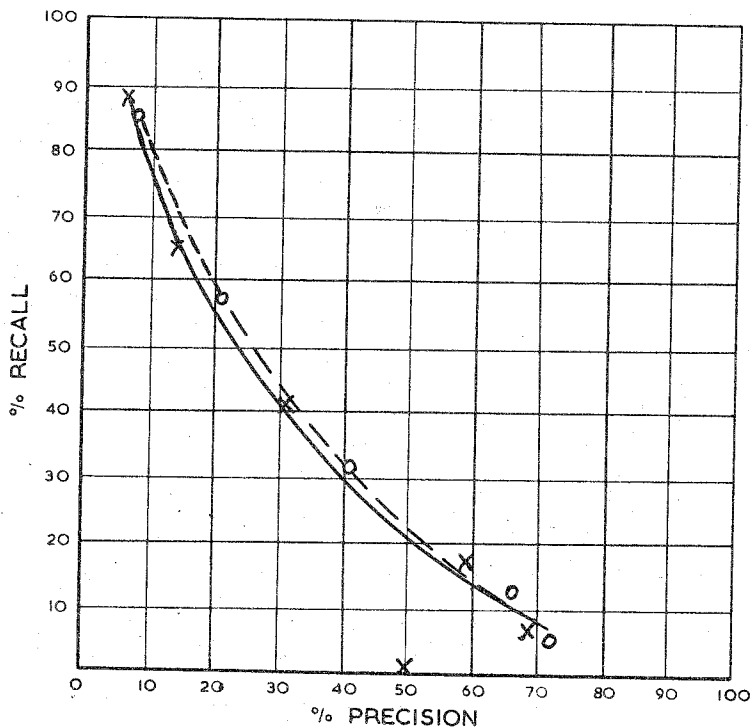


FIGURE 4.813P INDEX LANGUAGE III.4.a SEARCH E  
200 DOCUMENTS  
(Index Language III.1.a Broken line)

FIGURE 4.814T

Index Language III.5.a (Controlled terms. Related terms. Coordination)  
 Exhaustivity of Indexing 3  
 Search Rule E  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	194	4,339	98.0%	4.3%	52.914%	42	42	42
2	156	1,660	78.8%	8.6%	20.239%	42	42	42
3	104	628	52.5%	14.2%	7.657%	40	41	41
4	48	155	24.2%	23.6%	1.890%	27	34	34
5	27	47	13.6%	36.5%	0.573%	17	24	24
6	11	10	5.6%	52.4%	0.122%	5	13	13
7	2	2	1.0%	50.0%	0.024%	2	8	8
8	2	2	1.0%	50.0%	0.024%	2	4	4
9	2	1	1.0%	66.7%	0.012%	2	3	3

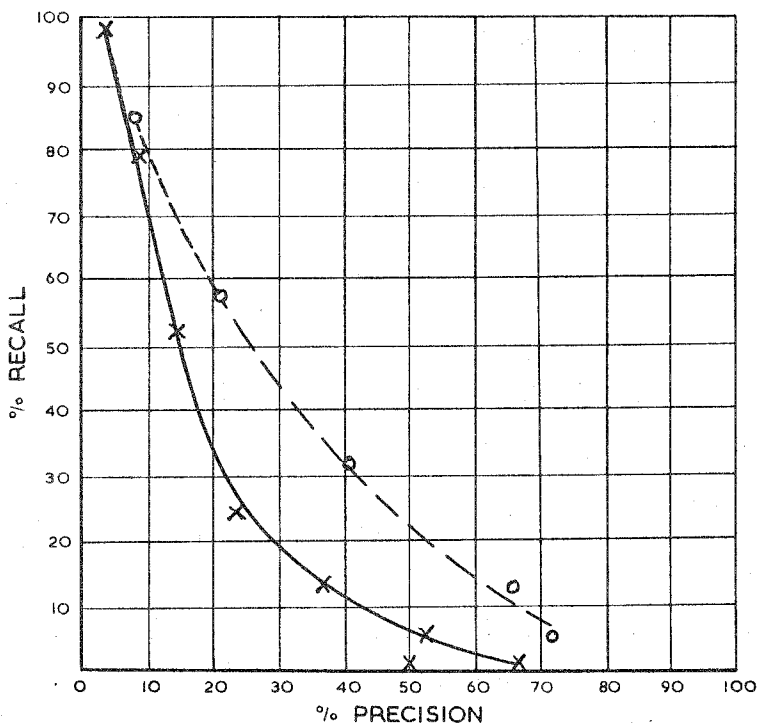


FIGURE 4.814P INDEX LANGUAGE III.5.a SEARCH E  
 200 DOCUMENTS  
 (Index Language III.1.a Broken line)

FIGURE 4.815T

Index Language III.6.a (Controlled terms. Narrower, broader and related terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	196	4,885	99.0%	4.0%	57.134%	42	42	42
2	169	2,092	85.4%	7.5%	25.506%	42	42	42
3	129	848	65.2%	13.2%	10.339%	42	41	41
4	60	210	30.3%	22.2%	2.560%	31	34	34
5	30	55	15.2%	35.3%	0.671%	18	24	24
6	11	10	5.6%	52.4%	0.122%	5	13	13
7	2	2	1.0%	50.0%	0.024%	2	8	8
8	2	2	1.0%	50.0%	0.024%	2	4	4
9	2	1	1.0%	66.7%	0.012%	2	3	3

8506

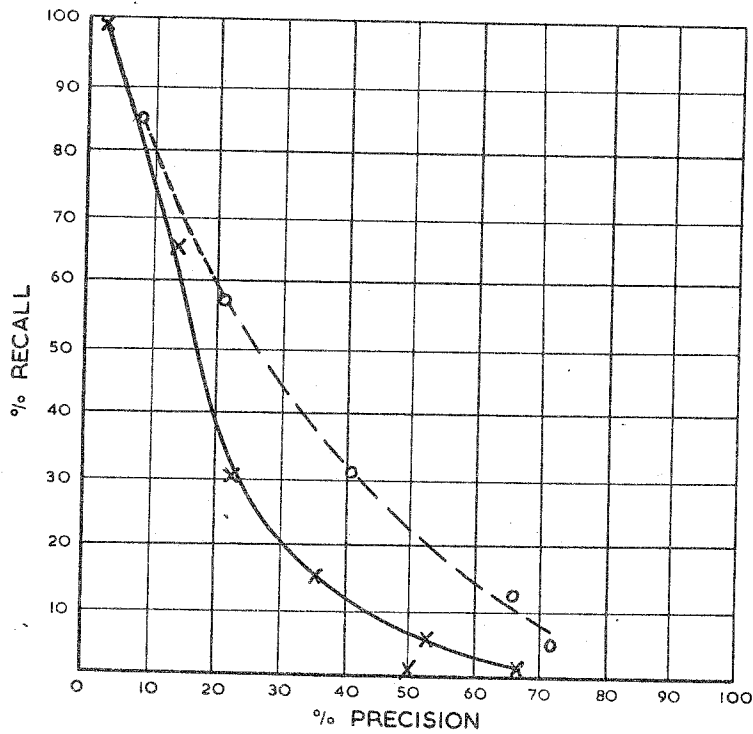


FIGURE 4.815P INDEX LANGUAGE III.6.a SEARCH E  
200 DOCUMENTS  
(Index Language III.1.a Broken line)

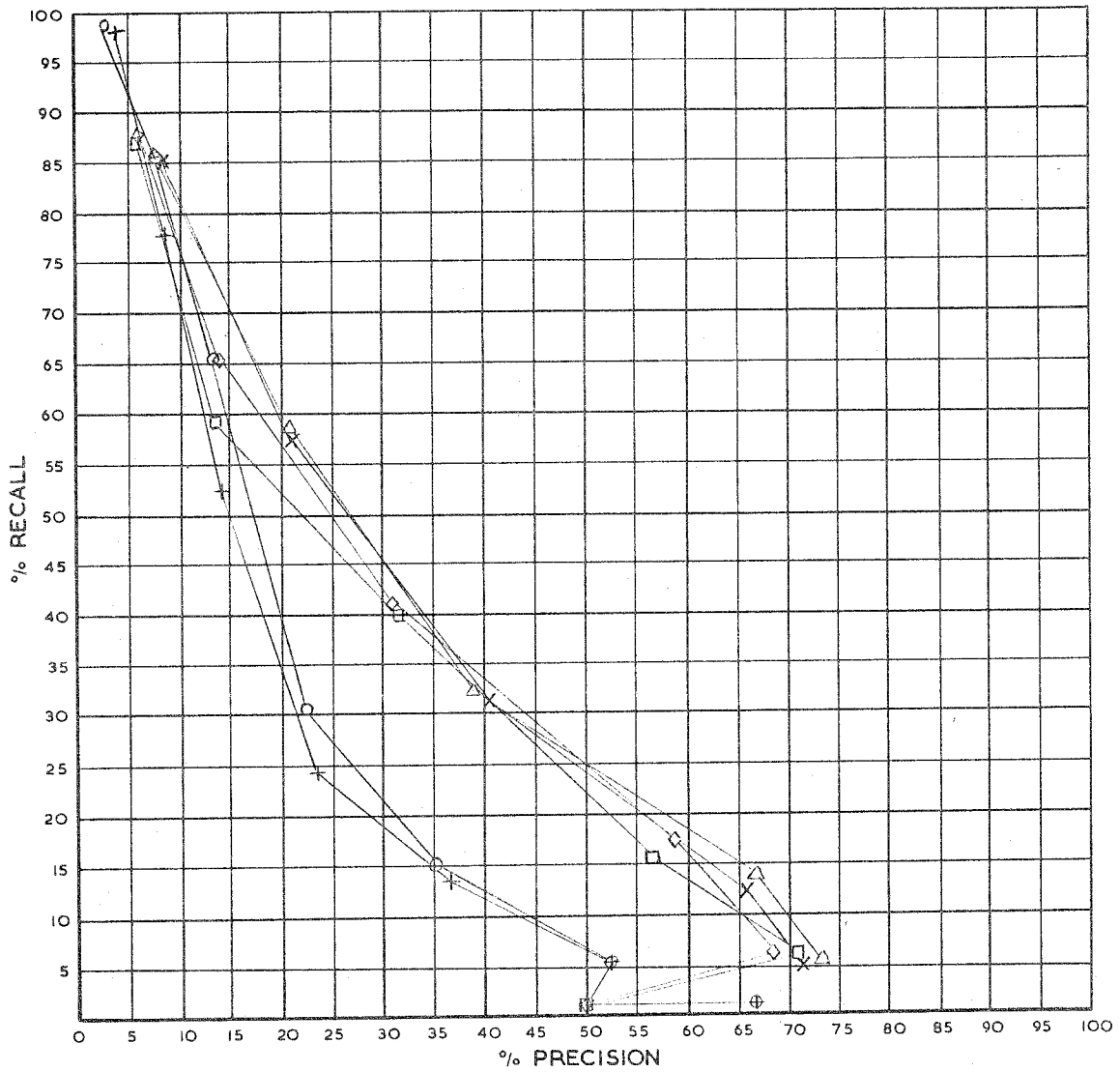


FIGURE 4.816P INDEX LANGUAGES III.1.a (X), III.2.a (A), III.3.a (□), III.4.a (◇), III.5.a (†), III.6.a (O). SEARCH E. 200 DOCUMENTS (Figures 4.810T - 4.815T)

FIGURE 4.820T

Index Language III.1.a (Controlled terms. Basic terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 350 (Subset 2)

Number of Questions 77 (Subset 7)

Number of Relevant Documents 454

Generality Number 16.8

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	361	6,320	79.5%	5.4%	23.849%	77	77	77
2	231	1,413	50.9%	14.1%	5.333%	73	77	77
3	109	289	24.0%	27.4%	1.091%	54	73	73
4	40	31	8.8%	56.3%	0.117%	23	63	63
5	13	6	2.9%	68.4%	0.023%	9	47	47
6	0	2	0.0%	0.0%	0.008%	1	28	28
7	0	1	0.0%	0.0%	0.004%	1	17	17
8	0	1	0.0%	0.0%	0.004%	1	11	11
9	0	0				0	6	6
10	0	0				0	1	1

FIGURE 4.821T

Index Language III.2.a (Controlled terms. Narrower terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 350 (Subset 2)

Number of Questions 77 (Subset 7)

Number of Relevant Documents 454

Generality Number 16.8

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	377	6,762	83.0%	5.3%	25.516%	77	77	77
2	245	1,554	54.0%	13.6%	5.865%	73	77	77
3	113	335	24.9%	25.2%	1.264%	54	73	73
4	46	35	10.1%	56.8%	0.132%	25	63	63
5	13	6	2.9%	68.4%	0.023%	9	47	47
6	0	2	0.0%	0.0%	0.008%	1	28	28
7	0	1	0.0%	0.0%	0.004%	1	17	17
8	0	1	0.0%	0.0%	0.004%	1	11	11
9	0	0				0	6	6
10	0	0				0	1	1

FIGURE 4.822T

Index Language III.3.a (Controlled terms. Broader terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 350 (Subset 2)

Number of Questions 77 (Subset 7)

Number of Relevant Documents 454

Generality Number 16.8

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	376	9,202	82.8	3.9%	37.724%	77	77	77
2	259	2,562	57.0%	9.2%	9.669%	73	77	77
3	133	496	30.4%	21.8%	1.872%	60	73	73
4	53	65	11.7%	44.9%	0.245%	31	63	63
5	18	9	4.0%	66.7%	0.034%	14	47	47
6	2	3	0.4%	40.0%	0.011%	3	28	28
7	0	1	0.0%	0.0%	0.004%	1	17	17
8	0	1	0.0%	0.0%	0.004%	1	11	11
9	0	0				0	6	6
10	0	0				0	1	1

FIGURE 4.823T

Index Language III.4.a (Controlled terms. Narrower and Broader terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 350 (Subset 2)

Number of Questions 77 (Subset 7)

Number of Relevant Documents 454

Generality Number 16.8

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	394	9,369	86.8%	4.0%	35.354%	77	77	77
2	276	2,757	60.8%	9.1%	10.405%	73	77	77
3	213	599	46.9%	26.2%	2.261%	60	73	73
4	58	73	12.8%	44.3%	0.276%	33	63	63
5	18	11	4.0%	62.1%	0.042%	14	47	47
6	2	3	0.4%	40.0%	0.011%	3	28	28
7	0	1	0.0%	0.0%	0.004%	1	17	17
8	0	1	0.0%	0.0%	0.004%	1	11	11
9	0	0				0	6	6
10	0	0				0	1	1

FIGURE 4.824T

Index Language III.5.a (Controlled terms. Related terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 350 (Subset 2)

Number of Questions 77 (Subset 7)

Number of Relevant Documents 454

Generality Number 16.8

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	430	13,965	94.7%	3.0%	52.698%	77	77	77
2	350	5,896	77.1%	5.6%	22.252%	77	77	77
3	211	1,943	46.5%	9.8%	7.333%	72	73	73
4	103	441	22.7%	18.9%	1.664%	53	63	63
5	40	132	8.8%	23.3%	0.498%	31	47	47
6	12	27	2.6%	30.8%	0.102%	9	28	28
7	2	5	0.4%	28.6%	0.019%	4	17	17
8	2	3	0.4%	40.0%	0.011%	2	11	11
9	2	1	0.4%	66.7%	0.004%	2	6	6
10	0	0				0	1	1

FIGURE 4.825T

Index Language III.6.a (Controlled terms. Narrower, broader and related terms. Coordination)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 350 (Subset 2)

Number of Questions 77 (Subset 7)

Number of Relevant Documents 454

Generality Number 16.8

Coord-ination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	440	15,172	96.9%	2.8%	57.252%	77	77	77
2	370	7,383	81.5%	4.8%	27.865%	77	77	77
3	252	2,664	55.5%	8.6%	10.054%	72	73	73
4	126	630	27.8%	16.7%	2.378%	59	63	63
5	45	167	9.9%	20.7%	0.630%	35	47	47
6	13	36	2.9%	26.5%	0.136%	11	28	28
7	2	7	0.4%	22.2%	0.026%	5	17	17
8	2	3	0.4%	40.0%	0.011%	2	11	11
9	2	1	0.4%	66.7%	0.004%	2	6	6
10	0	0				0	1	1

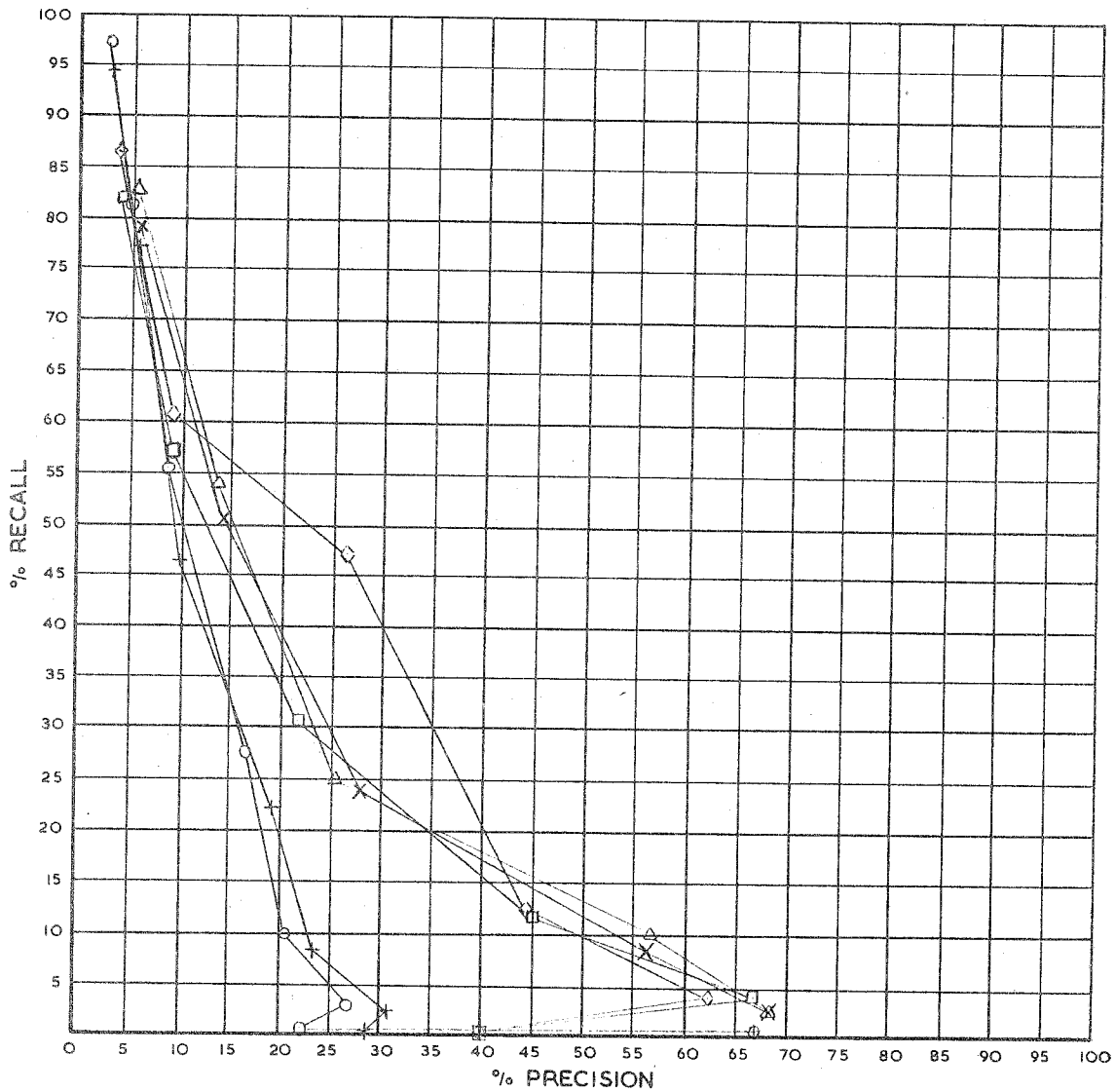


FIGURE 4.826P INDEX LANGUAGES III.1.a (X), III.2.a (Δ), III.3.a (□), III.4.a (◇), III.5.a (+), III.6.a (O), SEARCH E. 350 DOCUMENTS (Figures 4.820T - 4.825T)

FIGURE 4.830T

Index Language III.1.e (Controlled terms. Basic terms. Coordination, weighting)

Exhaustivity of Indexing 3

Search Rule E

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	155	(-)	78.3%	(-)	(-)	42	0	42
2	90	215	45.5%	29.5%	2.621%	36	42	42
3	39	27	19.7%	59.1%	0.329%	23	41	41
4	11	4	5.6%	73.3%	0.049%	7	34	34
5	6	0	3.0%	100.0%	0.000%	3	24	24
6	0	0				0	13	13
7	0	0				0	8	8
8	0	0				0	4	4
9	0	0				0	3	3

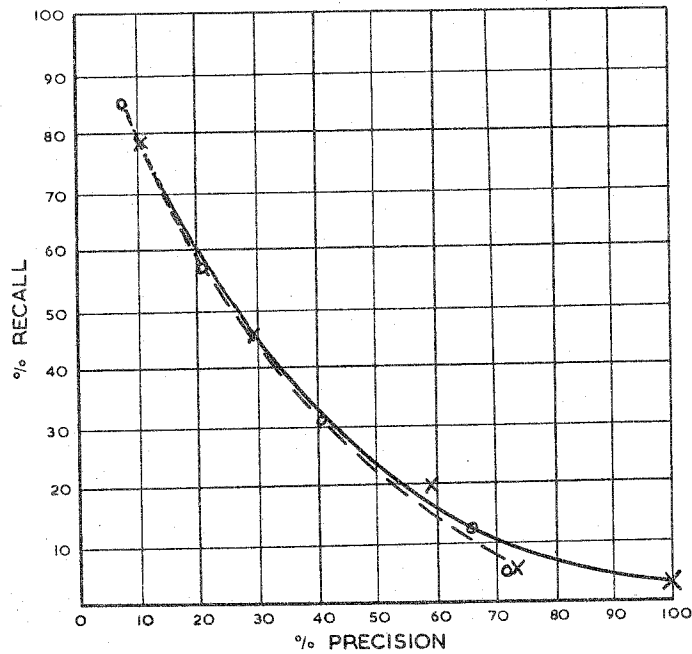


FIGURE 4.830P INDEX LANGUAGE III.1.e SEARCH E  
(Index Language III.1.a Broken Line)

FIGURE 4.831T

Index Language III.6.e (Controlled terms. Narrower, broader and related terms. Coordination weighting)  
 Exhaustivity of Indexing 3  
 Search Rule E  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	192	(-)	97.0%	(-)	(-)	42	0	42
2	137	1,134	69.2%	10.8%	13.826%	41	42	42
3	88	345	44.4%	20.3%	4.206%	38	41	41
4	39	56	19.7%	41.1%	0.683%	24	34	34
5	20	15	10.1%	57.1%	0.183%	11	24	24
6	8	2	4.0%	80.0%	0.024%	3	13	13
7	0	0				0	8	8
8	0	0				0	4	4
9	0	0				0	3	3

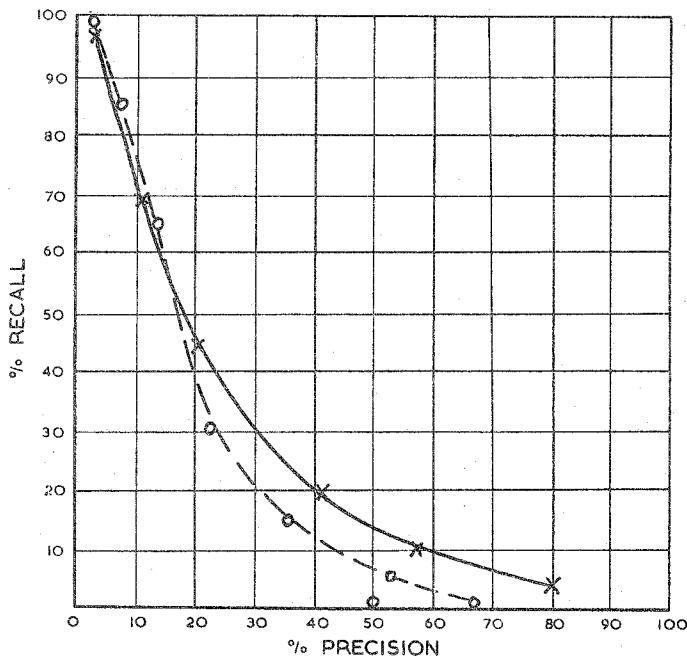


FIGURE 4.831P INDEX LANGUAGE III.6.e SEARCH E  
 (Index Language III.6.a Broken line)

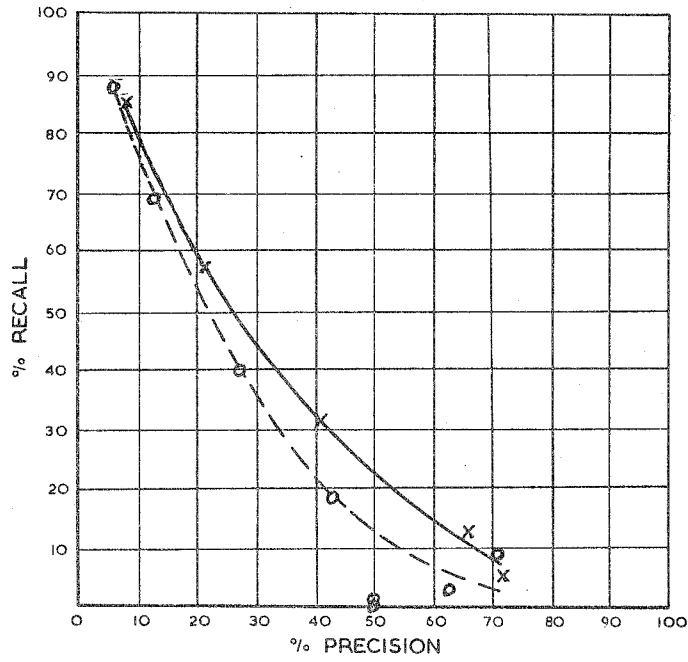


FIGURE 4.840P INDEX LANGUAGE III.1.a. SEARCH A (BROKEN LINE) AND SEARCH E. 200 DOCUMENTS (FIGURES 4.800T and 4.810T)

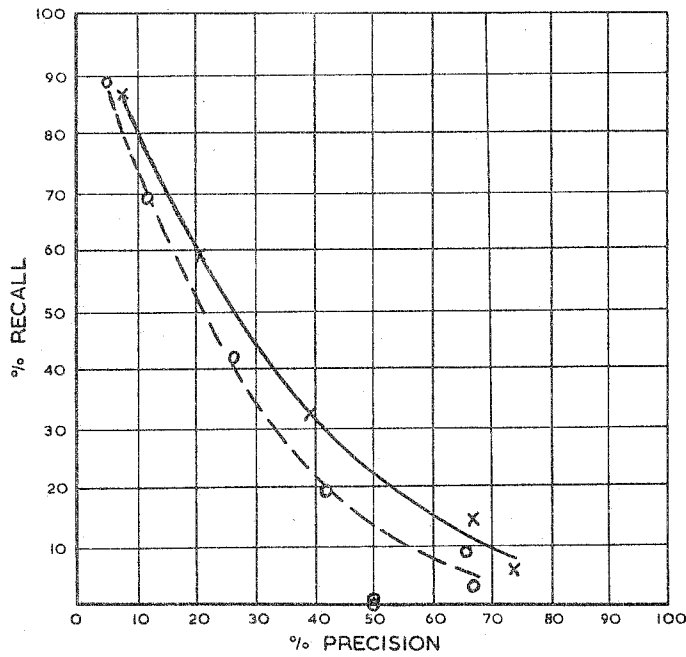


FIGURE 4.841P INDEX LANGUAGE III.2.a. SEARCH A (BROKEN LINE) AND SEARCH E. 200 DOCUMENTS (FIGURES 4.801T and 4.811T)

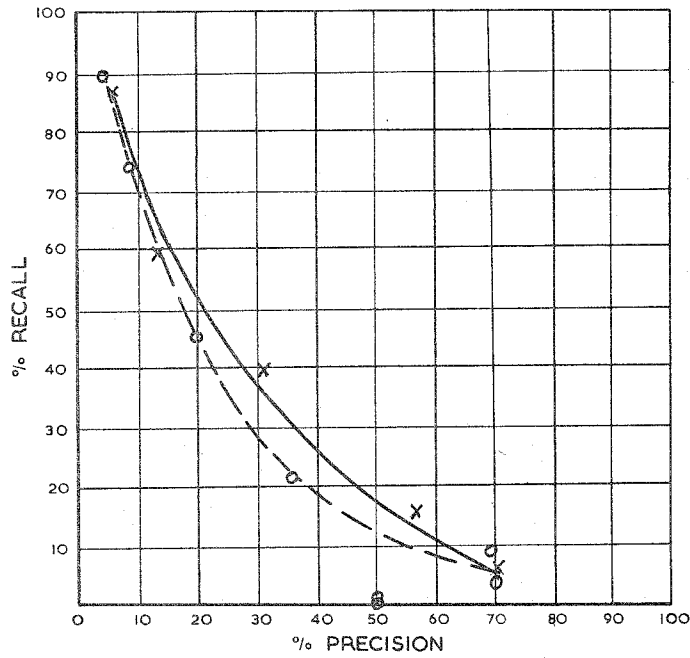


FIGURE 4.842P INDEX LANGUAGE III.3.a. SEARCH A (BROKEN LINE) AND SEARCH E. 200 DOCUMENTS (FIGURES 4.802T and 4.812T)

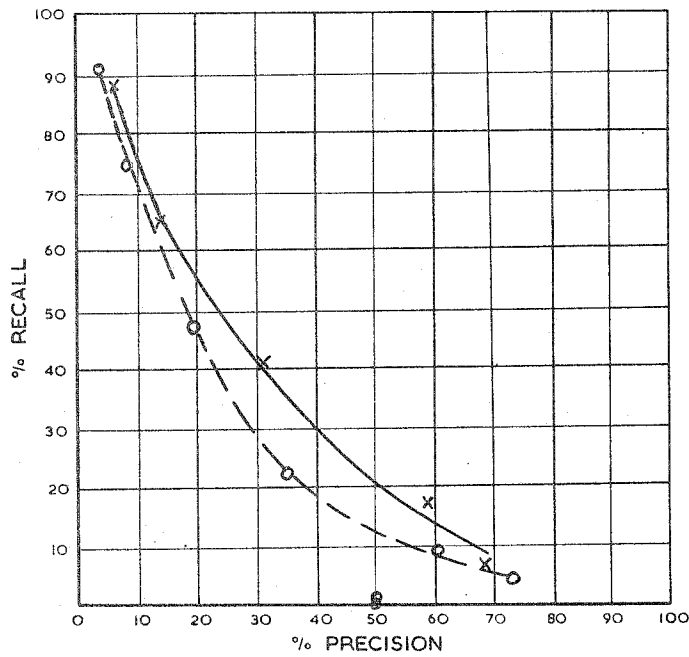


FIGURE 4.843P INDEX LANGUAGE III.4.a. SEARCH A (BROKEN LINE) AND SEARCH E. 200 DOCUMENTS (FIGURES 4.803T and 4.813T)

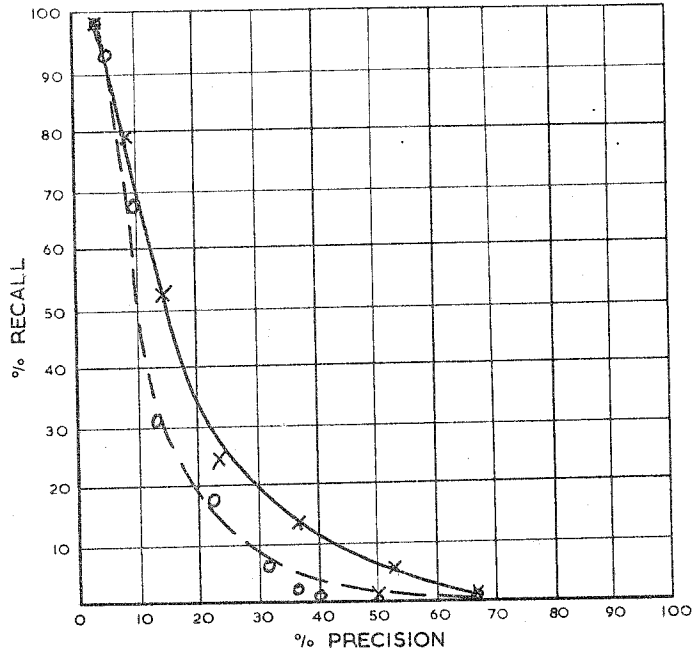


FIGURE 4.844P INDEX LANGUAGE III.5.a. SEARCH A (BROKEN LINE) AND SEARCH E. 200 DOCUMENTS (FIGURES 4.804T and 4.814T)

