

ED 025 288

52

LI 001 241

By- Schultz, Claire K.; And Others

Evaluation of Indexing by Group Consensus. Final Report.

Institute for the Advancement of Medical Communication, Philadelphia, Pa.

Spons Agency- Office of Education (DHEW), Washington, D.C. Bureau of Research.

Bureau No- BR-7-0622

Pub Date 30 Aug 68

Contract- OEC-1-7-070622-3890

Note- 46p.

EDRS Price MF-\$0.25 HC-\$2.40

Descriptors- Analysis of Variance, Comparative Analysis, *Evaluation Criteria, *Evaluation Techniques, *Groups, *Indexing, *Information Retrieval, Performance Factors, Reliability, Standards

Identifiers- *Criterion Group Method

This study was designed to explore the practicality, flexibility, reliability, and sensitivity of the Criterion Group Method of measuring the effectiveness and efficiency of indexing, a method using a criterion group to set the standard for "ideal" indexing. These major variables were examined: (1) size of document sample, (2) size of Criterion Group, (3) instructions to indexers and use of a vocabulary guide, (4) three methods of editing raw indicia to make terms comparable, and (5) two methods of weighting indexers' scores. Scores earned by a set of eight professional indexers, by individual authors of the test documents and, in some cases, scores for title sets or medical students' indexing were compared within selected treatments to measure the extent to which the detectability of differences was achieved by each treatment. A two-way analysis of variance was used to relate reliability of test scores to document sample size and criterion group size. From the results of these studies of the methodologic variables, it was concluded that the criterion group method of evaluating indexing can be a practical yardstick for a wide variety of managerial, research, and educational uses. Appendixes include the rationale of the method, a literature review, information on materials employed and subjects participating in study trials, and details on manual and computer implementation. (Author/JB)

BR 7-0622

PA-52

LI 001241



FINAL REPORT

PROJECT No. 7-0622

CONTRACT No. OEC 1-7-070622-3890

EVALUATION OF INDEXING BY GROUP CONSENSUS

AUGUST 30, 1968

ED025288

U. S. DEPARTMENT OF
HEALTH, EDUCATION, AND WELFARE

OFFICE OF EDUCATION
BUREAU OF RESEARCH

LI 001241

FINAL REPORT

PROJECT No. 7-0622

CONTRACT No. OEC 1-7-070622-3890

EVALUATION OF INDEXING BY GROUP CONSENSUS

CLAIRE K. SCHULTZ, PETER B. HENDERSON, AND RICHARD H. ORR

INSTITUTE FOR ADVANCEMENT OF MEDICAL COMMUNICATION

PHILADELPHIA, PA. 19104

AUGUST 30, 1968

THE RESEARCH REPORTED HEREIN WAS PERFORMED PURSUANT TO A CONTRACT OEC-1-7-070622-3890 WITH THE OFFICE OF EDUCATION, U. S. DEPARTMENT OF HEALTH, EDUCATION, AND WELFARE. CONTRACTORS UNDERTAKING SUCH PROJECTS UNDER GOVERNMENT SPONSORSHIP ARE ENCOURAGED TO EXPRESS FREELY THEIR PROFESSIONAL JUDGMENT IN THE CONDUCT OF THE PROJECT. POINTS OF VIEW OR OPINIONS STATED DO NOT, THEREFORE, NECESSARILY REPRESENT OFFICIAL OFFICE OF EDUCATION POSITION OR POLICY.

U.S. DEPARTMENT OF
HEALTH, EDUCATION, AND WELFARE

OFFICE OF EDUCATION
BUREAU OF RESEARCH

U.S. DEPARTMENT OF HEALTH, EDUCATION & WELFARE
OFFICE OF EDUCATION

THIS DOCUMENT HAS BEEN REPRODUCED EXACTLY AS RECEIVED FROM THE PERSON OR ORGANIZATION ORIGINATING IT. POINTS OF VIEW OR OPINIONS STATED DO NOT NECESSARILY REPRESENT OFFICIAL OFFICE OF EDUCATION POSITION OR POLICY.

ACKNOWLEDGEMENTS

MRS. RITA PORRECA, WORKING UNDER THE DIRECTION OF DR. MYRNA GOPNIK, WROTE THE COMPUTER PROGRAM FOR A CONTEXT-DEPENDENT METHOD OF THESAURUS PROCESSING. DR. ANDREW BAGGLEY, DEPT. OF EDUCATION, UNIVERSITY OF PENNSYLVANIA, SERVED AS STATISTICAL CONSULTANT TO THE PROJECT.

WE ARE MOST GRATEFUL TO THE SCIENTISTS AND PROFESSIONAL INDEXERS WHO GAVE THEIR TIME, AND TO THE STUDENTS WHO WORKED AT A NOMINAL FEE, TO GENERATE THE DATA THAT MADE THIS STUDY POSSIBLE.

TABLE OF CONTENTS

	PAGE
INTRODUCTION	1
AIM OF STUDY	1
METHODOLOGIC DESIDERATA	1
ORGANIZATION OF REPORT	2
ESSENTIALS OF METHOD	3
SELECTING THE DOCUMENT SAMPLE	3
SELECTING AND CONSTRUCTING THE CRITERION GROUP	3
INSTRUCTING TEST INDEXERS	4
ESTABLISHING CRITERION AND TEST SETS	4
WEIGHTING THE CRITERION SETS	5
SCORING THE TEST SETS	5
FINDINGS ON METHODOLOGIC VARIABLES	7
VARIABLE 1 - SIZE OF DOCUMENT SAMPLE	7
VARIABLE 2 - SIZE OF CRITERION GROUP	9
VARIABLE 3 - INSTRUCTIONS TO TEST INDEXERS	11
VARIABLE 4 - PROCEDURES FOR EDITING CRITERION AND TEST SETS	13
VARIABLE 5 - WEIGHTING SCHEME FOR SCORING	16
VARIABLE 6 - CONFOUNDING BEFORE SCORING	17
CONCLUSIONS	17
APPENDIX A	
RATIONALE OF METHOD AND LITERATURE REVIEW	A 1-12
APPENDIX B	
MATERIALS EMPLOYED IN STUDY TRIALS	B 1-3
APPENDIX C	
SUBJECTS PARTICIPATING IN STUDY TRIALS	C 1-2
APPENDIX D	
PROCEDURES FOR MANUAL IMPLEMENTATION	D 1-2
APPENDIX E	
COMPUTER IMPLEMENTATION	E 1-4

SUMMARY

THE CRITERION GROUP METHOD TESTS THE EFFECTIVENESS AND EFFICIENCY OF TEST INDEXING SETS, USING A CRITERION GROUP TO SET THE STANDARD FOR "IDEAL" INDEXING. THE CRITERION GROUP FOR A PARTICULAR APPLICATION IS CHOSEN BY THE TEST ADMINISTRATOR, CONSISTENT WITH HIS OWN CONCEPT OF WHO REPRESENTS THIS "IDEAL". MATCHING TEST SETS OF INDEXING TERMS WITH THE CRITERION SET YIELDS AS MANY DEGREES OF MATCH AS THERE ARE MEMBERS OF THE CRITERION GROUP (REFERRED TO AS CONSENSUS NUMBER).

THIS STUDY WAS DESIGNED TO EXPLORE THE PRACTICALITY, FLEXIBILITY, RELIABILITY, AND SENSITIVITY OF THE METHOD. TO DO THIS, IT EXAMINES THE MAJOR VARIABLES: (1) SIZE OF THE DOCUMENT SAMPLE, (2) SIZE OF THE CRITERION GROUP, (3) EFFECT OF VARIOUS INSTRUCTIONS TO INDEXERS AND USE OF A VOCABULARY GUIDE, (4) EFFECTS OF THREE METHODS OF EDITING RAW INDICIA TO MAKE TERMS COMPARABLE, AND (5) TWO ALTERNATIVE METHODS OF WEIGHTING INDEXERS' SCORES.

SCORES EARNED BY A SET OF EIGHT PROFESSIONAL INDEXERS, BY INDIVIDUAL AUTHORS OF THE TEST DOCUMENTS, AND IN SOME CASES SCORES FOR TITLE SETS OR MEDICAL STUDENTS' INDEXING, WERE COMPARED WITHIN SELECTED TREATMENTS TO MEASURE THE EXTENT TO WHICH THE DETECTABILITY OF DIFFERENCES WAS ACHIEVED BY EACH TREATMENT. A TWO-WAY ANALYSIS OF VARIANCE WAS USED TO RELATE RELIABILITY OF TEST SCORES TO DOCUMENT SAMPLE SIZE AND CRITERION GROUP SIZE.

RESULTS WITH REGARD TO PRACTICALITY SHOW THAT "INDICATIVE" TESTS (ALLOWING CONFIDENCE LIMITS OF ± 10 POINTS) AT THE 80% LEVEL OF CONFIDENCE CAN BE MADE WITH DOCUMENT SAMPLES AS SMALL AS 10 AND CRITERION GROUPS AS SMALL AS 4; 95% CONFIDENCE REQUIRES COMPARABLE VALUES OF 20 DOCUMENTS AND 9 CRITERION GROUP MEMBERS. IT IS POSSIBLE TO CONDUCT TESTS WITH ONLY A FEW "MECHANICAL" INSTRUCTIONS TO INDEXERS, NO VOCABULARY GUIDE, NO EDITING AND NO WEIGHTING DURING SCORING OTHER THAN USE OF THE CONSENSUS NUMBER EVEN THOUGH FROM THE STANDPOINT OF SENSITIVITY, THE METHOD CAN DETECT DIFFERENCES IN SCORES DUE TO EDITING METHOD, INSTRUCTIONS TO INDEXERS, USE OF A VOCABULARY GUIDE, OR WEIGHTING METHOD, SHOULD SUCH DETECTION BE DESIRABLE. THE METHOD IS FLEXIBLE IN THAT IT HAS BEEN SHOWN TO LEAVE AS OPTIONS VARIABLES SUCH AS METHOD OF INSTRUCTING INDEXERS, METHOD OF EDITING AND METHOD OF WEIGHTING. RELIABILITY IS PRIMARILY DEPENDENT ON THE SIZE OF DOCUMENT SAMPLE AND CRITERION GROUP, AS IS DEMONSTRATED GRAPHICALLY IN THE PAPER. FACE VALIDITY, OR INTUITIVE FEEL FOR THE MEANING OF TEST RESULTS, IS ENHANCED BY THE FACT THAT DIFFERENCES IN SCORES CAN BE EQUATED WITH DIFFERENCES IN THE INTERNATIONALLY KNOWN MEASURE OF "RECALL", AND SCORE DIVIDED BY THE NUMBER OF TERMS IN THE INDEXING SET YIELDS A RESULT SOMEWHAT ANALOGOUS TO "PRECISION". PERCENT MAXIMAL SCORE CAN ALSO BE EASIER TO ENVISION AS REFLECTING EFFECTIVENESS OF TEXT INDEXING SETS ON A 0-100 SCALE, WITH DIFFERENCES OF FROM 6-8 POINTS REPRESENTING SIGNIFICANCE, WHEN SIGNED-RANK TESTS ARE APPLIED.

INTRODUCTION

AIM OF STUDY

SEVERAL YEARS AGO, FOR A SPECIAL RESEARCH APPLICATION, WE DEVELOPED A METHOD FOR MEASURING THE "QUALITY" OR "EFFECTIVENESS", OF INDEXING. * AT THAT TIME WE FELT THIS METHOD COULD BE ADAPTED FOR A WIDE RANGE OF MANAGERIAL, EDUCATIONAL, AND RESEARCH APPLICATIONS AND, SINCE IT HAD SEVERAL IMPORTANT ADVANTAGES OVER OTHER METHODS, # IT MIGHT FILL THE CRITICAL NEED FOR A PRACTICAL YARDSTICK TO EVALUATE INDEXING AND SUBJECT CATALOGING. HOWEVER, THERE WERE A NUMBER OF QUESTIONS TO BE ANSWERED BEFORE ONE COULD BE CERTAIN THAT THE METHOD MET THE DEMANDING REQUIREMENTS FOR SUCH A YARDSTICK. THE PRESENT STUDY WAS UNDERTAKEN TO EXPLORE THESE QUESTIONS.

METHODOLOGIC DESIDERATA

FOR A TRULY GENERAL METHOD, APPLICABLE TO MANY TYPES OF INDEXING AND SUBJECT CATALOGING AND SUITABLE FOR SERVING A WIDE RANGE OF PURPOSES, CERTAIN METHODOLOGIC CHARACTERISTICS WOULD SEEM TO BE EITHER ESSENTIAL OR HIGHLY DESIRABLE. FIRST, THE METHOD SHOULD HAVE "FACE" VALIDITY IN THE EYES OF THOSE WHO WILL USE THE RESULTING MEASUREMENTS; AND SINCE INDIVIDUALS HAVE VARYING CONCEPTS OF WHAT CONSTITUTES "IDEAL" INDEXING, THE METHOD SHOULD ALLOW ONE THE OPTION OF CHOOSING A CRITERION CONCEPT THAT REFLECTS HIS OWN VALUES RATHER THAN BEGGING THE QUESTION OF WHAT THE "RIGHT" CONCEPT IS BY BUILDING IT INTO THE METHOD. SECOND, THE METHOD SHOULD BE PRACTICAL, IN TERMS OF TIME AND EFFORT REQUIRED; FOR ROUTINE OR EVERYDAY USE BY SMALL AND LARGE SERVICES AS WELL AS FOR ONE-TIME STUDIES AIMED AT OBTAINING "DEFINITIVE" MEASUREMENTS. THIRD, IF THE MEASUREMENTS OBTAINED ARE TO SERVE AS A BASIS FOR DECISIONS, ONE SHOULD KNOW HOW MUCH CONFIDENCE THEY MERIT--THAT IS, THEIR RELIABILITY, OR REPRODUCIBILITY, SHOULD BE STATISTICALLY DETERMINANT--AND THIS RELIABILITY SHOULD BE ADEQUATE TO WARRANT BASING IMPORTANT DECISIONS ON THE MEASUREMENTS. FOURTH, THE METHOD SHOULD BE FLEXIBLE IN THAT IT CAN ACCOMMODATE DIFFERENT TYPES OF INDEXING--FOR EXAMPLE, "KEYWORD" INDEXING WITH NO RESTRICTIONS ON ALLOWABLE TERMS, SUBJECT HEADINGS CONTROLLED BY AN AUTHORITY.