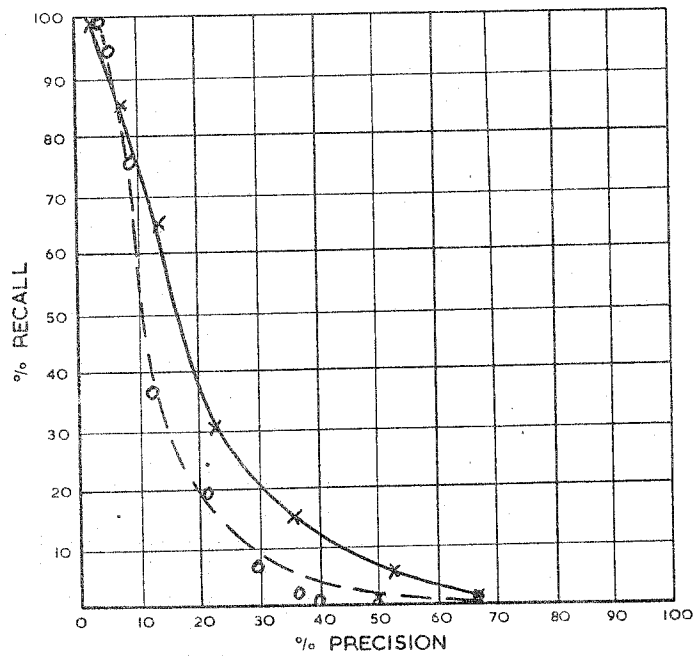


FIGURE 4.845P INDEX LANGUAGE III.6.a. SEARCH A (BROKEN LINE) AND SEARCH E. 200 DOCUMENTS (FIGURES 4.805T and 4.815T)

FIGURE 4.850T

Index Language III.5.a  
 Exhaustivity of indexing 3  
 Search Rule F  
 Document Relevance 1-4  
 Number of Documents in Collection 200 (subset 1)  
 Number of Questions 42 (subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Starting Terms	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
		Rel.	Non-rel.			
2	0	156	1,660	78.0%	8.6%	20.329%
	1	146	1,133	73.7%	11.4%	13.817%
3	2	114	428	57.6%	21.0%	5.218%
	0	104	628	52.5%	14.2%	7.657%
	1	100	534	50.5%	15.8%	6.512%
4	2	89	339	44.9%	21.1%	4.134%
	3	62	92	31.3%	40.3%	1.122%
	0	48	155	24.2%	23.6%	1.890%
	1	47	146	23.7%	24.3%	1.780%
5	2	47	100	23.7%	32.0%	1.220%
	3	44	55	22.2%	44.4%	0.671%
	4	25	13	12.6%	65.8%	0.158%
	0	27	47	13.6%	36.5%	0.573%
	1	27	46	13.6%	37.0%	0.561%
6	2	27	37	13.6%	42.2%	0.451%
	3	25	29	12.6%	46.3%	0.354%
	4	20	15	10.1%	57.1%	0.183%
	5	10	4	5.1%	71.4%	0.049%
	0	11	10	5.6%	52.4%	0.122%
	1	10	10	5.1%	50.0%	0.122%
	2	10	10	5.1%	50.0%	0.122%
3	10	6	5.1%	62.5%	0.073%	
4	10	4	5.1%	71.4%	0.049%	
5	6	2	3.1%	75.0%	0.024%	
6	0	0	2	0.0%	0.0%	0.024%

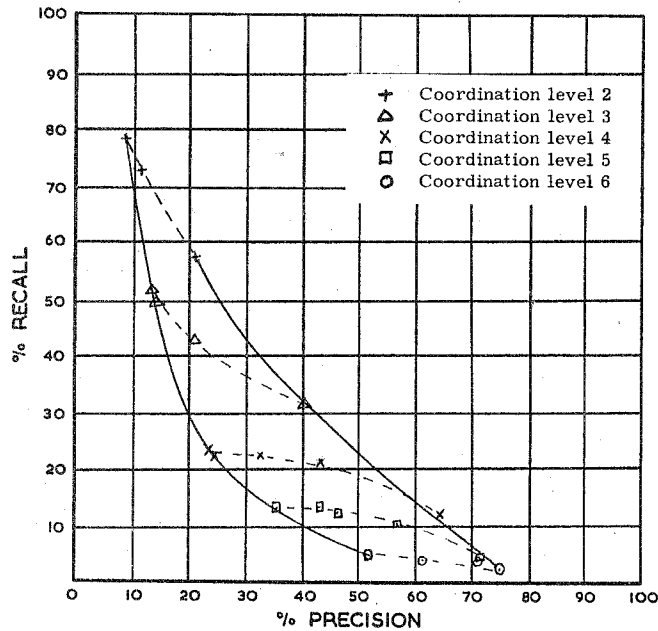


FIGURE 4.851T

Index Language III.6.a

Exhaustivity of indexing 3

Search Rule F

Document Relevance 1-4

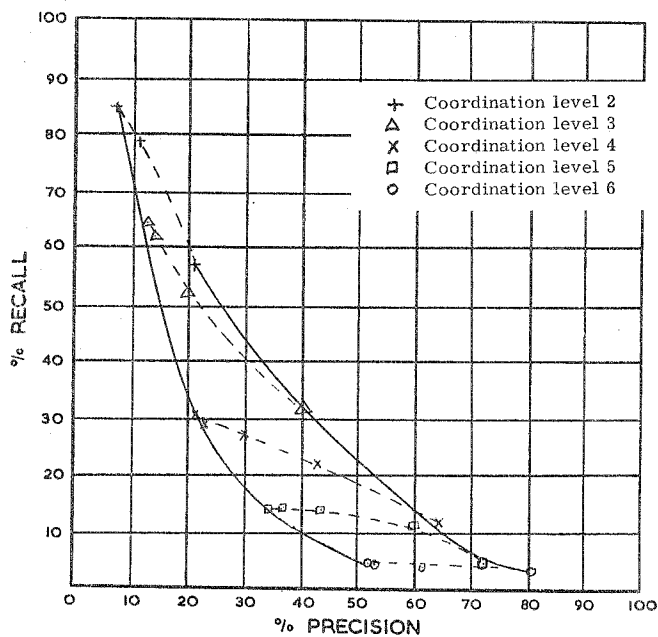
Number of Documents in Collection 200 (subset 1)

Number of Questions 42

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Starting Terms	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
		Rel.	Non-rel.			
2	0	169	2,092	85.4%	7.5%	25.506%
	1	158	1,321	79.3%	10.7%	16.110%
	2	114	428	57.6%	21.0%	5.218%
3	0	129	848	65.2%	13.2%	10.339%
	1	122	685	61.6%	15.1%	8.354%
	2	105	398	53.0%	20.9%	4.854%
4	3	62	92	31.3%	40.0%	1.122%
	0	60	210	30.3%	22.2%	2.560%
	1	59	192	29.8%	23.5%	2.341%
	2	57	128	28.8%	30.8%	1.561%
5	3	47	64	23.7%	42.3%	0.780%
	4	25	13	12.6%	65.8%	0.158%
	0	30	55	15.2%	35.3%	0.671%
	1	30	51	15.2%	37.0%	0.622%
	2	30	41	15.2%	42.3%	0.500%
6	3	30	30	15.2%	50.0%	0.366%
	4	23	15	11.6%	60.5%	0.183%
	5	10	4	5.1%	71.4%	0.049%
	0	11	10	5.6%	52.4%	0.122%
	1	10	10	5.1%	50.0%	0.122%
	2	10	9	5.1%	52.6%	0.110%
	3	10	6	5.1%	62.5%	0.073%
4	10	4	5.1%	71.4%	0.049%	
5	8	2	4.0%	80.0%	0.024%	
6	0	0	2	0.0%	0.0%	0.012%



Section ~~8~~<sup>9</sup> Abstracts and Titles

This section does not introduce any new index language device, since it deals with searches carried out on the abstracts and titles of the 200 document collection. As such it represents a variation in exhaustivity of indexing, and also, in some cases, different concept indexing to that done by the project staff. The searches are done first in single term natural language and then with word forms confounded.

Figures 4.900T and 4.901T present the results of searches on the titles only, and Figures 4.902T and 4.903T present the results obtained with titles and abstracts. Figure 4.904P is a plot of the four tables.

In order to make a comparison from the viewpoint of exhaustivity, Figures 4.910T and 4.911T present the results of Index Language I.3.a (word forms) at exhaustivity levels of 2 and 1, on the 200 document collection. Plot 4.912P now gives the performance curves for the three levels of exhaustivity with the project indexing, for the search on titles alone and for the search on titles and abstracts.

We were able to do this series of tests by having a print out of a concordance of the titles and abstracts of the 200 documents. This was prepared with the SMART programme, and we are indebted to Prof. Salton for making it available

LIST OF FIGURES

	Index Language	Exhaustivity	No. of Questions	Question Subset	Document Collection	Plots
4.900T	IV.1.a		42	2	200	
4.901T	IV.2.a		42	2	200	
4.902T	IV.3.a		42	2	200	
4.903T	IV.4.a		42	2	200	
4.904P	IV.1.a		42	2	200	4.900T
	IV.2.a					4.901T
	IV.3.a					4.902T
	IV.4.a					4.903T
4.910T	I.3.a	2	42	2	200	
4.911T	I.3.a	1	42	2	200	
4.912P	IV.2.a		42	2	200	4.901T
	IV.4.a					4.903T
	I.3.a	3				4.201T
	I.3.a	2				4.910T
	I.3.a	1				4.911T

FIGURE 4.900T

Index Language IV.1.a (Single terms. Natural language. Coordination)

Exhaustivity of Indexing Title only

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	155	1,902	78.3%	7.5%	23.189%	42	<del>42</del>	42
2	108	463	54.5%	18.9%	5.645%	40	42	42
3	58	76	29.3%	43.3%	0.927%	32	42	42
4	26	15	13.1%	63.4%	0.183%	16	41	41
5	15	3	7.6%	83.3%	0.037%	9	39	39
6	6	0	3.0%	100.0%	0.000%	3	33	33
7	3	0	1.5%	100.0%	0.000%	1	27	27
8	0	0				0	18	18
9	0	0				0	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1

2830

FIGURE 4.901T

Index Language IV.2.a (Single terms. Word forms. Coordination)

Exhaustivity of Indexing Title only

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	159	2,017	80.3%	7.3%	24.592%	42	<del>42</del>	42
2	111	512	56.1%	17.8%	6.242%	40	42	42
3	65	94	32.8%	40.9%	1.146%	34	42	42
4	29	18	14.6%	61.7%	0.219%	16	41	41
5	17	4	8.6%	81.0%	0.049%	12	39	39
6	7	0	3.5%	100.0%	0.000%	3	33	33
7	3	0	1.5%	100.0%	0.000%	1	27	27
8	0	0				0	18	18
9	0	0				0	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1

3036

148

306

FIGURE 4.902T

Index Language IV.3.a (Single terms. Natural language. Coordination)

Exhaustivity of Indexing Title and Abstract

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	185	6,261	93.4%	2.9%	76.353%	42	0	42
2	161	2,158	81.3%	6.9%	26.311%	42	42	42
3	120	854	60.6%	12.3%	10.412%	40	42	42
4	75	288	37.9%	20.7%	3.511%	32	41	41
5	44	97	22.2%	31.2%	1.183%	21	39	39
6	27	25	13.6%	51.9%	0.305%	13	33	33
7	12	2	6.1%	85.7%	0.024%	6	27	27
8	3	0	1.5%	100.0%	0.000%	2	18	18
9	2	0	1.0%	100.0%	0.000%	1	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1

FIGURE 4.903T

Index Language IV.4.a (Single terms. Word forms. Coordination)

Exhaustivity of Indexing Title and Abstract

Search Rule A

Document Relevance 1 - 4

Number of Documents in Collection 200 (Subset 1)

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	187	6,770	94.4%	2.7%	82.540%	42	0	42
2	166	2,654	83.8%	5.9%	32.358%	42	42	42
3	130	1,150	65.7%	10.2%	14.021%	41	42	42
4	92	432	46.5%	17.6%	5.267%	35	41	41
5	51	158	25.8%	24.4%	1.926%	28	39	39
6	32	53	16.2%	37.6%	0.646%	19	33	33
7	20	13	10.1%	60.6%	0.158%	10	27	27
8	9	4	4.5%	69.2%	0.049%	7	18	18
9	3	1	1.5%	75.0%	0.012%	4	11	11
10	1	0	0.5%	100.0%	0.000%	1	7	7
11	0	0				0	3	3
12	0	0				0	1	1

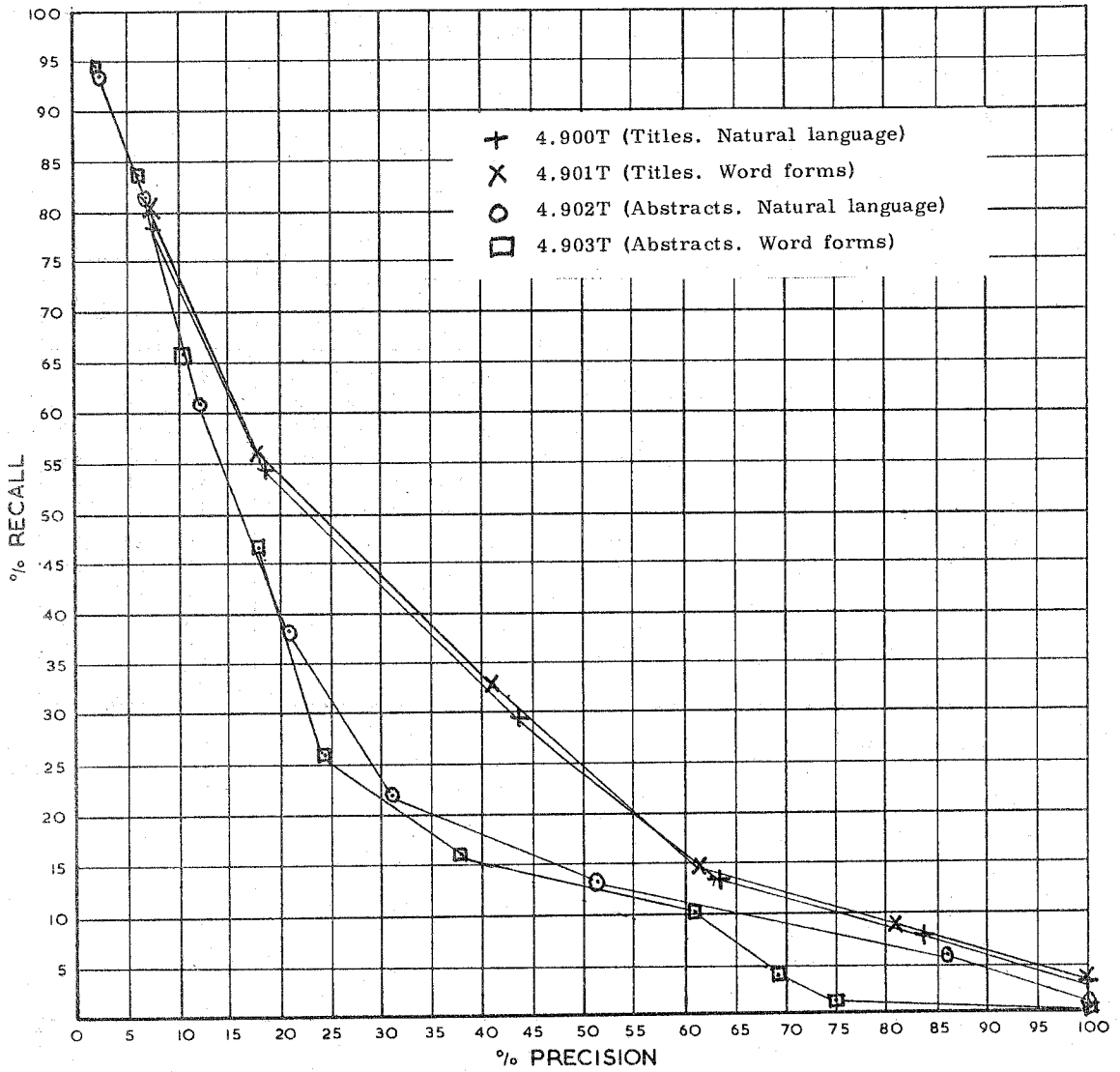


FIGURE 4.904P INDEX LANGUAGES IV.1.a, IV.2.a, IV.3.a and IV.4.a. (Figures 4.900T - 4.903T)

FIGURE 4.910T

Index Language I.3.a (Single terms. Word forms. Coordination)  
 Exhaustivity of Indexing 2  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	188	4,037	94.9%	4.4%	49.220%	42	42	42
2	162	1,867	81.8%	8.0%	22.763%	42	42	42
3	128	694	64.6%	15.6%	8.461%	40	42	42
4	87	208	43.9%	29.5%	2.536%	32	41	41
5	47	66	23.7%	41.6%	0.805%	22	39	39
6	30	17	15.2%	63.8%	0.207%	13	33	33
7	11	2	5.6%	84.6%	0.024%	6	27	27
8	3	0	1.5%	100.0%	0.000%	2	18	18
9	0	0				0	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1

FIGURE 4.911T

Index Language I.3.a (Single terms. Word forms. Coordination)  
 Exhaustivity of Indexing 1  
 Search Rule A  
 Document Relevance 1 - 4  
 Number of Documents in Collection 200 (Subset 1)  
 Number of Questions 42 (Subset 2)  
 Number of Relevant Documents 198  
 Generality Number 23.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	177	2,626	89.4%	6.3%	32.017%	42	42	42
2	138	1,101	69.7%	11.1%	13.424%	42	42	42
3	103	308	52.0%	25.1%	3.755%	36	42	42
4	60	80	30.3%	42.9%	0.975%	28	41	41
5	30	23	15.2%	56.6%	0.280%	14	39	39
6	17	3	8.6%	85.0%	0.037%	9	33	33
7	7	0	3.5%	100.0%	0.000%	3	27	27
8	1	0	0.5%	100.0%	0.000%	1	18	18
9	0	0				0	11	11
10	0	0				0	7	7
11	0	0				0	3	3
12	0	0				0	1	1

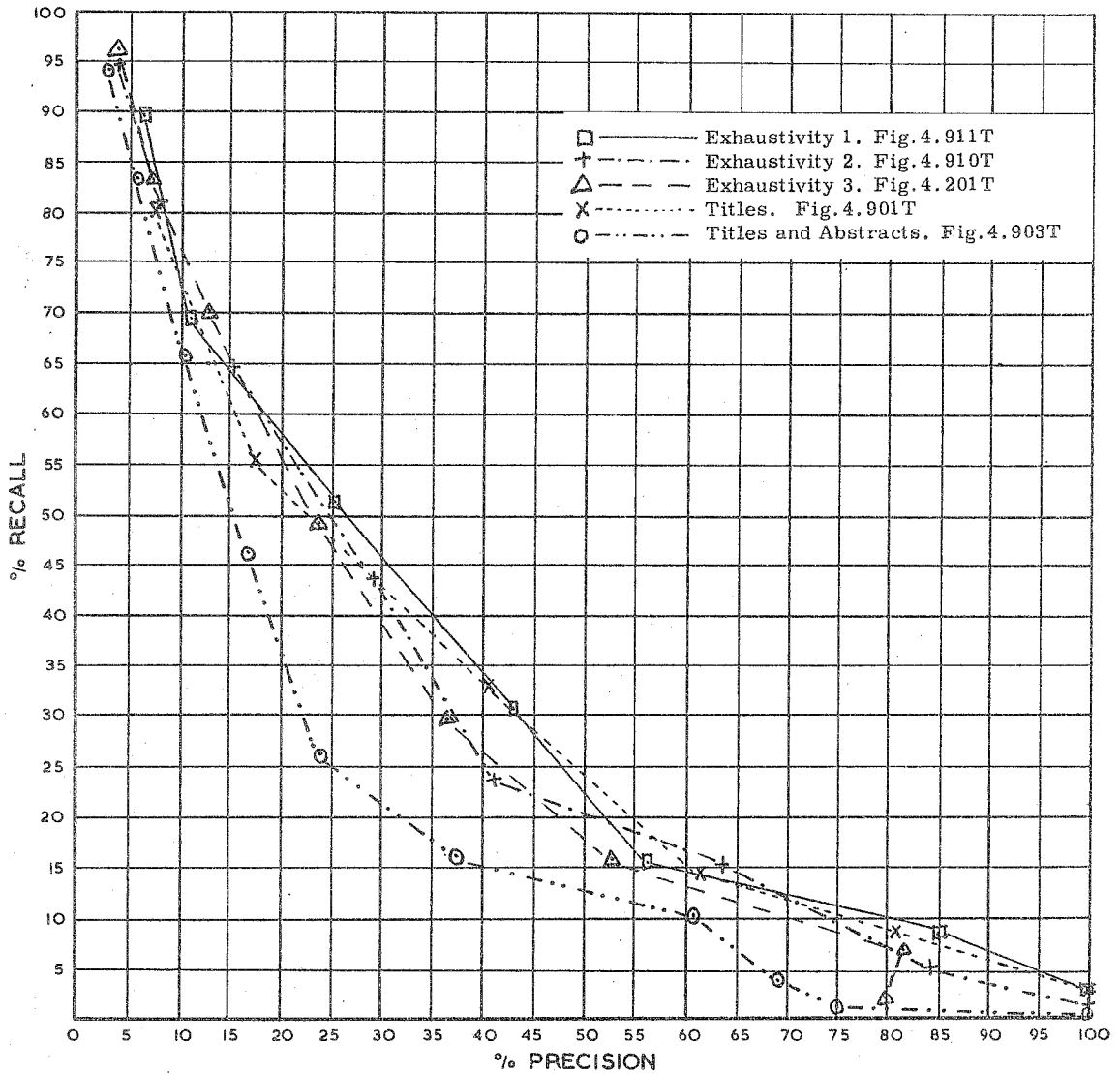


FIGURE 4.912P PLOT OF RESULTS FOR THREE LEVELS OF EXHAUSTIVITY OF INDEXING, FOR TITLES AND FOR ABSTRACTS, ON 42 QUESTIONS, 200 DOCUMENT COLLECTION, INDEX LANGUAGE I.3.a

*The merit of Exh 1 vs Exh 3 is reversed (Fig 5.23P)*

## CHAPTER 5

### Simulated ranking and document output cut-off

There is confusion of ends and means in this type of attack upon measurement in principle. Perhaps if medicine threw away the thermometer, the encephalograph, the X-ray, and all other technicalities, medicine would become much more human! How much more preferable the tender hand on the brow than a nasty piece of glass in the mouth - how inhuman! But is it sympathy and fellow-feeling that we want from the physician or a technical competence to identify the condition and give us the cure? The bedside manner still has a place in the cure, even although the hand on the brow has been replaced by the thermometer.

L.T. Wilkins: Social Deviance, page 9

With all the results so far given, the presentation has been on the basis of coordination level cut-offs. The reader is invited to consider the same test results, but now presented on the basis of a simulated ranking order and a document output cut-off. In Chapter 3, one of the main problems considered was that of totalling the results of a set of questions that was heterogenous in having different numbers of starting terms and matching terms. Several solutions were considered, but only brief mention was made of one possible method, namely document output cut-off. Although this method was recognised as having many advantages, it was decided not to use it for the main test results; this was partly because of the additional effort required to obtain the necessary prerequisite of a ranking order, but also because it would have involved a transformation of the test results as actually obtained by the co-ordination level cut-off. At a later date a simpler method of deriving a simulated ranking order was found and, in trying this out, it was shown that there was a possibility of obtaining an 'area measure' which could be used for producing an order of performance effectiveness for the different index languages. Therefore, the majority of the test searches were converted to a simulated ranking order, and in this chapter the results are presented by the document output cut-off method.

The influence of the SMART system was mainly responsible for our original investigation into attempting to obtain a ranked output for the Cranfield test searches. In the SMART system, the output of a search is arranged in an order of decreasing correlation with the search question; this is established by each document having a scoring that is obtained by calculations based on the match between the request terms and the document terms in the particular dictionary being tested. Thus every document in the collection is assigned a rank order number, the rank position reflecting the correlation with the search system. A sample output from the SMART system, showing the results for Question 147 searched on the Cranfield 200 document collection for fourteen different options, is given in Fig. 5.1. This output sheet shows, for each of the fourteen options, the file numbers of the fifteen highest ranked documents and also the rank numbers of the five documents which are relevant to this particular question. The heading at the top of each section refers to the particular option being tested, and it can be seen that, with 'ABSTR OLD QS', for instance, the five relevant documents, Nos. 708, 711, 713, 712 and 709 were ranked 21, 32, 68, 76 and 122 respectively.

In Fig. 5.2. are shown the conventional search results for 42 questions by Index Language I.l.a, and these are set out in coordination levels.

0147CONTROLS 5 RELEVANT

ABSTR OLD QS	ABSTR F NULL	ABSTR NEW QS	INDEX OLD QS	INDEX F NULL	INDEX NEW QS	CRAN CONCON	ABSTR F NULL	INDEX F NULL
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
1 792 21 708	1 792 13 712	1 792 7 712	1 792 18 712	1 792 14 712	1 792 3 712	1 597 4 708	1 792 4 708	1 792 20 712
2 33H 32 711	2 670 21 708	2 787 19 711	2 748 36 708	2 799 26 711	2 683 69 708	2 717 60 712	2 670 32 706	2 670 32 706
3 797 68 713	3 880 22 711	3 799 97 713	3 800 83 713	3 748 35 708	3 712 103 709	3 748 65 709	3 799 33 711	3 799 33 711
4 683 76 712	4 33H 41 713	4 324 101 708	4 437 117 711	4 316 147 709	4 797 139 713	4 708 98 713	4 880 56 713	4 880 56 713
5 574 122 709	5 331 76 709	5 683 149 709	5 799 137 709	5 11A 166 713	5 988 153 709	5 451 128 711	5 748 99 709	5 748 99 709
6 799 7 880	6 874 7 700	6 932 7 712	6 797 6 683	6 683 6 683	6 572 6 371	6 371 6 33H	6 33H 6 33H	6 33H 6 33H
7 880 7 700	7 700 7 700	7 712 7 712	7 572 7 572	7 670 7 670	7 437 7 437	7 10A 7 10A	7 316 7 316	7 316 7 316
8 655 8 451	8 451 8 451	8 707 8 707	8 683 8 683	8 683 8 683	8 799 8 799	8 916 8 916	8 331 8 331	8 331 8 331
9 995 9 704	9 704 9 704	9 416 9 416	9 070 9 070	9 666 9 666	9 682 9 682	9 509 9 509	9 11A 9 11A	9 11A 9 11A
10 569 10 919	10 919 10 919	10 321 10 321	10 509 10 509	10 701 10 701	10 675 10 675	10 982 10 982	10 874 10 874	10 874 10 874
11 800 11 748	11 748 11 748	11 655 11 655	11 988 11 988	11 33H 11 33H	11 670 11 670	11 700 11 700	11 683 11 683	11 683 11 683
12 793 12 988	12 988 12 988	12 53H 12 53H	12 656 12 656	12 451 12 451	12 316 12 316	12 792 12 792	12 700 12 700	12 700 12 700
13 416 13 712	13 712 13 712	13 874 13 874	13 108 13 108	13 572 13 572	13 681 13 681	13 331 13 331	13 451 13 451	13 451 13 451
14 700 14 415	14 415 14 415	14 316 14 316	14 983 14 983	14 712 14 712	14 701 14 701	14 797 14 797	14 080 14 080	14 080 14 080
15 34A 15 08D	15 08D 15 08D	15 919 15 919	15 701 15 701	15 705 15 705	15 707 15 707	15 794 15 794	15 704 15 704	15 704 15 704
RNK REC= 0.0470	RNK REC= 0.0867	RNK REC= 0.0402	RNK REC= 0.0304	RNK REC= 0.0387	RNK REC= 0.0338	RNK REC= 0.0423	RNK REC= 0.0625	RNK REC= 0.0625
LOG PRE= 0.2410	LOG PRE= 0.2859	LOG PRE= 0.2509	LOG PRE= 0.2327	LOG PRE= 0.2448	LOG PRE= 0.2433	LOG PRE= 0.2508	LOG PRE= 0.2577	LOG PRE= 0.2577
NOR REC= 0.6882	NOR REC= 0.8379	NOR REC= 0.6329	NOR REC= 0.6144	NOR REC= 0.6174	NOR REC= 0.5600	NOR REC= 0.6213	NOR REC= 0.7692	NOR REC= 0.7692
NOR PRE= 0.3029	NOR PRE= 0.4471	NOR PRE= 0.3390	NOR PRE= 0.2701	NOR PRE= 0.3172	NOR PRE= 0.3115	NOR PRE= 0.3386	NOR PRE= 0.3624	NOR PRE= 0.3624
OVERALL= 0.2680	OVERALL= 0.3726	OVERALL= 0.2911	OVERALL= 0.2710	OVERALL= 0.2635	OVERALL= 0.2771	OVERALL= 0.2930	OVERALL= 0.3202	OVERALL= 0.3202
NOR_OVR= 0.3029	NOR_OVR= 0.6369	NOR_OVR= 0.3390	NOR_OVR= 0.2701	NOR_OVR= 0.3172	NOR_OVR= 0.3115	NOR_OVR= 0.3386	NOR_OVR= 0.3624	NOR_OVR= 0.3624

ABSTR NEW QS	ABSTR F NULL	INDEX NEW QS	CRAN CONCON	ABSTR F NULL	INDEX F NULL
TOP 15 RELEVANT	TOP 15 RELEVANT	TOP 15 RELEVANT	TOP 15 RELEVANT	TOP 15 RELEVANT	TOP 15 RELEVANT
1 792 5 412	1 792 12 712	1 792 4 712	1 597 7 708	1 597 6 712	1 597 6 708
2 797 30 711	2 797 31 711	2 683 40 711	2 792 14 712	2 792 7 708	2 792 25 712
3 683 96 708	3 670 35 708	3 799 53 708	3 717 36 711	3 717 103 709	3 717 46 711
4 799 122 713	4 799 60 713	4 712 165 713	4 797 102 709	4 683 122 711	4 799 102 709
5 712 171 709	5 880 102 709	5 748 175 709	5 748 136 713	5 748 138 713	5 748 136 713
6 324 6 324	6 324 6 324	6 797 6 797	6 712 6 712	6 708 6 708	6 708 6 708
7 988 7 33H	7 33H 7 33H	7 316 7 316	7 708 7 708	7 316 7 316	7 316 7 316
8 992 8 683	8 683 8 683	8 988 8 988	8 324 8 324	8 797 8 797	8 451 8 451
9 572 9 331	9 331 9 331	9 11A 9 11A	9 451 9 451	9 451 9 451	9 11A 9 11A
10 437 10 992	10 992 10 992	10 572 10 572	10 683 10 683	10 988 10 988	10 371 10 371
11 707 11 874	11 874 11 874	11 437 11 437	11 371 11 371	11 683 11 683	11 683 11 683
12 416 12 712	12 712 12 712	12 670 12 670	12 992 12 992	12 572 12 572	12 10A 12 10A
13 682 13 700	13 700 13 700	13 080 13 080	13 10A 13 10A	13 10A 13 10A	13 670 13 670
14 321 14 707	14 707 14 707	14 682 14 682	14 712 14 712	14 437 14 437	14 916 14 916
15 675 15 451	15 451 15 451	15 666 15 666	15 916 15 916	15 916 15 916	15 080 15 080
LOG REC= 0.0354	LOG REC= 0.0625	LOG REC= 0.0343	LOG REC= 0.0508	LOG REC= 0.0399	LOG REC= 0.0476
LOG PRE= 0.2453	LOG PRE= 0.2621	LOG PRE= 0.2478	LOG PRE= 0.2704	LOG PRE= 0.2844	LOG PRE= 0.2605
NOR REC= 0.5805	NOR REC= 0.7692	NOR REC= 0.5672	NOR REC= 0.7128	NOR REC= 0.6297	NOR REC= 0.6923
NOR PRE= 0.3188	NOR PRE= 0.3802	NOR PRE= 0.3283	NOR PRE= 0.4027	NOR PRE= 0.3843	NOR PRE= 0.3717
OVERALL= 0.2806	OVERALL= 0.3256	OVERALL= 0.3212	OVERALL= 0.3212	OVERALL= 0.3043	OVERALL= 0.3081
NOR_OVR= 0.3188	NOR_OVR= 0.3802	NOR_OVR= 0.3283	NOR_OVR= 0.4027	NOR_OVR= 0.3843	NOR_OVR= 0.3717

FIGURE 5.1T EXAMPLE OF SMART OUTPUT WITH CRANFIELD 200 DOCUMENT

Q	1+		2+		3+		4+		5+		6+		7+		8+		9+	
	R	N	R	N	R	N	R	N	R	N	R	N	R	N	R	N	R	N
79	2	60	1	7	1	1												
100	4	167	3	71	3	50	1	2										
116	6	169	5	92	4	51	3	25	1	12	0	4						
118	5	123	5	49	5	31	3	17	3	10	1	2	1	0				
119	6	170	6	82	4	37	3	13	3	6	0	1						
121	3	30	3	5	3	0	3	0	1	0								
122	5	107	5	41	3	11	1	0										
123	3	92	3	24	3	3												
126	2	95	2	62	2	20	2	4	2	0								
130	4	153	4	28	1	0												
132	2	78	2	51	1	31	1	5										
136	6	63	6	23	6	8	6	4	5	4	5	3	3	2	2	1	0	1
137	6	147	6	60	6	22	6	10	2	4								
141	1	82	1	35	1	4	1	1										
145	12	168	12	102	11	47	7	23	6	10	4	1						
146	7	37	4	1	1	1												
147	3	97	2	35	1	13	0	7	0	3	0	1						
148	4	36	4	15	4	4	2	0	1	0								
167	4	182	4	105	3	51	3	21	3	9	1	1						
170	2	109	1	45	1	18	1	6	1	1								
181	2	164	2	42	1	7												
182	3	175	1	47	1	5												
189	2	64	0	12	0	1												
190	7	162	6	45	5	10	3	0										
223	2	148	2	75	2	38	2	19	2	3	2	1	2	0				
224	5	80	4	65	2	27	0	3										
225	6	158	4	91	4	43	2	17	0	5								
226	7	60	4	19	4	2	4	1	2	0								
227	2	83	2	35	2	8	2	3	1	0	1	0						
230	7	42	2	0	1	0												
250	8	162	8	54	8	25	8	7	5	4	3	0						
261	4	131	4	34	4	13	4	5	4	0	4	0	3	0				
264	2	104	2	29	2	5	2	1	1	0	1	0	1	0				
266	5	164	4	32	0	8	0	1										
268	5	23	5	1	4	0	2	0										
269	4	34	4	4	2	0	1	0										
272	4	183	4	123	4	66	4	22	3	4	3	1	2	0	2	0		
273	7	33	6	10	5	1	2	0										
274	5	177	4	81	3	28	2	8										
317	2	118	2	69	2	31	2	10	0	2	0	1						
323	5	162	5	69	4	26	0	2	0	1								
360	8	143	8	59	8	14	5	4	3	2	0	1						

FIGURE 5.2T SEARCH RESULTS BY COORDINATION LEVEL CUTOFF FOR SINGLE TERM INDEX LANGUAGE (I.1.a) WITH 42 QUESTIONS AND 200 DOCUMENT COLLECTION.

(R = Relevant documents retrieved  
N = Non-relevant documents retrieved)

By using these figures it was found possible to obtain a simulated ranking output. This is done by assigning a rank order number to each relevant document retrieved by means of the equations:-

$${}^cR_n = X_c + (n - Y_c) \left( \frac{x_c + 1}{y_c + 1} \right)$$

where  ${}^cR_n$  is the rank order number of the  $n^{\text{th}}$  relevant document to be retrieved

- $c$  is the coordination level at which the  $n^{\text{th}}$  relevant document is retrieved
- $x_c$  is the additional number of documents retrieved at coordination level  $c$ . (i.e. those not retrieved at a higher coordination level)
- $y_c$  is the additional number of relevant documents retrieved at coordination level  $c$ . (i.e. those not retrieved at a higher coordination level)
- $X_c$  is the total number of documents retrieved before searching at coordination level  $c$ . (i.e. at higher coordination levels)
- $Y_c$  is the total number of relevant documents retrieved before searching at coordination level  $c$ . (i.e. at higher coordination levels)

${}^cR_n$  is taken to the nearest whole number but if its value falls exactly between two whole numbers, it is taken to the lower whole number for odd numbered questions and to the higher whole number for even numbered questions. Two examples to illustrate the effect are taken from Fig. 5.2. With Question 100, no documents are retrieved at a coordination level higher than four, so for this question, the various values are as follows:

Question 100

At level  $c=4$ , then  $x_4 = 3$ ,  $y_4 = 1$ ,  $X_4 = 0$ ,  $Y_4 = 0$

At level  $c=3$ , then  $x_3 = 50$ ,  $y_3 = 2$ ,  $X_3 = 3$ ,  $Y_3 = 1$

At level  $c=2$ , then  $x_2 = 21$ ,  $y_2 = 0$ ,  $X_2 = 53$ ,  $Y_2 = 3$

At level  $c=1$ , then  $x_1 = 97$ ,  $y_1 = 1$ ,  $X_1 = 74$ ,  $Y_1 = 3$

∴ For Relevant Document 1, retrieved at level 4 :-

$${}^4R_1 = 0 + (1 - 0) \left( \frac{3 + 1}{1 + 1} \right) = 0 + 2 = 2$$

For Relevant Document 2, retrieved at level 3 :-

$${}^3R_2 = 3 + (2 - 1) \left( \frac{50 + 1}{2 + 1} \right) = 3 + 17 = 20$$

For Relevant Document 3, retrieved at level 3 :-

$${}^3R_3 = 3 + (3 - 1) \left( \frac{50 + 1}{2 + 1} \right) = 3 + 34 = 37$$

For Relevant Document 4 retrieved at level 1 :-

$${}^1R_4 = 74 + (4 - 3) \left( \frac{97 + 1}{1 + 1} \right) = 74 + 49 = 123$$

In the next example considered, Question 123, there are actually four relevant documents; no documents are retrieved at a coordination

Q	REL	1	2	3	4	5	6 -7	8 -10	11 -15	16 -20	21 -30	31 -50	51 -75	76 -100	101 -125	126 -150	151 -175	176 -200
79	3	x									x					x		
100	4		x							x	x					x		
116	6							x		x	x	x			x			
118	5	x					x	x			xx							
119	6					x	x				x		xx					
121	3	x	x	x							x	x						
122	5	x				x		x			x	x						
123	4	x		x		x											x	
126	2	x	x															
130	4	x						x	x		x							
132	4				x							x			x		x	
136	6		x	x		x	x	xx										
137	6		x			x		xx	xx									
141	1	x																
145	12	x	x	x	x			x	x	x		xxx	x	x				
146	9		x	x	x	x			x		x	x		x		x		
147	5							x			x		x			x	x	
148	4	x	x		x		x											
167	4	x				x		x						x				
170	2		x											x				
181	2				x						x							
182	4				x									x	x			x
189	2										x	x						
190	7	x	x	x			x		x			x			x			
223	2	x	x															
224	5								x		x	x	x	x				
225	6							x	x		x	x			x	x		
226	7	x	x	x		x						xx	x					
227	2	x		x														
230	7	x	x					x		x	x	xx						
250	8	x	x	x		x	x	x	xx									
261	4	x	x	x	x													
264	2	x	x															
266	5								x	x	x	x		x				
268	5	x	x	x	x		x											
269	4	x	x		x		x											
272	4	x	x		x					x								
273	7	x	x	x	x		x		x		x							
274	5				x		x			x			x			x		
317	2					x		x										
323	5					x			x	x	x	x						
360	8		x		x	x	x	x	x	xx								
Totals		23	21	13	13	12	11	16	14	10	18	17	8	7	5	6	3	1
Recall		12	22	29	35	41	47	56	62	67	76	85	89	92	95	98	99	100
Precision		55	51	45	42	39	32	26	20	16	12	8	6	4	4	3	3	2

FIGURE 5.3T DOCUMENT OUTPUT CUT-OFF SCORE SHEET FOR INDEX LANGUAGE I.1.a FOR 42 QUESTIONS WITH 200 DOCUMENT COLLECTION.

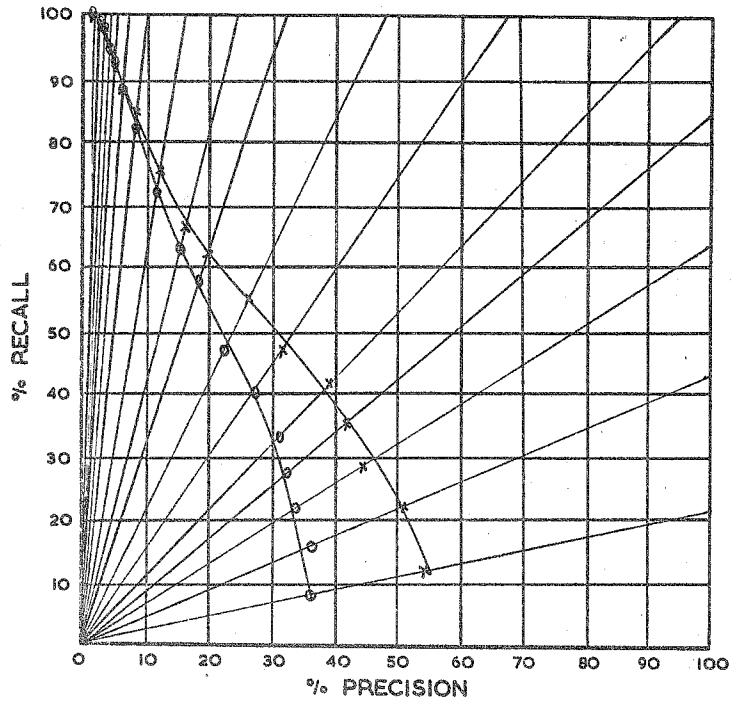


FIGURE 5.4P PLOT OF RESULTS FOR  
INDEX LANGUAGES I.1.a AND  
I.9.a BY DOCUMENTS OUTPUT  
CUT-OFF METHOD, SHOWING  
DOCUMENT OUTPUT CUT-OFF  
LINES.

- x INDEX LANGUAGE I.1.a
- o INDEX LANGUAGE I.9.a

level higher than three. It will be seen from Fig. 5.2. that at the single term level, only three of these documents have been found. The remaining relevant document can only be retrieved by searching through the remainder of the collection, namely 105 documents, and therefore at  $c=0$ ,  $x_0$  is taken to be 105. In addition the equations do not always produce whole numbers, so  ${}^cR_n$  has to be taken to the nearest whole number, or to the lower whole number where the value falls exactly between two whole numbers (since Q123 is an odd-numbered question).

Question 123

At level  $c=3$ , then  $x_3 = 6$ ,  $y_3 = 3$ ,  $X_3 = 0$ ,  $Y_3 = 0$

At level  $c=2$ , then  $x_2 = 21$ ,  $y_2 = 0$ ,  $X_2 = 6$ ,  $Y_2 = 3$

At level  $c=1$ , then  $x_1 = 68$ ,  $y_1 = 0$ ,  $X_1 = 27$ ,  $Y_1 = 3$

At level  $c=0$ , then  $x_0 = 105$ ,  $y_0 = 1$ ,  $X_0 = 95$ ,  $Y_0 = 3$

Then :-

$${}^3R_1 = 0 + (1 - 0) \left( \frac{6 + 1}{3 + 1} \right) = \frac{7}{4} = 2$$

$${}^3R_2 = 0 + (2 - 0) \left( \frac{6 + 1}{3 + 1} \right) = \frac{7}{2} = 3$$

$${}^3R_3 = 0 + (3 - 0) \left( \frac{6 + 1}{3 + 1} \right) = \frac{21}{4} = 5$$

$${}^0R_4 = 95 + (4 - 3) \left( \frac{105 + 1}{1 + 1} \right) = 95 + 53 = 148$$

The argument for this simulated ranking method is given in Appendix 5A.

When all such rankings have been calculated for the searches with a single index language, the results are entered on a score sheet as in Fig. 5.3T, which represents the results as given in Fig. 5.2T. Seventeen ranking groups were selected to have approximately the same number of documents falling in to each group; these were 1; 2; 3; 4; 5; 6-7; 8-10; 11-15; 16-20; 21-30; 31-50; 51-75; 76-100; 101-125; 126-150; 151-175; and 176-200. A cross is put in the appropriate column of the score sheet for every relevant document for the 42 questions. From the score-sheet, the total number of relevant documents retrieved at each of the seventeen cut-off levels can now be obtained. In Fig. 5.3T it is shown that, in the 42 searches, the first document retrieved was relevant on 23 occasions. As there were 198 documents relevant to the 42 questions, the recall ratio at this stage can be calculated as  $\frac{23}{198} \times 100 = 12\%$ ; the precision ratio is calculated on the basis of one document having been retrieved for each question, and is therefore  $\frac{23}{42} \times 100 = 55\%$ . In 21 of the searches, the second document retrieved was relevant, making a total of 44 relevant documents so far retrieved, so the recall ratio increases to 22%. The precision ratio is now calculated on the basis of 2 x 42 documents having been retrieved, and is therefore 51%. Recall and precision ratios are similarly calculated for each document output cut-off level; ultimately the recall ratio will reach 100%.

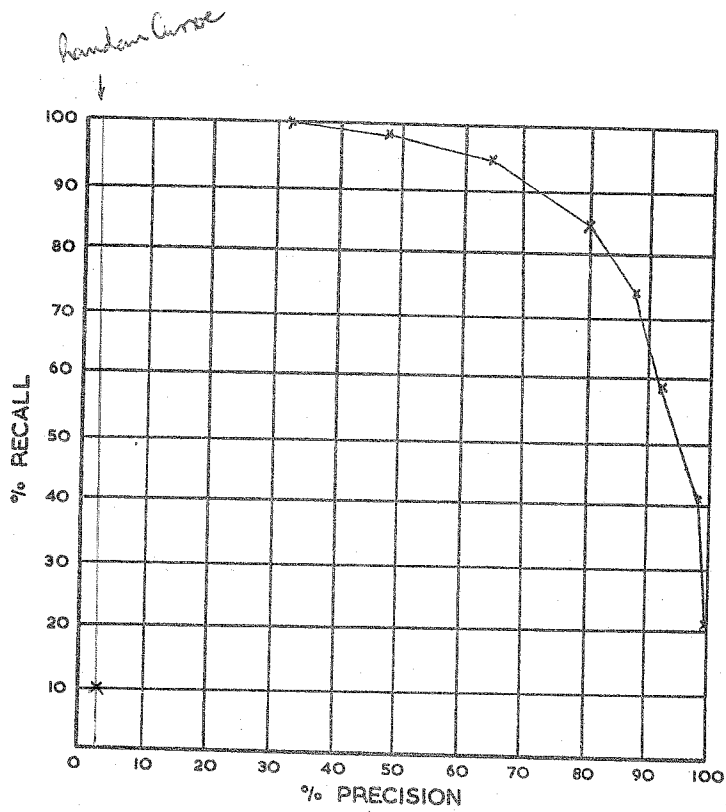


FIGURE 5.5P MAXIMUM POSSIBLE PERFORMANCE CURVE WITH DOCUMENT OUTPUT CUT-OFF FOR CRANFIELD TEST COLLECTION OF 200 DOCUMENTS AND 42 QUESTIONS.

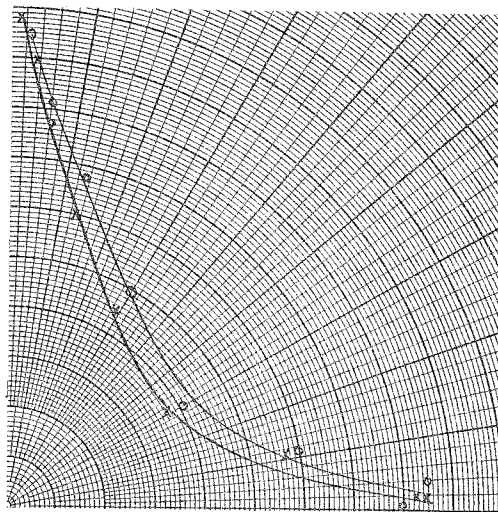


FIGURE 5.6P POLAR GRAPH OF RESULTS FOR INDEX LANGUAGE I.1.a (FIG. 4.140T) AND INDEX LANGUAGE I.6.a (FIG. 4.203T) (COORDINATION LEVEL CUTOFFS)

o I.1.a      x I.6.a

Such recall and precision figures can be plotted on a conventional graph as in Fig. 5.4P, which shows the results of index language I.1.a (as in Fig. 5.3T) and also index language I.9.a. These curves can be compared with Fig. 4.206P and show the same superiority of index language I.1.a over index language I.9.a.

There is, however, an important difference. The positions of the points in Fig. 4.206P were determined by coordination level cut-offs, and were therefore random in relation to each other. With Fig. 5.4P, if straight lines are drawn radiating from the point of origin, these will, as can be seen, pass through the corresponding points in each curve. This is due to the fact that the cut-off is based on document output, and recall and precision ratios are now interdependent. It is known that there are 198 documents relevant to the 42 questions, so, on average, 4.7 documents are relevant to each question. When only one document is retrieved for each question, even if every such document were relevant, the recall ratio could not possibly be higher than  $\frac{100 \times 42}{198} = 21.2\%$ , although it would, of course, represent a precision ratio of 100%. If any of the documents are not relevant, then the recall ratio will always fall on some point along the line which goes from the point of origin to a recall of 21.2% at 100% precision. Therefore at any given document output cut-off, a drop in recall ratio with any one system as against any other system must also involve a drop in the precision ratio. Similarly, when two documents are retrieved in each search, the maximum recall ratio is 42.4% and with this particular document/question set, 100% recall cannot possibly be reached until at least five documents are retrieved for each question. This would, however, represent a total of 210 documents. Since there are only 198 relevant documents in the collection, the theoretical maximum precision ratio would then be  $\frac{198}{210} \times 100 = 94.3\%$ . As more documents are retrieved, so the maximum possible precision ratio must drop, and these document output cut-off performance lines can be calculated as has been done in Fig. 5.4P.

Because of the fact that Question 141 had only one relevant document, it would not be possible in this collection to obtain the theoretically maximum figures for recall and precision beyond the single document cut-off level. Similarly, there are thirteen questions which have more than five relevant documents, and 100% recall could not possibly be obtained until twelve documents have been retrieved, this number representing the highest figure for documents relevant to a single question. This does not affect the position of the lines, which would be different, however, for other situations where there are more or less relevant documents per question.

As previously mentioned, it is not possible to obtain the theoretically maximum performance beyond the single document output cut-off, since Q141 has only one relevant document. As ten questions have only two relevant documents, there must be a further deviation from the theoretical maximum beyond this stage. In Fig. 5.5P is shown the actual possible maximum performance that could be obtained with this collection. Achieving this performance would imply that for each question all the relevant documents were retrieved before any non-relevant documents were retrieved.

In Fig. 5.4P the lines radiating from the point of origin have been based on the document output cut-off for this particular test situation, but the performance curves could be drawn on a polar coordinate graph with the lines radiating at regular intervals as in Fig. 5.6P. The original purpose of using this type of graph was to investigate the possibility that

↓

Q	REL	1	2	3	4	5	6 -7	8 -10	11 -15	16 -20	21 -30	31 -50	51 -75	76 -100	101 -125	126 -150	151 -175	176 -200
79	3											x	x	x				
100	4	x		x		x	x											
116	6			x				x	x	x	xx							
118	5				x					x	x	x	x					
119	6	x	x			x					x		x			x		
121	3	x	x	x								x	x					
122	5	x		x					x			x						
123	4		x		x			x				x						
126	2	x	x															
130	4								x			x				xx		
132	4		x		x				x							x		
136	6	x	x	x	x	x	x											
137	6			x				xx	x			x	x					
141	1	x																
145	12	x	x	x	x	x	x	x	xx				x	x				
146	9		x					xx	xx	x		xx	xx					x
147	5						xx									xx	x	
148	4	x	x				x								x			
167	4							xx		x	x							
170	2		x											x				
181	2				x			x										
182	4											x			x		xx	
189	2		x	x														
190	7	x	x	x	x		x	x	x									
223	2	x	x															
224	5								x	x	x	x	x					
225	6							x			x	x	xx			x		
226	7	x	x				x		x	xx				x				
227	2	x	x															
230	7		x		x				xx	x	x			x				
250	8	x	x	x	x	x		x	x		x							
261	4	x	x	x	x													
264	2		x		x													
266	5						x		xx	x	x							
268	5	x	x		x			xx										
269	4	x	x	x					x									
272	4		x						xx	x								
273	7	x		x	x	x	xx	x										
274	5		x					x	x						x	x		
317	2							x	x									
323	5							x	x		x			xx				
360	8	x	x		x		x	x	x	x	x							
Totals		19	24	13	14	6	12	20	23	11	14	10	13	6	8	5	-	-

198

FIGURE 5.7T

DOCUMENT OUTPUT CUTOFF SCORE SHEET  
 FOR SMART 'CRAN CON-CON, INDEX NEW QS'  
 FOR 42 DOCUMENTS WITH 200 DOCUMENT  
 COLLECTION

comparison could be made between different index languages by measuring the performance over the whole curve, and the polar coordinate graphs were first tried with the performance curves obtained by the conventional coordination level cut-off as given in Chapter 4, where there was no direct relationship between the various cut-offs. The intention was to calculate the area encompassed by the performance curve within certain limits; with Fig. 5.6P (which is similar to Fig. 4.203P) it was calculated that, in the area bounded by 95% recall and 85% precision, Index Language I.1.a had an area measure of 24.9 while Index Language I.6.a had an area measure of 21.1. It seemed to be unnecessary to do this with these new plots, since the document output cut-off automatically gave an exact match between systems. It was therefore hypothesised that obtaining a normalised recall ratio for all the systems tested would permit an 'order of effectiveness' to be determined. To obtain this normalised recall ratio, the recall ratio at each of the seventeen document cut-off levels would be summed and then divided by seventeen.

|| awc's  
Normalised  
Recall

It was possible to test this idea by using the output from the SMART searches on the same collection. As previously stated, Professor Salton had results for fourteen different options, and Fig. 5.1T shows the output for question 147. Having similar output sheets for all 42 questions, it was possible to prepare a score sheet for each option. As an example the score sheet for 'Cran. Con Con Index News QS' is shown in Fig. 5.7T. Reference to Fig. 5.1T will show that the five relevant documents for Question 147 were ranked at 6, 7, 103, 122 and 138, and it can be seen that this is shown in the appropriate columns of Fig. 5.7T. The recall and precision ratios based on this procedure were obtained for the fourteen SMART options and the results are shown in Fig. 5.8T. The normalised recall ratios for each option were then calculated and are shown in Fig. 5.9T. A normalised recall and normalised precision for each question is given in the output sheets of the SMART searches (see Fig. 6.1) and finally calculated for the complete set of questions; the figures so obtained are also given in Fig. 5.9T. In Fig. 5.10T these two sets of results are arranged in order of effectiveness the higher figures representing the better results. It will be seen that, with very minor variations, the order obtained by the Cranfield normalised recall is the same as that obtained with the SMART normalised recall, with a rank correlation of +.991. This would appear to validate the ranking groups used at Cranfield, and also the simple method we have used to obtain the normalised recall ratio.

To sum up what has been so far discussed, the document ranking method has two major advantages.

1. It enables a series of cut-offs to be applied with equal consistency (i.e. an equal cut-off ratio,  $\frac{100(a+b)}{N}$ ) between tests of different systems using the same document/question sets, and thus solves the problem of totalling sets of results which was discussed in Chapter 3.
2. It enables a series of recall ratios to be obtained which are directly comparable, and permits the calculation of a single measure of performance, normalised recall.

Regarding the measure itself, it was conceived (in a slightly different form) and originally used by Professor Salton. It is a method of representing performance over the whole of the operational range and therefore differs fundamentally from the 'single-point composite measures' which were discussed in Chapter 3. In experimental work of the nature described

DOCUMENTS OUTPUT CUT-OFF	S1		S2		S3		S4		S5		S6		S7		S8		S9		S10		S11		S12		S13		S14					
	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P		
1	8	36	10	45	10	45	12	57	10	45	10	45	10	45	10	45	10	45	10	45	10	45	12	57	10	45	10	45	10	45	10	45
2	15	36	19	44	18	43	24	57	18	43	18	43	24	57	18	43	18	43	20	48	17	40	22	51	18	43	22	51	22	51	19	44
3	21	33	26	40	24	38	28	44	24	37	25	40	28	44	27	43	25	40	27	43	25	40	27	43	24	38	28	44	28	44	25	40
4	26	31	29	35	31	36	27	32	30	35	33	39	29	34	33	39	33	39	34	40	31	36	32	38	31	36	31	36	33	39		
5	31	29	32	30	35	33	31	30	36	34	36	34	34	32	38	36	38	36	39	37	36	34	38	36	35	33	38	36	36	34		
6-7	36	24	40	27	41	28	37	25	43	29	45	30	43	29	45	30	45	30	45	31	43	29	47	32	42	28	44	30	44	30		
8-10	44	21	48	23	51	24	40	19	51	24	52	25	51	24	51	24	51	24	53	25	50	24	54	25	53	25	55	26	56	26		
11-15	52	16	58	18	61	19	51	16	58	18	63	20	59	18	59	19	59	19	62	20	62	20	63	20	65	20	66	21	65	20		
16-20	57	13	65	15	66	16	56	13	62	15	68	16	63	15	65	15	65	15	69	16	67	16	67	16	69	16	72	17	70	17		
21-30	66	10	71	11	76	12	62	10	73	12	73	12	68	11	75	12	75	12	79	12	74	12	74	12	76	12	79	12	76	12		
31-50	77	7	79	7	82	8	76	7	84	8	85	8	78	7	83	8	83	8	87	8	82	8	84	8	82	8	84	8	84	8		
51-75	87	5	87	5	89	6	86	5	92	6	92	6	88	6	91	6	91	6	92	6	89	6	93	6	89	6	90	6	90	6		
76-100	92	4	90	4	92	4	92	4	95	5	96	5	93	4	95	5	95	5	95	5	92	4	96	5	93	4	93	4	95	4		
101-125	94	4	92	3	95	4	96	4	96	4	98	4	96	4	98	4	98	4	98	4	96	4	98	4	97	4	97	4	97	4		
126-150	95	3	95	3	97	3	97	3	97	3	99	3	98	3	98	3	98	3	99	3	97	3	98	3	98	3	98	3	100	3	100	3
151-175	96	3	97	3	98	3	98	3	98	3	99	3	99	3	99	3	98	3	99	3	99	3	99	3	100	3	100	3	100	3	100	3
176-200	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2

FIG. 5-8T RECALL AND PRECISION RATIOS FOR 14 SMART LANGUAGES BASED ON CRANFIELD DOCUMENT OUTPUT CUTOFF.  
(R = Recall Ratio, P = Precision Ratio)

S1 Abstracts old qs. 2  
S2 Abstracts f null. 1  
S3 Abstracts new qs. 3  
S4 Indexing old qs. 4  
S5 Indexing f null. 5  
S6 Indexing new qs. 7  
S7 Cranfield concon. 4  
S8 Abstracts and indexing f null. 11  
S9 Abstracts and indexing new qs. 11  
S10 Abstracts new qs. and indexing f null. 8  
S11 Indexing new qs and f null. 7  
S12 Cranfield concon and abstracts new qs. 10  
S13 Cranfield concon and indexing new qs. 13  
S14 Cranfield concon and indexing f null. 14

<u>SMART LANGUAGE</u>	<u>CRANFIELD NORMALISED RECALL</u>	<u>SMART NORMALISED RECALL AND PRECISION</u>	NR	NP
S1	58.64	1.492	8602	6319
S2	61.06	1.546	8709	6754
S3	62.70	1.573	8867	6861
S4	58.58	1.495	8629	6335
S5	62.41	1.573	8897	6831
S6	64.88	1.609	8992	7094
S7	61.82	1.548	8776	6703
S8	63.64	1.594	8946	6991
S9	65.13	1.618	9061	7115
S10	62.94	1.579	8880	6913
S11	64.94	1.617	9019	7160
S12	63.64	1.593	8947	6988
S13	65.23	1.624	9036	7212
S14	64.82	1.612	9017	7114

FIGURE 5.9T PERFORMANCE FIGURES FOR SMART LANGUAGES

<u>ORDER</u>	<u>CRANFIELD</u>	<u>SMART</u>	NR	NP
1	S13	S13	S9	S13
2	S9	S9	S13	S11
3	S11	S11	S11	S9
4	S6	S14	S14	S14
5	S14	S6	S6	S6
6	S8	S8	S12	S8
7	S12	S12	S8	S12
8	S10	S10	S5	S10
9	S3	S3	S10	S3
10	S5	S5	S3	S5
11	S7	S7	S7	S2
12	S2	S2	S2	S7
13	S1	S4	S4	S4
14	S4	S1	S1	S1

(10) (3) (6)

FIGURE 5.10T COMPARISON OF RANKING OF SMART LANGUAGES BY CRANFIELD AND SMART NORMALISED MEASURES

DOCUMENT OUTPUT CUT-OFF	I-1		I-2		I-3		I-5		I-6		I-7		I-8		I-9	
	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P
47	12	55	12	57	12	57	13	60*	11	52	10	48	9	43	8	36
84	22	51	23	54	23	54	19	45	21	49	21	49	19	44	16	37
126	29	45	30	47	30	48	28	44	29	45	29	46	28	44	22	34
168	35	42	36	42	37	43	32	38	35	42	33	39	32	38	27	32
210	41	39	41	39	43	40	36	34	40	38	40	38	38	36	33	31
294	47	32	48	32	48	32	45	30	47	32	46	31	46	31	40	27
470	56	26	55	26	56	26	53	25	55	26	53	25	55	26	47	22
	62	20	63	20	64	20	59	19	62	19	63	20	62	20	58	18
16-20	67	16	67	16	70	17	65	15	66	16	67	16	68	16	63	15
21-30	76	12	76	12	76	12	73	12	73	12	76	12	76	12	72	11
31-50	85	8	85	8	86	8	82	8	83	8	86	8	85	8	82	8
51-75	89	6	89	6	89	6	88	6	89	6	91	6	91	6	89	6
76-100	92	4	92	4	93	4	91	4	92	4	93	4	93	4	93	4
101-125	95	4	95	4	95	4	94	4	95	4	95	4	96	4	95	4
126-150	98	3	98	3	98	3	96	3	97	3	97	3	98	3	97	3
151-175	99	3	99	3	99	3	98	3	99	3	99	3	99	3	98	3
176-200	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2
NORMALISED RECALL	65.00	65.23	65.82	63.05	64.47	64.05	64.41	61.17								

SINGLE TERM LANGUAGES

- I-1 Natural language
- I-2 Synonyms
- I-3 Word endings
- I-5 Synonyms and quasi-synonyms
- I-6 Synonyms, word endings, and quasi-synonyms
- I-7 Hierarchical reduction first stage
- I-8 Hierarchical reduction second stage
- I-9 Hierarchical reduction third stage

FIGURE 5.11T RECALL AND PRECISION RATIOS AND NORMALISED RECALL FOR SINGLE TERM INDEX LANGUAGES (AVERAGE OF NUMBERS)  
(R = Recall Ratio, P = Precision Ratio)

in this report, it appears to give a valid single measure for comparing the performance of different systems, and, without wishing to be overdogmatic, appears more suitable for this purpose than anything else that has been proposed.

Having - to our satisfaction - established the reasonableness both of the simulated ranking method and also the method for obtaining normalised recall, the procedure was used for the four main groups of index languages. Fig. 5.11T gives the recall and precision ratios for the eight single term languages, while Fig. 5.12T gives similar figures for the fifteen concept languages. The results of the six controlled languages are given in Fig. 5.13T and the searches of titles and abstracts are shown in Fig. 5.14T. These tables also show the normalised recall ratio for each index language. In Fig. 5.15T the index languages are rearranged into an order based on this normalised recall ratio, from which it can be seen that the highest score (65.82) is obtained by Index Language I.3.a (single terms, word forms), with the lowest score (44.64) for Index Language II.1.a (single concepts, natural language). It will be noted that this table also includes the fourteen SMART options.

The figures given so far have been based on what has earlier been described as the average of numbers, and it might be thought that the document ranking method would be particularly susceptible to aberrations which the average of numbers sometimes produces. The results have therefore been recalculated by the average of ratios. To do this, as can be seen from the example in Fig. 5.16T, the indication of a relevant document is replaced by the number representing the percentage of the total recall ratio for that particular question. Thus, with question 79, there were three relevant documents, each document therefore representing 33.3% of the total. With question 100, having four relevant documents, each relevant document is 25% of the total. Question 141 has only one relevant document, so the retrieval of this single document represents 100% recall. These figures are summed for each column, then aggregated and finally, of course, reach a total of 4200. Recall figures can then be obtained.

This process was carried out for all the index languages, and as can be seen from Fig. 5.17T this results in a general increase of two or three points in the normalised recall ratio; however, when placed in order, as in Fig. 5.18T, it can be seen that this order is virtually unchanged from that obtained with the average of numbers, with a positive rank correlation of +.992.

Fig. 5.19T shows the result of ranking documents on the complete collection of 1400 documents. It covers the 42 questions with Index Language I.1.a., and is therefore directly comparable with Fig. 5.3T which was based on the smaller collection of 200 documents. The first eleven ranking groups have been retained, after which they are enlarged to take in the greater number of documents. Fig. 5.20P gives the performance curves for the two situations, and shows that, as would be expected, the smaller generality number for the 1400 document collection adversely affects the performance.