* THE DEVELOPMENT OF THIS METHOD WAS DESCRIBED IN: SCHULTZ, CLAIRE K., SCHULTZ, WALLACE L., AND ORR, RICHARD H., "COMPARATIVE INDEXING: TERMS SUPPLIED BY BIOMEDICAL AUTHORS AND DOCUMENT TITLES." AMERICAN DOCUMENTATION 16, 4, (OCTOBER 1965), PP. 299-312.

THE RATIONALE UNDERLYING THE DEVELOPMENT OF THE METHOD IS GIVEN IN APPENDIX A.

LIST WITH OR WITHOUT HIERARCHICAL STRUCTURE, INDEXING DONE BY PEOPLE OR BY MACHINE, ETC. FIFTH, IT SHOULD BE SENSITIVE ENOUGH TO DETECT DIFFERENCES IN THE RELATIVE MERIT OF INDEXING PRODUCED BY THE ALTERNATIVE PROCEDURES OR AGENTS THOSE USING THE METHOD MAY WISH TO ASSESS. COLLECTIVELY, THESE FIVE GENERAL DESIDERATA--FACE VALIDITY, PRACTICABILITY, RELIABILITY, FLEXIBILITY, AND SENSITIVITY--REPRESENT A STRINGENT SET OF REQUIREMENTS A TRULY GENERAL METHOD SHOULD MEET. IN ANY PARTICULAR APPLICATION, OF COURSE, THERE MUST ALWAYS BE TRADE-OFFS BETWEEN VALIDITY AND PRACTICALITY, AND BETWEEN RELIABILITY AND PRACTICALITY; HOWEVER, IT SHOULD BE POSSIBLE TO ACHIEVE COMPROMISES THAT ARE ACCEPTABLE. THIS STUDY AIMED AT EXPLORING THE METHODOLOGIC VARIABLES THAT GOVERN THE TRADE-OFFS REQUIRED AND INFLUENCE THE METHOD'S FLEXIBILITY AND SENSITIVITY.

ORGANIZATION OF REPORT

IN THE SUCCEEDING SECTIONS OF THIS REPORT, WE WILL DESCRIBE THE BASIC OPERATIONS REQUIRED TO APPLY THE GENERAL METHOD; GIVE THE RESULTS OF TRIALS AND ANALYSES DESIGNED TO EXPLORE CRITICAL METHODOLOGIC VARIABLES, DISCUSS THE IMPLICATIONS OF THESE FINDINGS AS THEY RELATE TO THE DESIDERATA SET FORTH ABOVE, AND OFFER SOME CONCLUSIONS REGARDING THE METHOD'S POTENTIAL RANGE OF APPLICATIONS. FOR CLARITY OF PRESENTATION, ALL SUBSIDIARY DETAIL WILL BE RELEGATED TO THE APPENDICES.

ESSENTIALS OF METHOD

IN THE SIMPLEST TERMS, THE CRITERION-GROUP METHOD CAN BE DESCRIBED AS FOLLOWS: FOR EACH DOCUMENT IN THE TEST CORPUS A SET OF TERMS CHARACTERIZING THAT DOCUMENT IS FIRST ESTABLISHED BY MERGING ALL TERMS CHOSEN BY THE MEMBERS OF A CRITERION GROUP, EACH OF WHOM MAKES HIS CHOICES INDEPENDENTLY. THIS INDEXING SET IS THEN CONSIDERED THE STANDARD (CRITERION SET) AGAINST WHICH OTHER SETS OF INDICIA (TEST SETS) FOR THE SAME DOCUMENT ARE TESTED. IN THIS METHOD THE TERMS IN THE SETS TO BE TESTED ARE NOT SCORED ON A BLACK-OR-WHITE SCALE--THAT IS, THEY ARE NOT SIMPLY RATED AS "MATCHING" OR "NOT MATCHING" THE TERMS IN THE CRITERION SET; OUR SCALE ALLOWS FOR AS MANY SHADES OF GRAY AS THERE ARE MEMBERS OF THE CRITERION GROUP. CONDUCTING A TEST REQUIRES SIX BASIC OPERATIONS.

SELECTING THE DOCUMENT SAMPLE

IN ANY SPECIFIC APPLICATION OF THE METHOD, THE DOCUMENTS FOR WHICH INDEXING IS TO BE EVALUATED SHOULD BE A REPRESENTATIVE SAMPLE OF THE DOCUMENT UNIVERSE OF INTEREST. THIS SAMPLE MAY BE SELECTED FROM THIS UNIVERSE BY ANY OF THE USUAL SAMPLING PROCEDURES BASED ON RANDOM SELECTION. WHEN THE SAMPLE TO BE USED IS LARGE, A SIMPLE RANDOM SAMPLING PROCEDURE CAN BE USED; HOWEVER, FOR SMALL SAMPLES, A STRATIFIED RANDOM SAMPLE MAY BE PREFERABLE. FOR THIS OPERATION THE MOST IMPORTANT VARIABLE IS THE SIZE OF THE SAMPLE, WHICH SHOULD BE LARGE ENOUGH TO PROVIDE THE RELIABILITY NEEDED FOR THE PARTICULAR PURPOSE. ON THE OTHER HAND, SINCE THE NUMBER OF DOCUMENTS IS A MAJOR FACTOR IN DETERMINING THE EFFORT AND EXPENSE OF RUNNING A TEST, THIS NUMBER SHOULD BE NO LARGER THAN NECESSARY.

SELECTING AND INSTRUCTING THE CRITERION GROUP

WHAT TYPE OF INDIVIDUALS SHOULD CONSTITUTE THE CRITERION GROUP DEPENDS UPON ONE'S CONCEPT OF "IDEAL" OR "STANDARD" INDEXING AND THE PURPOSE TO BE SERVED. IN OUR ORIGINAL STUDY* THE AIM WAS TO TEST HOW WELL AUTHOR-SUPPLIED INDICIA MATCHED THE LANGUAGE OF POTENTIAL USERS; THEREFORE, A GROUP OF THE AUTHOR'S PEERS SERVED AS THE CRITERION GROUP. HOWEVER, IT MIGHT BE CONSIDERED APPROPRIATE

* SCHULTZ, CLAIRE K., WALLACE L. SCHULTZ, AND RICHARD H. ORR. COMPARATIVE INDEXING: TERMS SUPPLIED BY BIOMEDICAL AUTHORS AND DOCUMENT TITLES. AMERICAN DOCUMENTATION 16, 4, (OCTOBER, 1965). PP. 299-312.

FOR THE CRITERION GROUP TO CONSIST OF "EXPERT" INDEXERS SELECTED ON SOME BASIS FOR THE QUALITY OF THEIR WORK. IDEALLY, FROM WHATEVER UNIVERSE THE CRITERION GROUP IS DRAWN, THE SELECTION PROCEDURE SHOULD INSURE THAT THE GROUP IS REPRESENTATIVE OF THAT UNIVERSE; BUT PRACTICAL CONSTRAINTS MAY REQUIRE ONE TO SETTLE FOR SELECTING MEMBERS OF THE GROUP BY NON-RANDOM PROCEDURES. OTHER THINGS BEING EQUAL, THE LARGER THE GROUP THE MORE LIKELY IT WILL BE REPRESENTATIVE; AND A UNIVERSE THAT IS RELATIVELY HOMOGENEOUS CAN BE ADEQUATELY REPRESENTED BY A SMALLER CRITERION GROUP THAN A UNIVERSE THAT IS HETEROGENEOUS. THE SIZE OF THE CRITERION GROUP, LIKE THE SIZE OF THE DOCUMENT SAMPLE, AFFECTS THE COST OF USING THE METHOD; THEREFORE, THIS VARIABLE IS ALSO AN IMPORTANT DETERMINANT OF PRACTICALITY.

ANOTHER VARIABLE IN THIS OPERATION IS HOW THE GROUP IS INSTRUCTED TO CARRY OUT ITS TASK, INCLUDING WHETHER THEY ARE GIVEN ANY SORT OF A TERMINOLOGY "GUIDE" EXPLICITLY OR IMPLICITLY INTENDED TO STRUCTURE THEIR RESPONSES.

INSTRUCTING TEST INDEXERS

IN ANY APPLICATION WHERE AN INDIVIDUAL, A GROUP OF INDIVIDUALS, OR A MACHINE INDEXES DOCUMENTS FOR THE SPECIFIC PURPOSE OF TESTING THE RESULTING INDICIA, INSTRUCTIONS OR RULES ON HOW TO CARRY OUT THE TASK WILL HAVE TO BE GIVEN. THESE INSTRUCTIONS MAY OR MAY NOT BE EQUIVALENT TO THOSE GIVEN THE CRITERION GROUP. IN APPLICATIONS WHERE THE INDEXING TO BE TESTED HAS BEEN PRODUCED AS PART OF AN ONGOING SERVICE, THIS VARIABLE DOES NOT REPRESENT A TEST "OPTION". AGAIN, IF ONE DESIRES TO GENERALIZE FROM THE FINDINGS REGARDING THE QUALITY OF THE TESTED INDEXING TO SOME LARGER UNIVERSE, THE QUESTION OF REPRESENTATIVENESS ARISES; THEN, THE METHOD OF SELECTION AND SIZE OF THE GROUP REQUIRE CAREFUL CONSIDERATION.

ESTABLISHING CRITERION AND TEST SETS

IF EITHER THE CRITERION GROUP OR THE TEST INDEXERS ARE ALLOWED TO USE FREE LANGUAGE*, A DECISION IS REQUIRED ON WHETHER THEIR OUTPUT SHOULD BE EDITED, OR STANDARDIZED, BEFORE CRITERION AND TEST SETS ARE COMPARED; AND IF STANDARDIZING IS DONE, WHAT RULES SHOULD BE FOLLOWED. WITHOUT STANDARDIZATION, SYNONYMS AND TRIVIAL VARIATIONS--FOR EXAMPLE, SINGULAR AND PLURAL FORMS OF THE SAME TERM--WILL BE

* OR IF MACHINE INDEXING IS TO BE TESTED

COUNTED AS DIFFERENT TERMS. HOWEVER, ANY EDITING INCREASES THE COST OF A TEST; AND ALL HUMAN EDITING IS PRONE TO INCONSISTENCIES AND BIASES THAT MAY AFFECT THE RELIABILITY AND THE VALIDITY OF TEST RESULTS.

WEIGHTING THE CRITERION SETS

IN THIS METHOD, SOME SCHEME IS REQUIRED FOR WEIGHTING THE TERMS USED FOR INDEXING A DOCUMENT TO REFLECT THE CONSENSUS THAT EXISTS AMONG THE CRITERION GROUP WITH RESPECT TO APPROPRIATE INDEXING TERMS FOR THAT DOCUMENT. MANY SCHEMES COULD BE EMPLOYED, BUT PERHAPS THE SIMPLEST IS TO WEIGHT EACH TERM IN THE CRITERION SET BY THE NUMBER OF CRITERION GROUP MEMBERS WHO USED IT TO CHARACTERIZE THE DOCUMENT AND TO GIVE ANY TERM NOT USED BY AT LEAST ONE MEMBER OF THE CRITERION GROUP (THAT IS, ANY TERM NOT IN THE CRITERION SET) A WEIGHT OF ZERO TO INDICATE ITS "UNDESIRABILITY". ALTERNATIVE SCHEMES CAN BE DEVISED THAT WILL INCREASE OR DECREASE THE EFFECT OF CONSENSUS AND WILL CHANGE THE "PENALTY" FOR USING TERMS THAT ARE NOT IN THE CRITERION SET. (SEE APPENDIX D FOR DETAILS ON WEIGHTING AND AN EXAMPLE OF AN ALTERNATIVE SCHEME.)

SCORING THE TEST SETS

THE WEIGHTS THUS ESTABLISHED ARE EMPLOYED TO SCORE EACH TEST SET BY ADDING THE WEIGHTS FOR EACH TERM IN THE SET. THE "RAW SCORE" FOR A TEST SET IS THEN STANDARDIZED BY EXPRESSING IT AS A PERCENTAGE OF THE HIGHEST SCORE POSSIBLE FOR THAT SET, OR THE "VARIABLE SCORE", WHICH IS DETERMINED BY THE SUM OF THE WEIGHTS FOR ALL TERMS IN THE CRITERION SET. THUS IF A TEST SET SCORES 0%, IT MEANS THAT NO TERM IN THE SET WAS USED BY ANY MEMBER OF THE CRITERION GROUP; AND A SCORE OF 100% MEANS THAT THE TEST SET CONTAINS ALL THE TERMS USED BY THE CRITERION GROUP COLLECTIVELY.

WHEN THE CRITERION GROUP CONSISTS OF POTENTIAL USERS, THE PERCENT MAXIMAL SCORE IS ANALOGOUS TO CLEVERDON'S "RECALL" MEASURE; AND IF DESIRED, A SUPPLEMENTARY FIGURE OF MERIT ANALOGOUS TO HIS "PRECISION" MEASURE MAY ALSO BE CALCULATED BY TAKING INTO CONSIDERATION THE FREQUENCY WITH WHICH TERMS NOT IN THE CRITERION SETS (NON-SCORING OR "ZERO TERMS") APPEAR IN THE TEST SETS. * (SEE APPENDIX D FOR DETAILS ON SCORING.)

* THE RELATION OF MEASURES DERIVED BY THIS METHOD TO OTHER MEASURES OF INDEXING PERFORMANCE ARE SUGGESTED IN APPENDIX A. A FULL DISCUSSION OF THESE RELATIONS IS OUTSIDE THE SCOPE OF THIS REPORT.

IN THIS OPERATION, ONE MAY WISH TO GIVE SOME CREDIT FOR TEST SET TERMS THAT, ALTHOUGH NOT IDENTICAL TO TERMS IN THE CRITERION SET, ARE SUBSUMED BY CRITERION SET TERMS IN A GIVEN INDEXING VOCABULARY, THE METHOD ALLOWS THE OPTION OF DEALING WITH SUCH MISMATCHES BY "CONFOUNDING" OR "GENERIC POSTING" BEFORE SCORING THE TERM SETS. * THIS COMPLICATES SCORING AND HENCE INCREASES THE COST OF A TEST; BUT IT MAY BE APPROPRIATE IN SOME APPLICATIONS.

* ALTERNATIVELY, GENERIC-SPECIFIC TRANSFORMATIONS AS WELL AS STANDARDIZATION OF SYNONYMS MAY BE DONE IN THE EDITING OPERATION.