In Chapter 4, Section 8, were given the performance figures for the controlled term languages with Search E, which required some intellect to be applied to the search formulation. The result of ranking the output from these searches is given in Fig. 5.21T, and the

DOCUMENTS OUTPUT CUT-OFF	II - 1		II - 2		II - 3		II - 4		II - 5		II - 6		II - 7		II - 8		II - 9		II - 10		II - 11		II - 12		II - 13		II - 14		II - 15					
	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P	R	P		
1	7	31	7	31	8	38	8	36	6	29	7	31	10	48	10	48	10	48	9	43	8	38	8	38	6	29	5	21	6	29				
2	12	29	13	30	12	29	15	36	11	26	13	30	19	45	17	39	19	45	16	38	18	42	18	42	14	32	9	21	11	26				
3	15	24	18	29	17	26	21	33	18	29	21	33	18	29	17	26	26	41	27	42	27	42	25	39	24	38	17	27	12	18	17	26		
4	18	21	20	24	28	33	21	24	22	26	25	29	21	24	22	26	22	26	29	34	30	36	30	35	30	35	24	28	19	23	21	25		
5	21	20	21	20	31	29	24	23	29	27	29	27	22	21	26	24	26	24	33	31	36	34	38	36	32	30	27	26	21	20	25	23		
6-7	24	16	28	19	35	23	29	20	35	23	31	21	30	20	32	21	41	28	43	29	44	30	36	24	36	24	32	22	27	18	33	22		
8-10	29	14	32	15	41	20	37	17	39	19	36	17	39	18	39	18	45	21	45	21	51	24	41	20	39	19	33	16	41	19				
11-15	31	10	36	11	45	14	42	13	47	15	43	13	48	15	47	15	52	16	52	16	60	19	48	15	48	15	43	13	50	16				
16-20	34	8	39	9	49	12	48	11	53	13	48	11	55	13	56	13	56	13	56	13	64	15	66	15	52	12	54	13	47	11	55	13		
21-30	37	6	43	7	54	8	57	9	60	9	56	9	60	9	66	10	66	10	60	9	75	12	76	12	59	9	64	10	57	9	68	10		
31-50	46	4	52	5	60	6	65	6	68	6	62	6	67	6	74	7	67	6	67	6	84	8	83	8	65	6	74	7	70	7	80	8		
51-75	55	3	59	4	67	4	71	4	74	5	70	4	74	5	80	5	74	5	74	5	89	6	88	6	73	5	81	5	79	5	88	6		
76-100	67	3	70	3	75	4	77	4	80	4	78	4	81	4	85	4	80	4	80	4	93	4	92	4	78	4	85	4	84	4	91	4		
101-125	77	3	80	3	84	3	85	3	88	3	85	3	87	3	90	3	87	3	87	3	95	4	96	4	87	3	91	3	91	3	95	4		
126-150	87	3	89	3	91	3	93	3	94	3	92	3	93	3	93	3	94	3	93	3	97	3	98	3	92	3	95	3	95	3	96	3	98	3
151-175	99	3	99	3	99	3	99	3	99	3	99	3	99	3	99	3	100	3	99	3	100	3	100	3	99	3	99	3	99	3	99	3		
176-200	100	2	100	3	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2	100	2
NORMALISED RECALL	44.64		47.41		53.52		52.05		55.05		51.82		53.88		55.76		57.11		62.88		63.05		55.41		55.88		52.47		57.41					

SIMPLE CONCEPT INDEX LANGUAGES

II-1	Natural language	II-6	II-2 + coordinate (selection)	II-11	II-8 + II-10
II-2	Synonyms	II-7	II-6 + collateral (selection)	II-12	II-2 + complete species
II-3	II-2 + species (selection)	II-8	II-5 + II-7	II-13	II-12 + superordinate
II-4	II-2 + superordinate	II-9	II-2 + alphabetical collateral first stage (selection)	II-14	II-13 + complete collateral
II-5	II-3 + II-4	II-10	II-9 + alphabetical collateral second stage (selection)	II-15	II-11 + II-14

FIGURE 5.12T RECALL AND PRECISION RATIOS AND NORMALISED RECALL FOR SIMPLE CONCEPT INDEX LANGUAGES (AVERAGE OF NUMBERS)  
(R = Recall Ratio, P = Precision Ratio)

DOCUMENT OUTPUT CUT-OFF	III-1			III-2			III-3			III-4			III-5			III-6			IV-1			IV-2			IV-3			IV-4				
	R	P	R	R	P	R	R	P	R	R	P	R	R	P	R	R	P	R	R	P	R	R	P	R	R	P	R	R	P	R		
1	12	55	12	57	11	52	11	52	10	48	10	45	8	38		10	48	10	48	10	48	11	52		10	48	10	48	11	52		
2	17	40	19	44	16	38	13	31	12	29	12	29		17	38	17	39	18	43	18	43		17	38	17	39	18	43	18	43		
3	26	40	25	40	24	37	23	36	19	30	18	28		25	39	24	38	23	39	24	38		25	39	24	38	23	39	24	38		
4	31	36	32	38	28	33	28	33	25	29	24	28		29	34	28	34	29	34	29	34		29	34	28	34	29	34	30	35		
5	36	34	35	33	34	32	33	31	27	25	27	25		32	30	32	30	32	30	32	30		32	30	32	30	35	32	35	32		
6-7	43	29	43	29	41	28	42	28	36	24	34	23		39	25	41	28	39	25	39	25		39	25	39	25	39	25	39	25		
8-10	50	24	49	23	47	22	46	22	46	22	42	20		45	21	47	22	47	22	47	22		45	21	47	22	47	22	46	21		
11-15	58	18	58	18	55	17	54	17	52	16	53	17		55	17	55	17	56	17	54	17		55	17	55	17	56	17	54	17		
16-20	61	14	62	15	60	14	60	14	59	14	60	14		60	14	62	15	64	15	64	15		60	14	62	15	64	15	64	15		
21-30	71	11	69	11	68	11	68	11	68	11	69	11		64	10	67	10	70	11	69	11		64	10	67	10	70	11	69	11		
31-50	80	8	81	8	77	7	76	7	79	7	81	8		76	7	78	7	80	8	80	8		76	7	78	7	80	8	80	8		
51-75	85	5	85	5	86	5	85	5	88	6	88	6		82	5	83	5	87	6	86	6		82	5	83	5	87	6	86	6		
76-100	88	4	90	4	88	4	89	4	93	4	93	4		84	4	85	4	91	4	92	4		84	4	85	4	91	4	92	4		
101-125	95	4	94	4	93	4	94	4	98	4	98	4		89	3	91	3	94	4	95	4		89	3	91	3	94	4	95	4		
126-150	97	3	97	3	95	3	95	3	100	3	99	3		95	3	96	3	97	3	97	3		95	3	96	3	97	3	97	3		
151-175	100	3	99	3	99	3	99	3	100	3	100	3		100	3	100	3	99	3	100	3		100	3	100	3	99	3	100	3		
176-200	100	2	100	2	100	2	100	2	100	2	100	2		100	2	100	2	100	2	100	2		100	2	100	2	100	2	100	2		
<b>NORMALISED RECALL</b>	61.76	61.76	61.76	60.11	59.70	59.58	59.17																									

**CONTROLLED TERM INDEX LANGUAGES**

- III-1 Controlled terms
- III-2 Narrower terms
- III-3 Broader terms
- III-4 Related terms
- III-5 Narrower and broader terms
- III-6 III-4 + III-5

**ABSTRACTS AND TITLES**

- IV-1 Titles natural language
- IV-2 Titles word forms
- IV-3 Abstracts and titles natural languages
- IV-4 Abstracts and titles word forms

FIGURE 5.13T RECALL AND PRECISION RATIOS AND NORMALISED RECALL FOR CONTROLLED TERM INDEX LANGUAGES (AVERAGE OF NUMBERS)  
(R = Recall Ratio, P = Precision Ratio)

FIGURE 5.14T RECALL AND PRECISION RATIOS AND NORMALISED RECALL FOR SINGLE TERM SEARCHES ON ABSTRACTS AND TITLES (AVERAGE OF NUMBERS)  
(R = Recall Ratio, P = Precision Ratio)

<u>ORDER</u>	<u>NORMALISED RECALL</u>	<u>INDEXING LANGUAGE</u>
1	65.82	I-3 Single terms. Word forms
2=	65.23	I-2 Single terms. Synonyms
2=	65.23	S-13 SMART Concon and indexing new qs.
4	65.13	S-9 SMART Abstract and indexing new qs.
5	65.00	I-1 Single terms. Natural language
6	64.94	S-11 SMART Indexing new qs. and f null
7	64.88	S-6 SMART Indexing new qs.
8	64.82	S-14 SMART Concon and indexing f null
9	64.47	I-6 Single terms. Synonyms, word forms, quasi-synonyms
10	64.41	I-8 Single terms. Hierarchy second stage
11	64.05	I-7 Single terms. Hierarchy first stage
12	63.64	S-8 SMART Abstracts and indexing f null
12=	63.64	S-12 SMART Indexing new qs. and f null
14	63.05	I-5 Single terms. Synonyms. Quasi-synonyms
14=	63.05	II-11 Simple concepts. Hierarchical and alphabetical selection
16	62.94	S-10 SMART Abstracts new qs. and indexing f null
17	62.88	II-10 Simple concepts. Alphabetical second stage selection
18	62.70	S-3 SMART Abstracts new qs.
19	62.41	S-5 SMART Indexing f null
20	61.82	S-7 SMART Concon
21	61.76	III-1 Controlled terms
21=	61.76	III-2 Controlled terms. Narrower terms
23	61.17	I-9 Single terms. Hierarchy third stage
24	61.06	S-2 SMART Abstracts f null
25	60.94	IV-3 Abstracts. Natural language
26	60.82	IV-4 Abstracts. Word forms
27	60.11	III-3 Controlled terms. Broader terms
28	59.76	IV-2 Titles. Word forms
29	59.70	III-4 Controlled terms. <i>Related terms Narrower + broader terms.</i>
30	59.58	III-5 Controlled terms. <i>Narrower and broader terms Related terms.</i>
31	59.17	III-6 Controlled terms. Narrower, broader and related terms
32	58.94	IV-1 Titles. Natural language
33	58.64	S-1 SMART Abstracts old qs.
34	58.58	S-4 SMART Indexing old qs.
35	57.41	II-15 Simple concepts. Complete combination
36	57.11	II-9 Simple concepts. Alphabetical first stage selection
37	55.88	II-13 Simple concepts. Complete species and superordinate
38	55.76	II-8 Simple concepts. Hierarchical selection
39	55.41	II-12 Simple concepts. Complete species
40	55.05	II-5 Simple concepts. Selected species and super ordinate
41	53.88	II-7 Simple concepts. Selected coordinate and collateral
42	53.52	II-3 Simple concepts. Selected species
43	52.47	II-14 Simple concepts. Complete collateral
44	52.05	II-4 Simple concepts. Superordinate
45	51.82	II-6 Simple concepts. Selected coordinate
46	47.41	II-2 Simple concepts. Synonyms
47	44.64	II-1 Simple concepts. Natural language

FIGURE 5.15T ORDER OF EFFECTIVENESS BASED ON NORMALISED RECALL FOR 33 CRANFIELD AND 14 SMART INDEX LANGUAGES (AVERAGE OF NUMBERS)

Q	REL	1	2	3	4	5	6 -7	8 -10	11 -15	16 -20	21 -30	31 -50	51 -75	76 -100	101 -125	126 -150	151 -175	176 -200
79	3	33										34					33	
100	4		25							25	25						25	
116	6							16		17	16	17	16		17			
118	5	20					20	20			40							
119	6			16		17	16				17		34					
121	3	33	34	33														
122	5	20				20		20			20	20						
123	4	25		25		25												25
126	2	50	50															
130	4	25						25	25		25							
132	4				25								25		25		25	
136	6		16	17		16	17	34										
137	6		16			17		33	34									
141	1	100																
145	12	8	9	8	9			8	9	8		25	9	8				
146	9		11	11	11	11			11		11	11		11		12		
147	5							20			20		20			20	20	
148	4	25	25		25		25											
167	4	25				25		25						25				
170	2				50									50				
181	2		50								50							
182	4				25									25	25			25
189	2										50	50						
190	7	14	14	14			15		14			14			15			
223	2	50	50															
224	5								20		20	20	20	20				
225	6							16	17		16	17			17	17		
226	7	14	14	15		14						28	15					
227	2	50		50														
230	7	14	14					15		14	14	29						
250	8	12	13	12		13	13	12	25									
261	4	25	25	25	25													
264	2	50	50															
266	5								20	20	20	20		20				
268	5	20	20	20	20		20											
269	4	25	25		25		25											
272	4	25	25		25					25								
273	7	14	14	14	15		14				14			15				
274	5				20		20			20			20			20		
317	2					50		50										
323	5					20			20	20	20	20						
360	8		12		13	12	13	12	13	25								
Totals		677	512	260	288	240	198	306	208	174	378	330	134	174	99	127	70	25

Aggregate total  
48,051

FIGURE 5.16T DOCUMENT OUTPUT CUT-OFF SCORE SHEET AS FIGURE 5.3T CONVERTED TO AVERAGE OF RATIOS.

∴ 42 x 17  
= 67.298

Index Language	Normalised Recall		Index Language	Normalised Recall
I-1	67.2	67.3	III-1	64.2
I-2	67.7	67.6	III-2	64.5
I-3	68.5	67.9	III-3	62.6
I-5	65.6	64.4	III-4	62.4
I-6	66.9		III-5	61.7
I-7	67.4		III-6	61.7
I-8	67.1			
I-9	63.5			
II-1	45.6		IV-1	61.5
II-2	49.0		IV-2	62.4
II-3	55.2		IV-3	62.7
II-4	53.5		IV-4	63.1
II-5	56.3			
II-6	53.8			
II-7	55.6			
II-8	56.8			
II-9	59.3			
II-10	64.9			
II-11	65.1			
II-12	57.2			
II-13	58.4			
II-14	55.0			
II-15	59.8			

FIGURE 5.17T NORMALISED RECALL FOR CRANFIELD INDEX LANGUAGES BASED ON AVERAGE OF RATIOS.

<u>ORDER</u>	<u>INDEX</u> <u>LANGUAGE</u>		<u>ORDER</u>	<u>INDEX</u> <u>LANGUAGE</u>	
1	I-3	(1)	25	IV-4	(26)
2	S-13	(2)	26	IV-3	(25)
3	I-2	(2)	27	III-3	(27)
4	I-7	(11)	28	III-4	(29)
5	S-11	(6)	29	IV-2	(28)
6	I-1	(5)	30	III-5	(30)
7	S-14	(8)	31	III-6	(31)
8	S-9	(4)	32	IV-1	(32)
9	S-8 ?	(10)	33	S-1	(33)
10	I-6	(9)	34	S-4	(34)
11	S-6	(7)	35	II-15	(35)
12	S-8 ?	(12)	36	II-9	(36)
13	S-12	(13)	37	II-13	(37)
14	S-10	(16)	38	II-12	(39)
15	I-5	(14)	39	II-8	(38)
16	II-11	(15)	40	II-5	(40)
17	S-3	(18)	41	II-7	(41)
18	II-10	(17)	42	II-3	(42)
19	III-2	(22)	43	II-14	(43)
20	S-5	(19)	44	II-6	(45)
21	S-2	(24)	45	II-4	(44)
22	III-1	(21)	46	II-2	(46)
23	I-9	(23)	47	II-1	(47)
24	S-7	(20)			

FIGURE 5.18T ORDER OF EFFICIENCY BASED ON NORMALISED RECALL FOR CRANFIELD AND SMART INDEX LANGUAGES CALCULATED BY AVERAGE OF RATIOS (FIGURES IN BRACKETS REPRESENT ORDER WHEN CALCULATED BY AVERAGE OF NUMBERS AS IN FIG. 5.14T)

Q	REL	1	2	3	4	5	6	8	11	16	21	31	51	101	201	401	601	801	1101	
79	3			x																
100	4										x			x	x				x	
116	6										x	x	x	x	x				x	
118	5		x							xx				xx						
119	6							x	x	x			x	x	x					
121	3	x	x	x																
122	5					x						x	x	x	x					
123	4						x		x	x									x	
126	2	x	x																	
130	4	x											x	x	x					
132								x						x					x	x
136	6		x	x		x	xx		x											
137	6			x			x			x	x	xx								
141	1			x																
145	12	x		x		x	xx			x	x		x	xxxx						
146	9	x			x				x	x	x		x	xx					x	
147	5						x						x			x			x	x
148	4	x				x			x	x									x	
167	4						x			x		x				x				
170	2		x																	
181	2								x					x						
182	4										x								xx	x
189	2													x	x					
190	7		x	x		x					x	x		x					x	
223	2	x		x																
224	5											x	xxx						x	
225	6										x	x		xx			x		x	
226	7	x	x			x		x								xxx				
227	2	x						x												
230	7	x					x						xx	xx	x					
250	8	x		x	x		x	x			xx	x								
261	4	x	x	x	x															
264	2	x					x													
266	5											x	xx	x			x			
268	5	x	x		x	x			x											
269	4	x		x					x		x									
272	4	x	x			x							x							
273	7	x	x	x		x	x		x				x							
274	5						x		x				x			x			x	
317	2											x	x							
323	5										x	x	xx	x						
360	8					x	x	x						x	xx					
Totals		17	11	12	5	9	14	6	10	9	13	13	21	25	14	2	1	13	3	

FIGURE 5.19T DOCUMENT OUTPUT CUT-OFF SCORE SHEET FOR INDEX LANGUAGE I.1.a FOR 42 QUESTIONS WITH 1,400 DOCUMENT COLLECTION.

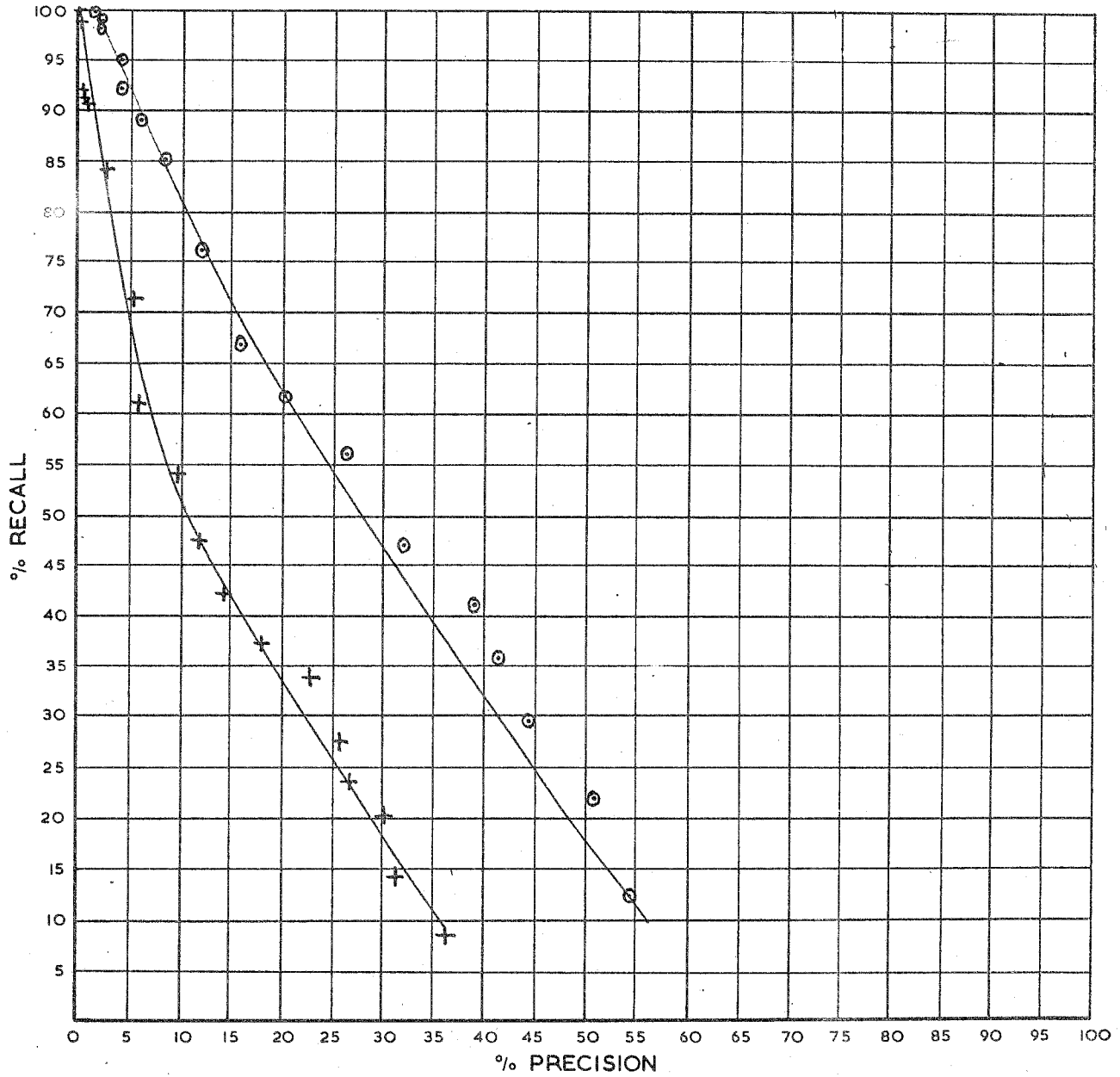


FIGURE 5.20P COMPARISON OF PERFORMANCE, BASED ON DOCUMENT OUTPUT CUT-OFF, FOR COLLECTIONS OF 1,400 AND 200 DOCUMENTS WITH INDEX LANGUAGE I.1.a AND 42 QUESTIONS.  
+ 1,400 Document collection. o 200 Document Collection

normalised recall ratio is shown for each index language by Search E and Search A. It will be seen that there is an improvement with each language of from 1 to 2 points.

Fig. 5.22T shows the ranking score sheet for Index Language I.1.a. with the 42 questions on the 200 document collection, but with the lowest level of exhaustivity of indexing. Fig. 5.23P compares these results with those obtained under similar conditions except that exhaustivity was at its highest level (as Fig. 5.3T).

Four grades of document relevance were used in the tests, and the effect on performance of each of the relevance grades has been considered in Section 6 of Chapter 4. An alternative method of scoring performance from that so far used would be to take account of these relevance gradings by giving each document a weighting related to its relevance grading. The use of the document output cut-off method and normalised recall permits this to be done in what might be considered to be a meaningful manner. A simple form of weighting is to give a score of 4 to those documents rated relevance 1, a score of 3 for documents of relevance 2, a score of 2 for documents of relevance 3 and a score of 1 for documents rated relevance 4. The effect of this would be that question 119, for instance, which has two documents (1378 and 1667) rated relevance 2 and four documents (1324, 1666, 1670 and 2391) rated relevance 3 would now have a total "retrieval score" of  $(2 \times 3) + (4 \times 2) = 14$ .

Referring back to Fig. 5.3T, the score sheet for this question would be amended to show the weighting of each relevant document according to the order in which the documents of the two levels of relevance were retrieved. This was done for the 42 questions by Index Language I.1.a and the amended score sheet is given as Fig. 5.24T. The recall ratio is now determined on the total "points" score for the set of questions, which is 421. At a document cut-off of 1, the recall ratio is therefore shown to be  $\frac{58 \times 100}{421} = 14\%$  and the recall ratios are similarly calculated for the other sixteen cut-off groups. The normalised recall ratio is then calculated as being 67.12.

This procedure was repeated for five other index languages to find whether the effect of a weighting score made any difference to their comparative performance. As can be seen from Fig. 5.25T, there was for each case an increase of approximately two points in the normalised recall, so it does not appear that this method of weighting makes any significant difference to the overall comparison.

The exercise was repeated using different weightings, with a score of 10 for documents rated relevance 1, a score of 5 for documents rated relevance 2, a score of 3 for documents rated relevance 3 and a score of 1 for documents rated relevance 4. This resulted in a further small increase in the normalised recall ratios, but made no significant difference in the comparison between systems. It would be incorrect to state that some form of weighting might not be useful in certain circumstances, but it would seem that it does not have any particular value in this test.

In connection with the normalised recall ratio, it is obvious that there is what could be considered a minimum figure which is based on the random retrieval of the whole collection for every question. For instance, the three relevant documents of Question 79 would, with random retrieval, be ranked 50, 100 and 150, while the seven relevant documents of Question 190 would be ranked 25, 50, 75, 100, 125, 150 and 175. With this particular document/question set, the normalised recall ratio based on this random

DOCUMENT OUTPUT CUT-OFF	III-1		III-2		III-3		III-4		III-5		III-6	
	R	P	R	P	R	P	R	P	R	P	R	P
1	13	62	13	62	11	50	10	48	10	48	8	38
2	21	50	21	50	16	37	16	37	15	36	13	30
3	31	48	31	48	24	38	25	40	23	36	22	34
4	36	42	35	42	29	36	29	36	28	33	28	33
5	39	37	38	37	35	33	35	33	31	29	33	31
6-7	46	31	45	31	43	29	43	29	38	26	39	27
8-10	53	25	52	25	49	23	51	24	46	22	46	22
11-15	62	20	63	20	59	18	60	19	54	17	55	17
16-20	67	16	68	16	63	14	65	15	62	14	65	15
21-30	72	11	74	12	70	10	71	11	70	10	74	12
31-50	79	7	79	7	78	7	79	7	80	7	82	8
51-75	85	5	86	5	85	5	86	5	87	5	89	6
76-100	88	4	88	4	89	4	90	4	91	4	93	4
101-125	93	3	93	3	92	3	93	3	95	3	95	3
126-150	96	3	95	3	95	3	95	3	98	3	98	3
151-175	100	2	100	2	99	2	99	2	99	2	100	2
176-200	100	2	100	2	100	2	100	2	100	2	100	2
NORMALISED RECALL												
Search E	63.58		63.64		61.00		61.58		60.41		61.17	
Search A	61.76		61.76		60.11		59.70		59.58		59.17	

FIGURE 5.21T RECALL AND PRECISION FIGURES FOR INDEX LANGUAGES III.1 - III.6 FOR SEARCH E BY DOCUMENT OUTPUT CUT-OFF METHOD, TOGETHER WITH NORMALISED RECALL FOR SEARCH E AND SEARCH A.  
(R = Recall Ratio, P = Precision Ratio)

Q	REL	1	2	3	4	5	6	8	11	16	21	31	51	76	101	126	151	176
							-7	-10	-15	-20	-30	-50	-75	-100	-125	-150	-175	-200
79	3									x				x		x		
100	4							x			x		xx					
116	6			x					x	x		x	x			x		
118	5			x			x	x	x			x						
119	6	x					x	x	x		x		x					
121	3	x	x		x													
122	5		x		x			x	x			x						
123	4	x	x					x							x			
126	2	x									x							
130	4			x			x		xx									
132	4						x				x				x		x	
136	6	x	x	x		x	x	x										
137	6		x		x		x	x	x		x							
141	1	x																
145	12	x	x	x		x		x	x	x	xx	x		x		x		
146	9					x		x	x	xx			x	x		x	x	
147	5		x								x	x		x		x		
148	4		x	x										x		x		
167	4	x		x								x		x				
170	2	x										x						
181	2		x						x									
182	4		x										x			x	x	
189	2									x					x			
190	7	x	x		x		x	x		x			x					
223	2	x	x															
224	5						x					xx	xx					
225	6					x		x	x		x		x			x		
226	7	x	x	x	x						x	x			x			
227	2	x		x														
230	7	x					x	x	x	x	x				x			
250	8	x		x	x		xx	xx									x	
261	6	x	x	x	x													
264	2	x		x														
266	5							x			x	x	x			x		
268	5	x	x	x	x		x											
269	4		x	x		x	x											
272	4		x			x	x			x								
273	7	x	x	x	x			x	x	x								
274	5			x				x	x				x		x			
317	2			x								x						
323	5						x		x	x		x	x					
360	8			x		x	x	x		x	x	x	x					
Totals		19	18	18	9	6	16	19	15	12	13	15	13	7	5	10	3	

FIGURE 5.22T DOCUMENT OUTPUT CUT-OFF SCORE SHEET ON INDEX LANGUAGE I.1.a FOR 42 QUESTIONS WITH 200 DOCUMENT COLLECTION FOR INDEXING AT LEVEL OF EXHAUSTIVITY 1.

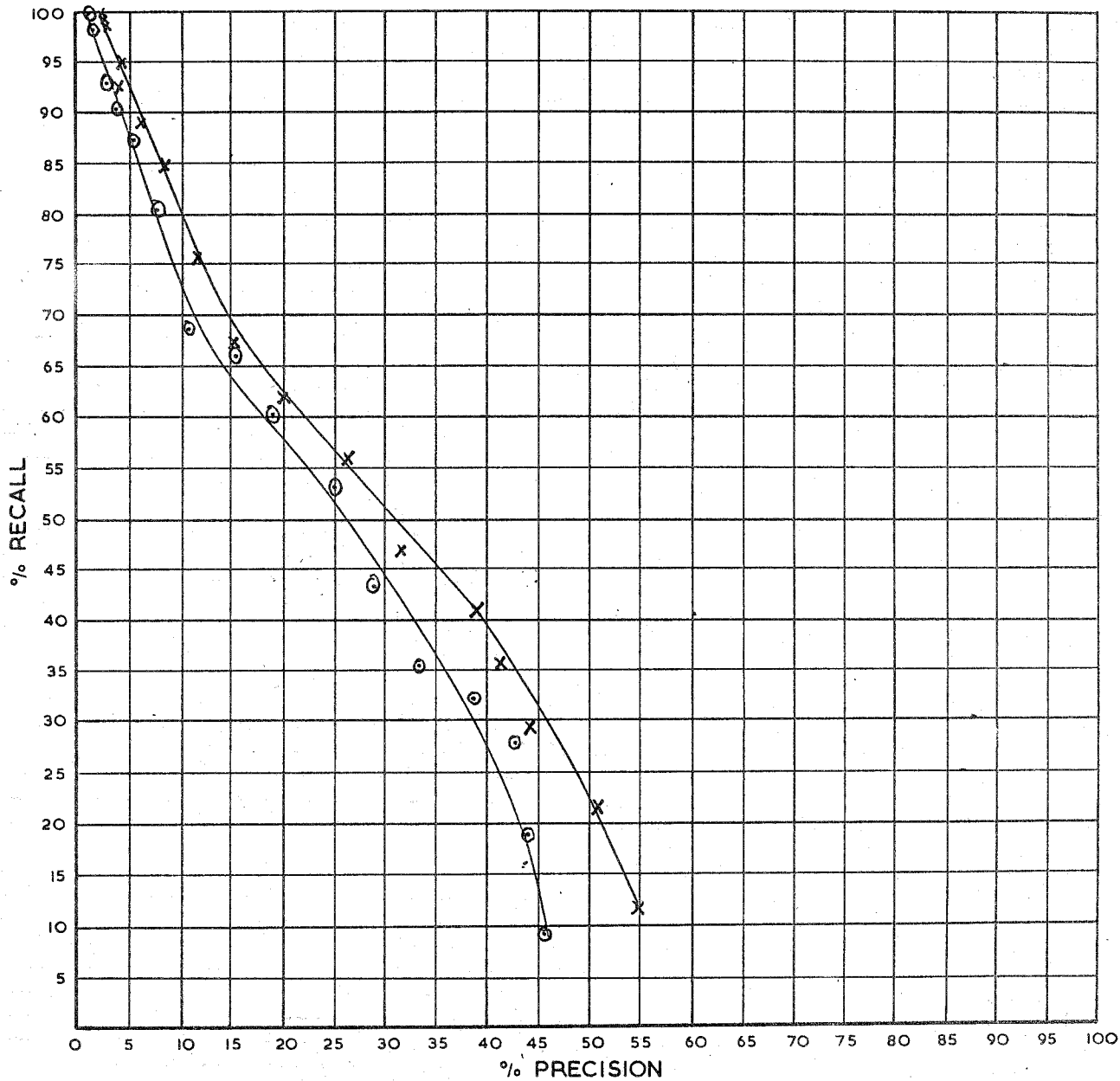


FIGURE 5.23P COMPARISON OF PERFORMANCE, BASED ON DOCUMENT OUTPUT CUT-OFF, FOR EXHAUSTIVITY LEVEL 1 (Fig. 5.22T) AND EXHAUSTIVITY LEVEL 3 (Fig. 5.3T).  
o Exhaustivity Level 1      x Exhaustivity Level 3

*(The next figure is a complete reversal of Fig 4.912P)*  
*See Fig 4.912P*

Q	REL	1	2	3	4	5	6	8	11	16	21	31	51	76	101	126	151	176
79	5	2									2					1		
100	9		3							2	2					2		
116	11							3		2	3	1	1		1			
118	11	3					2	2			4							
119	14			2		3	2				2		5					
121	7	2	3	2														
122	7	2				1		1			1	2						
123	12	4		2		4											2	
126	5	4	1															
130	11	3						3	2		3							
132	11				2							4			2		3	
136	14		2	3		2	3	4										
137	11		3			2		4	2									
141	2	2																
145	16	2	1	2	1			2	1	1		4	1	1				
146	16		2	2	2	2			2		1	2		1		2		
147	7							1			1		1			1	3	
148	10	3	3		2		2											
167	11	4				3		2						2				
170	5		3											2				
181	5				2						3							
182	10				3									2	2			3
189	4										2	2						
190	15	2	3	2			2		3			2			1			
223	4	2	2															
224	10								2		3	1	1	3				
225	13							2	2		4	2			2	1		
226	11	1	1	2		2						4	1					
227	5	2		3														
230	17	2	2					3		2	4	4						
250	15	2	3	2		2	1	2	3									
261	12	4	4	2	2													
264	8	4	4															
266	13								3	2	3	2		3				
268	11	1	3	2	3		2											
269	8	2	1		4		1											
272	8	2	2		2					2								
273	12	3	1	2	1		2		2		1							
274	12				2		2			3			2			3		
317	6					3		3										
323	12					3			3	1	3	2						
360	15		2		2	2	2	2	1	4								
Totals	421	58	49	28	28	29	21	34	26	19	42	32	12	14	8	10	8	3

FIGURE 5.24T RESULTS AS Fig. 5.3T ADJUSTED FOR WEIGHTING BASED ON RELEVANCE GRADES.

Index Language	Normalised Recall Ratio (basic)	Normalised Recall Ratio (weighted)
I.1.a	65.00	67.12
I.7.a	64.05	65.94
III.1.a	61.76	63.64
III.6.a	59.17	61.06
II.9.a	57.11	58.94
II.5.a	55.05	57.11

FIGURE 5.25T COMPARISON OF NORMALISED RECALL RATIOS BY BASIC SCORING METHOD (as Fig. 5.15T) AND BY WEIGHTED SCORING METHOD FOR SIX INDEX LANGUAGES.

retrieval would be 26%. On the other hand, as was discussed earlier in this chapter, the theoretical maximum performance cannot be achieved due to the different numbers of relevant documents for each question, so the highest possible normalised recall ratio would be 86.70%.

It should also be emphasised that the normalised recall ratio only has meaning within the context of the manner in which it has been calculated. In this particular case it was by averaging the results of seventeen cut-off groups as given on page 198. Assume that the number of groups had been reduced to thirteen by combining the first six groups into two larger groups covering documents ranked 51 - 100 and documents ranked 101 - 200. The effect of doing this would be to reduce the normalised recall ratio for index language I.1.a from 65% to 55.7%. On the other hand, if the original groups were broken down so that no groups had more than ten rankings, the normalised recall ratio based upon the resulting twenty-seven groups would be 75.1%. At the same time, the effect of either of these actions would be to change, as considered in the previous paragraph, the minimum figure based on random retrieval and the maximum possible figure.

## CHAPTER 6

### Supplementary tests and results

Any social agency has a duty to study and evaluate its effectiveness and to seek continuously to improve the methods it employs to achieve its objectives. It is not enough to believe, however sincerely, that we are doing good. It is not enough to invoke 'experience' or to collect meaningless and misleading information... It is not enough to rely upon the support of colleagues and those in the same professional group and to accept their endorsement of our work as proof of its effectiveness. Professional in-group support does not measure effectiveness and does not absolve us from accountability for our decisions. The effectiveness of social agencies, it is claimed, is a question to be determined empirically by methods which can be repeated and verified by others.

L.T. Wilkins: Social Deviance, pages 5 and 6

Whereas in the preceding chapter, the main test results were considered on the basis of the document output cut-off method, with normalised recall ratios, we now return to the basic method used in Chapter 4, and present a series of mainly disconnected notes on various supplementary matters. In some cases, new data are presented; in other cases data which have already been given in Chapter 4 is brought together in different ways in order to illustrate more effectively certain points.

### Comparative Results

It is difficult to make direct comparison between the main index languages, because of the inevitable variations created by different numbers of starting terms. However, Fig. 6.1P shows the performance curves for Single Term Natural Language (I.1.a), Simple Concept Natural Language (II.1.a) and Controlled Term, Basic Terms (III.1.a). These might be considered to be comparable since they are all concerned with the basic terms in the particular vocabulary, but the inability of the Simple Concept Index Language to obtain a higher recall figure than 36.9% is due to the severe restrictions which interfixing imposes. That the Controlled Term Index Language also suffers a drop, as compared to Single Term Index Languages, of 7.6% in maximum recall is for the same reason, but the effect is not so severe in this case, since fewer single terms are interfixed. In general the Single Term Natural Language appears to give the best performance.

More reasonable is to make comparison between the index language which have the highest normalised recall ratios in each of the three main groups. \* These would appear to be Index Languages I.3.a (Single term. Word forms), II.10.a, (Simple Concept. Second alphabetical collateral selected), and III.2.a, (Controlled term. Narrower terms). The results are given in Fig. 6.2P, and show that the Simple Concept index language has made a large increase in maximum recall, but again the Single Term index language appears to give the best performance over the whole curve, thus bearing out the figures presented in Chapter 5.

\* II 11a

is best

+ this will

III 1a

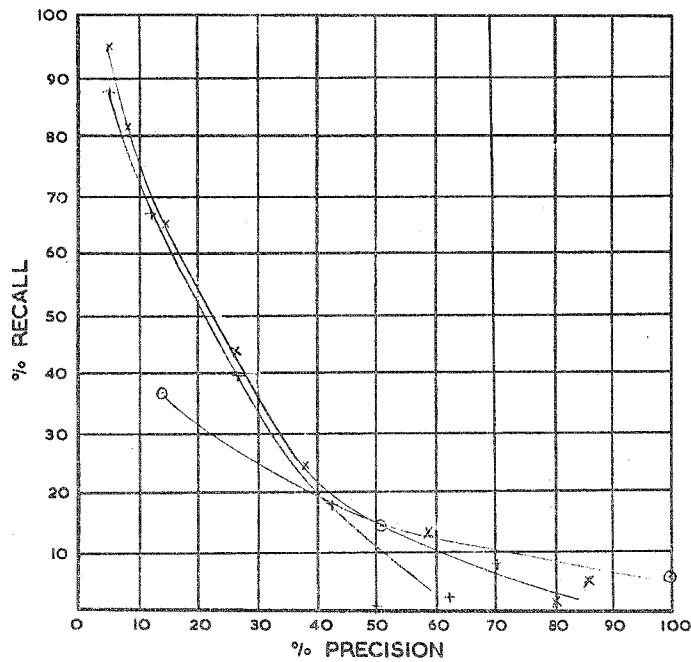


FIGURE 6.1P COMPARISON OF BASIC LANGUAGES IN THE THREE MAIN GROUPS

- x I.1.a. (Fig. 4.140T)
- o II.1.a. (Fig. 4.700T)
- + III.1.a. (Fig. 4.800T)

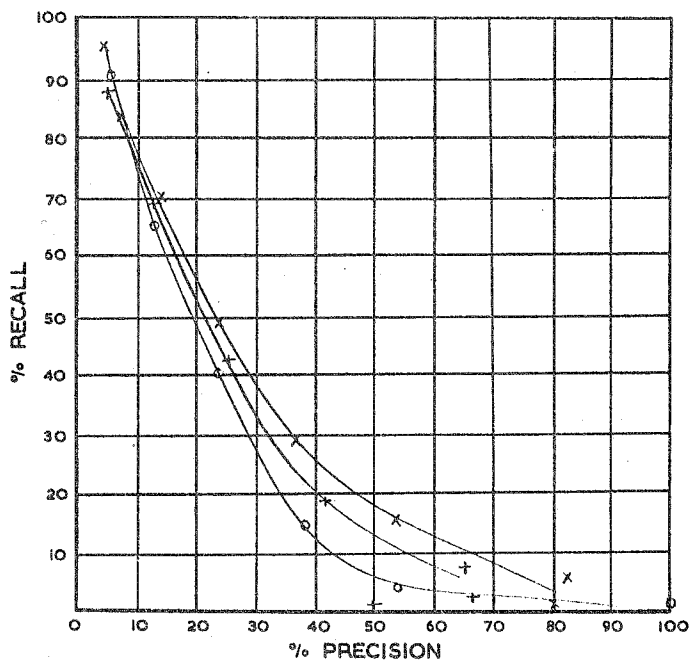


FIGURE 6.2P COMPARISON OF INDEX LANGUAGES GIVING BEST PERFORMANCE IN THREE MAIN GROUPS

- x I.3.a. (Fig. 4.201T)
- o II.10.a. (Fig. 4.709T)
- + III.2.a. (Fig. 4.801T)

Document Relevance

In Chapter 4, Section 6, the effect of relevance was considered, and the results were presented and plotted for documents of different grades of relevance according to the coordination level cut-off. Fig. 6.3T shows the same results as are given in Figs. 4.610T to 4.613T, but now grouped according to relevance grade for each coordination level.

Coordination Level	Relevance Grade	Recall Ratio	Precision Ratio
1	1	94.7	0.3
	1-2	94.2	0.5
	1-3	93.2	0.9
	1-4	94.2	1.1
2	1	85.3	0.7
	1-2	80.6	1.0
	1-3	79.1	2.0
	1-4	77.8	2.4
3	1	60.0	1.2
	1-2	56.1	1.8
	1-3	54.5	3.4
	1-4	48.8	3.6
4	1	42.1	2.2
	1-2	37.4	3.2
	1-3	32.7	5.3
	1-4	29.6	5.8
5	1	25.3	3.1
	1-2	21.3	4.3
	1-3	16.5	6.4
	1-4	16.3	7.6
6	1	14.7	4.5
	1-2	13.5	6.9
	1-3	9.8	9.2
	1-4	9.7	11.4
7	1	7.4	7.1
	1-2	6.5	10.2
	1-3	5.3	16.2
	1-4	5.3	19.2
8	1	3.2	14.3
	1-2	3.9	25.0
	1-3	2.0	25.0
	1-4	1.9	29.2

FIGURE 6.3T RESULTS FOR INDEX LANGUAGE I.1.a FOR 42 QUESTIONS WITH 1400 DOCUMENTS FOR FOUR GRADES OF RELEVANCE.

Plotted as a series of short graphs in Fig. 6.3P, these illustrate yet again the inverse relationship of recall and precision.

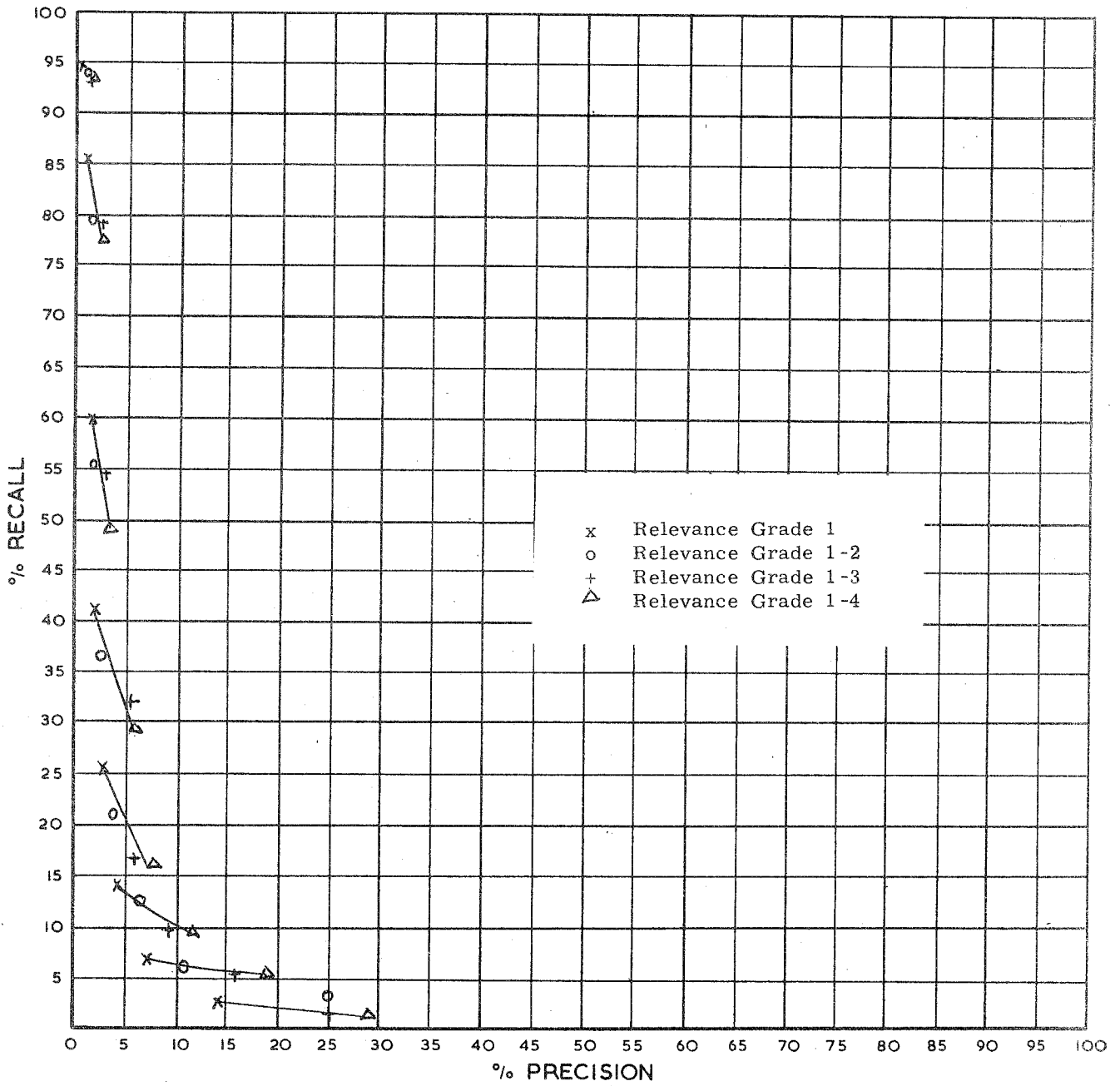


FIGURE 6.3P PLOTS OF EFFECT OF DOCUMENT RELEVANCE AT EIGHT COORDINATION LEVELS

### Basic and supplementary questions

Figs. 6.4T and 6.5T present the results on Index Language I.1.a for the 221 questions when these are divided into the 94 basic questions and 127 supplementary questions (see Vol. 1, Appendix 3G). The basic questions have a generally superior performance, particularly in the middle range of coordination levels and this can be partly accounted for by the higher generality number for this group. On the other hand, documents relevant to the supplementary questions have an average relevance grading that is higher than that for the basic questions (2.7 as against 3.0), and this would have been expected to more than counter the previous effect. It might be suspected that the difference in performance is due to a stronger artificial match between the basic questions and, say, the document titles than exists with the supplementary questions. While analysis does not bear this out, no other adequate explanation can be offered, and the matter is considered again in Chapter 8.

### Average of ratios

On pages 51 to 56, the matter of averaging sets of results was considered, the discussion being on the question of using the average of ratios or the average of numbers. To go into this in more detail, the subset of 35 seven-starting-term questions with Index Language I.1.a on the 1400 document collection is used to demonstrate some difficulties that arise with the average of ratios. Numerical results for the 35 questions can be found in Appendix 4A and the results are presented (by the average of numbers) in Fig. 4.110T.

It can be seen from Fig. 6.6T that, when ratios are obtained for each individual question, three different situations arise. Firstly, there are those questions (e.g. Q82) where it is possible to include recall and precision ratios at all coordination levels to the maximum of 7 (since these are all seven-starting-term questions). Secondly, there are those questions (e.g. Q294) where no documents are retrieved at the higher coordination levels, so no ratios can be included. Thirdly, there are those questions (e.g. Q40) where at the higher coordination levels no relevant documents are retrieved although some non-relevant documents are retrieved. This latter situation is indicated in Fig. 6.6T by an asterisk in the appropriate column. Because of these three different situations, it is a matter for argument as to the figure which should be used for obtaining the average ratios. As an example, at the coordination level of four, the sum of the precision ratios is 561.7. In order to obtain the average precision ratio for the whole set of questions, this figure could be divided by 35, this representing the total number of questions. Alternatively it could be divided by 28, representing the questions which, at this particular coordination level, retrieved some documents, either relevant or non-relevant. Finally it could be divided by 23, representing the number of questions which, at this particular coordination level, retrieved relevant documents. With the results by the average of numbers for comparison, the precision ratios obtained by these three methods are given in Fig. 6.7T.

The first method is clearly unsatisfactory; it would appear to be relatively immaterial as to whether method 2 or 3 should be used, but it is obviously important that when results are presented by the average of ratios, it should be made quite clear as to which procedure has been adopted. The complexity involved in presenting results by the average of ratios is an additional reason why, in this report, we have preferred to

FIGURE 6.4T

Index Language I.1.a  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1-4  
 Number of Documents in Collection 1400  
 Number of Questions 94 Basic Questions  
 Number of Relevant Documents 737  
 Generality Figure 5.6

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d
	Rel.	Non-rel.			
1	694	69037	94.2%	1.0%	51.946%
2	604	26777	82.0%	2.2%	20.148%
3	474	10425	64.3%	4.3%	7.844%
4	324	3874	44.0%	7.7%	2.914%
5	179	1286	24.3%	12.2%	0.967%
6	92	333	12.5%	21.6%	0.251%
7	49	112	6.6%	30.8%	0.082%
8	14	24	1.9%	36.8%	0.018%
9	3	3	0.4%	50.0%	0.002%

FIGURE 6.5T

Index Language I.1.a  
 Exhaustivity of Indexing 3  
 Search Rule A  
 Document Relevance 1-4  
 Number of Documents in Collection 1400  
 Number of Questions 127 Supplementary Questions  
 Number of Relevant Documents 853  
 Generality Figure 4.7

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d
	Rel.	Non-rel.			
1	816	90085	95.6%	0.9%	48.174%
2	679	31345	79.6%	2.1%	16.762%
3	472	11508	55.3%	3.9%	6.154%
4	282	3485	33.1%	7.5%	1.864%
5	135	1094	15.6%	10.8%	0.585%
6	62	366	7.3%	14.5%	0.195%
7	25	94	2.9%	20.7%	0.049%
8	8	19	0.8%	28.8%	0.010%
9	5	2	0.6%	71.4%	0.001%
10	1	0	0.1%	100.0%	0.000%

QUESTION NUMBER	COORDINATION LEVELS													
	1		2		3		4		5		6		7	
	R	P	R	P	R	P	R	P	R	P	R	P	R	P
2	66.7	3.3	37.5	7.6	12.5	12.0	8.3	100	4.2	100				
9	100	0.8	75.0	6.4	25.0	20.0								
34	88.9	1.4	11.1	0.7	-	*	-	*						
40	100	0.2	100	0.5	-	*	-	*		*				
49	100	1.6	100	13.9	66.7	40.0								
67	100	1.7	60.0	2.4	30.0	4.1	-	*						
81	100	1.4	93.3	2.4	86.7	3.7	73.3	7.8	40.0	15.0	20.0	42.9		
82	100	2.0	100	4.2	92.9	8.0	85.7	26.7	64.3	60	42.9	66.7	7.1	100
87	100	0.3	100	0.7	100	1.1	100	8.0	25	33.3				
95	100	1.4	100	2.6	70.0	4.2	50	38.5						
113	100	1.4	88.2	3.0	47.1	5.4	23.5	13.8	5.9	100	<del>5.9</del>	<del>100</del>		
122	100	0.6	100	1.5	60.0	3.6	20	11.1						
131	100	1.3	90.9	3.1	54.5	7.9	27.3	27.3	9.1	100				
132	50.0	0.4	50.0	0.8	25.0	1.4	25.0	5.3						
142	100	1.9	100	3.8	91.7	6.4	83.3	12.5	75.0	14.3	33.3	13.3	33.3	33.3
145	100	1.1	100	2.7	92.3	6.6	61.5	9.3	53.8	21.9	38.5	62.5		
157	100	2.2	71.4	8.3	57.1	16.7								
160	100	0.7	100	1.4	80.0	6.3	60.0	16.7	40.0	50.0				
165	100	0.5	75.0	1.2	75.0	3.6	50.0	11.8	25.0	25.0				
170	100	0.3	50.0	0.5	50.0	1.4	50.0	6.3	50.0	33.3				
171	100	0.6	100	2.1	100	4.4	66.7	11.8	-	*	-	*		
177	100	1.5	88.9	7.1	77.8	11.3	55.6	41.7	-	*				
181	100	0.2	100	1.0	50.0	3.8								
189	100	0.6	-	*	-	*	-	*						
205	100	0.4	100	1.3	66.7	2.8	33.3	20.0						
211	100	0.8	100	4.1	83.3	11.1	16.7	20.0						
219	100	2.5	81.8	7.6	72.7	17.0	27.3	23.1	-	*				
261	100	0.5	100	2.3	100	5.8	100	22.2	100	80	100	100	75	100
285	93.8	1.3	87.5	3.2	62.5	4.8	50.0	16.0	18.8	27.3				
292	77.8	1.3	33.3	3.6	22.2	20.0	11.1	100	<del>11.1</del>	<del>100</del>				
293	100	1.4	60.0	2.9	40.0	10.0								
294	92.3	10.7	76.9	40.0										
299	75.0	0.8	41.7	3.2	8.3	11.1								
315	100	1.6	85.7	2.4	71.4	11.8	14.3	11.8						
338	75.0	0.8	75.0	3.2	75.0	13.0	-	*						
Totals	3319.5	49.5	2733.2	151.7	1946.4	279.3	1092.9	561.7	511.1	660.1	285.4	115.4	233.3	

FIGURE 6.6T. PERFORMANCE RESULTS FOR 35 SEVEN-STARTING-TERM QUESTIONS WITH INDEX LANGUAGE I.1.a CALCULATED BY THE AVERAGE OF RATIOS.

1400 coll see Fig 7.110 T page 90  
 For data in figures subtract from page 289-295.  
 - \* These are counted as retained levels in MRTCL as of ratios.

Coordination Level	Average of numbers	Average of Ratios			Recall	
		1 (Total divided by 35) <i>also = potential with this case</i>	2 (Total divided by figure shown in brackets)	3 (Total divided by figure shown in brackets)	$\div 35$ (= potential)	$\div$ actual
1	1.1%	1.4%	1.4%(35)	1.4%(35)	94.8%	94.8% (35)
2	2.7%	4.3%	4.3%(35)	4.4%(34)	78.1%	78.1% (35)
3	5.2%	8.0%	8.2%(34)	9.0%(31)	55.6%	57.2% (34)
4	13.5%	16.0%	20.1%(28)	24.4%(23)	31.2%	39.0% (28)
5	23.8%	<del>21.7%</del> 18.9%	<del>42.2%(18)</del> 38.8%(17)	54.3%(14)	14.6%	30.1% (17)
6	37.2%	<del>11.0%</del> 8.2%	<del>55.0%(7)</del> 47.6%(6)	64.2%(6)	6.7%	39.1% (6)
7	50.0%	6.7%	77.8%(3)	77.8%(3)	3.3%	38.5% (3)

FIGURE 6.7T. PRECISION RATIOS OBTAINED BY THREE DIFFERENT AVERAGE OF RATIOS PROCEDURES.

use the average of numbers.

Comparison of documents dealing with aerodynamics and structures

The main sets of test results in Chapter 4 were concerned with a subset of 42 questions all of which dealt with aerodynamics rather than structures. For comparison purposes, a set of 42 questions on structures was prepared. Searched on the 1400 document collection, with index language I.1.a, the tests results are given in Fig. 6.8T. Comparison is made in Fig. 6.9P with the results as given in Fig. 4.120T for the 42 aerodynamic questions under the same conditions. This plot shows an unusual characteristic, in that at the higher recall levels, the structure questions have superior precision, but at a recall ratio of about 25%, the curves cross over, and the aerodynamic questions have the better performance.

There are two reasons why one would expect the structure questions to do better. Firstly there are more relevant documents, and therefore the generality number is higher, namely 4.3 as against 3.4. Secondly, although to calculate the generality number N is presumed to be 1400, real N must (as argued on pages 71 - 76) be considerably less than this number.

If the position at a coordination level of 3 is considered, the performance figures are as follows:

Aerodynamics (As Fig. 4.120T)			Structures (As Fig. 6.9T)		
Recall Ratio	Precision Ratio	Fallout Ratio	Recall Ratio	Precision Ratio	Fallout Ratio
66.7%	3.2%	6.790%	67.5%	8.6%	1.732%

To allow for the difference in the generality number, the precision ratio for the aerodynamic questions can be adjusted by the equation given on page 73 and this would result in a new precision ratio of 4.1% which continues to be well below the comparable figure for the structures questions.

FIGURE 6.8T

Index Language I.1.a

Exhaustivity of Indexing 3

Search Rule A

Document Relevance 1-4

Number of Documents in Collection 1400

Number of Questions 42 Structures Questions

Number of Relevant Documents 255

Generality Number 4.3

Coordination Level	Documents Retrieved		Recall Ratio a/a+c	Precision Ratio a/a+b	Fallout Ratio b/b+d	x	y	z
	Rel.	Non-rel.						
1	251	22629	98.4%	1.1%	38.652%	42	42	42
2	216	6081	84.7%	3.4%	10.387%	42	42	42
3	172	1822	67.5%	8.6%	3.112%	41	42	42
4	92	602	36.1%	13.3%	1.028%	34	41	41
5	43	182	16.9%	19.1%	0.311%	23	37	37
6	23	65	9.0%	26.1%	0.111%	15	30	30
7	12	21	4.7%	36.4%	0.036%	10	25	25
8	2	4	0.8%	33.3%	0.007%	2	19	19
9	1	2	0.4%	33.3%	0.003%	2	14	14
10	0	0				0	10	10
11	0	0				0	4	4

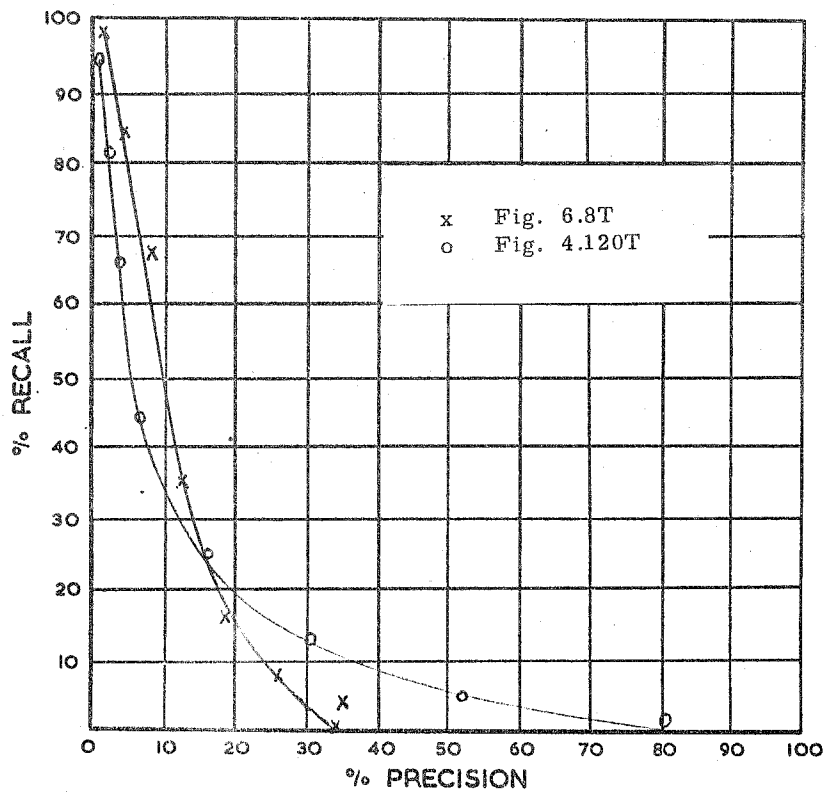


Fig 6.9P

Since this direct adjustment on the basis of generality does not equate performance, it is therefore necessary to consider whether N should be revised for the structures questions. It has already been established (on pages 71 to 76) that real N for the aerodynamics questions is not 1400, but in the region of 1027, at which figure the true generality number is 4.6. One might first hypothesise that the remaining 373 documents represent N for the structures questions; the corresponding generality number would be 16.3. To match this new figure, the adjusted precision ratio for the aerodynamics questions would now be 10.6%, which is higher than the figure (8.6%) for the structures questions. It therefore appears that in the collection of 1400 documents, there must be a subset which is common to both. Using the methods given on pages 71 to 76, N for the structures questions is shown to be probably at least 474, which gives a generality number of 12.8. Adjusted to this generality number the precision ratio for the aerodynamics questions is now 8.6%, the same as for the structures questions. The fallout ratios also now match; for the aerodynamics questions, where N = 1027, the fallout ratio is 9.2%; for the structures questions, where N = 474, the fallout ratio is 9.270%.

The phrase "probably at least 474" was used because no account has been taken of the possibility that the performance figures will be affected by the comparative firmness of the terminology of aerodynamics and structures. The phrase, in fact, implied a belief that aerodynamics has the mushier or more imprecise language, and that for this reason, one would expect the set of structures questions to provide the better performance.

However, the matter is complicated even if this latter point is ignored. At a coordination level of five, the structures questions have a performance of 16.1% recall and 18.1% precision. No exact matching figures can be obtained from Fig. 4.120T, but reference to 4.125P shows that, for the aerodynamic questions, at 16% recall, precision would be approximately 25%. Adjusted for generality on the basis worked out earlier, this would increase the precision ratio to 62%, which is far in advance of the figures for the structures questions. On the other hand at a single term level, it is found that the 42 structures questions have retrieved a total of 22,929 documents, which is an average of 538 documents for each question. This is a figure larger than the 474 documents earlier hypothesised as representing N.

The above discussion is neither clear nor conclusive, and offers no explanation for the crossover in the performance figures of the two sets of questions (which is probably an aberration caused by the relatively small number of results). Rather it serves to point up some of the difficulties which are involved in attempting to compare performance in different subject areas by the coordination level cut-off, and emphasises the necessity for further research in this and related fields.

#### Performance comparison by coordination levels

In Chapter 4, all the tables of results and accompanying performance curves were based on the variation of coordination level. From these tables, sets of figures are extracted where the coordination level is held constant while the variable is the index language. Figs. 6.10T and 6.11T deal with the Single Term index languages at a coordination level of 3 and 6. Figs. 6.12T and 6.13T present the results at coordination levels of 2 and 4 for the Simple Concept index languages, while Figs. 6.14T and 6.15T present results at the same coordination levels for the Controlled Term index languages.

FIGURE 6.10T

Coordination Level 3

Exhaustivity 3

Document Relevance 1-4

Search Rule A

Number of Documents in Collection 200

Number of Questions 42

Number of Relevant Documents 198

Generality ~~Figure~~ 23.6

Index Language	Documents Retrieved		Recall Ratio	Precision Ratio	Fall-out Ratio	x	y	z
	Rel.	Non-rel.						
I.1.a	132	761	66.7%	14.8%	9.278%	42	42	42
I.2.a	134	837	67.7%	13.8%	10.205%	42	42	42
I.3.a	139	913	70.2%	13.2%	11.131%	42	42	42
I.5.a	144	1613	72.7%	8.2%	19.666%	42	42	42
I.6.a	151	1785	76.3%	7.8%	21.763%	42	42	42
I.7.a	143	989	72.2%	12.6%	12.058%	42	42	42
I.8.a	146	1081	73.7%	11.9%	13.180%	42	42	42
I.9.a	163	2034	82.3%	7.4%	24.799%	42	42	42

FIGURE 6.11T

Coordination Level 6

Exhaustivity 3

Document Relevance 1-4

Search Rule A

Number of Documents in Collection 200

Number of Questions 42

Number of Relevant Documents 198

Generality ~~Figure~~ 23.6

Index Language	Documents Retrieved		Recall Ratio	Precision Ratio	Fall-out Ratio	x	y	z
	Rel.	Non-rel.						
I.1.a	25	17	12.6%	59.5%	0.207%	15	33	33
I.2.a	25	23	12.6%	52.1%	0.280%	15	33	33
I.3.a	32	28	16.2%	53.3%	0.341%	17	33	33
I.5.a	41	69	20.7%	37.3%	0.841%	23	33	33
I.6.a	47	87	23.7%	35.1%	1.061%	23	33	33
I.7.a	31	31	15.7%	50.0%	0.378%	20	33	33
I.8.a	33	37	16.7%	47.1%	0.451%	21	33	33
I.9.a	46	171	23.2%	21.2%	2.085%	25	33	33

FIGURE 6.12T

Coordination Level 2

Exhaustivity 3

Document Relevance 1-4

Search Rule A

Number of Documents in Collection 200

Number of Questions 42

Number of Relevant Documents 198

Generality Figure 23.6

Index Language	Documents Retrieved		Recall Ratio	Precision Ratio	Fall-out Ratio	x	y	z
	Rel.	Non-rel.						
II.1.a	28	27	14.1%	50.9%	0.329%	18	41	41
II.2.a	44	65	22.2%	40.4%	0.792%	25	41	41
II.3.a	60	115	30.3%	27.9%	1.890%	32	41	41
II.4.a	67	255	33.8%	20.8%	3.109%	36	41	41
II.5.a	81	359	40.9%	18.4%	4.377%	39	41	41
II.6.a	57	155	28.8%	26.9%	1.890%	33	41	41
II.7.a	75	318	37.9%	19.1%	3.877%	33	41	41
II.8.a	99	605	50.0%	14.1%	7.376%	40	41	41
II.9.a	75	296	37.9%	20.2%	3.609%	33	41	41
II.10.a	129	942	65.2%	12.0%	11.485%	32	41	41
II.11.a	146	1259	73.7%	10.4%	15.350%	41	41	41
II.12.a.	72	287	36.4%	20.1%	3.500%	34	41	41
II.13.a	117	937	59.1%	11.1%	11.424%	40	41	41
II.14.a	143	2047	72.2%	6.5%	24.957%	40	41	41
II.15.a	168	2590	84.8%	6.1%	31.578%	41	41	41

FIGURE 6.13T

Coordination Level 4

Exhaustivity 3

Document Relevance 1-4

Search Rule A

Number of Documents in Collection 200

Number of Questions 42

Number of Relevant Documents 198

Generality Figure 23.6

Index Language	Documents Retrieved		Recall Ratio	Precision Ratio	Fall-out Ratio	x	y	z
	Rel.	Non-rel.						
II.1.a	0	0				0	32	32
II.2.a	2	0	1.0%	100%	0.000%	2	32	32
II.3.a	4	0	2.0%	100%	0.000%	3	32	32
II.4.a	2	8	1.0%	20.0%	0.098%	4	32	32
II.5.a	6	10	3.0%	37.5%	0.122%	5	32	32
II.6.a	4	5	2.0%	44.4%	0.061%	4	32	32
II.7.a	6	9	3.0%	40.0%	0.110%	4	32	32
II.8.a	14	33	7.1%	29.8%	0.402%	10	32	32
II.9.a	13	6	6.6%	68.4%	0.073%	8	32	32
II.10.a	27	44	13.6%	38.0%	0.536%	18	32	32
II.11.a	34	99	17.2%	25.6%	1.207%	19	32	32
II.12.a	9	3	4.5%	75.0%	0.037%	5	32	32
II.13.a	20	54	10.1%	27.0%	0.658%	11	32	32
II.14.a	26	179	13.1%	12.7%	2.182%	19	32	32
II.15.a	51	325	25.8%	13.6%	3.962%	26	32	32

FIGURE 6.14T

Coordination Level 2

Exhaustivity 3

Document Relevance 1-4

Search Rule A

Number of Documents in Collection 200

Number of Questions 42

Number of Relevant Documents 198

Generality Figure 23.6

Index Language	Documents Retrieved		Recall Ratio	Precision Ratio	Fall-out Ratio	x	y	z
	Rel.	Non-rel.						
III.1	136	946	68.7%	12.6%	11.534%	42	42	42
III.2	137	1024	69.2%	11.8%	12.485%	42	42	42
III.3	147	1575	74.2%	8.5%	19.203%	42	42	42
III.4	148	1661	74.7%	8.2%	20.251%	42	42	42
III.5	184	3044	92.9%	5.7%	37.113%	42	42	42
III.6	187	3482	94.4%	5.1%	42.465%	42	42	42

FIGURE 6.15T

Coordination Level 4

Exhaustivity 3

Document Relevance 1-4

Search Rule A

Number of Documents in Collection 200

Number of Questions 42

Number of Relevant Documents 198

Generality Figure 23.6

Index Language	Documents Retrieved		Recall Ratio	Precision Ratio	Fall-out Ratio	x	y	z
	Rel.	Non-rel.						
III.1	36	49	18.2%	42.4%	0.597%	20	34	34
III.2	38	53	19.2%	41.8%	0.646%	20	34	34
III.3	43	78	21.7%	35.5%	0.951%	22	34	34
III.4	44	82	22.2%	34.9%	1.000%	23	34	34
III.5	62	420	31.3%	12.9%	5.121%	19	34	34
III.6	72	527	36.4%	12.0%	6.425%	31	34	34

If one makes the assumption that the coordination level of 3 for the Single Term index languages is approximately equal to a coordination level of 2 for the Simple Concept and Controlled Term index languages, then it is possible to present in a bar chart a representation of what happens in regard to recall and precision ratios when moving from one index language to another. Index Language II.1.a has the lowest recall ratio and highest precision ratio so this is taken as the starting point in Fig. 6.16T.