FINDINGS ON METHODOLOGIC VARIABLES

THIS STUDY FOCUSED ON SIX OF THE METHODOLOGIC VARIABLES SELECTED FROM THOSE IDENTIFIED ABOVE. THESE SIX VARIABLES WERE SELECTED BECAUSE, FOR A PRIORI REASONS, WE FELT THEY COULD BE MAJOR DETERMINANTS OF THE METHOD'S PRACTICALITY, FLEXIBILITY, RELIABILITY, AND SENSITIVITY, AND BECAUSE THEY COULD BE INVESTIGATED WITHOUT ESTABLISHING A NEW DOCUMENT CORPUS. TO EXPLORE THE EFFECTS OF THESE VARIABLES, WE CARRIED OUT SPECIAL ANALYSES OF THE DATA OBTAINED IN THE ORIGINAL APPLICATION OF THE METHOD AND ALSO CONDUCTED TRIALS TO OBTAIN NEW DATA BEARING ON THESE VARIABLES. THE MAJOR FINDINGS ARE SUMMARIZED AND DISCUSSED BELOW. DETAILS ON THE MATERIALS, SUBJECTS, AND MANUAL AND COMPUTER PROCEDURES REFERRED TO ARE GIVEN IN THE APPENDICES.

VARIABLE 1. -- SIZE OF DOCUMENT SAMPLE

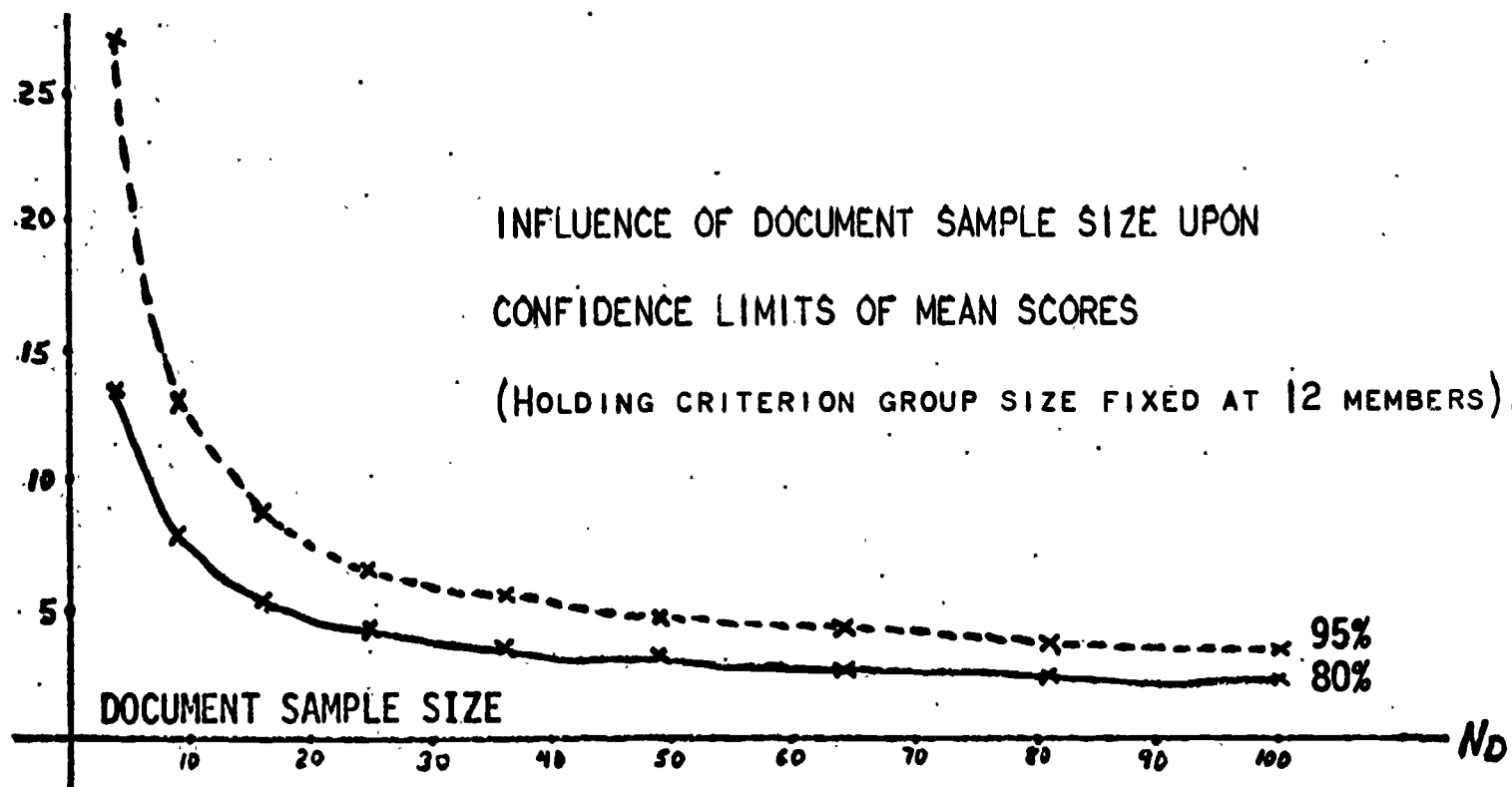
AN INDICATION OF HOW THE RELIABILITY OF TEST SCORES DEPENDS ON THE SIZE OF THE DOCUMENT SAMPLE USED FOR A TEST IS PROVIDED BY THE STANDARD DEVIATION OF THE % MAXIMAL SCORES FOR INDIVIDUAL DOCUMENTS FROM THE MEAN SCORE FOR ALL DOCUMENTS IN THE TEST SAMPLE. WE FOUND THAT THE SAMPLE STANDARD DEVIATION IS MODERATELY AFFECTED BY OTHER METHODOLOGIC VARIABLES. ASSESSING THE EFFECTS OF EACH VARIABLE ON RELIABILITY SINGLY AND IN COMBINATION WITH OTHER VARIABLES WAS NOT FEASIBLE; HOWEVER, THE EFFECTS OF THE TECHNIQUE USED FOR EDITING TERM SETS (VARIABLE 4) AND OF THE SCHEME EMPLOYED FOR WEIGHTING BEFORE SCORING (VARIABLE 5) WERE EXPLORED AND WILL BE DISCUSSED LATER IN CONNECTION WITH THESE VARIABLES. FOR THE STUDIES REPORTED IN THIS SECTION AND IN THE SECTION DEVOTED TO CRITERION GROUP SIZE (VARIABLE 2), THE EDITING TECHNIQUE AND WEIGHTING SCHEME REMAINED CONSTANT.*

WHEN SETS OF TERMS PRODUCED BY 8 PROFESSIONAL INDEXERS FOR EACH DOCUMENT WERE SCORED AGAINST THE SET OF TERMS SUPPLIED BY THE CRITERION GROUP OF 12 COLLECTIVELY, THE STANDARD DEVIATION OF SCORES FOR TERM SETS AVERAGED OVER THE 8 INDEXERS FROM THE GRAND MEAN FOR A SAMPLE OF 128 DOCUMENTS WAS 17 POINTS (% OF MAXIMAL SCORE). IN A SAMPLE OF 32 DOCUMENTS, THE CORRESPONDING STANDARD DEVIATIONS FOR SCORES OF TERM SETS PRODUCED BY INDIVIDUAL INDEXERS RANGED FROM 16 TO 20. FOR TERM SETS SUPPLIED BY AUTHORS, THE STANDARD DEVIATION OF INDIVIDUAL TERM SET SCORES FROM THE MEAN FOR 256 DOCUMENTS WAS 17 POINTS. FIGURE 1 GIVES THE CONFIDENCE LIMITS FOR MEAN SCORES BASED ON DIFFERENT SAMPLE SIZES WHEN THE OBSERVED SAMPLE STANDARD DEVIATION OF 17 POINTS IS TAKEN AS AN ESTIMATE OF THE STANDARD DEVIATION FOR THE DOCUMENT POPULATION FROM WHICH THE SAMPLES WERE DRAWN.

* COMPUTER EDITING AND WEIGHTING SCHEME #1 WERE EMPLOYED THROUGHOUT.

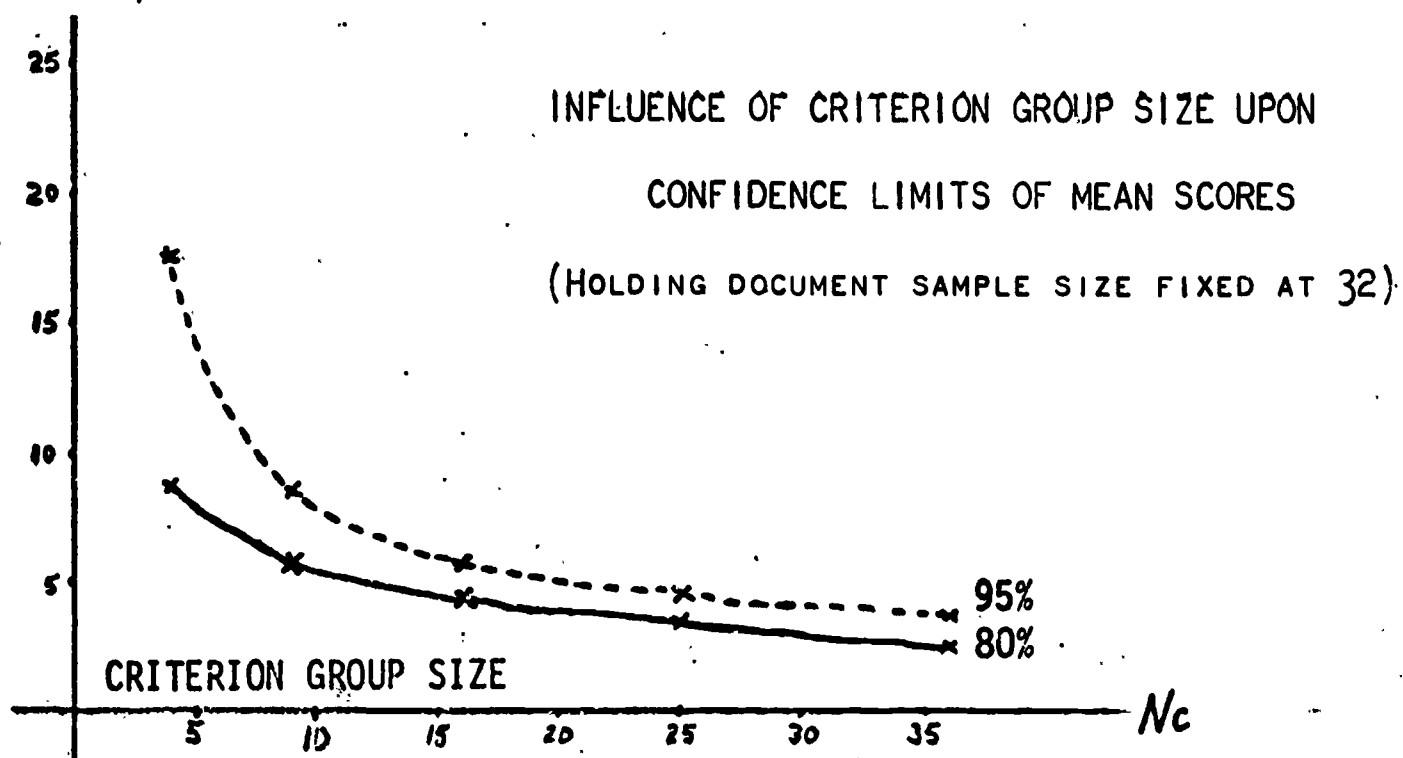
CONFIDENCE
LIMITS
(% MAXIMAL SCORE)

Figure 1



CONFIDENCE
LIMITS
(% MAXIMAL
SCORE)

Figure 2



THE RANDOM VARIATION IN SCORES ATTRIBUTABLE TO DOCUMENTS WILL, OF COURSE, DEPEND UPON THE HETEROGENEITY OF THE DOCUMENT POPULATION FROM WHICH THE SAMPLE WAS DRAWN; AND SINCE THE PRESENT STUDIES WERE LIMITED TO OUR DOCUMENT CORPUS, ONE CANNOT SAY THE OBSERVED SAMPLE VARIATION IS A GOOD ESTIMATE OF THE VARIATION THAT WILL BE ENCOUNTERED IN APPLICATIONS OF THE METHOD WITH OTHER DOCUMENT POPULATIONS. HOWEVER, THESE FINDINGS SHOULD PROVIDE AT LEAST A ROUGH IDEA OF THE GENERAL SIZE OF DOCUMENT SAMPLE REQUIRED IN APPLICATIONS WHERE IT IS IMPORTANT TO MEASURE THE EFFECTIVENESS OF A GIVEN INDEXING "TREATMENT" WITHIN SPECIFIED CONFIDENCE LIMITS. IT CAN BE SEEN THAT WHERE THERE IS A NEED FOR RELATIVELY PRECISE MEASUREMENTS, E.G., WITHIN ± 5 POINTS AT THE 95% CONFIDENCE LEVEL, SAMPLES OF 50 TO 100 DOCUMENTS WILL PROBABLY SUFFICE UNLESS VARIATION IN THE DOCUMENT POPULATION IS CONSIDERABLY GREATER THAN IN OUR CORPUS. FOR MANY APPLICATIONS, THIS DEGREE OF PRECISION WILL NOT BE NECESSARY AND USEFUL RESULTS CAN BE OBTAINED WITH CONSIDERABLY SMALLER SAMPLES, --FOR EXAMPLE, WHERE A ROUGH ESTIMATE (± 10 POINTS WITH 80% CONFIDENCE) CAN BE USEFUL, SAMPLES OF 10 DOCUMENTS MAY SUFFICE.

MANAGERS OF INDEXING SERVICES AND RESEARCHERS ATTEMPTING TO DEVELOP INDEXING "SYSTEMS" OFTEN NEED TESTS TO INDICATE WHETHER TWO INDEXING TREATMENTS GIVE SIGNIFICANTLY AND MATERIALLY DIFFERENT RESULTS. FOR SUCH USES, TESTS WITH SMALL SAMPLES SHOULD PROVIDE AN ADEQUATE BASIS FOR WORKING DECISIONS ON MATTERS WHERE THE COST OF BEING WRONG IS NOT GREAT. THE USE OF SMALL SAMPLE TESTS THAT TAKE ADVANTAGE OF THE REDUCED VARIABILITY ACHIEVED BY EMPLOYING THE SAME SAMPLE TO TEST TWO DIFFERENT TREATMENTS WILL BE ILLUSTRATED LATER.

VARIABLE 2 --SIZE OF CRITERION GROUP

ONE COULD ASSESS THE EFFECT OF THIS VARIABLE DIRECTLY BY SEEING HOW THE SCORES OF A GIVEN INDEXING TREATMENT FOR A GIVEN DOCUMENT SAMPLE CHANGE AS THE NUMBER OF INDIVIDUALS IN THE CRITERION GROUP INCREASES. HOWEVER, WHEN SCORING IS DONE MANUALLY AND THE DOCUMENT SAMPLE IS OF ANY SIZE, THE WORK REQUIRED FOR EACH INCREMENT IN THE SIZE OF THE CRITERION GROUP IMPOSES SEVERE LIMITATIONS ON THIS APPROACH. FOR THIS REASON, IN OUR ORIGINAL PROJECT, WE WERE ONLY ABLE TO ASSESS THIS VARIABLE CRUDELY BY GROUPING SCORES BASED ON HALF OF OUR CRITERION GROUP OF 12 SCIENTISTS WITH SCORES BASED ON THE OTHER HALF. WITH THE DEVELOPMENT OF A COMPUTER PROGRAM FOR SCORING, SYSTEMATIC ASSESSMENT OF THE EFFECT OF CRITERION GROUP SIZE BECAME FEASIBLE; HOWEVER, THE COST OF A DEFINITIVE STUDY WAS STILL MATERIAL SO WE CONSIDERED ALTERNATIVE APPROACHES THAT WOULD BE MORE ECONOMICAL AND ALSO BE USEFUL FOR UNFINISHED STUDY OF THE CHARACTERISTICS OF INDIVIDUAL SCIENTISTS THAT MAY INFLUENCE HOW EFFECTIVE INDEXING IS FOR THEM. ALTHOUGH IN THIS METHOD DEFINITIVE SCORING OF A TEST SET OF INDEXING TERMS IS BASED ON COMPARISONS WITH A "COMPOSITE" CRITERION SET ESTABLISHED BY MERGING THE TERMS USED BY EACH MEMBER OF THE CRITERION GROUP TO DESCRIBE A

GIVEN DOCUMENT, WE HAVE DEMONSTRATED EMPIRICALLY THAT THE SCORE BASED ON A COMPOSITE CRITERION SET CAN BE USEFULLY APPROXIMATED UNDER CERTAIN CONDITIONS BY AVERAGING SCORES FOR A TEST SET BASED ON INDIVIDUAL CRITERION SETS, CONSISTING OF THE TERMS USED BY EACH CRITERION GROUP MEMBER INDIVIDUALLY.* THIS SUGGESTED ANOTHER APPROACH TO ASSESSING THE EFFECT OF CRITERION GROUP SIZE UTILIZING ANALYSIS OF VARIANCE TECHNIQUES. DETAILS OF THESE ANALYSES WOULD BE INAPPROPRIATE HERE, BUT THE MAJOR FINDINGS RELATING TO THE EFFECT OF CRITERION GROUP SIZE WILL BE SUMMARIZED VERY BRIEFLY.

THESE ANALYSES INDICATE THAT AN APPROPRIATE MODEL FOR PRESENT PURPOSES IS ONE IN WHICH THE TOTAL VARIANCE IN SCORES IS PARTITIONED INTO 3 ADDITIVE COMPONENTS ATTRIBUTABLE TO DOCUMENT VARIANCE, CRITERION GROUP VARIANCE, AND RESIDUAL ERROR. WHEN CRITERION GROUP VARIANCE IS HELD CONSTANT, THIS MODEL GIVES THE SAME ESTIMATE FOR DOCUMENT VARIANCE AS THAT OBTAINED BY "EXPERIMENTAL" OR DIRECT, DETERMINATION OF DOCUMENT SAMPLE STANDARD DEVIATION REPORTED EARLIER. WHEN DOCUMENT VARIANCE IS HELD CONSTANT, THE MODEL GIVES AN ESTIMATE FOR CRITERION GROUP VARIANCE CENTERED AROUND 121 POINTS (STANDARD DEVIATION, 11 POINTS). THE EFFECT OF SAMPLING ERROR ATTRIBUTABLE TO THIS SOURCE ON TEST SCORE RELIABILITY IS SHOWN IN FIGURE 2 WHERE THE CONFIDENCE LIMITS ARE CALCULATED FROM THIS ESTIMATED VARIANCE.

ON A A PRIORI BASIS, ONE WOULD EXPECT CRITERION GROUP VARIANCE TO DEPEND UPON THE HETEROGENEITY OF THE POPULATION THE GROUP REPRESENTS. IT HAS NOT BEEN FEASIBLE TO TEST THIS HYPOTHESIS SYSTEMATICALLY; HOWEVER, WE HAVE SCORED AUTHOR-INDEXER TEST SETS AGAINST ANOTHER CRITERION GROUP--THE 8 PROFESSIONAL INDEXERS. RATHER SURPRISINGLY, THE SAME ESTIMATE OF CRITERION GROUP VARIANCE WAS OBTAINED. THESE INDEXERS ALSO CONSTITUTE A RELATIVELY HETEROGENEOUS GROUP IN THAT THEIR APPROACHES TO INDEXING REFLECT A VARIETY OF DIFFERENT INDEXING SERVICES.

FIGURE 2 INDICATES THAT, WHEN PRECISE ESTIMATES OF INDEXING EFFECTIVENESS ARE CRITICAL, THE CRITERION GROUP WILL PROBABLY HAVE TO BE SIZABLE IF ONE IS TO HAVE MUCH CONFIDENCE THAT THEY ARE ADEQUATELY REPRESENTATIVE OF SOME LARGER POPULATION. FOR OUR ORIGINAL APPLICATION OF THIS METHOD, IT WAS IMPORTANT TO INCLUDE

*THE CONDITIONS UNDER WHICH SCORING BASED ON INDIVIDUAL CRITERION SETS APPROXIMATES SCORING BASED ON COMPOSITE CRITERION SETS ARE COMPLEX AND HAVE NOT BEEN COMPLETELY EXPLORED; HOWEVER, NUMEROUS TRIALS HAVE SHOWN THAT, WHEN WEIGHTING SCHEME #1 IS EMPLOYED THE APPROXIMATION IS GOOD AT LEAST FOR TERM SETS SUPPLIED BY OUR ORIGINAL CRITERION GROUP OF SCIENTISTS.

ENOUGH PEOPLE IN THE CRITERION GROUP THAT WE COULD BE REASONABLY CERTAIN ANOTHER SAMPLE FROM THE USER POPULATION THEY REPRESENTED WOULD NOT GIVE MATERIALLY DIFFERENT SCORES FOR THE INDEXING TREATMENTS WE WANTED TO ASSESS. WITHOUT A GUIDE AS TO HOW MANY WOULD BE "ENOUGH", WE THEREFORE MADE THE SAMPLE AS LARGE AS WE COULD WITHIN PRACTICAL CONSTRAINTS. IN ANY APPLICATION WHERE MEMBERS OF THE CRITERION GROUP ARE SUPPOSED TO REPRESENT SOME LARGE POPULATION, HOW LARGE THE GROUP SHOULD BE IS A CRITICAL CONSIDERATION SINCE THIS VARIABLE IS A MAJOR DETERMINANT OF THE OVERALL COST OF EMPLOYING THE METHOD. FOR OTHER APPLICATIONS, HOWEVER, THE REPRESENTATIVENESS OF THE CRITERION GROUP IS IRRELEVANT--FOR EXAMPLE, WHERE ONE CAN IDENTIFY A FEW "EXPERT" INDEXERS AND CONSIDER THEIR "OUTPUT" AS A VALID STANDARD. IF A CRITERION GROUP WERE SELECTED FROM THE "BEST" INDEXERS WORKING FOR A SINGLE SERVICE, IT SEEMS REASONABLE TO PREDICT THAT THEIR VARIANCE WILL BE MATERIALLY SMALLER THAN THAT FOUND IN THE TWO GROUPS WE STUDIED AND THAT A GROUP OF 3 OR 4 WILL PROBABLY BE OPTIMAL. EVEN WHERE THE CRITERION GROUP IS SUPPOSED TO REPRESENT SOME LARGER POPULATIONS, THERE ARE NUMEROUS POTENTIAL APPLICATIONS WHERE HIGH PRECISION IS NOT ESSENTIAL AND A CRITERION GROUP OF LESS THAN 10 MEMBERS WILL PROBABLY SUFFICE--WHERE ONLY ROUGH ESTIMATES ARE REQUIRED OR THE NEED IS FOR A QUICK TEST TO GUIDE THE KIND OF WORKING DECISIONS DISCUSSED IN CONNECTION WITH DOCUMENT SAMPLE SIZE.

VARIABLE 3--INSTRUCTIONS TO TEST INDEXERS

WHETHER THE METHOD COULD ACCOMMODATE INDEXING DONE WITHOUT ANY VOCABULARY GUIDE, SUCH AS THE AUTHOR-INDEXING FORM EMPLOYED IN THE ORIGINAL APPLICATION, WAS AN IMPORTANT QUESTION CONCERNING THE METHOD'S FLEXIBILITY; AND WHETHER THE METHOD COULD DETECT DIFFERENCES IN INDEXING PRODUCED BY ASKING INDEXERS TO FOLLOW DIFFERENT RULES HAD A BEARING ON ITS SENSITIVITY. BOTH OF THESE QUESTIONS WERE EXPLORED IN NUMEROUS SMALL-SCALE EXPERIMENTS, IN WHICH DIFFERENT TYPES OF SUBJECTS--INDIVIDUALS WITH AND WITHOUT INDEXING EXPERIENCE, AND WITH AND WITHOUT BIOMEDICAL KNOWLEDGE--WERE ASKED TO INDEX DOCUMENT SAMPLES UNDER TRIAL CONDITIONS. THE KIND OF EVIDENCE THESE EXPERIMENTS PROVIDED RELATING TO THE TWO QUESTIONS CAN BE ILLUSTRATED BY THE RESULTS OF ONE SERIES OF EXPERIMENTS, WHICH IS SUMMARIZED IN TABLE I. WITH NO GUIDE AND NO EXPLICIT RULES, THE MEAN SCORES FOR GROUP A AND GROUP B ON THE 10 DOCUMENTS IN SUBSAMPLE X (32% vs. 25%) WERE, AS ONE WOULD EXPECT, NOT SIGNIFICANTLY DIFFERENT.* THERE WERE

*THE SIGNED-RANK (WILCOXON) TEST WAS EMPLOYED TO TEST THE SIGNIFICANCE OF THE OBSERVED DIFFERENCE. HEREAFTER, ALL STATEMENTS CONCERNING THE SIGNIFICANCE OF DIFFERENCES ARE BASED ON THE SIGNED-RANK TEST IF THE SAME SUBSAMPLE OF DOCUMENTS WAS EMPLOYED FOR BOTH INDEXING "TREATMENTS," AND THE RANK TEST (VARIOUSLY CALLED WILCOXON T TEST OR THE MANN-WHITNEY U TEST) WAS USED WHEN THE DOCUMENT SUBSAMPLES DIFFERED.

Table I. Trials of Indexing Rules and Aids with Medical Student Subjects

Test Indexers	Document Subsample X			Document Subsample Y			Document Subsample Z		
	No Rules			Rule 1			Rule 2		
	Without Guide	With Guide	% Max. Score	Without Guide	With Guide	% Max. Score	Without Guide	With Guide	% Max. Score
Group A of Non-Professional Indexers (10 medical students)	32%	46%		--	36%		--	49%	
	25%	--		28%	--		41%	--	
Group B of Non-Professional Indexers (9 medical students)									

All test sets were weighted by Scheme #1. Each of the 3 document subsamples contained 10 documents.

Rule 1 "Please index this sample of documents the way you think their authors would index them."

Rule 2 "After studying the enclosed record of how authors actually indexed the last document sample you were given, please try again to apply Rule 1."

NO PROBLEMS IN STANDARDIZING AND SCORING TEST SETS UNDER SUCH UNSTRUCTURED CONDITIONS, AND THE MEAN SCORE FOR GROUP B REMAINED STABLE WHEN RULE 1 WAS IMPOSED FOR THEIR SECOND SUBSAMPLE (Y)-- THIS RULE MAY BE CONSIDERED A CONTROL IN THAT IT WAS NOT EXPECTED TO MAKE A DIFFERENCE. HOWEVER, WHEN RULE 2 WAS IMPOSED FOR THEIR THIRD SUBSAMPLE, A SIGNIFICANT DIFFERENCE (99% CONFIDENCE) IN MEAN SCORES RESULTED. THE IMPLICATIONS OF TRIALS WITH GROUP A ARE LESS CLEAR CUT SINCE, FOR COMPARISONS OF INTEREST IN THE PRESENT CONTEXT, THE VARIABLES ARE CONFOUNDED. THESE TRIALS, IN CONJUNCTION WITH EVIDENCE PROVIDED BY SIMILAR EXPERIMENTS WITH OTHER TYPES OF SUBJECTS, INDICATE THAT THE METHOD IS INDEED FLEXIBLE ENOUGH TO ACCOMMODATE INDEXING DONE WITHOUT A VOCABULARY GUIDE AND THAT IS SENSITIVE ENOUGH TO DETECT THE EFFECT OF INDEXING RULES AND INDEXING AIDS. IN ADDITION, THESE ILLUSTRATE HOW SMALL DOCUMENT SAMPLES MAY SUFFICE FOR SOME APPLICATIONS.

VARIABLE 4--PROCEDURES FOR EDITING CRITERION AND TEST SETS.

THREE DIFFERENT PROCEDURES FOR EDITING WERE ASSESSED FOR THEIR GENERAL EFFECTS ON TEXT SENSITIVITY. THE PROCEDURES WERE AS FOLLOWS:

NO EDITING COMPLETELY UNEDITED TEST SETS WERE COMPARED WITH AND SCORED AGAINST THE UNEDITED TERM SETS OF THE CRITERION GROUP ON A WORD-BY-WORD BASIS. THIS MEANT THAT WHERE "NONSUBSTANTIVE" WORDS, SUCH AS "IN" AND "OF", WHICH WERE PRESENT IN AN UNEDITED CRITERION SET, MATCHED WORDS IN A TEST SET SCORING CREDIT WAS GIVEN. ON THE OTHER HAND, NO SCORING CREDIT WAS GIVEN IF A TEST SET WORD FAILED TO MATCH A "SUBSTANTIVE" WORD IN THE CRITERION SET BECAUSE OF A SLIGHT ORTHOGRAPHIC DIFFERENCE, E.G., A WORD ENDING.