#### Effects of precision devices

In Chapter 4, Section 3, the results of tests on the Single Term index languages with interfixing and partitioning were presented. Figure 6.17T and 6.18T make extracts from these tables of the figures at the coordination level of 4.

#### Effects of question generality

The individual results for each of the 221 questions with the 1400 document collection and Index Language I.1.a are given in Appendix 4A, and the figures for this particular set of results are given in Fig. 4.100T. As discussed in Chapter 3, this set of questions was a heterogeneous group in a number of respects; various breakdowns have now been made.

First the questions have been grouped according to the number of documents relevant to each question, and table 6.19T shows the recall and precision ratios for each of the groups.

There appears to be a general trend towards a lower recall ratio at any given coordination level for those searches where there are increased numbers of relevant documents; as usual this is matched by a higher precision ratio. If the questions having 1-4 relevant documents and the questions having 16 or more relevant documents are grouped, then this change becomes more apparent, as is shown in Fig. 6.20T.

However, the marked increase in the precision ratio at any given recall ratio is obviously due to the large increase in the generality number of the questions having 16 or more relevant documents. If one considers the fallout ratio, it can be seen from Fig. 6.21P that when recall is plotted against fallout, those questions which have four or less relevant documents have markedly superior performance.

It would probably be correct to say that, as a rule, a question having few relevant documents is a specific query, while a question having a very large number of relevant documents is likely to be a general question. From this it is reasonable to hypothesise that a specific question should present a simpler retrieval problem, a general question. Without suggesting that the results presented above prove this hypothesis, they can certainly be said to support it.

#### Effect of number of postings

The next breakdown of the 221 questions was made by grouping the questions according to the numbers of total postings of the question search terms; information on this is included with the set of results in Appendix 4A. For instance, as can be checked from Appendix 5.1 of Volume I, the three search terms of Question 295 (i.e. 'uniformly', 'loaded', 'sectors'.) have a total of only 46 postings, while for Question 106, the nine search terms have a total of 3,474 postings. Ten groups were formed on this basis, each

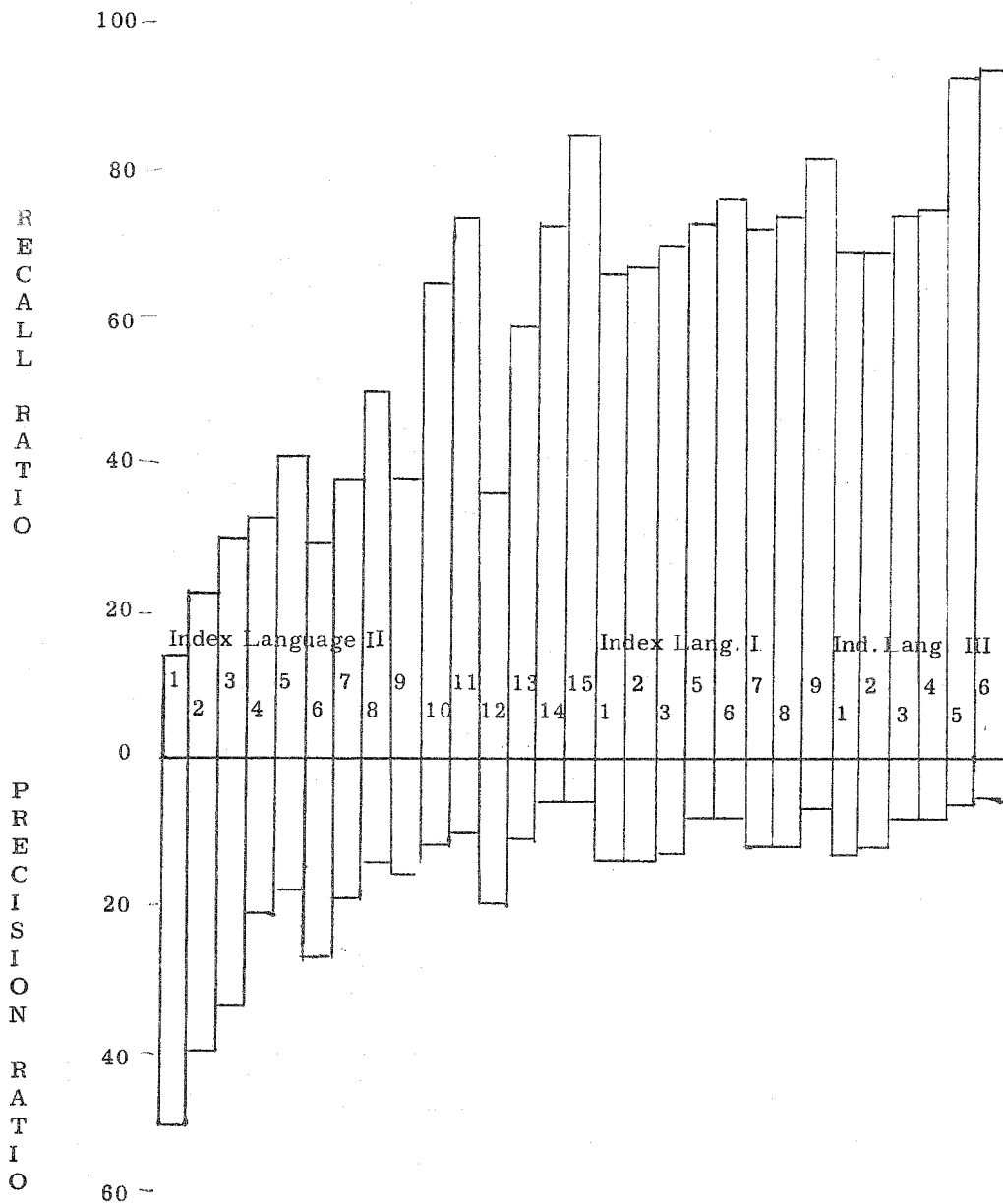


FIGURE 6.16T RECALL AND PRECISION RATIOS WITH 200 DOCUMENTS AND 42 QUESTIONS AT COORDINATION LEVEL OF 2 FOR SIMPLE CONCEPT AND CONTROLLED TERM INDEX LANGUAGES; AND AT COORDINATION LEVEL OF 3 FOR SINGLE TERM INDEX LANGUAGES.

FIGURE 6.17T

Exhaustivity of Indexing 3  
 Search Rule B  
 Document Relevance 1-4

Number of Documents 1400  
 Number of Questions 19 (Subset 4)  
 Number of Relevant Documents 131  
 Generality Number 4.9

Coordination Level 4

Figure	Index Language	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio	x	y	z
		Rel.	Non-rel.						
4.300T	I.1.a	37	127	28.2%	22.6%	0.480%	15	19	19
4.310T	I.1.b (Partitioning)	26	75	19.8%	25.7%	0.283%	12	19	19
4.320T	I.1.c (Interfixing)	25	54	19.1%	31.6%	0.204%	12	19	19
4.330T	I.1.d (Interfixing & Partitioning)	16	36	12.2%	30.8%	0.136%	10	19	19

FIGURE 6.18T

Exhaustivity of Indexing 3  
 Search Rule B  
 Document Relevance 1-4

Number of Documents 1400  
 Number of Questions 19 (Subset 4)  
 Number of Relevant Documents 131  
 Generality Number 4.9

Coordination Level 4

Figure	Index Language	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio	x	y	z
		Rel.	Non-rel.						
4.302T	I.6.a	58	765	44.3%	7.0%	2.890%	19	19	19
4.312T	I.6.b (Partitioning)	39	321	29.8%	10.8%	1.213%	16	19	19
4.322T	I.6.c (Interfixing)	38	312	29.0%	10.9%	1.179%	17	19	19
4.332T	I.6.d (Interfixing & Partitioning)	27	148	20.6%	15.4%	0.559%	12	19	19

*These figures don't match with page 23*

No. of relevant documents per question	No. of questions in group	COORDINATION LEVEL									
		1	2	3	4	5	6	7	8	9	10
1	6	100 0.1	83.3 0.2	66.7 0.3	50.0 1.5	29.0 9.0	14.5 18.8	12.9 29.6	3.2 28.6	1.6 33.7	
2	31	95.2 0.3	78.0 0.7	69.4 1.7	45.2 4.0	29.0 9.0	14.5 18.8	12.9 29.6	3.2 28.6	1.6 33.7	
3	18	98.1 0.4	94.4 1.3	66.7 2.8	46.3 5.5	24.1 7.4	18.5 14.5	9.3 38.5	3.7 100		
4	25	92.0 0.5	76.0 1.1	58.0 2.5	37.0 6.5	23.0 17.0	11.0 36.7	6.0 66.7	2.0 100		
5	26	100 0.6	90.0 1.5	66.2 3.0	39.2 4.8	18.5 6.9	7.7 9.9	4.6 25.0	2.3 100	0.8 100	
6	20	97.4 0.7	90.4 1.7	78.9 3.5	62.3 6.0	41.2 9.1	26.3 15.1	13.2 15.2	5.3 24.0	3.5 57.1	0.9 100
7	16	98.2 1.2	76.8 3.4	52.7 6.2	25.0 9.9	12.5 18.2	3.6 17.4	2.7 50.0	0.9 50.0		
8	14	97.3 1.2	85.7 2.7	59.8 3.9	35.7 6.6	16.1 7.4	6.3 7.4	0.9 11.1			
9	15	97.8 1.3	83.7 3.7	59.3 8.6	35.6 14.9	16.3 19.5	8.1 36.7	5.2 53.8	0.7 50.0	0.7 100	
10	7	94.3 1.2	84.3 2.3	57.1 4.0	35.7 9.5	11.4 14.5	8.6 25.0	2.9 50.0	2.9 66.7	1.4 100	
11-15	30	93.3 3.2	82.6 3.7	61.0 6.0	38.2 9.7	20.1 13.4	10.4 22.1	3.5 21.9	0.5 11.8		
16-20	8	99.3 1.6	87.2 3.1	61.7 6.2	44.0 11.7	18.4 19.7	4.3 30.0	2.1 100			
21-25	2	81.3 3.1	45.8 4.3	16.7 4.5	10.4 10.9	4.2 15.4					
26-30	1	78.6 5.5	50.0 9.3	32.1 31.0	10.7 100	7.1 100	3.6 100				
31-40	2	75.0 4.5	51.4 5.6	40.3 9.2	33.3 14.2	15.3 25.6	9.7 38.9	4.2 75.0			

FIGURE 6.19T RESULTS OF 221 SEARCHES ON 1400 DOCUMENTS WITH INDEX LANGUAGE I.1.a GROUPED ACCORDING TO NUMBER OF DOCUMENTS RELEVANT TO QUESTION.

COORDINATION LEVEL	1-4 RELEVANT DOCUMENTS			16 OR MORE RELEVANT DOCUMENTS		
	RECALL RATIO	PRECISION RATIO	FALLOUT RATIO	RECALL RATIO	PRECISION RATIO	FALLOUT RATIO
1	96.7%	0.3%	50.181%	77.7%	6.2%	19.670%
2	86.1%	0.9%	15.635%	49.3%	5.5%	14.042%
3	68.1%	1.7%	6.382%	31.1%	8.8%	5.378%
4	45.9%	4.1%	1.687%	21.6%	14.7%	2.101%
5	25.4%	8.1%	0.473%	10.1%	25.9%	0.486%
6	15.6%	16.1%	0.128%	5.4%	38.1%	0.147%
7	10.7%	32.5%	0.035%	2.0%	75.0%	0.010%
8	1.6%	28.6%	0.006%			
9	0.8%	33.3%	0.002%			

FIGURE 6.20T RESULTS AS FIGURE 6.19T GROUPED FOR QUESTIONS WITH FEW AND MANY RELEVANT DOCUMENTS

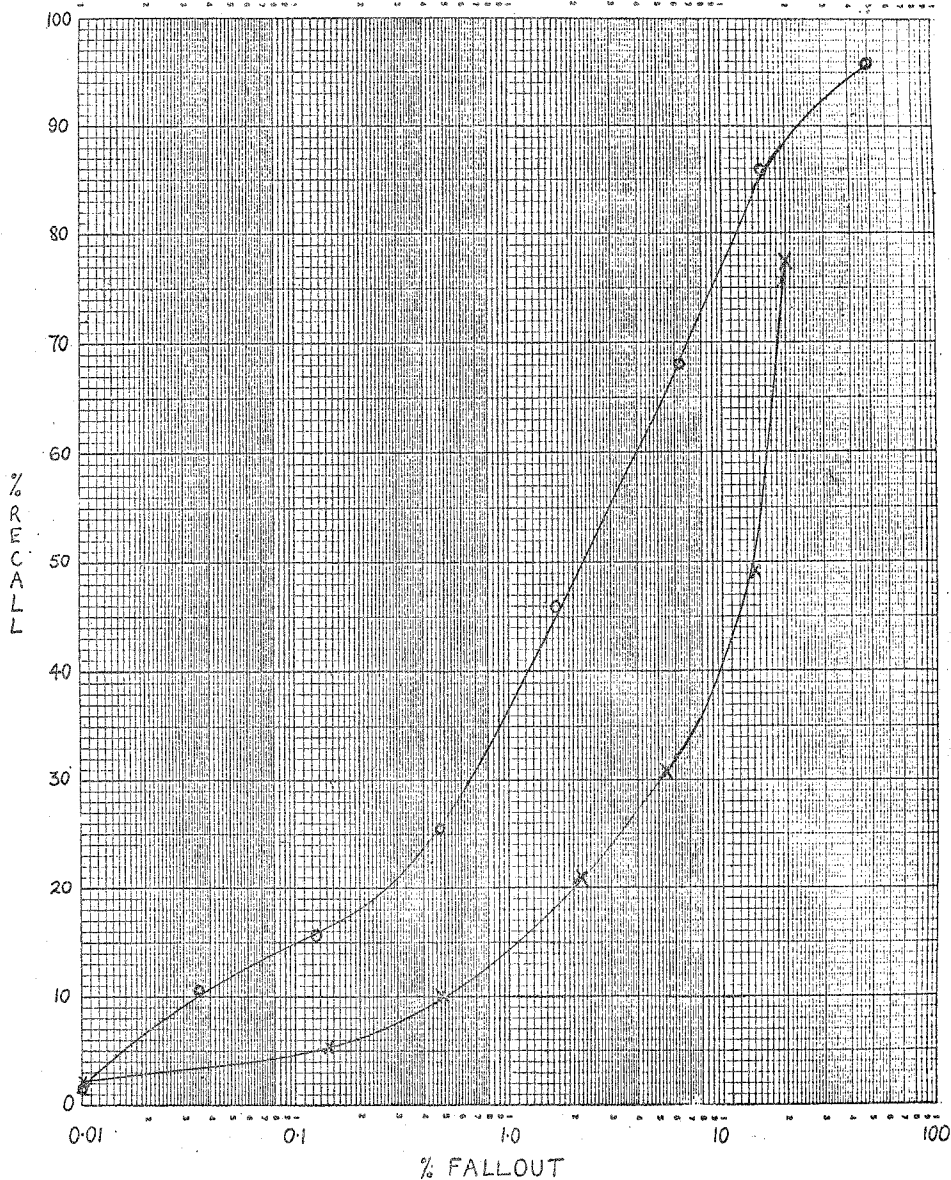


FIGURE 6.21P RECALL/FALLOUT PLOT OF RESULTS PRESENTED IN FIGURE 6.20T

containing approximately similar numbers of questions; the results are presented in Fig. 6.22T.

With some minor aberrations, these show, for any given coordination level, that the increase in total term postings results in a regular increase in recall ratio accompanied by a corresponding decrease in precision ratio. The next stage was to group the questions in relation to the average number of postings for each term. However, the preliminary stage of making up the groups of questions by this method showed that the groups differed little from those used in Fig. 6.22T, so no further work on this was done.

#### Order of retrieval of relevant documents

An analysis was made of the effect on retrieval of individual relevant documents in moving from one index language to another. The results in Chapter 4 show, for instance, that with Index Language I.1.a, at a coordination level of 5, there were 94 relevant documents retrieved (see Fig. 4.200T). With Index Language I.6.a, at a coordination level of 6, there were 87 relevant documents retrieved (see Fig. 4.203T). The question is whether the change of index language and the increase in coordination levels resulted in a different or a similar set of retrieved documents.

To investigate this point, nine index languages were selected namely I.1.a, I.5.a, I.8.a, II.1.a, II.5.a, II.10.a, II.15.a, III.1.a and III.6.a. For the 42 questions, the records were checked to find the order of retrieval of the relevant documents for each of the nine languages. Some examples are given in Fig. 6.23T, which shows, for Questions 118, 170 and 250, the coordination levels at which the relevant documents were retrieved, and in Fig. 6.24T a ranked order of retrieval. From this type of data for the 42 questions, it was too involved to sort out what happened to each individual relevant document, but an analysis was made for each of the three main groups of index languages to find what happened to the relevant documents ranked first and last in the basic languages (i.e. I.1.a, II.1.a, and III.1.a). While it was not possible to make a clear cut decision every time, Fig. 6.25T shows that in the very large majority of cases, the change from one language to another did not alter the retrieval rank of the first and last documents retrieved.

TOTAL NUMBER OF TERM POSTINGS	COORDINATION LEVEL																				
	1		2		3		4		5		6		7		8		9		10		
	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	
0/300	85.1	3.9	56.5	19.8	23.6	49.4	7.1	65.0	1.0	100.0											
301/600	93.7	1.9	70.8	7.2	45.7	20.6	10.2	33.9	4.0	77.7	1.1	66.6	.5	100.0							
601/300	91.3	1.1	71.4	3.6	51.2	9.0	29.9	23.3	11.1	31.0	8.6	64.2	3.8	66.6	1.4	75.0	-	-	-	100.0	
801/1000	97.3	.7	84.3	2.1	67.8	5.3	48.6	13.7	27.8	32.6	13.0	60.0	6.9	80.0	.8	100.0	.8	100.0			
1001/1200	95.0	.7	84.2	2.4	54.2	5.2	32.1	11.8	15.7	23.1	6.4	40.9	4.2	100.0							
1201/1400	96.9	.4	84.3	2.3	61.6	5.0	43.4	9.7	25.7	14.8	14.1	21.8	7.5	45.4	1.0	100.0					
1401/1600	98.4	.7	86.4	2.8	66.9	3.9	45.8	9.7	19.5	14.0	6.7	18.3	3.7	31.2	.7	33.3	.7	100.0			
1601/1900	98.7	.5	94.1	1.3	70.7	2.7	47.4	6.1	24.0	12.9	5.8	12.8	3.2	20.0	1.2	33.3	-	-	-	-	
1901/2500	98.9	.8	91.5	1.5	79.3	3.0	60.3	6.3	33.8	11.8	15.8	23.0	5.8	28.9	.5	33.3	.5	100.0			
2501/3500	100.0	.8	100.0	1.2	94.6	1.8	82.9	3.3	50.0	4.4	36.1	8.5	14.8	10.4	10.6	23.8	4.2	66.6	1.0	100.0	

FIGURE 6.22T RESULTS OF 221 SEARCHES ON 1400 DOCUMENTS WITH INDEX LANGUAGE I.1.2 GROUPED ACCORDING TO TOTAL NUMBER OF STARTING TERM POSTINGS

QUESTION NUMBER	INDEX LANGUAGE	COORDINATION LEVEL								
		1	2	3	4	5	6	7	8	
118	I.1			1324 1378		1667 1666 1670				
	I.5			1378		1324		1667 1670	1666	
	I.8			1324	1378	1666	1667 1670			
	II.1	1667 1324 1378		1666 1670						
	II.5	1667 1378	1324		1666 1670					
	II.10		1324 1378		1666 1670	1667				
	II.15			1324 1378	1666 1670	1667				
	III.1		1324 1378		1667 1666 1670					
	III.6			1324 1378	1667 1666 1670					
170	I.1	1605				1360				
	I.5			1605			1360			
	I.8		1605		1360					
	II.1	1360								
	II.5	1605			1360					
	II.10	1605			1360					
	II.15		1605		1360					
	III.1		1605		1360					
250	I.1				2364 2367 1335	1798 1316	1311 1415 1416			
		I.5			2367 1335	1798 2364 1316	1311 1415 1416			
		I.8			2367 1335	1798 2364 1316	1311 1415 1416			
	II.1	1415 2364 1335	1311 1798 1316							
		II.5		1416 1798 2364 2367 1335	1311 1415 1316					
	II.10	2367	1798 2364 1316 1335	1311 1415 1416						
		II.15		1798 2364 2367 1335	1311 1415 1416 1316					
	III.1		1335 2367	2364	1316 1798	1311 1415 1416				
		III.6		1335 2367	2364 1798	1316 1415 1416				

FIGURE 6.23T COORDINATION LEVEL OF RETRIEVAL OF RELEVANT DOCUMENTS FOR QUESTIONS 118, 170 AND 250 BY NINE INDEX LANGUAGES

QUESTION NUMBER	RELEVANT DOCUMENTS	INDEX LANGUAGE								
		I.1	I.5	I.8	II.1	II.5	II.10	II.15	III.1	III.6
118	1324	4=	4	5	3=	3	4=	4=	4=	4=
	1378	4=	5	4	3=	4	4=	4=	4=	4=
	1666	1=	1	3	1=	1=	2=	2=	1=	1=
	1667	1=	2=	1=	3=	4=	1	1	1=	1=
	1670	1=	2=	1=	1=	1=	2=	2=	1=	1=
170	1360	1	1	1	1	1	1	1	1	1
	1605	2	2	2	2	2	2	2	2	2
250	1311	1=	1=	1=	1=	1=	1=	1=	1=	1=
	1316	4=	4=	4=	1=	1=	4=	1=	4=	4=
	1335	6=	7=	7=	4=	4=	4=	5=	7=	8=
	1415	1=	1=	1=	4=	1=	1=	1=	1=	1=
	1416	1=	1=	1=	7=	4=	1=	1=	1=	1=
	1798	4=	4=	4=	1=	4=	4=	5=	4=	4=
	2364	6=	4=	4=	4=	4=	4=	5=	6	6=
2367	6=	7=	7=	7=	4=	8	5=	7=	6=	

FIGURE 6.24T RANKED ORDER OUTPUT FOR RELEVANT DOCUMENTS OF QUESTIONS 118, 170 AND 250 BY NINE INDEX LANGUAGES

INDEX LANGUAGE GROUP	HIGHEST RANKED DOCUMENTS		LOWEST RANKED DOCUMENTS	
	Maintained Position	Changed Position	Maintained Position	Changed Position
I SINGLE TERMS	33	3	33	1
II SIMPLE CONCEPTS	26	7	28	2
III CONTROLLED TERMS	30	6	29	4

FIGURE 6.25T EFFECT ON RANK OF HIGHEST AND LOWEST RANKED DOCUMENTS IN MOVING TO DIFFERENT INDEX LANGUAGE FOR 42 QUESTIONS

## CHAPTER 7

### CITATION INDEXING AND BIBLIOGRAPHIC COUPLING

It is true that librarians do an almost religious job of storing information; it is placed on record, but without evaluation, and much of it is not worth its rental of space. Each operational unit tends to use its own special language, and translators are very few. In warfare, the questioning of returning forces is regarded as a highly skilled speciality. Would social action (and perhaps social research) not gain much if an analogous speciality could be created to assess and consolidate relevant information? Such an organization could not confine its attention to information retrieval, no matter how efficient such a retrieval process might be. Storage and retrieval systems do not represent evaluation and consolidation of information .

L.T. Wilkins: Social Deviance Page 111

In Chapter 7 of Vol. I, an account was given of the compilation of a citation index, and the subsequent preparation of bibliographic coupling groups, with a view to an evaluation being made of this index.

As stated in this earlier account, it is a matter for some argument as to how this type of index can be tested in an experimental situation. In carrying out a test in an operational environment, there would be no difficulty beyond the effort required, but although several different ways of presenting the results of this test have been considered, there does not appear to be any procedure which can be considered entirely satisfactory. For this reason two sets of figures are being given; the first method probably results in a performance which is better (in comparison with the results obtained with conventional systems) than should be the case, but it has the major advantage that it does not involve any manipulation of the test results and therefore permits direct comparison to be made between different subsets of questions, different relevance decisions etc. If an evaluation of a citation index can be carried out in an experimental environment, then the second method of presenting results is probably nearer the real performance of the system, and is used for comparison with the results of the conventional index languages.

As described in some detail in Vol. I (page 110 and Fig. 7.5) the score sheets for each question gave the results with coupling strength from 1 to 7. The basic scoring at the seven coupling levels for the 42 aerodynamics questions with the 1400 document collection is shown in Fig. 7.1T and the results for this set of questions are presented in Fig. 7.2T. Fig. 7.3T presents the results for the 42 questions dealing with structures, while Fig. 7.4T gives the results for the 35 questions having 7 starting terms. A comparative plot of these three question sets is given in Fig. 7.5P.

All the results so far shown are obtained with documents of relevance grades 1-4; Figs. 7.6T, 7.7T and 7.8T show the results for documents of relevance grades 1-3, relevance grades 1-2 and relevance grade 1, with the 42 aerodynamics questions. Fig. 7.9P plots the results of the four grades of relevance.

QUES- TION	COUPLING LEVEL													
	R <sup>1</sup>	NR <sup>1</sup>	R <sup>2</sup>	NR <sup>2</sup>	R <sup>3</sup>	NR <sup>3</sup>	R <sup>4</sup>	NR <sup>4</sup>	R <sup>5</sup>	NR <sup>5</sup>	R <sup>6</sup>	NR <sup>6</sup>	R <sup>7</sup>	NR <sup>7</sup>
79	2	136	2	53	0	15	0	3						
100	4	76	4	19	4	10	2	7	2	2	2	0		
116	6	204	6	77	3	33	3	16	3	7	3	5	2	2
118	3	262	3	64	2	23	2	7	0	4				
119	6	317	6	83	4	31	3	10	0	5				
121	-	-	-	-	-	-	-	-	-	-	-	-	-	-
122	2	201	0	45	-	18	-	7	-	5	-	0	-	1
123	-	73	-	15	-	8	-	3	-	3	-			
126	2	65	2	22	-	8	-	6	-	3	-	1		
130	4	136	2	23	2	10	2	3	-	2	-	1	-	1
132	2	397	-	108	-	40	-	20	-	10	-	6	-	4
136	6	48	5	8	3	5	2	3	2	3	-	2	-	1
137	6	125	4	33	4	11	3	5	3	2	3	-	2	-
141	-	6	-	2	-	2								
145	12	314	12	76	10	38	9	19	9	9	5	4	5	3
146	9	363	9	82	8	45	7	19	7	8	4	4	4	3
147	5	24	3	2	3	-	3	-	2	-	2	-		
148	-	40	-	11	-	3	-	1						
167	4	145	2	36	-	25	-	12	-	3				
170	-	79	-	17	-	3								
181	-	103	-	26	-	12	-	4	-	1	-	1		
182	3	47	2	9	-	4								
189	-	19	-	3	-	1								
190	5	42	5	6	4	1								
223	2	63	2	13	-	1								
224	3	221	2	56	-	8	-	1						
225	4	636	4	199	3	95	3	56	3	39	3	30	3	17
226	7	74	7	15	5	3	3	0	3	0				
227	2	8	2	1	2	1								
230	2	274	-	65	-	27	-	3						
250	7	213	7	62	6	37	5	20	2	11	2	5	2	2
261	-	30	-	8	-	2	-	1	-	1				
264	2	9	2	1	2	1	2	1	2	-	2	-	2	-
266	4	141	3	35	-	14	-	3	-	1	-	1		
268	-	2	-	-										
269	-	-	-	-										
272	3	94	3	23	2	10	2	3	-	2	-	1	-	1
273	6	43	6	7	2	4	2	3	-	3	-	2	-	2
274	-	49	-	12	-	1								
317	2	13	2	4	2	-								
323	2	66	2	13	2	8	-	2	-	1	-	1	-	1
360	4	337	4	112	3	34	3	19	3	10	3	5	3	-

FIGURE 7.1T RESULTS FOR 42 SEARCHES WITH 1400 DOCUMENTS BY BIBLIOGRAPHIC COUPLING

FIGURE 7.2T

Index Language Citation Indexing and Bibliographic Coupling  
Document Relevance 1-4

Number of Documents in Collection 1400  
Number of Questions 42 (Subset 2 Aerodynamics)  
Number of Relevant Documents 198  
Generality Number 3.4

Coupling Strength	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
	Rel.	Non-rel.			
1	131	5495*	66.1%	2.3%*	9.345%*
2	113	1446	56.3%	7.2%	2.459%
3	76	592	38.2%	11.4%	1.007%
4	56	257	28.1%	17.9%	0.437%
5	41	135	20.6%	23.3%	0.230%
6	29	70	14.6%	29.3%	0.119%
7+	23	38	11.6%	37.7%	0.065%

FIGURE 7.3T

Index Language Citation Indexing and Bibliographic Coupling  
Document Relevance 1-4

Number of Documents in Collection 1400  
Number of Questions 42 (Structures)  
Number of Relevant Questions 252  
Generality Number 4.3

Coupling Strength	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
	Rel.	Non-rel.			
1	179	5962*	71.0%	2.9%*	10.684%*
2	138	1569	54.8%	8.1%	2.668%
3	90	541	35.7%	14.3%	0.920%
4	51	140	20.2%	26.7%	0.238%
5	36	93	14.3%	27.9%	0.193%
6	31	54	12.3%	36.5%	0.092%
7+	25	28	9.9%	47.2%	0.047%

FIGURE 7.4T

Index Language Citation Indexing and Bibliographic Coupling  
Document Relevance 1-4

Number of Documents in Collection 1400  
Number of Questions 35 (Subset 1)  
Number of Relevant Documents 237  
Generality Number 5.9

Coupling Strength	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
	Rel.	Non-rel.			
1	195	11658	67.9%*	1.6%*	23.792%*
2	156	3068	54.3%	4.8%	6.261%
3	115	1152	40.1%	9.1%	2.351%
4	70	464	27.5%	14.5%	0.947%
5	56	299	19.5%	15.7%	0.610%
6	43	174	15.0%	19.8%	0.355%
7+	40	101	13.9%	28.4%	0.206%

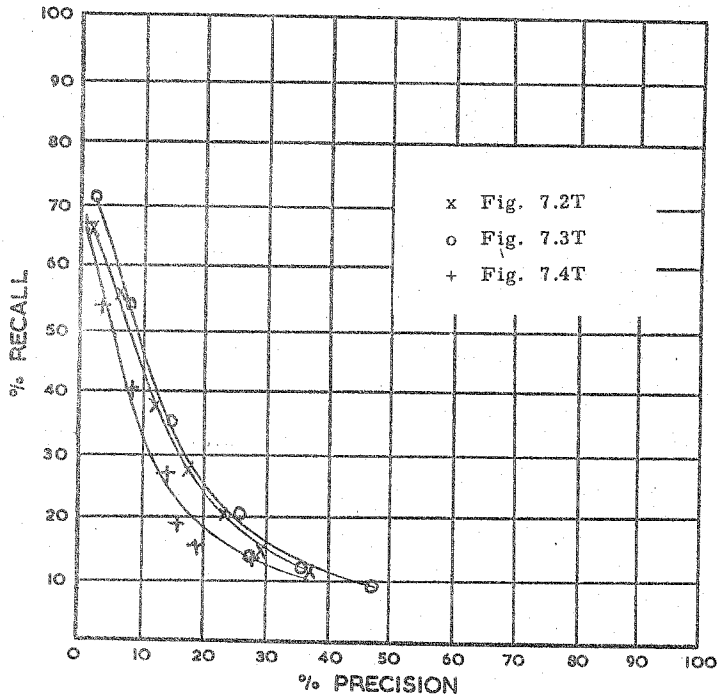


FIGURE 7.5P PLOT OF PERFORMANCE WITH BIBLIOGRAPHIC COUPLING FOR SETS OF 42 AERODYNAMIC QUESTIONS, 42 STRUCTURES QUESTIONS and 35 SEVEN-STARTING-TERM QUESTIONS.