COMPUTER EDITING THE COMPUTER FILE OF THESAURUS RULES EDITED BOTH THE CRITERION SETS AND TEST SETS TO ELIMINATE MOST "NONSUBSTANTIVE" WORDS AND TO STANDARDIZE WORD ENDINGS.

MANUAL EDITING A HUMAN EDITOR ATTEMPTED TO APPLY TO CRITERION AND TEST SETS THE THESAURUS RULES INCORPORATED IN THE COMPUTER EDITING PROGRAM; HOWEVER, THIS WAS DONE LARGELY BY MEMORY AND THE EDITOR UNDOUBTEDLY CONSIDERED A WIDER RANGE OF CONTEXTS THAN WAS AVAILABLE TO THE COMPUTER. FOR EXAMPLE, TERMS THAT DID NOT MATCH BECAUSE OF MISSPELLING WERE CREDITED BY THE HUMAN EDITOR.

INITIAL EXPLORATORY TRIALS WITH SAMPLES OF 8 DOCUMENTS SUGGESTED THAT AS COMPARED TO NO EDITING, BOTH COMPUTER AND MANUAL EDITING INCREASED THE METHOD'S ABILITY TO PICK UP DIFFERENCES IN SCORES GIVEN BY DIFFERENT INDEXING TREATMENTS OF THE SAME DOCUMENTS--FOR EXAMPLE, PROFESSIONAL INDEXERS VS. AUTHOR INDEXERS. HOWEVER, LATER TRIALS CONDUCTED WITH SAMPLES OF 32 DOCUMENTS INDICATED THAT, IN THIS REGARD, ANY ADVANTAGE OF THESE PROCEDURES OVER NO EDITING WAS RELATIVELY SMALL. SOME OF THE CRITICAL COMPARISONS IN THE LATER TRIALS ARE SHOWN IN TABLE 11.* THE PRINCIPAL EFFECT OF EDITING IS TO INCREASE SCORES FOR ALL INDEXING TREATMENTS AND THIS INCREASE IS SOMEWHAT MORE MARKED WITH MANUAL EDITING THAN COMPUTER EDITING; HOWEVER, THE DIFFERENCES BETWEEN MEAN SCORES FOR TWO DIFFERENT INDEXING TREATMENTS IS NOT UNIFORMLY INCREASED. IN ADDITION, THE STANDARD DEVIATIONS, WHICH ALSO AFFECT SENSITIVITY, ARE GENERALLY INCREASED BY BOTH COMPUTER AND HUMAN EDITING PROCEDURES. IT IS OF SOME INTEREST TO NOTE THAT THE COMPUTER PROGRAM QUITE SUCCESSFULLY SIMULATED A HUMAN EDITOR; THE MEAN SCORE OF ALL PROFESSIONAL INDEXER TEST SETS OVER 32 DOCUMENTS WAS 34 (STANDARD ERROR, 1.1) WHEN EDITED BY COMPUTER, AS COMPARED TO 36 (STANDARD ERROR, 1.5) WHEN THE SAME TEST SETS WERE MANUALLY EDITED.

ALTHOUGH THE FACT THAT WITHOUT EDITING NON-SUBSTANTIVE WORDS PRESENT IN THE CRITERION SET ARE COUNTED IN SCORING MAY OFFEND ONE'S INTUITIVE SENSE OF TEST VALIDITY, THE FINDINGS SEEM TO INDICATE THAT EDITING MAKES A RELATIVELY SMALL CONTR!BUTION TO TEST SENSITIVITY. BOTH HUMAN AND COMPUTER EDITING IS RELATIVELY COSTLY; THE FORMER SHOULD BE DONE BY EXPERIENCED INDEXERS, AND THE LATTER IS DEFINITELY UNECONOMIC UNLESS LARGE VOLUMES OF TEST SETS ARE TO BE PROCESSED OR A SUITABLE THESAURUS PROGRAM HAS ALREADY BEEN WRITTEN. IF EDITING IS OMITTED, THE REMAINING OPERATIONS CAN BE CARRIED OUT BY CLERICAL PERSONNEL. THIS IS A PRACTICAL CONSIDERATION THAT MAY BE IMPORTANT FOR SOME APPLICATIONS. HAVING THE OPTIONS OF NO EDITING, COMPUTER EDITING, OR HUMAN EDITING INCREASES THE METHOD'S FLEX!BILITY AND RANGE OF POTENTIAL APPLICATIONS.

* EACH OF THE CONTRASTS SHOWN WERE ALREADY KNOWN TO BE SIGNIFICANT FROM LARGE SAMPLE TESTS WITH 128 TO 282 DOCUMENTS BUT THE DIFFERENCES WERE OF AN ORDER THAT MIGHT POSE A "CHALLENGE" FOR SMALL SAMPLE TESTS.

TABLE II COMPARISONS TO ASSESS EFFECT OF EDITING PROCEDURES
ON TEST SENSITIVITY

	NONE		COMPUTER		MANUAL	
	MEAN	S.D.	MEAN	S.D.	MEAN	S.D.
<u>PROFESSIONAL INDEXERS</u> <u>VS. AUTHOR SETS</u>						
P.I.	26	8	34	10	36	9
AUTHOR	34	11	37	12	54	16
DIFFERENCE	S*		S		S*	
<u>AUTHOR VS. TITLE SETS</u>						
AUTHOR	(NOT DONE)		(NOT DONE)		54	16
TITLE					34	17
DIFFERENCE					S*	

*S = SIGNIFICANT AT 80% LEVEL OR HIGHER; WHERE THERE IS AN ASTERISK THE DIFFERENCE WAS ALSO SIGNIFICANT AT THE 95% LEVEL OR HIGHER. ALL SCORES ARE % MAXIMAL SCORES (WEIGHTING SCHEME 1) ON THE SAME 32 DOCUMENTS, AND THE STANDARD DEVIATIONS IN PERCENTAGE POINTS ARE GIVEN IN PARENTHESES BELOW EACH SCORE. STANDARD DEVIATIONS WERE CALCULATED AS DESCRIBED EARLIER IN THE SECTION DEVOTED TO THE EFFECT OF DOCUMENT SAMPLE SIZE.

VARIABLE 5--WEIGHTING SCHEME FOR SCORING

TRIALS OF TWO DIFFERENT WEIGHTING SCHEMES WERE CONDUCTED PRIMARILY TO DETERMINE WHETHER TEST SENSITIVITY WAS AFFECTED BY THIS VARIABLE. IN SCHEME 1, TERMS IN TEST SETS ARE WEIGHTED BY THE FREQUENCY WITH WHICH THEY WERE USED BY THE MEMBERS OF THE CRITERION GROUP IN DESCRIBING THE GIVEN DOCUMENT; WHEREAS, IN SCHEME 2, WHICH WAS THE ONE EMPLOYED IN THE ORIGINAL APPLICATION, THE SQUARE OF THIS FREQUENCY IS USED FOR WEIGHTING. IT CAN BE SEEN THAT THE LATTER SCHEME PLACES MUCH GREATER EMPHASIS ON "POPULAR" CRITERION GROUP RESPONSES. LIKE EDITING, SCHEME 2 HAS THE EFFECT OF RAISING THE SCORES OF MOST TEST SETS AND GENERALLY INCREASES THE STANDARD DEVIATIONS IN COMPARISON WITH SCHEME 1; HOWEVER IT ALSO COMMONLY INCREASES THE DIFFERENCES BETWEEN MEAN SCORES FOR DIFFERENT INDEXING TREATMENTS. THE RESULTING EFFECT ON TEST SENSITIVITY IS COMPLEX. AS ONE EXAMPLE, WITH SCHEME 1 MEAN SCORES FOR COMPUTER-EDITED PROFESSIONAL INDEXER SETS VS. AUTHOR SETS ARE 34 (s.d., 10) VS. 37 (s.d., 12); WHEREAS, WITH SCHEME 2 THE CORRESPONDING VALUES ARE 49 (s.d., 16) VS. 59 (s.d., 18). FOR THIS CONTRAST, THE ADVANTAGE OF SCHEME 2 IS APPARENT, BUT NOT MARKED. ON THE OTHER HAND, FOR THE CONTRAST BETWEEN MANUAL-EDITED AUTHOR SETS VS. TITLE SETS, SCHEME 2 IS GREATLY SUPERIOR--54 (s.d., 16) VS. 34 (s.d., 17) AS COMPARED TO 74 (s.d., 16) VS. 35 (s.d., 26). SINCE WEIGHTING BY SCHEME 2 ENTAILS A RELATIVELY SMALL INCREMENT IN EFFORT OVER WHAT IS REQUIRED WITH SCHEME 1, IT MAY BE A USEFUL OPTION IN SOME CIRCUMSTANCES.

WE CONSIDERED WEIGHTING SCHEMES THAT WOULD "PENALIZE" OVERASSIGNMENT OF TERMS MORE HEAVILY THAN EITHER SCHEME 1 OR SCHEME 2; FOR EXAMPLE, BY GIVING A NEGATIVE WEIGHT TO TERMS IN TEST SETS THAT WERE NOT USED BY ANY NUMBER OF THE CRITERION GROUP. HOWEVER, THE SCHEMES CONSIDERED HAD NUMEROUS TECHNICAL DISADVANTAGES; AND SINCE THE % MAXIMAL SCORE DIVIDED BY THE TOTAL NUMBER OF TERMS IN THE TEST SET CAN SERVE AS A MEASURE OF INDEXING EFFICIENCY, AS CONTRASTED WITH EFFECTIVENESS, THIS MATTER HAS NOT BEEN PURSUED FURTHER.

VARIABLE 6--CONFOUNDING BEFORE SCORING

AS AN EXPLORATORY TRIAL OF THE EFFECT OF CONFOUNDING, I.E., GENERIC POSTING, ALL PROFESSIONAL INDEXER AND AUTHOR SETS FOR 32 DOCUMENTS WERE RESCORED AFTER ADDING TO EACH TEST SET ANY TERMS SHOWN BY THE VOCABULARY GUIDE AS GENERIC TO TERMS IN THE ORIGINAL TEST SET. SCORING CREDIT WAS THEN GIVEN TO SUCH ADDED TERMS WHEN THEY MATCHED CRITERION SET TERMS. CONFOUNDING INCREASED THE GRAND MEAN FOR THE PROFESSIONAL INDEXER SETS BY 6 POINTS, AND THE MEAN FOR AUTHOR SETS WAS ALSO INCREASED BY 6 POINTS; IN BOTH CASES, THE STANDARD DEVIATION WAS UNCHANGED. THESE FINDINGS SUGGEST THAT CONFOUNDING HAS LITTLE OR NO EFFECT ON TEST SENSITIVITY, WHICH WAS THE MAIN QUESTION PROMPTING THE TRIAL. CONFOUNDING, HOWEVER, MAY HAVE AN ADVANTAGE FOR CERTAIN APPLICATIONS, E.G., IN TESTS IN THE CONTEXT OF AN INDEXING SYSTEM THAT EMPLOYS HIERARCHICAL STRUCTURE, AND WHERE IT MAY BE DESIRABLE TO MAKE MORE COMPARABLE INDEXING DONE AT DIFFERENT LEVELS OF SPECIFICITY. IF AN HIERARCHICAL ORGANIZATION OF INDEXING TERMINOLOGY HAS BEEN CREATED PRIOR TO SCORING, THE PROCESS CAN BE CARRIED OUT DURING EITHER MANUAL OR COMPUTER SCORING AT A RELATIVELY LOW COST. CONFOUNDING THEREFORE REPRESENTS A USEFUL OPTION THAT INCREASES THE METHOD'S FLEXIBILITY.

CONCLUSIONS

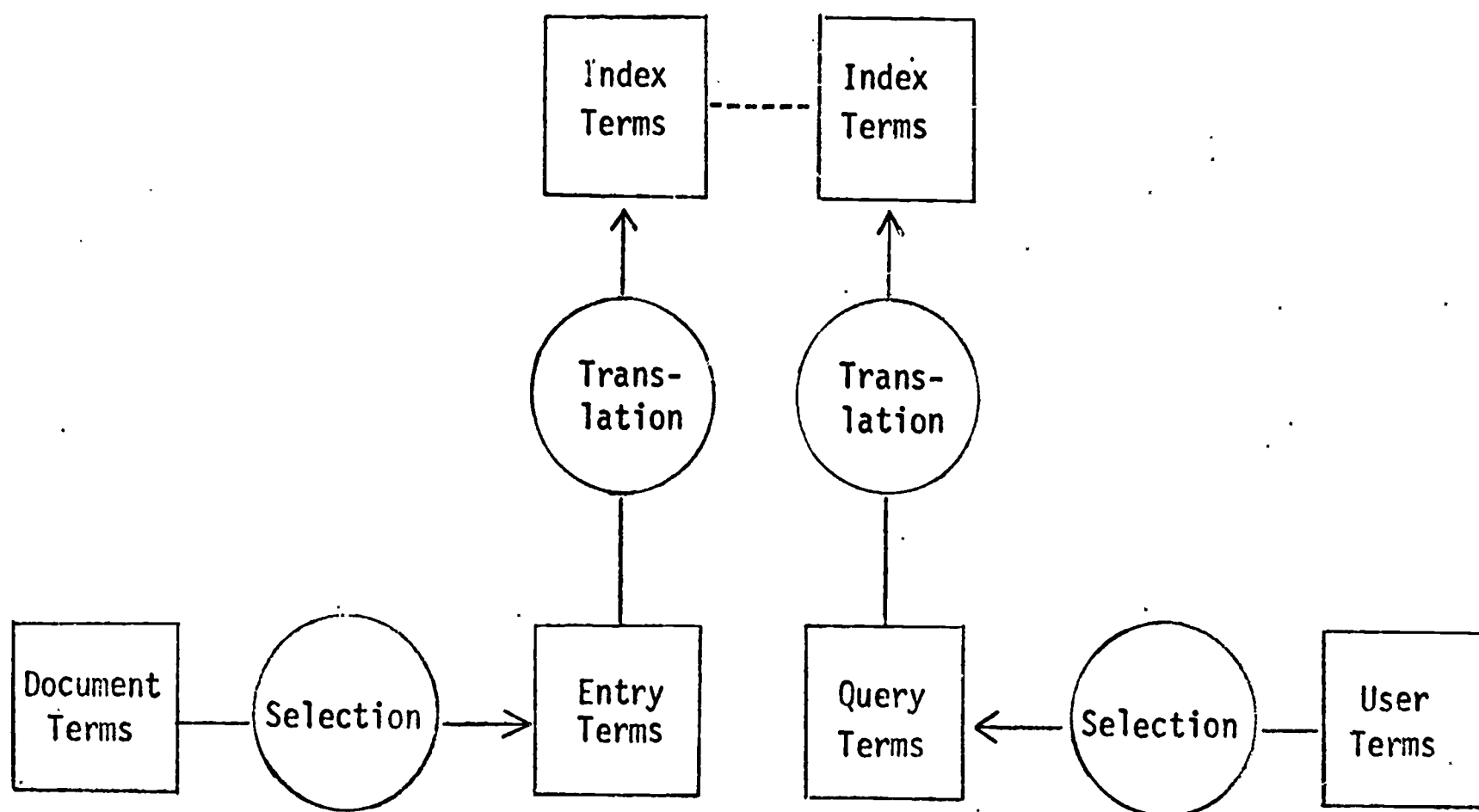
FROM THE RESULTS OF THESE STUDIES OF THE METHODOLOGIC VARIABLES, WE HAVE CONCLUDED THAT THE CONSENSUS-GROUP METHOD OF EVALUATING INDEXING CAN BE A PRACTICAL YARDSTICK FOR A WIDE VARIETY OF MANAGERIAL, RESEARCH, AND EDUCATIONAL USES.