FIGURE 7.6T

Index Language Citation Indexing and Bibliographic Coupling  
Document Relevance 1-3

Number of Documents in Collection 1400  
Number of Questions 35 (Subset 1)  
Number of Relevant Documents 212  
Generality Number 4.3

Coupling Strength	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
	Rel.	Non-rel.			
1	139	-	65.5%	-	-
2	106	2021	50.0%	5.0%	4.124%
3	75	759	35.4%	9.0%	1.549%
4	48	370	22.6%	11.5%	0.755%
5	38	187	17.9%	16.9%	0.382%
6	29	102	13.7%	22.1%	0.208%
7+	28	58	13.2%	32.6%	0.118%

FIGURE 7.7T

Index Language Citation Indexing and Bibliograph  
Document Relevance 1-2

Number of Documents in Collection 1400  
Number of Questions 35  
Number of Relevant Documents 79  
Generality Number 1.6

Coupling Strength	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
	Rel.	Non-rel.			
1	51	-	64.5%	-	-
2	35	843	44.3%	4.0%	1.720%
3	26	412	32.9%	5.9%	0.841%
4	12	197	15.2%	5.7%	0.402%
5	11	105	13.9%	9.5%	0.214%
6	9	66	11.4%	12.0%	0.135%
7+	8	41	10.1%	16.3%	0.084%

FIGURE 7.8T

Index Language Citation Indexing and Bibliographic Coupling  
Document Relevance 1

Number of Documents in Collection 1400  
Number of Questions 35  
Number of Relevant Documents 18  
Generality Number 0.4

Coupling Strength	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
	Rel.	Non-rel.			
1	10	-	55.5%	-	-
2	7	194	38.9%	3.5%	0.396%
3	5	84	27.8%	5.6%	0.171%
4	2	41	11.1%	4.7%	0.084%
5	2	28	11.1%	6.7%	0.057%
6	2	18	11.1%	10.0%	0.037%
7+	2	12	11.1%	14.3%	0.024%

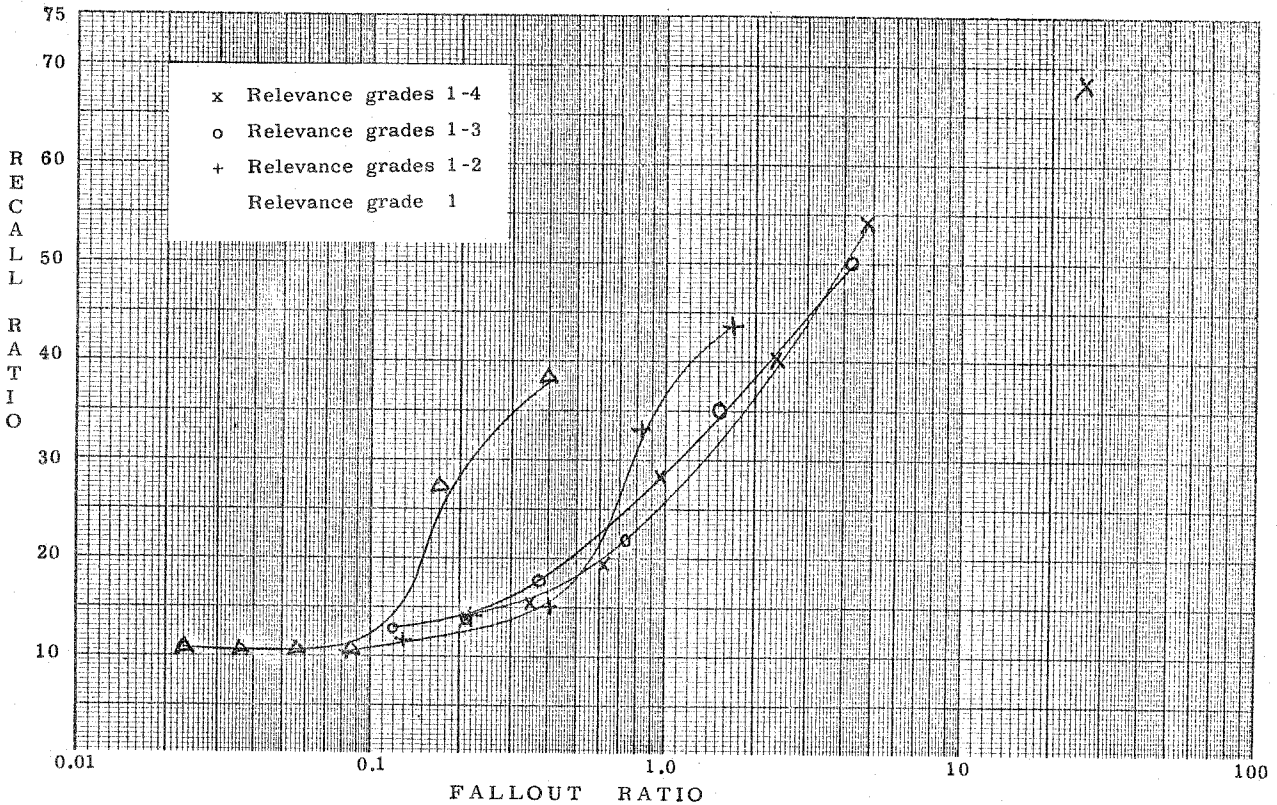


FIGURE 7.9P RECALL/FALLOUT GRAPH FOR PERFORMANCE WITH BIBLIOGRAPHIC COUPLING FOR DOCUMENTS OF FOUR GRADES OF RELEVANCE, WITH 1400 DOCUMENT COLLECTION AND 35 SEVEN-STARTING-TERM QUESTIONS.

In the method used for compiling the scores in the above results, what might be described as the "entry document" was, in the scoring, also counted as a successfully retrieved document; in other words, a previously known relevant document was scored as being a successfully retrieved relevant document. To put it at its simplest, Q227 (as can be seen from Fig. 7.1T) has two relevant documents, the numbers of which were 2087 and 2088. When document 2087 was used as an "entry document", it was found that it had three references in common with document 2088, and therefore both documents were entered as being retrieved at a coupling strength of 3. To take another example Q100 has four relevant documents, numbers 1785, 1786, 1787 and 1788. In the test search, documents 1787 and 1788 were found to have a coupling strength of 6, and documents 1785 and 1786 had a coupling strength of 3. However, there were no references that were common to the pair of documents 1785 and 1786 on the one hand or the pair of documents 1787 and 1788 on the other hand. In spite of this, it would be scored as all four relevant documents having been retrieved at a coupling strength of 3 and lower. As a third example, for Q116 there were six relevant documents, numbers 1317, 1574, 1575, 1576, 1578 and 1656. In the search, document 1576 had a coupling strength of 6 with documents 1574 and 1578, and a coupling strength of 2 with documents 1575 and 1317. In addition document 1317 had a coupling strength of 2 with document 1656. Therefore, at this coupling level this would be recorded as a successful retrieval of all six relevant documents.

By the second method of presenting the results, allowance would be made for these various situations. With Q227, the "entry document" would be eliminated from the scoring; it would be considered that there was only one relevant document, and that this was retrieved. With Q100, however, the first "entry document" would be eliminated from the scoring, but since there was no link between the two pairs of documents, it would be considered that of the three remaining relevant documents, two had been retrieved. With Q116, the "entry document" would be eliminated from the scoring, but since the other five documents were linked either directly or indirectly with the "entry document", all these five documents would be included in the scoring.

On the other hand, with those questions such as Q.122 or Q.132, where no relevant documents were retrieved, the total of relevant documents would in each case be reduced by one.

The result of this exercise is to produce a new set of performance figures where there are now only 156 relevant documents, and the results are presented in Fig. 7.10T. In doing this, it is only the recall and precision ratios that are changed, for the fallout ratio remains the same as in Fig. 7.2T.

It was earlier suggested that it would be reasonable to compare the results by this method with those obtained by the coordination level cut-off. However, as the generality number has been changed, by eliminating 42 relevant documents, it is necessary for this to be done on a recall/fallout graph as in Fig. 7.11P where comparison is made with the Single Term index languages which gave the best and worst performance.

Further tests were carried out where account was taken of the proportional match between documents, this being based on the number of references in the documents concerned. The procedure for doing this was described on pages 110 and 112 of Vol. I. It can make no difference to the

FIGURE 7.10T

Index Language Citation Indexing and Bibliographic Coupling  
Document Relevance 1-4

Number of Documents in Collection 1400  
Number of Questions 42 (Subset 2)  
Number of Relevant Documents 156  
Generality Number 2.7

Coupling Strength	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
	Rel.	Non-rel.			
1	95	5495*	60.8%	1.7%*	9.345%*
2	79	1446	50.5%	5.1%	2.459%
3	51	592	32.6%	7.9%	1.007%
4	36	257	23.0%	12.2%	0.437%
5	25	135	16.0%	15.6%	0.230%
6	18	70	11.5%	20.4%	0.119%
7+	12	38	7.6%	24.0%	0.065%

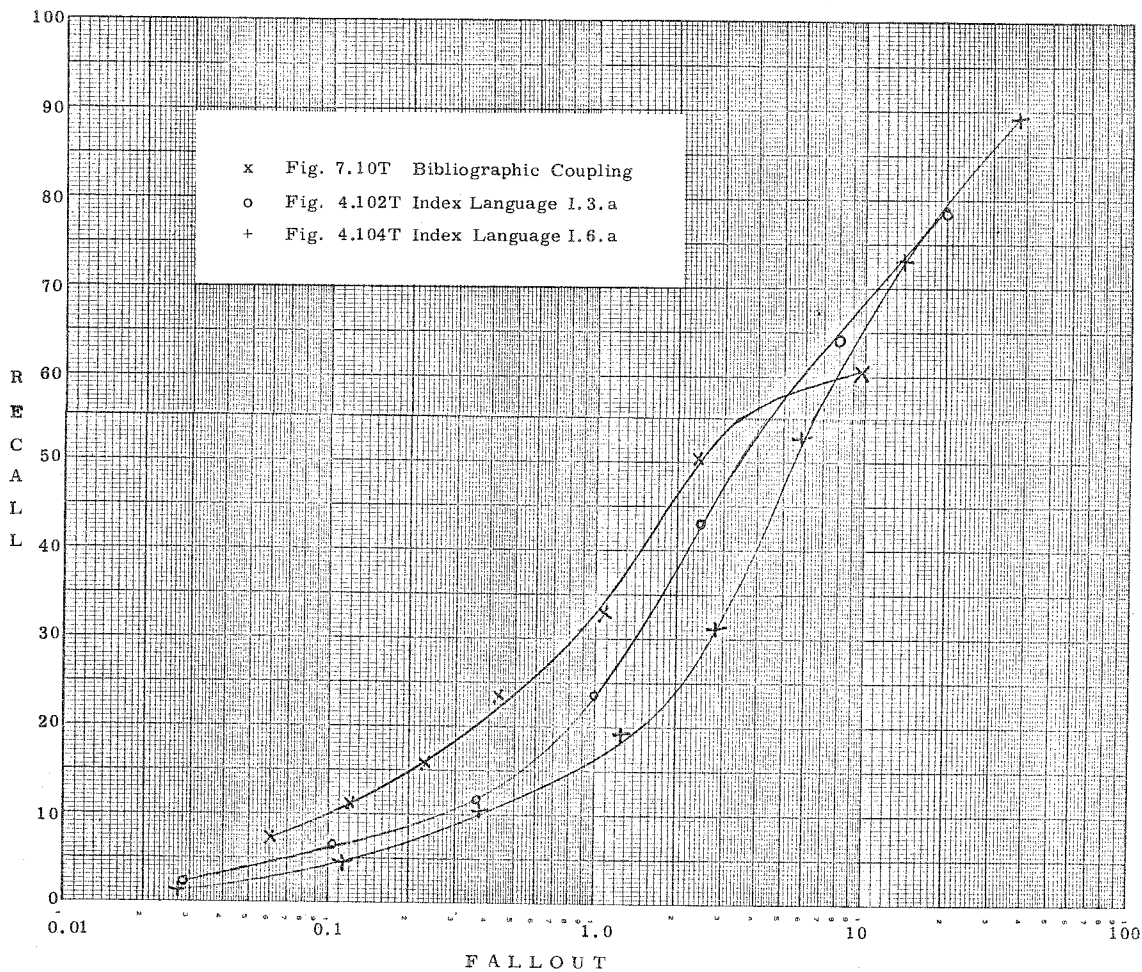


FIGURE 7.11P RECALL FALLOT PLOT FOR BIBLIOGRAPHICAL COUPLING AND SINGLE TERM INDEX LANGUAGES I.3.a AND I.6.a WITH 1400 DOCUMENTS AND 42 AERODYNAMIC QUESTIONS

FIGURE 7.12T

Index Language Citation Indexing and Bibliographic Coupling (Weighted)  
Documents Relevance 1-4

Number of Documents in Collection 1400

Number of Questions 42 (Subset 2)

Number of Relevant Documents 198

Generality Number 3.4

Weighted Coupling Strength	Documents Retrieved		Recall Ratio	Precision Ratio	Fallout Ratio
	Rel.	Non-rel.			
150+	131	5495*	66.1%	2.3%*	9.377%*
81-150	122	1758	61.6%	6.5%	3.000%
51-80	111	1387	56.0%	7.4%	2.367%
31-50	90	854	45.5%	9.5%	1.457%
21-30	67	432	33.8%	13.4%	0.737%
16-20	51	207	25.7%	19.7%	0.353%
11-15	38	126	19.1%	23.1%	0.221%
6-10	23	59	11.6%	28.0%	0.101%
3-5	12	20	6.0%	37.5%	0.034%
1-2	2	1	1.0%	66.7%	0.002%

final figure of relevant and non-relevant documents retrieved, but replaces the groups formed at the various coupling levels as given in earlier totals with new groups based on the weighted scores. The results on the 42 aerodynamic questions are shown in Fig. 7.12T; although different groups are formed, there appears to be little variation from the performance for the same document/question set presented in Fig. 7.2T.

As stated in the opening chapter of this volume, we have considerable reservations in presenting these results, in particular when it comes to attempting to make comparison with the performance obtained by conventional methods. One thing that can be stated positively is that the same inverse relationship exists; bibliographic coupling is a precision device which has very much the same effect as coordination in a conventional system.

Since approximately 12% of the documents did not contain any references, it was inevitable that the maximum recall ratio should fall well short of 100%. In the event it appears that, with this collection, something around 70% recall might be expected; for any recall ratio lower than this, the performance appears to compare quite favourably with conventional indexing.

## CHAPTER 8

### CONCLUSIONS

The first step in testing a theory (qua theory) is to examine it to see what deductions can be made from it - to set up postulates which may be tested either experimentally or by observations of the 'real-life' situation. That is to say, the first step in testing a theory is to state the practical consequences of it. If the deduced practical consequences (operational definitions) are proved to be unsustainable, the theory is discredited. No theory can ever be proved to be true; it is held for so long as no better theory can be found.

L.T. Wilkins: Social Deviance Page 36

Although the results presented in this volume inevitably represent only a condensation of the tens of thousands of individual results which have been obtained, it is hoped that they are in sufficient detail for anyone interested to make their own interpretation. It might, therefore, be argued that much of this final chapter is redundant, and that it would be better to leave readers to reach their own conclusions. However, the following comments are offered as a personal contribution, with the hope - and expectation - that others will feel free to deduce and argue.

The results have been presented in three main ways. Firstly, there are the details of the search results for the various index languages, recall and precision devices and search rules as obtained with the conventional coordination level cut-off. Secondly, some of these results have been regrouped to illustrate various aspects of the test and thirdly, many of the test results have been re-calculated by the document output cut-off method based on simulated ranking. While the opinions presented in this chapter may be illustrated by referring to a particular set of figures, they are not usually based on a single result.

Within the definition as given in Chapter 2 of Volume 1, every set of figures supports the original hypothesis of an inverse relationship between recall and precision. It is immaterial which variable is changed to give a new system; it may be the coordination level (e.g. Fig. 4.100T), the exhaustivity of indexing (e.g. Fig.4.912P), the recall devices (e.g. Fig. 6.10T), the precision devices (e.g. Fig.6.17T), the search programmes (e.g. Fig.4.850T), or the relevance decisions (e.g. Fig. 6.3P); it has been impossible to find any exception to what can be claimed as a basic rule.

Quite the most astonishing and seemingly inexplicable conclusion that arises from the project is that the single term index languages are superior to any other type. This is mainly evidenced by the results based on the normalised recall ratios of Fig. 5.15T, but also, although less obviously, by the comparison of different systems using the conventional coordination level cut-off (see Fig. 6.2P). This conclusion is so controversial and so unexpected that it is bound to throw considerable doubt on the methods which have been used to obtain these results, and our own first reaction was to doubt the evidence. A complete recheck has failed to reveal any discrepancies, and unless one is prepared to say that the whole test conception is so much at fault that the results are completely distorted, then there is no other course except to attempt to explain the results which seem to offend against every canon on which we were trained as librarians.

<u>ORDER</u>	<u>NORMALISED RECALL</u>		<u>INDEXING LANGUAGE</u>
1	65.82	I-3	Single terms. Word forms
2	65.23	I-2	Single terms. Synonyms
3	65.00	I-1	Single terms. Natural Language
4	64.47	I-6	Single terms. Synonyms, word forms, quasi-synonyms
5	64.41	I-8	Single terms. Hierarchy second stage
6	64.05	I-7	Single terms. Hierarchy first stage
7=	63.05	I-5	Single terms. Synonyms. Quasi-synonyms
7=	63.05	II-11	Simple concepts. Hierarchical and alphabetical selection
9	62.88	II-10	Simple concepts. Alphabetical second stage selection
10=	61.76	III-1	Controlled terms. Basic terms
10=	61.76	III-2	Controlled terms. Narrower terms
12	61.17	I-9	Single terms. Hierarchy third stage
13	60.94	IV-3	Abstracts. Natural language
14	60.82	IV-4	Abstracts. Word forms
15	60.11	III-3	Controlled terms. Broader terms
16	59.76	IV-2	Titles. Word forms
17	59.70	III-4	Controlled terms. <del>Related terms</del> <i>Narrower and broader terms</i>
18	59.58	III-5	Controlled terms. <del>Narrower and broader terms</del> <i>Related terms</i>
19	59.17	III-6	Controlled terms. Narrower, broader and related terms
20	58.94	IV-1	Titles. Natural language
21	57.41	II-15	Simple concepts. Complete combination
22	57.11	II-9	Simple concepts. Alphabetical first stage selection
23	55.88	II-13	Simple concepts. Complete species and superordinate
24	55.76	II-8	Simple concepts. Hierarchical selection
25	55.41	II-12	Simple concepts. Complete species
26	55.05	II-5	Simple concepts. Selected species and superordinate
27	53.88	II-7	Simple concepts. Selected coordinate and collateral
28	53.52	II-3	Simple concepts. Selected species
29	52.47	II-14	Simple concepts. Complete collateral
30	52.05	II-4	Simple concepts. Superordinate
31	51.82	II-6	Simple concepts. Selected coordinate
32	47.41	II-2	Simple concepts. Synonyms
33	44.64	II-1	Simple concepts. Natural language

FIGURE 8.1T ORDER OF EFFECTIVENESS BASED ON NORMALISED  
RECALL FOR 33 CRANFIELD INDEX LANGUAGES  
(AVERAGE OF NUMBERS)

Before considering some of the particularly striking aspects of the ranked order of effectiveness as given in Fig. 5.15T, there are certain points to be noted about this table. The normalised recall ratios range from 65.82% to 44.64% and this range encompasses some 33 different index languages plus 14 languages (or options) of the SMART system. It is impossible to state here what is a significant difference; most people who have been consulted agree that anything less than 1% is probably of doubtful significance, but that a difference of 3% or 4% almost certainly represents a significant change in performance. Rather than try to postulate on this point, we would prefer to rely on the consistency with which certain actions have certain effects.

For convenience of discussion, the normalised recall table, with the SMART results deleted, is reprinted as Fig. 8.1. It can be seen that the Single Term index languages rank 1, 2, 3, 4, 5, 6, 7 and 12 with the normalised recall ratio ranging from 65.82% and 61.17%. Starting from the base of natural language (with a score of 65.00%), the use of synonyms and word forms shows a slight improvement, whereas an enlargement of the classes by quasi-synonyms and hierarchical grouping detracts from the performance.

Of the six Controlled Term index languages, that using only the basic terms gave the best performance, with a ranking of 10 and a normalised recall ratio of 61.76%, this being a slight improvement on the lowest score with a Single Term index language. As narrower, broader and related terms are brought in, ranking orders for the other five Controlled Term index languages are 10, 15, 17, 18 and 19, with the lowest score being 59.17%.

The searches on abstracts and titles gave four languages which ranked 13, 14, 16 and 20, the range being from 60.94% to 58.94%. The abstracts (which included titles) seem to be marginally better than the titles on their own. It is interesting that, with the abstracts, the confounding of word forms results in a slightly lower score, whereas the reverse is true with the titles.

The highest rank of the Simple Concept index languages is 7, with a normalised recall ratio of 63.05%. Another language in this group is ranked 9, but the other thirteen Simple Concept index languages occupy the final ranks from 21 to 33. The two Simple Concept index languages which perform reasonably well are - surprisingly - those where the selection of additional related terms is based not on the classification schedules but on the rotated alphabetical index (see Vol. 1, Appendix 5.5).

In Fig. 8.1 it is significant that Single Term Natural Language I.1.a has a score of 65.00%, while Simple Concept Natural Language II.1.a has the lowest score of 44.64%. There is only one difference between these two index languages. In the former, the single terms are free; in the latter exactly the same single terms are interfixed into concepts. Index Language II.1.a represents the concept taken directly from the terminology of the document, e.g. 'conical afterbody', 'centrifugal compressor'; Index Language I.1.a uses exactly the same words, but they are broken down to the single terms, i.e. 'conical', 'afterbody', 'centrifugal', 'compressor'. It would therefore seem that interfixing is such a powerful device that it can severely depress the performance when calculated by the normalised recall ratio. Even when one considers the performance by coordination level cut-off, it can be seen from Fig. 4.700T and from the composite graph in Fig. 4.715P, that the Simple Concept Natural Language II.1.a has a very low maximum recall ratio, which is not compensated for by a particularly good precision ratio. Because it is so relatively inefficient, one finds that, for the Simple Concept index languages, the broadening of

*but not so  
with average  
of ratios*

classes by the use of various recall devices results in a considerable improvement in performance, which is contrary to the effect observed with the Single Term index languages. This leads to the following conclusions.

There was in this test an optimum level of specificity in the terms which were used. The conceptual terms of the Simple Concept index languages were over-specific when used in natural language, this high level of specificity being related to the strength of interfixing between the single terms of the natural language. Because of this, the broadening of the natural language concepts into more general classes resulted in a significant improvement in performance, in that it helped to overcome the high specificity. On the other hand, the Single Terms in natural language appear to have been near to the correct level of specificity; only to the relatively small extent of grouping true synonyms and word forms could any improvement in performance be obtained. Contrary to the experience of Simple Concepts, the broadening of the classes by the use of quasi-synonyms or hierarchical grouping resulted in a significant loss of performance. In between these two extremes of Single Term and Simple Concepts came the Controlled Terms. Less specific than the Concepts but more specific than the Single Terms, the effect of broadening the classes from the Controlled Terms Basic Terms (Index Language III.1.a) was to depress the performance, although not to the same extent as single terms.

While the evidence is not so easy to interpret from the tables and plots of the main test results as given in Chapter 4, it is quite obvious that within the various groups of index languages - where a direct comparison can be made - there is a difference between systems, and that these substantiate the rankings which are given in Chapter 5.

To restate the main conclusions more precisely

1. In the environment of this test, it was shown that the best performance was obtained by the use of Single Term index languages.
2. With these Single Term index languages, the formation of groups of terms or classes beyond the stage of true synonyms or word forms resulted in a drop of performance.
3. The use of precision devices such as interfixing and partitioning was not as effective as the basic precision device of coordination.

In the light of these unexpected conclusions, it is necessary to consider very carefully the test environment and to see whether there is any factor which could have distorted the results.

The subject field is a matter on which it is difficult to argue. There has in the past been a tendency to assume that, with an imprecise (mushy) subject language, where the same notion can be expressed in several different ways, there is the necessity for broad grouping of terms in the index language. Yet it seems possible that this imprecision is such that it is virtually impossible to make any logical practical grouping or class which can improve overall performance. To form a single class of two vague, imprecise terms may merely add confusion to confusion, so that any resulting improvement in the retrieval of relevant documents is more than outweighed by the increase in the retrieval of non-relevant documents.

In Chapter 6, the results were given for a set of questions dealing with aircraft structures, where, it has been earlier suggested, the subject language is less mushy. The results are not easy to interpret, but it appears probable that the assumption that aerodynamics represents the mushier language was unjustified. In the final chapter of Volume 1, we said that "It would seem, that next to the question of relevance assessments, the determination of the effect of subject language precision is the most important problem to be tackled". This opinion still holds, and we find it impossible to say categorically that the subject area of the test collection did not have an influence on the comparative test results.

Undoubtedly the size of the test collection (on which the normalised recall ratios are based) is smaller than one would have liked. The test results presented in Chapter 4, Section 1, show that the smaller sets of documents and questions were representative of the complete document collection and question set, but these tests were only concerned with the Single Term index languages, and it will be necessary to await confirmation on this point from the tests being carried out using the complete collection with the SMART system. However, there appears to be no justification for suggesting that the size of the test collection could have significantly affected the comparison between systems.

A matter that has already been raised in reviews of Volume 1 (e.g. Ref.14), and will undoubtedly be argued again is the matter of relevance decisions used in this test. It was in fact considered in the earlier volume, and the reader is referred in particular to the table on page 14 of Volume I. However, since that section was written, the matter of relevance has become the object of research and investigation in its own right, and it may be worth reopening and expanding the argument in the hope that some of the complexities introduced by psychological overtones might be clarified.

Consider first the matter of the evaluation of an operational information retrieval system, which we have earlier described as covering all stages from the first receipt of an enquiry to the stage of supplying the requester with the references to the set of documents (or, if the system is so designed, to an actual set of documents) which represent the system's answer to his enquiry. It is particularly stressed that the process starts with the first receipt of an enquiry. This enquiry is expressed in the form of a "stated requirement"; anyone with practical experience of information work will know that quite often the stated requirement is far removed from the real needs of the questioner. The greater the expertise of the information staff concerned, the greater the probability that it will be possible, before commencing a search, to reduce the gap between the real and stated needs of the enquirer.

However, in such a situation, namely the evaluation of an operational system, it is essential that the relevance assessments should be based on the real needs of the questioner; it therefore follows that the questioner must make the relevance judgements. Only if this is done can it be found whether there are any errors (i.e. the retrieval of non-relevant documents, or the non-retrieval of relevant documents) which are due to a failure to bridge the gap between the real and the stated needs. At the same time, however, it is necessary to determine the relevance of documents in relation to the stated needs. With these two sets of relevance judgements, it is possible to pinpoint the reasons for the failures in the complete system.

These two types of relevance are called "user relevance" and "stated relevance". The former can only be decided by the questioner himself, but "stated relevance" can be determined (as has been argued in the table on

page 14 of Volume 1) by anybody with reasonable knowledge of the subject field.

On the other hand, if the evaluation is only intended to cover a sub-system of the complete operational system, such as the index language, then there is not the same necessity of having "user relevance" decisions; in fact, such decisions could introduce an additional variable which might mitigate against the interpretation of the test results, and a set of "stated relevance" decisions could be more satisfactory.

So far the argument has been concerned with the evaluation of operational systems. All the tests of experimental systems have been or are being conducted in artificial, created environments. Under such circumstances, "user relevance" decisions cannot be obtained, and in the few tests so far carried out, "stated relevance" decisions of one kind or another have been used. However, in this particular project, as explained in the first Volume (pages 21 - 23) an endeavour was made to simulate "user relevance" decisions. At the same time (and contrary to what was done in Cranfield I), we deliberately eschewed any effort to interpret the stated needs; in all cases the search terms were based solely on the terminology of the question. Whether the original decision to simulate user relevance decisions was correct has already been considered (Vol. 1, page 114) and tentatively the conclusion was there reached that it might have assisted the interpretation of the test results if, instead, stated relevance decisions had been used. On the whole, this is a view to which we would still subscribe but for one fact. If stated relevance decisions had been used, and assuming the test results had shown the similar superiority of Single Term Natural Language, then it would have been virtually impossible to refute an argument that the results were unduly influenced by the relevance decisions.

In the artificial situation, a person - or a group of persons - is presented with a search question (which may have been devised by someone else) and a set of documents (or their surrogates in the form of titles or abstracts) and told to make a series of decisions as to which documents are relevant. He can be given specific instructions, such as the type of person that he is supposed to be or the purpose for which he is supposed to require the information. Whatever such instructions he may receive, he is ultimately faced with a sequence of words which make up the question, and other sequence of words which make up the documents, and by the intensity with which the words and the meaning of the question appear to match the words and the meaning of a document, he must decide that a given document is or is not relevant to a given question. In this artificial situation it seems reasonable to assume - and such experimental evidence as is available bears out the assumption - that there will be a closer direct match between the actual words of a question and a relevant document, than is the case in the natural situation of a questioner making user relevance decisions. Conversely, and just as important, there will, in the artificial situation, be a lower match between the question and a non-relevant document than will often be the case with user relevance judgements.

Under such circumstances, it is highly probable that system performance will be better with stated relevance decisions, than with user relevance decisions, since a source of possible error in the complete system has been eliminated. This is not an important factor in the present investigation, since the objective is not to obtain maximum performance per se, but is concerned with the comparison between the performance of different index languages. The important point is that stated relevance decisions which can

only be based on a match between words in the document and the question, might be expected to favour systems using precise natural language, while user relevance decisions might logically be expected to favour systems which bring in groups of related terms. The conclusion is therefore reached that the method of obtaining relevance decisions in this test could not have been responsible for the unexpected results, since any influence it might have had would have tended to work in the opposite direction.

Without going against the above argument, Vickery (Ref.15) rightly points out that "There are still verbal links between source document and question; the questions supplied by the author - some time after doing the research - were formulated after the cited papers had been read and possibly influenced the wording of his question." This raises two separate questions; firstly, is it very much different to what happens in a real life situation, and secondly, is the effect serious enough to distort the test results? To consider the first point, experience in the evaluation test of Medlars at the National Library of Medicine has shown that the majority of questioners are already aware of certain relevant documents before asking for a search to be carried out. It therefore seems likely that, in real life, search questions must often be influenced by the terminology of relevant documents, and therefore the procedure which was adopted in this test for obtaining questions is not far removed from what normally happens. If, however, the actions of those who prepared the search questions were significantly different from what happens in real life, then it is necessary to consider whether the results are likely to have been distorted. To determine whether this is so would require a far deeper analysis of the individual searches than has so far been done. Our own opinion is that if such an analysis were made, it would show that in the large majority of cases there had been no serious distortion, and it is difficult to believe that the few cases where it might have occurred would have been sufficient to produce the significant - and consistent - variations in performance.

The concept indexing was done by selecting from each document those concepts which appeared to be of importance. This being an intellectual task it is not possible to argue that it was done correctly. Readers of the reports on Cranfield I will recollect that the errors of the indexers were the cause of a significant number of failures to retrieve relevant documents, but that considered as a percentage of total indexing, it represented a very low "error rate". Usually in that test the errors were errors of omission. The higher level of indexing exhaustivity, and the longer time devoted to indexing each document made it less likely that these would occur in this project, and some analysis of the failures to retrieve relevant documents has not revealed any significant errors in this respect. Certainly it does not seem plausible to suggest that any such errors could have influenced the comparative results.

While the complete indexing was more exhaustive than would normally be the case, the assignment of an indexing weight to each concept permitted the testing of various levels of indexing exhaustivity. The test results are given in Chapter 4, Section 4, and again show that whatever the level of indexing exhaustivity might be, the effect of moving from Index Language I.1.a to Index Language I.6.a is consistent, and there is no evidence to suggest that the exhaustivity of indexing affected the comparison between different index languages.

Concerning this level of exhaustivity of indexing, it again becomes obvious that there was an optimum in regard to this particular document/question set. The lowest level of exhaustivity of indexing investigated was the search on titles only; the highest level of exhaustivity occurred with the

search on abstracts. Intermediary were the three levels of indexing done by the project staff. Figure 8.2T shows the normalised recall ratios obtained in these five cases, all using natural language terms.

Index Language	Average No. of Terms	Normalised Recall Ratio
Titles	7	59.76%
Level 1 Single Term Natural Language	14	62.88%
" 2 Single Term Natural Language	22	63.57%
" 3 Single Term Natural Language	33	65.00%
Abstracts	Approx 60	60.94%

FIGURE 8.2T NORMALISED RECALL RATIOS FOR FIVE LEVELS OF EXHAUSTIVITY

There is the possibility that the selection of terms by the indexer was more descriptive of the document content than those terms used for the titles and the abstracts, but the main variable in these five results concerns the level of indexing exhaustivity. It would seem that while the titles were at too low a level of exhaustivity, the gradual increase in the level, up to an average of 33 terms, brought about an improvement in performance. However, the higher level of exhaustivity represented by the abstracts (probably about 60 terms per document) was too high, resulting in the retrieval of large numbers of additional non-relevant documents, so that the performance only represented a slight improvement on that obtained with titles. This hypothesis is supported by the effect with titles and abstracts of enlarging the classes by the use of word forms. With titles, where it has been shown that the level of exhaustivity is too low, the use of word forms improves the normalised recall ratio from 58.94% to 59.76%. With abstracts, however, no such improvement is noted; already there are too many terms and the use of word forms results in a fall from 60.94% to 60.82%. Admittedly this in itself cannot be considered a significant change, but taken in the context of the other results, appears to be of some importance.

The compilation of the dictionaries or schedules was done, in the main, by Mr. Jack Mills. Although there can be few people more competent in such work, there can obviously be no guarantee but that different classes in the Single Term index languages might have given an improved performance as compared to natural language. However, it seems unlikely that the classes prepared for the Simple Concept index languages could have been solely responsible for the relatively poor performance as compared to the Single Term index languages. With the Controlled Term index languages, the classes of terms were formed on the basis of groupings given in the Thesaurus of Engineering Terms of the Engineers Joint Council, yet the use of any groupings except Narrower Terms (Index Language III.2.a) resulted in a loss of performance.