APPENDIX A

RATIONALE OF METHOD AND LITERATURE REVIEW

A MODEL OF INFORMATION RETRIEVAL

THE FOLLOWING¹⁵ DIAGRAM WHICH HAS BEEN FREELY ADAPTED FROM KYLE^{18*} AND HYSLOP,⁵ REPRESENTS A SIMPLIFIED[#] MODEL OF THE CHAIN OF PROCESSES IN AN "INFORMATION RETRIEVAL" SYSTEM. THIS MODEL CAN ACCOMODATE ANY SYSTEM IN WHICH DOCUMENTS ARE INDEXED PRIOR TO THE RECEIPT OF QUERIES, WHETHER THE INDEXING IS DONE BY PEOPLE, MACHINES, OR MAN-MACHINE COMBINATIONS.



* IN THE FOLLOWING PAGES, ONLY ONE REFERENCE IS CITED ON MOST POINTS; OFTEN SEVERAL OTHER REFERENCES WOULD BE EQUALLY APPROPRIATE. SELECTION OF THE ONE USED AS AN EXAMPLE WAS USUALLY FORTUITOUS AND IS NOT MEANT TO IMPLY ATTRIBUTION OF PRIORITY OR NOVELTY.

THE OPERATIONS OF CODING, FILING, AND MATCHING ARE OMITTED HERE; THE "SHORT-CIRCUIT" IS SYMBOLIZED BY THE DOTTED LINE BETWEEN THE TOPMOST BOXES.

IN THIS MODEL, INDEXING IS DEPICTED AS A TWO-STEP PROCESS (CIRCLES LABELLED 1 AND 2). THE OPERATIONS PERFORMED DURING THE FIRST STEP, SELECTION, DETERMINE WHICH "ASPECTS" OF THE DOCUMENT WILL BE REPRESENTED IN THE INDEX. (THIS STEP IS ALSO CALLED "CONCEPT ANALYSIS" ⁹, "DOCUMENT ANALYSIS" ³⁵, "DETECTION" ³, AND VARIOUS OTHER NAMES.) THE OUTPUT OF SELECTION IS A SET OF ENTRY TERMS. (SYNONYMS FOR ENTRY TERMS INCLUDE, AMONG OTHERS, "ENTRY EXPRESSIONS" ²³, "DETECTION TERMS" ¹, "CLUE WORDS" ¹⁹, "INDICATOR WORDS" ¹², AND "CANDIDATE TERMS" ⁴⁰.) IN THE SECOND STEP, THE ENTRY TERMS ARE TRANSLATED INTO A SET OF INDEX TERMS. (THIS TRANSLATION IS COMMONLY REFERRED TO AS "STANDARDIZATION" OR "VOCABULARY CONTROL".) THE DISTINCTION BETWEEN THE TWO STEPS WAS APTLY DESCRIBED BY KYLE AS THE DIFFERENCE BETWEEN "WHAT TO INDEX" AND "HOW TO INDEX IT". ¹⁸

IN SYSTEMS WHERE NO ATTEMPT IS MADE TO CONTROL THE NUMBER OF DIFFERENT TERMS THAT MAY APPEAR IN THE INDEX (FOR EXAMPLE, SYSTEMS USING KWIC INDEXING, OR "PURE" UNITERMS ³⁷, THE TRANSLATION STEP IS, OF COURSE, MISSING; ENTRY TERMS ARE INDEX TERMS. SINCE FEW SYSTEMS REQUIRE INDEXERS TO RECORD ENTRY TERMS ROUTINELY, INDEXING MAY ALSO APPEAR TO BE A ONE-STEP PROCESS IN MANY SYSTEMS WHERE INDEX TERMS ARE CONTROLLED. IN SUCH CASES, HOWEVER, IT IS REASONABLE TO POSTULATE THAT THE TWO STEPS OCCUR IN THE INDEXER'S MIND, EVEN THOUGH THERE IS EVIDENCE TO SUGGEST THAT PROFESSIONAL INDEXERS MAY SOMETIMES THINK DIRECTLY IN CONTROLLED INDEXING LANGUAGE WHEN DECIDING WHICH ASPECTS OF A DOCUMENT SHOULD BE REPRESENTED IN THE INDEX. ²⁴ DESPITE THE FACT THAT IT IS OFTEN DIFFICULT TO SEPARATE CLEANLY THE SELECTION AND TRANSLATION STEPS, THE DISTINCTION IS VERY USEFUL IN ANALYZING THE INDEXING PROCESS BECAUSE SELECTION POSES THEORETICAL AND PRACTICAL PROBLEMS OF A DIFFERENT ORDER OF DIFFICULTY THAN THOSE OF TRANSLATION.

THE IMPORTANCE OF THE SELECTION STEP

CLEVERDON ⁸ AND OTHERS ²⁰ HAVE POSTULATED THAT, GIVEN A WELL-DEVELOPED "INDEXING LANGUAGE," * THE TRANSLATION STEP CAN BE REDUCED TO A CLERICAL OR MACHINE-LIKE ROUTINE; WHEREAS, THE SELECTION STEP IS AN INTELLECTUAL TASK. THE FACT THAT TRANSLATION WAS SUCCESSFULLY AUTOMATED IN 1963, ³⁰ AND THAT COMPUTER PROGRAMS TO ACCOMPLISH THE TRANSLATION STEP HAVE SINCE BEEN INTEGRATED INTO SEVERAL OPERATING SERVICES, ^{22,33} AS WELL AS BEING DEMONSTRATED (AS CONTRASTED TO SIMULATED) IN EXPERIMENTAL TRIALS, SUCH AS, ARTANDI'S, ¹ INDICATED THE VALIDITY OF THE POSITION THAT THE PROBLEMS OF SELECTION ARE OF A DIFFERENT ORDER THAN THOSE OF TRANSLATION.

* AS A MINIMUM, AN INDEXING LANGUAGE INCLUDES A SET, OR VOCABULARY, OF ENTRY TERMS; A SET OF INDEX TERMS; AND RULES FOR TRANSLATING FROM ONE SET TO THE OTHER. INDEXING LANGUAGES MAY HAVE VARIOUS OTHER ELEMENTS, BUT THESE THREE ARE ESSENTIAL.

IN ADDITION TO BEING A MORE CHALLENGINGLY DIFFICULT STEP, THE EFFECTIVENESS WITH WHICH SELECTION IS CARRIED OUT FIXES THE UPPER LIMITS ON THE PERFORMANCE OF THE ENTIRE CHAIN OF PROCESSES IN AN INFORMATION RETRIEVAL SYSTEM. AGAIN, IT WAS CLEVERDON WHO, WHETHER HE WAS THE FIRST TO DO SO OR NOT, CAN BE CREDITED FOR EMPHASIZING THIS IMPORTANT POINT AND MAKING IT CONVINCING. HE POINTED OUT THAT THE MAXIMUM PERFORMANCE ANY GIVEN SYSTEM IS CAPABLE OF, WITH REGARD TO "RECALL" AND "PRECISION" # DEPENDS UPON HOW COMPLETELY AND SPECIFICALLY ALL THE "CONCEPTS" IN THE DOCUMENTS HAVE BEEN INDEXED. ⁰ SINCE THE COMPLETENESS (OR "EXHAUSTIVITY") AND THE SPECIFICITY OF THE INDEX TERMS FOR A DOCUMENT CAN BE LESS, BUT NO GREATER, THAN THE COMPLETENESS AND SPECIFICITY OF THE ENTRY TERMS SELECTED FOR THAT DOCUMENT, IT FOLLOWS THAT HOW THE SELECTION STEP IS DONE DETERMINES THE HIGHEST LEVEL OF PERFORMANCE A GIVEN SYSTEM CAN PROVIDE -- VARIATIONS IN THE EFFECTIVENESS OF THE TRANSLATION STEP, OF THE INDEXING LANGUAGE ITSELF, AND OF ALL OTHER PROCESSES AND COMPONENTS IN THE SYSTEM CAN ONLY LOWER SYSTEM PERFORMANCE BELOW THIS THEORETICALLY ATTAINABLE LEVEL. IN OTHER WORDS, GOOD SELECTION IS A NECESSARY BUT NOT SUFFICIENT CONDITION FOR GOOD PERFORMANCE.

WE HAVE LEARNED HOW TO USE MACHINES IN OPERATING SYSTEMS TO EXECUTE, TIRELESSLY AND WITHOUT ERROR, THE TRANSLATIONS SPECIFIED BY AN INDEXING LANGUAGE; WE ARE BEGINNING TO LEARN HOW TO DESIGN AND USE INDEXING LANGUAGES SO THAT EITHER RECALL OR PRECISION CAN BE EMPHASIZED, DEPENDING ON WHAT THE REQUESTOR WANTS; AND MARKED PROGRESS HAS BEEN MADE IN IMPROVING CODING AND FILING, THE FINAL PROCESSES ON THE INDEXING SIDE OF THE INFORMATION RETRIEVAL "CHAIN" (SEE FIGURE, PAGE A-1). RELATIVE TO THE THEORETICAL AND PRACTICAL ADVANCES IN ALL THESE AREAS, PROGRESS APPEARS TO HAVE BEEN MUCH SLOWER IN UNDERSTANDING AND IMPROVING THE SELECTION STEP OF INDEXING. ONE CAN ARGUE THAT TODAY, AT THE PRESENT STATE-OF-THE-ART, SELECTION IS THE CRITICAL PROBLEM IN INDEXING, BOTH THEORETICALLY AND PRACTICALLY. FOR THESE REASONS WE WANTED TO DEVELOP AN EVALUATION METHOD THAT FOCUSED SPECIFICALLY ON THE SELECTION STEP AND COULD MEASURE ITS EFFECTIVENESS INDEPENDENT OF THE TRANSLATION STEP.

$$\# \text{ RECALL} = \frac{\text{NUMBER OF RELEVANT DOCUMENTS RETRIEVED IN RESPONSE TO A QUERY}}{\text{TOTAL NUMBER OF DOCUMENTS IN THE SYSTEM THAT ARE RELEVANT TO THE QUERY}}$$

$$\text{PRECISION} = \frac{\text{NUMBER OF RELEVANT DOCUMENTS RETRIEVED IN RESPONSE TO A QUERY}}{\text{TOTAL NUMBER OF DOCUMENTS RETRIEVED}}$$

CLEVERDON ORIGINALLY CALLED THE LATTER, "RELEVANCE RATIO" BUT LATER ACCEPTED THE SUGGESTION OF OTHERS AND CHANGED IT TO "PRECISION RATIO".⁹

REVIEW OF THE CRITERION PROBLEM IN INDEXING *

WHEN WE WERE WRITING UP THE FIRST USE OF OUR METHOD FOR EVALUATING INDEXING, THAT IS, OUR CRITERION MEASURE, CLEVERDON'S "RECALL" AND "PRECISION" RATIOS STILL HAD ALMOST THE STATUS OF AN INTERNATIONAL STANDARD. AT THAT TIME WE WERE RATHER APOLOGETIC ABOUT INTRODUCING A NEW CRITERION MEASURE, PARTICULARLY ONE THAT HAD NOT YET BEEN VALIDATED AGAINST THE "ULTIMATE" CRITERION CONCEPT, WHICH IN ITS FULLEST EXPLANATION RUNS SOMETHING LIKE THIS: PERFORMANCE OF A REAL SYSTEM, IN A REAL ENVIRONMENT, SUPPLYING REAL DOCUMENTS, FROM A REAL COLLECTION, IN RESPONSE TO REAL QUERIES, PROMPTED BY REAL PROBLEMS OF REAL USERS-- WITH PERFORMANCE RATED OBJECTIVELY ON THE BASIS OF HOW COMPLETELY THE SYSTEM RETRIEVES EVERY DOCUMENT IN THE COLLECTION THAT THE USER JUDGES AS "RELEVANT" TO HIS QUERY, AND HOW COMPLETELY IT RELIEVES THE USER OF THE CHORE OF WEEDING OUT DOCUMENTS HE FINDS IRRELEVANT. A NUMBER OF TRENDS THAT BEGAN SEVERAL YEARS AGO HAVE RECENTLY ACCELERATED, AND THE CRITERION "PROBLEM" HAS CHANGED MARKEDLY SINCE OUR HESITANT INTRODUCTION OF A NEW MEASURE. THESE TRENDS CAN BE SUMMARIZED AS FOLLOWS:

(1) THERE IS GROWING RECOGNITION OF THE NEED FOR, AND LEGITIMACY OF, PROXIMATE CRITERION MEASURES.

(2) IN ADDITION TO GOOD RETRIEVAL PERFORMANCE, AS MEASURED BY RECALL AND PRECISION, OTHER SYSTEM DESIDERATA ARE RECEIVING MORE EMPHASIS.

(3) THE UNIVERSAL APPROPRIATENESS AND GENERAL UTILITY OF RECALL AND PRECISION AS MEASURES OF THE ULTIMATE CRITERION CONCEPT IS BEING QUESTIONED MORE FREQUENTLY. ^{36, 38}

(4) THE CONCEPT OF "RELEVANCE", WHICH IS CENTRAL TO RECALL AND PRECISION MEASURES, IS UNDERGOING A RAPID AND DRASTIC METAMORPHOSIS. ^{13, 2, 11, 10}

(5) ON THE MOST FUNDAMENTAL LEVEL, THE IMPLICIT ASSUMPTION, BEHIND THE OLD "ULTIMATE" CRITERION CONCEPT IS BEING CHALLENGED; ⁴ THIS ASSUMPTION IMPLIES THAT EXHAUSTIVE SEARCH IS THE FUNCTION TO BE SERVED, WHEREAS IT IS ONLY ONE OF THE SEVERAL FUNCTIONS OF IR SYSTEMS (E.G., ALERTING, BROWSING, SEARCHING FOR "ENOUGH" DOCUMENTS TO MAKE A DECISION, ETC.), AND NOT NECESSARILY THE MOST IMPORTANT ONE.

* SNYDER'S DISTINCTIONS AMONG "CRITERION CONCEPT", "CRITERION MEASURE", AND "CRITERION VALUE" ³⁵ ARE USEFUL, AND WILL BE OBSERVED IN THE FOLLOWING DISCUSSION.

SINCE THE FIRST TWO TRENDS ARE ESPECIALLY PERTINENT TO THE MEASURES OF PREFERREDNESS TO BE USED IN THE PROPOSED STUDY, THEY WARRANT SOME DISCUSSION.

ALTHOUGH PROXIMATE CRITERIA CONCEPTS HAVE LONG BEEN USED FOR DAY-TO-DAY QUALITY CONTROL (E.G., ACCURACY AS JUDGED BY INDEXING SUPERVISORS), AND AS A BASIS FOR MANAGEMENT DECISIONS (E.G., QUALITY OF INDEXING AS JUDGED BY EXPERTS IN IR, OR BY EXPERTS IN THE SUBJECT-MATTER OF THE COLLECTION), THEY WERE CONSIDERED A KIND OF SECOND-CLASS MEASUREMENT AFTER RECALL AND PRECISION GAINED WIDE ACCEPTANCE IN THE IR COMMUNITY AROUND 1963. MORE RECENTLY, HOWEVER, BASED ON CONSIDERATIONS OF EXPERIMENTAL DESIGN, SNYDER³⁵ HAS OFFERED CONVINCING ARGUMENTS FOR THE UTILITY AND LEGITIMACY OF PROXIMATE, OR INTERMEDIATE, CRITERION CONCEPTS AND MEASURES. HE POINTS OUT THE NEED TO STUDY SEPARATELY THE DIFFERENT COMPONENTS IN AN IR SYSTEM USING SENSITIVE MEASURES SPECIFIC FOR THE COMPONENT BEING STUDIED. APPARENTLY HE IS NOT READY TO ABANDON THE OLD ULTIMATE CRITERION CONCEPT, HOWEVER, FOR HE ADDED THE PROVISIO THAT ANY PROXIMATE CRITERIA SHOULD BE VALIDATED AGAINST RETRIEVAL PERFORMANCE, PRESUMABLY MEASURED IN TERMS OF RECALL AND PRECISION. THE SAME NEEDS WERE EXPRESSED IN DIFFERENT WORDS BY A STUDY CONFERENCE SPONSORED BY NSF IN FEBRUARY 1965, WHERE THE CONSENSUS WAS THAT

"FOR THE TIME BEING, IN VIEW OF THE PRESENT STATE OF THE ART, EFFORTS TO DEVELOP AND TEST EVALUATION METHODS AND TO CONDUCT TESTS SHOULD BE CONCENTRATED ON SELECTED FEATURES OF DOCUMENT SEARCHING SYSTEMS IN SYSTEMS CONTEXT, RATHER THAN ON TOTAL SYSTEMS".²¹

IF ONE STILL HOLDS TO THE OLD ULTIMATE CRITERION CONCEPT, FOR WHICH THE PROPER MEASURES OF SYSTEM PERFORMANCE ARE EXPRESSED SOLELY IN TERMS OF TWO OR MORE RATIOS BASED ON THE FOUR-WAY PARTITION CREATED BY THE TWO DICHOTOMIES, RETRIEVED-NOT RETRIEVED AND RELEVANT-IRRELEVANT (OR AN N-WAY PARTITION, IF RELEVANCE IS RATED ON SOME SCALE), ALL OTHER MEASURES CAN BE ONLY PROXIMATE. THE VALIDITY OF SUCH MEASURES MUST, THEREFORE, DEPEND ON THEIR HAVING SOME DEPENDABLE RELATION TO THE ULTIMATE CRITERION MEASURE. IDEALLY, THIS RELATION IS DEMONSTRATED EMPIRICALLY; BUT IN PRACTICE IT IS OFTEN ASSUMED TO EXIST, EITHER BECAUSE THERE IS CONSENSUS THAT IT "SHOULD" EXIST OR BECAUSE IT FOLLOWS FROM AN ACCEPTED THEORY. SOONER OR LATER, HOWEVER, MOST PROXIMATE CRITERION MEASURES TEND TO ACQUIRE "FACE" VALIDITY AND ACHIEVE AN INDEPENDENT "STATUS" THAT IS ACCEPTED IN ALL EXCEPT THE MOST FORMAL USAGE. ANOTHER WAY PROXIMATE MEASURES BECOME INDEPENDENT OF THEIR ORIGINAL REFERENCE STANDARD IS BY BEING INCORPORATED INTO A NEW THEORY, OR BY A REDEFINING OF CONCEPTS CENTRAL TO THE OLD ULTIMATE CRITERION CONCEPT OF MEASURE. ALL OF THESE MECHANISMS ARE APPARENTLY

AT WORK IN THE IR FIELD TODAY, AND A HOST OF CRITERION MEASURES FOR INDEXING QUALITY ARE ACQUIRING STATUS.

MOST OF THESE "NEW" * MEASURES EMPLOY A GROUP OF INDIVIDUALS WHOSE COLLECTIVE RESPONSES ESTABLISH A CRITERION STANDARD AGAINST WHICH ANY "UNKNOWN" SAMPLE OF INDEXING IS MEASURED. THESE MEASURES FALL INTO THE FOLLOWING FOUR CATEGORIES, BASED ON THE TYPES OF INDIVIDUALS THAT COMPOSE THE "CRITERION GROUP". THE FOLLOWING TABLE CLASSIFIES SOME REPRESENTATIVE EXAMPLES OF SUCH MEASURES BUT DOES NOT INCLUDE THEM ALL:

<u>COMPOSITION OF CRITERION GROUP</u>	<u>SIZE OF GROUP</u>	<u>NAME OF CRITERION MEASURE</u>
I. <u>EXPERTS</u> (AUTHORITIES, ETC.)	AS FEW AS 1	"ACCURACY" 5
II. <u>INDEXERS</u>	2 OR MORE	"CONSISTENCY" 25 "PRECISION OF MEANING" 39
III. <u>AUTHORS</u>	NO. REPRESENTED IN DOCUMENT CORPUS	"RELEVANCE" 14
IV. <u>USERS</u> (SIMULATED QUERIES)	AS FEW AS 1	"RELEVANCE" 9 "REPRESENTATIVENESS" 16, 17, 26

CATEGORY IV INCLUDES ALL MEASURES IN WHICH THE CRITERION GROUP MEMBERS ARE NOT ACTUAL SYSTEM USERS JUDGING THE RELEVANCE OF SYSTEM RESPONSES TO THEIR OWN QUERIES, WHICH WERE GENERATED IN THE COURSE OF THEIR REGULAR WORK. THEREFORE, IT INCLUDES MEASURES IN WHICH THE CRITERION GROUP CONSISTS OF INDIVIDUALS WHO ARE ENTITLED TO, OR MIGHT BE EXPECTED TO, USE THE GIVEN SERVICE (POTENTIAL USERS), E.G., THE QUERISTS IN CLEVERDON'S LAST STUDY,⁹ OR INDIVIDUALS FROM A POPULATION CONSIDERED COMPARABLE TO THE SYSTEM'S CLIENTELE (SIMULATED USERS). THE "FACE" VALIDITY OF THESE MEASURES DEPENDS UPON ONE'S OPINION ON HOW CLOSELY THEY APPROACH REALITY.

THE METHODS USED TO CALCULATE CRITERION VALUES, AND THE SIZE OF CRITERION GROUPS, VARY WIDELY WITHIN A GIVEN CATEGORY, AS DOES THE PROBABLE RELIABILITY OF THE VALUES OBTAINED. IN A FEW CASES THESE MEASURES HAVE BEEN EMPIRICALLY VALIDATED AGAINST RETRIEVAL PERFORMANCE IN AN OPERATING SYSTEM. AS AN EXAMPLE, IN ONE SYSTEM, WHERE AN EXPERT'S JUDGEMENT COULD BE TAKEN AS FINAL AND "CORRECT", BRYANT DEMONSTRATED THAT "ACCURACY" VALUES CORRELATED HIGHLY WITH ACTUAL RETRIEVAL PERFORMANCE, AND THAT CONSISTENCY VALUES CORRELATED HIGHLY WITH ACCURACY. 5

* SOME ARE ACTUALLY OLD; BUT THE CONFIDENCE WITH WHICH THEY ARE USED SEEMS NEW.

ANOTHER WAY TO DEVELOP WHAT MIGHT BE CONSIDERED A SPECIAL TYPE OF A PROXIMATE CRITERION MEASURE IS TO "ADAPT" THE OVER-ALL SYSTEM CRITERION MEASURE (I.E., THE ULTIMATE CRITERION MEASURE) FOR ASSESSING SEPARATELY THE PERFORMANCE OF SOME SINGLE SYSTEM COMPONENT, SUCH AS INDEXING. IN A SPECIAL TEST SYSTEM, SUCH AS CLEVERDON'S OR THE COMPARATIVE SYSTEMS LABORATORY OF WESTERN RESERVE,²⁷ THIS CAN BE DONE BY ACTUALLY PERFORMING ALL THE OPERATIONS CONCERNED IN INFORMATION RETRIEVAL, KEEPING EVERYTHING CONSTANT EXCEPT THE COMPONENT UNDER STUDY.* HOWEVER, THIS PROCEDURE IS VERY DIFFICULT AND EXPENSIVE; AND ONE ALTERNATIVE IS TO "SIMULATE" THE CONDITION OF "ALL OTHER THINGS BEING EQUAL" BY SIMPLY SHORT-CIRCUITING PART OF THE CHAIN. THIS IS WHAT CLEVERDON DID IN HIS LAST SERIES OF STUDIES.⁹ HIS "SEARCHES" WERE PERFORMED BY PAPER-AND-PENCIL SIMULATION ON DOCUMENT-TERM MATRICES; THIS SIMULATION, ALTHOUGH IT WAS CARRIED OUT BY PEOPLE IN THIS STUDY WAS A CLERICAL OPERATION HE STATES COULD HAVE BEEN AUTOMATED. AN ESPECIALLY INTERESTING EXAMPLE OF A SHORT-CIRCUIT STRATEGY IS AN INGENIOUS METHOD KATTER HAS DEVELOPED TO ANALYZE AND COMPARE DOCUMENTS OR DOCUMENT REPRESENTATIONS, SUCH AS, INDEX TERMS, ABSTRACTS, ETC.¹⁶

THE FINAL STEP IN STREAMLINING THE EVALUATION OF INDEXING, OF COURSE, IS TO REPLACE THE QUALITATIVE SYSTEM MODEL WITH A MATHEMATICAL FORMULATION THAT PERMITS QUANTITATIVE PREDICTIONS OF RETRIEVAL PERFORMANCE GIVEN NUMERICAL VALUES FOR THE VARIABLES. THEN ONE CAN SIMPLY "PLUG IN" THE CRITERION VALUE FOR THE COMPONENT OR PROCESS ONE IS STUDYING AND CALCULATE THE ABSOLUTE VALUE FOR RETRIEVAL PERFORMANCE PREDICTED BY THE MATHEMATICAL MODEL OR THEORY. A NUMBER OF SUCH MODELS HAVE BEEN ADVANCED AS APPROPRIATE FOR AT LEAST PART OF A TOTAL SYSTEM (E.G., SALTON'S²⁶ AND BRYANT'S⁶, AND SOME OF THESE HAVE BEEN TESTED WITH VARYING DEGREES OF RIGOR. THIS ELEGANT WAY OF EVALUATING INDEXING SEEMS VERY ATTRACTIVE, ONCE THE UNDERLYING THEORY HAS BEEN WELL TESTED; BUT FOR ANY OF THE CURRENT MODELS, TESTING HAS THUS FAR BEEN LIMITED AND/OR CONFINED TO SPECIAL CASES WHERE SOME OF THE VARIABLES CAN BE SAFELY IGNORED.

* SNYDER EXPRESSES THE OPINION THAT, EVEN IF ONE COULD CONTROL ALL COMPONENTS OTHER THAN THE ONE BEING STUDIED, AND THEN SEE HOW CHANGES IN THIS ONE COMPONENT AFFECT RETRIEVAL, THE USUAL CRITERION MEASURES OF OVER-ALL SYSTEM PERFORMANCE WOULD BE TOO CRUDE AND INSENSITIVE TO ANSWER SOME IMPORTANT QUESTIONS ABOUT FACTORS INFLUENCING INDIVIDUAL COMPONENTS.^{35, 7} WHETHER THIS OPINION IS BASED ON THE RESULTS OF THE EARLY CRANFIELD STUDIES, WHICH SEEMED TO SHOW THAT THE OVER-ALL SYSTEM CRITERION MEASURE WAS REMARKABLY INSENSITIVE, ON THE FINDINGS OF OTHER STUDIES, OR ON THEORETICAL CONSIDERATIONS IS NOT CLEAR.

THE CRITERION-GROUP METHOD

OUR METHOD TESTS INDEXING AS A SUBSYSTEM OF THE INFORMATION STORAGE AND RETRIEVAL PROCESS. IT EMPLOYS A GROUP TO SET THE STANDARD AGAINST WHICH QUALITY IS TESTED, AS OPPOSED TO A SINGLE INDIVIDUAL'S JUDGMENT.