In Chapter 3, the statement was made that for any given question, the total number of postings of the search terms of that question must be equal to the total number of retrievals at the various coordination levels. To explain this point with a simple example, assume the search programme is

made up of four terms A, B, C and D, each of which have been used five times in the indexing of a set of documents as follows (x represents any other term or terms also used in indexing the documents):

<u>Document Number</u>	<u>Index Terms</u>
1	ADx
2	x
3	ACx
4	x
5	BCDx
6	x
7	Bx
8	x
9	ABCDx
10	x
11	BDx
12	x
13	Ax
14	x
15	BCDx
16	x
17	Cx
18	x
19	Ax
20	x

Searches for any combination of A, B, C and D would result in retrieval at various coordination levels as follows:

<u>Coordination Level</u>	<u>No. of Documents Retrieved</u>
4	1 (Document 9)
3	3 (Document 9, 5, 15)
2	6 (Document 9, 5, 15, 1, 3, 11)
1	10 (Document 9, 5, 15, 1, 3, 11, 7, 13, 17, 19)

Thus the sum of the retrievals (1+3+6+10=20) is the same as the total number of postings for the four terms.

The particular significance of this point is the effect on retrieval performance of enlarging the classes. Assume that the search terms are broadened by being grouped with a related term, A<sub>1</sub>, B<sub>1</sub>, C<sub>1</sub> or D<sub>1</sub>, and that these related terms have also each been used five times in the same set of 20 documents, the indexing being as follows:

<u>Document Number</u>	<u>Index Terms</u>
1	ADx
2	A <sub>1</sub> B <sub>1</sub> D <sub>1</sub> x
3	ACD <sub>1</sub> x
4	B <sub>1</sub> x
5	BCDA <sub>1</sub> x
6	C <sub>1</sub> x
7	BC <sub>1</sub> x
8	D <sub>1</sub> x
9	ABCDx
10	A <sub>1</sub> x
11	BDA <sub>1</sub> x
12	C <sub>1</sub> x
13	AB <sub>1</sub> C <sub>1</sub> x
14	D <sub>1</sub> x
15	BCDx
16	D <sub>1</sub> x
17	CA <sub>1</sub> B <sub>1</sub> x
18	C <sub>1</sub> x
19	Ax
20	B <sub>1</sub> x

Assuming the search now is for any coordination of (A + A<sub>1</sub>), (B + B<sub>1</sub>), (C + C<sub>1</sub>) and (D + D<sub>1</sub>), the retrieval at different coordination levels will be as follows

<u>Coordination Level</u>	<u>No. of Documents Retrieved</u>
4	2 (Document 5, 9)
3	8 (Document 5, 9, 2, 3, 11, 13, 15, 17)
2	10 (Document 5, 9, 2, 3, 11, 13, 15, 17, 1, 7)
1	20 (Document 5, 9, 2, 3, 11, 13, 15, 17, 1, 7, 4, 6, 8, 10, 12, 14, 16, 18, 20)

Again it is shown that the sum of the retrievals (40) equals the total postings for the four groups of terms. Assume now that there were four relevant documents, numbers 3, 7, 9 and 15. The performance in the two cases would then be as follows

<u>Coordination Level</u>	<u>Case A</u>				<u>Case B</u>			
	R	N-R	Recall Ratio	Precision Ratio	R	N-R	Recall	Precision Ratio
4	1	0	25%	100%	1	1	25%	50%
3	2	1	50%	66%	3	5	75%	38%
2	3	3	75%	50%	4	6	100%	40%
1	4	6	100%	40%	4	16	100%	20%

While this particular result has obviously been prepared to illustrate the point, it would seem that this is an example of what has been consistently happening in the test searches with the Single Term and Controlled Term index languages. Whereas the broadening of the term classes has increased the recall of relevant documents at higher coordination levels, the effect of doing this has been more than offset by the increased number of non-relevant documents. Only when the index terms being used are too precise, as in the case of the Simple Concept Natural Language, can the formation of broad classes of terms bring about an improvement.

Finally, it is necessary to consider the measures which have been used in this test, and to ask whether it is possible that some other measures would have brought about a change in the comparative results. Obviously suspect is the normalised recall ratio, based on a simulated rank output. While at first it might seem that such a measure is likely to weigh in favour of systems having high recall ratios, it is in fact mainly influenced by the first two ranked documents. At this stage, the recall ratios, as can be seen from Figures 5.11T - 5.14T, are as follows

Recall Ratio at Document Output Cut-off of 2	Index Language		
23%	I.2, I.3	12,	
22%	I.1	3	
21%	I.6, I.7	4, 5	
20%			
19%	I.5, I.8, II.9, III.2	6-9	7.5
18%	II.3, II.12, IV.3, IV.4	10-13	11.5
17%	II.10, III.1, IV.1, IV.2	14-17	15.5
16%	I.9, II.11, III.3, III.4	18-21	19.5
15%	II.5	22	
14%	II.13	23	
13%	II.2, II.8, III.5	24-26	25
12%	II.1, II.4, II.6, III.6	27-30	28.5
11%	II.7, II.15	31-32	31.5
10%			
9%	II.14	33	

*perhaps I - yes, other weak*

It will be seen that with the exception of Index Language II.3, which (at 18%) rises from 28 to 10-, there is a strong correlation between this ordering and the final ordering as given in Table 8.1. With the document output cut-off method, recall and precision are, as we explained earlier, completely interdependent, and therefore it would appear to be a measure that is quite impartial as between recall and precision. It is known that others are investigating different measures, and most of those that have been proposed have already been considered in Chapter 3. Now that the results of this test are available, it is to be hoped that proponents of new measures will be able to demonstrate any superiority over those used in this report. Until such time, there appears to be no reason to suggest that the measures have affected the comparative results.

With the possible doubtful exception of the subject field, there appears to be nothing in the test environment which could be held responsible for serious distortion of the results as between one system and another. Therefore it is necessary to proceed on the assumption that the results are

correct, and attempt to find the reasons why they should be as they are.

It would be quite incorrect to suggest that no-one has previously argued in favour of single terms, natural language and coordination, for these were the bedrock of the Uniterm System of coordinate indexing as originally propounded by the late Dr. Taube in 1951. But while the device of coordination - or, as we would now term it, post-coordination - continues in favour, there are few who now accept (for Information Retrieval Systems) uncontrolled vocabularies, and some who insist additionally on the use of links and roles. Even Dr. Taube himself was, within a couple of years of the inception of the Uniterm System, to start devising associated maps, and there is no indication, in the writings at that time of the group at Documentation Inc., of any awareness that the resultant increased recall would be more than offset by the lower precision.

There are doubtless indexes in existence which follow the original Uniterm principles, but one of the few persons who has consistently, in print, advocated the use of natural language and coordination is Mr. Th. te Nuyl with his L'Unité System (Ref. 16). Even so, for most people L'Unité System will be associated mainly with the ingenious coding system rather than the use of natural language. It is of interest to note that the clustering of the natural language terms into broad alphabetical groups (as in L'Unité) brings about the confounding of word forms, so, possibly unintentionally, te Nuyl did adopt a coding device which was, it would appear from the results of this test, the only way to improve performance over natural language.

Then there are, of course, permuted title indexes, which use the natural language of the title, but these can hardly be considered in the same light, since they do not have the facilities of post-coordination.

Therefore it is against these few that are ranged, for instance, the activities over the last fifty years of the Universal Decimal Classification, which is probably now more widely used than ever before. At the same time, a large number of national and international organisations are engaged in constructing thesauri, while many groups in the research field are endeavouring to develop computer methods for the formation of classes of terms (e.g. Ref. 17).

The effort that is put into these activities, by whichever process the classes may be formed, is presumably influenced by the widely held belief that it is only by such means that a high recall ratio is obtainable. Yet even in Cranfield I we reported that a recall ratio of 97% was possible merely by using the words in the titles. There was no way of knowing in that experiment the corresponding precision ratio, but it was not only assumed (correctly) that it would be very low, but it was also assumed that it would be lower than would have been the case if such a recall ratio had been obtained with a conventional index language.

As far as this test is concerned, the latter assumption would be unjustified; is it now reasonable to assume that the grouping of natural language terms to form controlled vocabularies, or the broadening of search strategies, must inevitably result in a loss in overall performance?

We would certainly not make such a statement on the basis of this single test; however, it would be surprising if the comparative test results were peculiar to the particular environment of this test, and it does seem

that the results are sufficiently convincing to justify a fresh look at firmly held beliefs.

The present position is that, in the very large majority of cases, the manager of an information retrieval system employs indexers who apply their intelligence to the documents which are to be entered into the system. The indexers select the important concepts which they then translate into the terms of a controlled vocabulary (e.g. a thesaurus or classification schedule). This has possibly involved a considerable amount of intelligence in its compilation, and requires more intelligence for its maintenance. At the stage of a question being received, the search staff will apply their intelligence to deciding the exact meaning of the question and to preparing a suitable search programme, using the terminology of the system. In doing this they will take advantage of the intelligence that has been applied to denoting the relationships between the index terms, either in the arrangement of a classification schedule or by the visual display of a thesaurus. Normally the search is then made, and the questioner receives the output.

It would appear to be a reasonable assumption that the more intelligence that is applied to any of these three stages (i.e. the indexing, the compilation and maintenance of a controlled vocabulary, and the determination of the search strategy) the better the result should be in terms of recall and precision. For example, the most direct measurement (which can be isolated) of the effect of using intelligence in this project is given in the series of results presented in Figures 4.840P - 4.845P, and again in Figure 5.21T. From the latter it can be seen that Search E (where intelligence was used in deciding the acceptable combinations of search terms) resulted in a 1% - 2% increase in normalised recall ratio as compared to Search A (where any combination of terms was accepted).

However, the mere use of intelligence is not enough, for in all cases it is necessary that the intelligence should be applied intelligently in relation to the needs of the system. An example of this relates to the level of exhaustivity in indexing. One cannot say categorically that the selection of seven terms to index a document indicates more or less intelligence than the selection of sixty terms, for it could be argued that, while the latter certainly requires more clerical effort, the former requires more intelligence in selecting the most important terms. However, in the environment of this test, the results show that intelligence was more effectively applied in selecting an average of some thirty-three terms for indexing. (It should be emphasised that this is in no way intended to imply that this level of exhaustivity would be the optimum in a different environment.)

Intelligence is a valuable - and relatively expensive - commodity, and should be used in the most efficient manner. An interpretation which could be placed on the results of this test is that there may well be operational situations in which one should take advantage of the intelligence that has already been applied by the author of a paper by accepting as index terms the key-words in the title or abstract. There is the folk-lore that titles do not represent a correct indication of the content of the document or that authors cannot write reasonable abstracts, and everyone can quote examples where this is the case. Such examples are, however, comparatively rare; for instance, of the many thousands of research papers issued by the National Aeronautics and Space Administration, it would be very difficult to find a single paper where the title did not present an adequate representation of the main subject matter or the summary did not cover all the items of importance. We would therefore argue that, in many operational systems,

*It's really a bridge between the two systems*

a case could be made out for dispensing with indexers within the system and for using the persons thus displaced to screen the search output. The indications are that information staff, merely on the basis of reading titles, can quickly and reliably screen out between 50% and 80% of the non-relevant documents which are retrieved in an average search, without any loss of relevant documents. Such a process would, in many cases, result in a far better service to the user, by giving an operational performance higher than that now obtained.

Additionally it can be argued that there are situations where the intellectual effort involved in the construction and maintenance of controlled vocabularies is unjustified. It is with very strict qualifications that this viewpoint is advanced; in certain subject fields it is almost certainly not true. One critical factor (there are certainly others) could be the occurrence of real synonyms as opposed to quasi-synonyms, or near-synonyms. To illustrate the difference between subject fields, it has been said that there are twenty-one synonyms for the term 'aspirin' (apart from trade names), any of which may be found in the literature, whereas in the subject field used in this project the number of true synonyms (in contrast to quasi-synonyms) was relatively small, the improvement in performance by grouping synonyms was equally small and was hardly sufficient to justify moving from natural language terms. It is difficult to believe that a controlled vocabulary should be less efficient than natural language, even though the evidence of this test points to such a conclusion. Apart from the theoretical reasons already advanced for this being so, there could be a more fundamental reason, and the answer may again lie in the intelligent application of intelligence. No one could deny that a large number of highly intelligent people have given a considerable amount of time to the maintenance of the Universal Decimal Classification or to the preparation of the Thesaurus of Engineering Terms of the Engineers Joint Council. It can, however, legitimately be asked whether these activities represent intelligent applications of intelligence. It may, in fact, be not possible to generate an efficient controlled vocabulary without the applied and close attention, over a relatively long period, of the professional staff of the operating group.

This test has shown that natural language, with the slight modifications of confounding synonyms and word forms, combined with simple coordination, can give a reasonable performance. This means that, based on such practice, a norm could be established for operational performance in any subject field, and it would then be for those who proposed new thesauri, new relational groups, links or roles, to show how the use of their techniques would improve on the norm. The availability of a computer programme, such as a simplified version of the SMART programme of Professor Salton, would make this relatively inexpensive.

Every quotation that has been taken from the book by Professor Wilkins is relevant to our final argument. We make no forecasts that a coordinate system will break down when it reaches a certain size, or any other speculations of this kind, for there is nothing that has been done in this project - or in any other experimental project recently completed or under way - which can justify categorical statements of this nature. As Cranfield I gave indications of the situation over the general field of information retrieval systems, so this project has shown, in a more specialised area, some of the basic problems which beset any and every operator of an information retrieval system. The results can be taken as an indication of what might be done to improve efficiency, but the application of the results to any given situation can only be on the basis of an evaluation of the operational system concerned.

In conclusion, we make no apology for repeating the quotation given at the beginning of this chapter

"The first step in testing a theory (qua theory) is to examine it to see what deductions can be made from it - to set up postulates which may be tested either experimentally or by observations of the 'real-life' situation. This is to say, the first step in testing a theory is to state the practical consequences of it. If the deduced practical consequences (operational definitions) are proved to be unsustainable, the theory is discredited. No theory can ever be proved to be true; it is held for so long as no better theory can be found."

To put it more colloquially "the proof of the pudding is in the eating". It must remain so with the results and conclusions of this project.

REFERENCES

1. STEVENS, N.D. Review of test reports of Cranfield I. Library Resources and Technical Services, 8, 1964, pp. 87-90.
2. AITCHISON, J. and CLEVERDON, C.W. Report of a test on the index of metallurgical literature of Western Reserve University. Cranfield, 1963.
3. VICKERY, B.C. On retrieval system theory. 2nd ed. London, Butterworths, 1965.
4. SWETS, J.A. Information retrieval systems. Science, 141, 1963, pp. 245 - 250.
5. FAIRTHORNE, R.A. Unpublished notes.
6. REES, A.M. The evaluation of retrieval systems. Comparative Systems Laboratory Technical Report No. 5. Western Reserve University, 1965.
7. SINNETT, J.D. An evaluation of links and roles used in information retrieval. Dayton, Air Force Materials Laboratory, 1964.
8. GOFFMAN, W. and NEWILL, V.A. Methodology for test and evaluation of information retrieval systems. Comparative Systems Laboratory Technical Report No. 2. Western Reserve University, 1964.
9. VERHOEFF, J. GOFFMAN, W. and BELZER, J. Inefficiency of the use of Boolean functions for information retrieval systems. Communications of the Association for Computing Machinery, 4, 1961, pp. 557-558, 594.
10. KENDALL, M.G. and STUART, A. The advanced theory of statistics. London, Griffin, 1961.
11. MORONEY, M.J. Facts from figures, 3rd ed. London, Penguin, 1965.
12. CROXTON, F.E. and COWDEN, D.J. Applied general statistics. New York, Prentice-Hall, 1939.
13. SALTON, G. Progress in automatic information retrieval. I.E.E.E. Spectrum, August 1965, pp. 90-103.
14. Review of 'Factors determining the performance of indexing systems', Vol. I. Times Literary Supplement, Aug. 11, 1966.

15. VICKERY, B.C. Review of 'Factors determining the performance of indexing systems', Vol. I.  
Jnl. of Documentation, 22, 1966, pp. 247-249.
16. te NUYL, Th. W. The L'Unité mechanised documentation system.  
Revue de Documentation, 28, 1962, pp. 140-147.
17. NEEDHAM, R.M. and SPARCK JONES, K. Keywords and clumps.  
Jnl. of Documentation, 20, 1964, pp. 5-15.

APPENDIX 3A

TABLES OF GENERALITY NUMBER, AND FALLOUT, RECALL AND  
PRECISION RATIOS

These tables provide a method of conversion between the recall, precision and fallout ratios and the generality number. The equations on which the tables are based, are given on pages 41 - 42.

The method of using the tables is as follows. Assume three separate search results where in the first case the generality number is 2, in the second case it is 5, and in the third case it is 25. At 60% recall ratio, the precision ratio is 15% in the first case, 25% in the second case, and 60% in the third case. It is desired to equate the performance in regard to recall and precision.

The correction is done by converting the precision ratios to what they would be at a generality number of 2. Reference to the tables for a generality number of 5 (page 270) shows that for a recall ratio of 60% and a precision ratio of 25%, the fallout ratio is 0.905%. Reference to the tables for a generality number of 2 (page 267) shows that for a recall ratio of 60% and the fallout ratio of 0.905%, the precision ratio would be approximately 11%.

For the third case, where the generality number was 25, reference to the appropriate table on page 281 gives a fallout ratio of 1.026% for a recall ratio of 60% and a precision ratio of 60%. From page 267 it is found that this fallout ratio of 1.026% at 60% recall occurs with a precision ratio of approximately 10%.

The original and corrected ratios are then as follows:

	Generality Number	Recall Ratio	Precision Ratio	Fallout Ratio	Corrected Precision Ratio
1.	2	60%	15%	0.681%	15%
2.	5	60%	25%	0.905%	11%
3.	25	60%	60%	1.026%	10%

This shows that the first case represents the best performance.













































APPENDIX 4A

This appendix presents one set of search results, namely those obtained with the 1400 document collection and 221 questions using the Single Term Natural Language I.1.a, with Exhaustivity of Indexing 3, and Document Relevance 1-4. It shows, therefore, the individual question results which are summarised in Figure 4.100T.

Against each question number is shown first the number of documents relevant to that particular question, the number of search terms in the question, together with the total number of postings for the search terms. At each coordination level is then shown the number of relevant and non-relevant documents which were retrieved. It will be noted that the sum of the documents retrieved equals the total number of postings. For instance, with question 1, the retrievals at the various coordination levels of relevant and non-relevant documents are (reading from left to right), 22, 381, 14, 136, 9, 20, 3, 2, 1, and the sum of those is 588 which is the same as the figure given for the total number of postings. The significance of this point is discussed in Chapter 8.

Question	No. of relevant documents	No. of search terms	Total no. of postings	COORDINATION LEVELS																			
				1	2	3	4	5	6	7	8	9	10										
1	28	9	588	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R		
2	24	7	631	22	381	14	136	9	20	3	0	2	0	0	1	0							
4	8	6	719	16	468	9	110	3	22	2	2	0	1	0									
8	2	15	1454	2	603	8	73	3	16	4	4												
9	4	7	569	4	912	2	347	2	163	2	17	2	4		1								
10	4	4	1115	4	513	3	44	1	4														
12	5	9	1773	4	988	4	114	1	4														
13	11	5	1261	5	1085	4	397	3	180	3	72	2	21		1								
15	3	5	1535	10	835	9	246	8	90	4	50	1	8										
18	8	8	756	3	1007	3	332	3	174	2	11	2	11										
22	7	9	797	7	509	6	149	4	42	4	20	2	8	2	2	1	0						
23	5	8	1611	6	609	5	139	2	31		5												
26	2	5	108	5	1158	4	346	4	76	2	16												
27	2	4	673	2	94		10		2														
29	2	3	84	2	409	2	203	2	54	1	0												
31	3	8	1686	2	75	2	4	1	0														
32	2	9	1716	3	1165	3	394	2	94	2	20	1	2										
33	3	8	2023	2	1173	2	362	1	135	1	30	1	9										
34	9	7	743	3	1257	3	457	2	167	2	77	1	37	1	16								
35	9	10	1284	8	572	1	140		21		1												
39	4	8	699	9	1000	9	193	8	46	7	11	1	0										
40	1	7	775	4	517	2	153	1	18	0	3	0	1										
41	32	3	233	1	430	1	209	0	120	0	10	0	4										
49	3	7	234	15	216	1	1																
50	9	10	1863	3	203	3	20	2	3														
51	6	11	3176	9	1200	9	548	8	65	6	13	4	1										
52	3	8	1004	6	1263	6	754	6	473	6	293	5	170	5	100	4	60	3	16	3	2	1	0
53	2	5	804	3	805	3	165	2	25	1	0												
54	9	9	2032	2	664	1	116	1	18	0	2	7	40	4	9	4	4	1	1	1	0		
55	7	5	1364	9	1195	9	485	9	172	8	74	7	40	4	9	4	4	1	1	1	0		
56	1	12	2351	7	1104	6	193	4	44	0	6	6											
57	6	6	1080	1	1136	0	780	0	297	0	125	0	11	0	1								
58	3	6	1356	3	778	3	207	3	77	2	7	0	22	0	3	10	2	1	1	0			
59	6	4	718	3	704	3	426	3	123	3	49	3	22	3	10	1	1	0					
61	3	5	453	6	370	6	256	5	64	4	7	1	6	4	7	1	1	0					
	3	5	453	3	399	3	40	1	6	0	0	1	6	4	7	1	1	0					









Question	No. of relevant documents	No. of search terms	Total no. of postings	COORDINATION LEVELS																	
				1	2	3	4	5	6	7	8	9									
273	7	5	176	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R		
274	5	5	1828	7	139	6	15	5	2	2	0	0	0	0	0	0	0	0	0		
275	2	11	867	5	1120	4	513	3	162	2	19	2	2	2	2	2	2	2	2		
277	13	6	173	2	514	2	223	2	76	2	24	2	2	2	2	2	2	2	2		
283	7	4	275	13	146	8	4	2	0	0	0	0	0	0	0	0	0	0	0		
284	9	4	257	6	241	4	24	2	0	0	0	0	0	0	0	0	0	0	0		
285	16	7	1847	9	219	8	20	1	0	0	0	0	0	0	0	0	0	0	0		
288	6	10	1695	15	1116	14	425	10	206	8	42	3	8	3	8	3	8	3	8		
291	11	5	140	6	1011	5	451	4	122	3	71	1	19	1	19	1	19	1	1		
292	9	7	635	11	71	8	32	6	7	3	2	2	2	2	2	2	2	2	2		
293	5	7	494	7	535	3	79	2	8	1	0	0	0	0	0	0	0	0	0		
294	13	7	191	5	360	3	106	2	18	2	18	2	18	2	18	2	18	2	18		
295	4	3	46	12	154	10	15	4	0	0	0	0	0	0	0	0	0	0	0		
296	9	4	655	3	37	1	4	1	0	0	0	0	0	0	0	0	0	0	0		
297	6	13	1300	9	524	9	91	6	12	5	87	4	45	1	14	0	3	0	3		
298	3	13	1287	6	585	6	342	3	196	5	3	4	45	1	14	0	3	0	3		
299	12	7	1226	3	664	3	330	3	160	22	75	12	32	1	12	12	12	12	12		
300	3	3	61	9	1053	5	150	1	8	8	0	0	0	0	0	0	0	0	0		
301	4	6	834	3	50	2	5	0	1	1	3	6	2	0	0	0	0	0	0		
303	8	6	852	4	575	3	177	3	61	3	6	2	0	0	0	0	0	0	0		
304	3	10	680	8	647	6	134	4	40	1	11	0	1	0	0	0	0	0	0		
306	16	10	1138	3	559	3	85	2	22	0	5	0	1	0	0	0	0	0	0		
314	14	10	2217	16	741	13	247	6	80	4	22	3	3	3	3	3	3	3	3		
315	14	7	1498	14	1277	13	530	11	268	6	73	4	20	1	0	0	0	0	0		
316	14	4	1414	14	886	12	484	10	75	2	15	15	15	15	15	15	15	15	15		
317	2	8	1811	14	1922	10	405	5	52	2	4	4	4	4	4	4	4	4	4		
321	2	9	951	2	916	2	514	2	242	2	117	0	13	0	1	1	1	1	1		
323	5	9	1416	2	632	2	189	2	88	2	27	0	6	0	1	1	1	1	1		
327	6	15	1199	5	974	5	327	4	83	0	14	0	4	0	0	0	0	0	0		
331	12	5	2064	6	704	6	262	6	123	5	61	3	17	1	4	4	4	4	4		
332	6	15	1368	12	1139	11	504	9	368	4	17	3	17	1	4	4	4	4	4		
333	11	5	1106	6	790	5	386	5	93	4	46	4	15	4	6	6	6	6	6		
335	14	9	873	11	771	9	201	7	74	4	28	1	0	0	0	0	0	0	0		
336	11	9	864	14	569	14	150	12	61	11	25	6	3	4	1	0	0	0	0		
338	4	7	497	11	559	10	161	10	71	8	26	3	3	1	0	0	0	0	0		
				3	371	3	92	3	20	0	5	5	20	3	1	0	0	0	0		

Question	No. of relevant documents	No. of search terms	Total no. of postings	COORDINATION LEVELS													
				1		2		3		4		5		6		7	
				R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R	R	N-R
339	2	10	1869	2	1154	2	433	2	194	0	60	0	20	0	2		
340	1	9	1887	1	1267	1	430	1	154	0	33	0	33	0			
347	15	9	1382	15	748	15	324	15	149	12	58	8	22	3	10	1	2
348	15	6	1225	15	798	14	305	11	48	5	24	1	4				
349	18	5	1444	18	1001	11	356	2	53	0	3						
352	19	10	2603	19	1177	19	654	19	386	16	217	10	75	1	10		
353	18	8	2400	18	1282	15	716	14	240	11	89	1	12	0	2		
355	9	4	572	9	436	9	95	7	14	1	1						
356	4	5	612	4	486	4	91	4	19	2	2						
360	8	9	1728	8	1068	8	413	8	167	5	40	3	6	0	2		
365	24	8	1373	23	752	13	374	5	148	3	41	1	11	0	2		

APPENDIX 5A

FORMULA FOR DOCUMENT RANKING BASED ON  
PROBABILITY CONSIDERATIONS

by

G.H. STEARMAN

If a particular question at a particular coordination level results in the retrieval of a total of N documents, of which R documents are relevant, then the average time taken to find each of the relevant documents when a large number of searches is made can be determined on the basis of the following assumptions:-

- (a) Each successive document is selected at random.
- (b) The same time is taken to inspect each document for relevancy so that, for example, if a relevant document is found at the 3rd choice, three units of time are taken and if at the 7th choice, seven units and so on.

If one unit of time is assigned to each choice, then the value of the average as defined above can be taken as the rank of the relevant document in a simulated ordering of the N documents.

Let:-

Total number of documents retrieved be N

Total number of relevant documents be R

Order of N be S (S = 1, 2 ... N)

Order of R be K (K = 1, 2 ... R)

Then the problem is to find an expression for  $P_{K,S}$ , the probability that the Kth relevant document will be found at the Sth inspection, where N and R are given. Then if  $Q_K$  is the simulated ranking, its value is given by the weighted sum

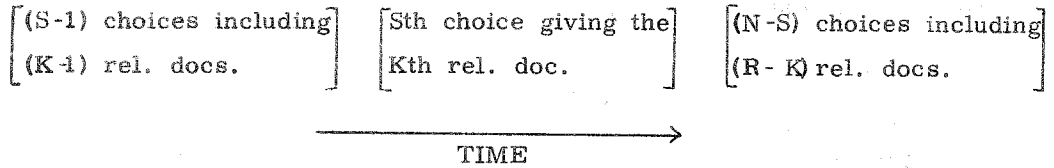
$$Q_K = \sum_{S=1,2 \dots N} S \cdot P_{K,S}$$

$A_C B$  means the number of ways of choosing B items from A and is expressed as

$$\frac{A!}{B! (A-B)!}$$

The probability  $P_{K,S}$  may be determined as the ratio of the number of configurations in which the Kth relevant document appears at the Sth inspection and the total number of ways in which R relevant documents may be arranged in N positions.

The typical layout would then be as shown:-



The number of ways in which this layout can be formed is the numerator of the required ratio and is given by

$${}^{(S-1)}C_{(K-1)} \cdot {}^{(N-S)}C_{(R-K)}$$

The denominator is simply  ${}^N C_R$

$$\text{Thus } P_{K,S} = \frac{{}^{(S-1)}C_{(K-1)} \cdot {}^{(N-S)}C_{(R-K)}}{{}^N C_R}$$

$Q_K$  can now be evaluated as indicated above for each value of K from 1 to R

$$\text{Thus } Q_K = \sum_{S=1}^N S \cdot \frac{{}^{(S-1)}C_{(K-1)} \cdot {}^{(N-S)}C_{(R-K)}}{{}^N C_R}$$

$$= \sum_{S=1}^N \frac{S(S-1)!}{(K-1)!(S-K)!} \cdot \frac{{}^{(N-S)}C_{(R-K)}}{{}^N C_R}$$

$$= \frac{K}{N} C_R \sum_{S=1}^N \frac{S!}{K!(S-K)!} \cdot {}^{(N-S)}C_{(R-K)}$$

$$= \frac{K}{N} C_R \sum_{S=1}^N \left[ {}^S C_K \cdot {}^{(N-S)}C_{(R-K)} \right]$$

The terms of the series vanish for  $K > S > (K+N-R)$  so that the limits of S may be changed to give

$$Q_K = \frac{K}{N} C_R \sum_{S=K}^{K+N-R} \left[ {}^S C_K \cdot {}^{(N-S)}C_{(R-K)} \right] \quad \text{Note that } K \leq R$$

The summation of this series is given by Schwatt ["Operations with series" - Chelsea, 2nd edn. page 47] in the form:-

$$\sum_{k=n}^{p-m} b C_n \cdot p-k C_m = p+1 C_{p-n-m}$$

Putting  $b=S$ ,  $p=N$ ,  $m=R-K$ ,  $n=K$  we obtain

$$\sum_{S=K}^{N-R+K} S C_K \cdot (N-S) C_{(R-K)} = (N+1) C_{(N-R)} = (N+1) C_{(N+1-N+R)} = (N+1) C_{(R+1)}$$

$$\therefore Q_K = \frac{K \cdot (N+1) C_{(R+1)}}{N C_R} = \underline{\underline{\frac{K(N+1)}{(R+1)}}}$$

This simple expression is the basis of the method used in Chapter 5; the above formal analysis shows it to be soundly based upon probability considerations.

SUBJECT INDEX

78

- Abstracts. 186, 254, 258
- Adjusted precision ratio. 73
- Aerodynamic questions. 228, 255
- Average of numbers. 53, 205
- Average of ratios. 53, 205, 225
- Averaging sets of results. 51
- Basic questions. 225
- Bibliographic coupling. 243
- Citation indexing. 243
- Collection size. 4, 256
- Comparative results. 221
- 'Composite' measures. 43, 49
- 'Composite' table. 16
- Concept-indexing. 7, 258
- Confounding of word forms. 9
- Contingency table. 33, 49
- Controlled term index languages.  
11, 159, 254, 262
- Coordination level. 56, 80, 202,  
230
- Dictionary compilation. 259
- Distillation ratio. 39
- Document output cut-off. 69, 200,  
262
- Document relevance. 15, 80, 140,  
215, 223
- Effectiveness measure. 44
- Exhaustivity of indexing. 6, 7, 78,  
121, 186, 215, 258-9
- Environmental factors. 4, 21, 255
- Evaluation of operational systems.  
256
- Fallout ratio. 36, 39, 81
- Generality number. 71, 80
- Hardware factors. 4, 6
- Index language groups. 78
- Information retrieval system  
defined. 4
- Interfixing. 12, 109
- Linear composite measures. 44
- 'Measure of merit'. 48
- Noise factor. 36
- Non-linear composite measures.  
48
- Normalised recall ratio. 202,  
205, 215, 254, 259
- Number of postings. 234, 259
- Operational factors. 4, 6
- Order of effectiveness. 209, 254
- Order of retrieval of relevant  
documents. 239
- Partitioning. 12, 109
- Performance measures. 32
- Precision devices. 11, 78, 109,  
159, 234, 255
- Precision ratio. 35, 81
- Presentation of results. 31
- Q factor. 48
- Quasi-synonyms. 9, 11
- Question generality. 234
- Question sets. 21, 80, 83, 225,  
228
- Rank order number. 195
- Recall devices. 11, 100, 148, 159
- Recall ratio. 35, 81
- Relevance decisions. 256
- Relevance of the documents. 15,  
80, 140, 215, 223
- Search rules. 12, 78, 109, 131,  
159, 205
- Simple concept index languages. 11,  
148, 254
- Simulated ranking method. 195, 205
- Single term index language. 9, 83,  
221, 254, 262
- SMART searches. 202
- Snobbery ratio. 35
- Software factors. 4, 6
- Specificity. 6, 255
- Starting term coordination level. 61,  
71
- Stated relevance. 256
- Structure questions. 228, 255
- Subject field. 4, 228, 255
- Supplementary questions. 225
- Term clusters. 9
- Term usage. 7
- Titles. 186, 254, 258
- Totalling results. 56
- Twin variable measures. 36
- User relevance. 256
- Weighting. 12, 78, 121
- Weighting related to relevance  
grading. 215



# Factors determining the performance of indexing systems. Volume 2, Test results

Cleverdon, Cyril W.

1966

---

Cleverdon CW, Keen M. (1966) Factors determining the performance of indexing systems; Volume 2, Test results. Aslib Cranfield research project

<http://hdl.handle.net/1826/863>

*Downloaded from CERES Research Repository, Cranfield University*