THIS CRITERION GROUP CAN BE MADE UP OF PROFESSIONAL INDEXERS, OF AUTHOR INDEXERS, OR OF DOCUMENT USERS. THE CHOICE IS LEFT TO THE PERSON USING THE METHOD, ACCORDING TO HIS JUDGMENT OF WHAT CONSTITUTES IDEAL INDEXING. IN ITS ORIGINAL APPLICATION, USERS CONSTITUTED THE CRITERION GROUP. 32

REFERENCES

1. ARTANDI, SUSAN S. BOOK INDEXING BY COMPUTER. THESIS. ANN ARBOR, MICHIGAN, UNIVERSITY MICROFILMS, INC., 1963.
2. ASSORIO, PETER G. CLASSIFICATION SPACE ANALYSIS. BOULDER, COLORADO, UNIVERSITY OF COLORADO. OCTOBER, 1964. AD 608 034.
3. BERNIER, CHARLES L. AND E. J. CRANE, "CORRELATIVE INDEXES. VIII. SUBJECT-INDEXING VS. WORD-INDEXING" JOURNAL OF CHEMICAL DOCUMENTATION, VOL. 2, No. 2, APRIL, 1962. P.117-122.
4. BRYANT, EDWARD C., "PROGRESS TOWARD EVALUATION OF INFORMATION RETRIEVAL SYSTEMS". IN COMMITTEE FOR INTERNATIONAL COOPERATION IN INFORMATION RETRIEVAL AMONG EXAMINING PATENT OFFICES. ANNUAL MEETING, PROCEEDINGS. WASHINGTON, D.C., SPARTAN BOOKS, 1965. P. 362-377.
5. BRYANT, EDWARD C., D. W. KING, AND P. J. TERRAGNO. SOME TECHNICAL NOTES ON CODING ERRORS. DENVER, COLORADO, WESTAT RESEARCH ANALYSTS, INC., JULY, 1963. WRA PO 7.
6. BRYANT, EDWARD C., DONALD T. SEARLES, AND ROBERT H. SHUMWAY. SOME THEORETICAL ASPECTS OF THE IMPROVEMENT OF DOCUMENT SCREENING BY ASSOCIATIVE TRANSFORMATIONS. DENVER, COLORADO, WESTAT RESEARCH ANALYSTS, INC., NOVEMBER 30, 1963. AD 628 191.
7. CLEVERDON, CYRIL W. REPORT ON THE TESTING AND ANALYSIS OF AN INVESTIGATION INTO THE COMPARATIVE EFFICIENCY OF INDEXING SYSTEMS. CRANFIELD, ENGLAND, ASLIB CRANFIELD RESEARCH PROJECT, OCTOBER, 1962.
8. CLEVERDON, CYRIL W. AND J. MILLS. THE TESTING OF INDEX LANGUAGE DEVICES. ASLIB PROCEEDINGS, VOL. 15, No. 4, APRIL, 1963. P. 106-130.
9. CLEVERDON, CYRIL W., JACK MILLS, AND MICHAEL KEEN. FACTORS DETERMINING THE PERFORMANCE OF INDEXING SYSTEMS, VOL. 1: DESIGN. 2 PARTS: TEST AND APPENDICES. CRANFIELD, ENGLAND, ASLIB CRANFIELD RESEARCH PROJECT, 1966.
10. CUADRA, CARLOS, R. V. KATTER, E. H. HOLMES AND E. M. WALLACE. EXPERIMENTAL STUDIES OF RELEVANCE JUDGMENTS: FINAL REPORT. 3 VOLUMES. NSFC-424. SANTA MONICA, CALIF., SYSTEMS DEVELOPMENT CORP. JUNE 30, 1967. TM-3520.

11. DOYLE, LAUREN B. IS RELEVANCE AN ADEQUATE CRITERION IN RETRIEVAL SYSTEM EVALUATION? IN H. P. LUHN, ED., AUTOMATION AND SCIENTIFIC COMMUNICATION, PROCEEDINGS OF THE ANNUAL MEETING OF THE AMERICAN DOCUMENTATION INSTITUTE, CHICAGO, ILL., OCTOBER, 1963. WASHINGTON, D.C., AMERICAN DOCUMENTATION INSTITUTE, 1963. PART 2, P.199-200.
12. HARLOW, JACQUES AND PAUL W. ABRAHAMS. RESEARCH IN INFORMATION RETRIEVAL. FINAL REPORT: AN INVESTIGATION OF THE TECHNIQUES AND CONCEPTS OF INFORMATION RETRIEVAL. PARAMUS, NEW JERSEY. ITT DATA AND INFORMATION SYSTEMS DIVISION, 31 JULY, 1964. AD 461 099.
13. HILLMAN, DONALD J. DOCUMENT RETRIEVAL THEORY, RELEVANCE AND THE METHODOLOGY OF EVALUATION. REPORT NO. 1: CHARACTERIZATION AND CONNECTIVITY. BETHLEHEM, PA., LEHIGH UNIVERSITY, CENTER FOR THE INFORMATION SCIENCES. 24 MAY, 1966.
14. HILLMAN, DONALD J. MATHEMATICAL THEORIES OF RELEVANCE WITH RESPECT TO SYSTEMS OF AUTOMATIC AND MANUAL INDEXING. IN H. P. LUHN, ED., AUTOMATION AND SCIENTIFIC COMMUNICATION, PROCEEDINGS OF THE AMERICAN DOCUMENTATION INSTITUTE ANNUAL MEETING, CHICAGO, ILL., OCTOBER, 1963. WASHINGTON, D.C., AMERICAN DOCUMENTATION INSTITUTE, 1963. PART 2, P. 323-324.
15. HYSLOP, MARJORIE R. THE ASM INFORMATION RETRIEVAL SYSTEM: AFTER CRANFIELD. JOURNAL OF DOCUMENTATION VOL. 21, No. 1, MARCH, 1965. P. 27-42.
16. KATTER, ROBERT V., EMORY H. HOLMES, AND RICHARD L. WEISS. EXPERIMENTAL INVESTIGATIONS OF A METHOD FOR ANALYZING DOCUMENT REPRESENTATION. TM-3090/000/00. SANTA MONICA, CALIFORNIA, SYSTEMS DEVELOPMENT CORPORATION. 9 AUGUST, 1966.
17. KATTER, ROBERT V. STUDY OF DOCUMENT REPRESENTATIONS; MULTIDIMENSIONAL SCALING OF INDEXING TERMS. FINAL REPORT. AUGUST 31, 1967. NSF GN-544 SANTA MONICA, CALIF. SYSTEMS DEVELOPMENT CORP. TM 3627.
18. KYLE, BARBARA R. F. INFORMATION RETRIEVAL AND SUBJECT INDEXING: CRANFIELD AND AFTER. JOURNAL OF DOCUMENTATION, Vol. 20, No. 2, JUNE, 1964. P. 55-69.
19. MARON, M. E. AUTOMATIC INDEXING: AN EXPERIMENTAL INQUIRY. JOURNAL OF THE ASSOCIATION FOR COMPUTING MACHINERY, Vol. 8, No. 3, JULY 1961. P. 404-417.
20. MONTGOMERY, CHRISTINE AND DON R. SWANSON. MACHINELIKE INDEXING BY PEOPLE. AMERICAN DOCUMENTATION, Vol. 13, No. 4, OCTOBER, 1962. P. 359-366.

21. NATIONAL SCIENCE FOUNDATION. SUMMARY OF STUDY CONFERENCE ON EVALUATION OF DOCUMENT SEARCHING SYSTEMS AND PROCEDURES. WASHINGTON, D.C., OCTOBER 2-3, 1964. WASHINGTON, D. C., NATIONAL SCIENCE FOUNDATION, 10 FEBRUARY, 1965.
22. NEWBAKER, H. R. AND T. R. SAVAGE. SELECTED WORDS IN FULL TITLE (SWIFT): A NEW PROGRAM FOR COMPUTER INDEXING. IN H. P. LUHN, ED., AUTOMATION AND SCIENTIFIC COMMUNICATION, PROCEEDINGS OF THE ANNUAL MEETING OF THE AMERICAN DOCUMENTATION INSTITUTE, CHICAGO, ILL. OCTOBER, 1963. WASHINGTON, D.C., AMERICAN DOCUMENTATION INSTITUTE, 1963. PART 2, P. 87-88.
23. O'CONNOR, JOHN J. MECHANIZED INDEXING METHODS AND THEIR TESTING. JOURNAL OF THE ASSOCIATION OF COMPUTING MACHINERY 11, 4, (OCTOBER 1964) PP. 437-449.
24. OLIVER, LAWRENCE H., CLAUDE MITCHELL, EDMUND W. FITZPATRICK, AND HERBERT JACOBSON. AN INVESTIGATION OF THE BASIC PROCESSES INVOLVED IN THE MANUAL INDEXING OF SCIENTIFIC DOCUMENTS. BETHESDA, MARYLAND. GENERAL ELECTRIC CORPORATION, INFORMATION SYSTEMS OPERATION. FEBRUARY 11, 1966.
25. RODGERS, DOROTHY J. A STUDY OF INTER-INDEXER CONSISTENCY. WASHINGTON, D.C., GENERAL ELECTRIC COMPANY, INFORMATION SYSTEMS SECTION, 29 SEPTEMBER, 1961.
26. SALTON, GERARD. INFORMATION STORAGE AND RETRIEVAL. SCIENTIFIC REPORT NO. ISR-13. DEPT. OF COMPUTER SCIENCE, CORNELL UNIVERSITY, ITHICA, N.Y. JANUARY, 1968.
27. SARACEVIC, TEFKO AND ALAN M. REES. TOWARDS THE IDENTIFICATION AND CONTROL OF VARIABLES IN INFORMATION RETRIEVAL EXPERIMENTATION. CLEVELAND, OHIO, WESTERN RESERVE UNIVERSITY, CENTER FOR DOCUMENTATION AND COMMUNICATION RESEARCH, JANUARY, 1966.
28. SCHULTZ, CLAIRE K., COMPILER AND EDITOR, GUIDE TO CURRENT TERMINOLOGY IN BIOMEDICAL RESEARCH. FEDERATION PROCEEDINGS VOL. 24, NO. 4, JULY-AUGUST, 1965.
30. SCHULTZ, CLAIRE K., EDITING AUTHOR-PRODUCED INDEXING TERMS AND PHRASES VIA A MAGNETIC TAPE THESAURUS AND A COMPUTER PROGRAM. IN LUHN, HANS P., ED., AUTOMATION AND SCIENTIFIC COMMUNICATION, WASHINGTON, D.C., AMERICAN DOCUMENTATION INSTITUTE, 1963. P. 9.
31. SCHULTZ, CLAIRE K. AND RICHARD H. ORR. EVALUATING INDEXING BY REFERENCE TO TERM-CHOICE PATTERNS OF A CRITERION GROUP. INTERNATIONAL FEDERATION FOR DOCUMENTATION 1965 CONGRESS, ABSTRACTS OF PAPERS, WASHINGTON, D.C., OCTOBER 10-15, 1965. P.86.

32. SCHULTZ, CLAIRE K., WALLACE L. SCHULTZ, AND RICHARD H. ORR. COMPARATIVE INDEXING: TERMS SUPPLIED BY BIOMEDICAL AUTHORS AND DOCUMENT TITLES. AMERICAN DOCUMENTATION 16, 4, (OCTOBER, 1965). PP. 299-312.
33. SHEPHERD, CLAYTON A. DESIGN AND IMPLEMENTATION OF THE AMERICAN SOCIETY FOR METALS MARK II DOCUMENTATION SYSTEM. IN PARAMETERS OF INFORMATION SCIENCE, PROCEEDINGS OF THE ANNUAL MEETING OF THE AMERICAN DOCUMENTATION INSTITUTE, PHILADELPHIA, PA., OCTOBER 5-8, 1964. WASHINGTON, D.C., SPARTAN BOOKS, 1964. P. 347-354.
34. SLAMECKA, V. AND J. JACOBY. EFFECT OF INDEXING AIDS ON THE RELIABILITY OF INDEXERS. FINAL TECHNICAL NOTE AF30 (602)-2616 DOCUMENTATION INC., WASHINGTON, D.C. JUNE, 1963.
35. SNYDER, MONROE B., ANNE W. SCHUMACHER, STEVEN E. MAYER AND M. DEAN HAVRON. METHODOLOGY FOR TEST AND EVALUATION OF DOCUMENT RETRIEVAL SYSTEMS: A CRITICAL REVIEW AND RECOMMENDATIONS. MCLEAN, VA., HUMAN SCIENCES RESEARCH, INC., JANUARY, 1966.
36. SWANSON, DON R. THE EVIDENCE UNDERLYING THE CRANFIELD RESULTS. LIBRARY QUARTERLY, VOL. 35 No. 1, JANUARY, 1965. P. 1-20.
37. TAUBE, MORTIMER. INSTALLATION MANUAL FOR THE UNITERM SYSTEM OF COORDINATE INDEXING. ARLINGTON HALL, VIRGINIA, ARMED SERVICES TECHNICAL INFORMATION AGENCY. OCTOBER, 1953.
38. TAUBE, MORTIMER. A NOTE ON THE PSUEDO-MATHEMATICS OF RELEVANCE. AMERICAN DOCUMENTATION, VOL. 16, No. 2, APRIL, 1965. P. 69-72.
39. TINKER, JOHN F. IMPRECISION IN MEANING MEASURED BY INCONSISTENCY OF INDEXING. AMERICAN DOCUMENTATION, VOL. 17, No. 2, APRIL, 1966. P. 96-102.
40. WALL, EUGENE. A UNIFIED ENGINEERING VOCABULARY FOR USE IN INFORMATION DISSEMINATION, INDEXING, STORAGE AND RETRIEVAL. IN H. P. LUHN, ED., AUTOMATION AND SCIENTIFIC COMMUNICATION, PROCEEDINGS OF THE ANNUAL MEETING OF THE AMERICAN DOCUMENTATION INSTITUTE, CHICAGO, ILL., OCTOBER, 1963. WASHINGTON, D.C., AMERICAN DOCUMENTATION INSTITUTE, 1963. PART 2, P. 37-38.

APPENDIX B

MATERIALS EMPLOYED IN STUDY TRIALS

DOCUMENTS

THE DOCUMENTS EMPLOYED IN ALL TRIALS CAME FROM THE SAME CORPUS, WHICH CONSISTED OF 285 BRIEF PRELIMINARY REPORTS OF RESEARCH. THE SOURCE OF THESE DOCUMENTS AND THE SELECTION OF THE CORPUS ARE DESCRIBED IN THE FOLLOWING PARAGRAPHS, WHICH ARE ADAPTED FROM THE ORIGINAL REPORT ON THE CRITERION GROUP METHOD.

EACH YEAR SEVERAL THOUSAND 10-MINUTE ORAL PAPERS REPORTING CURRENT BIOMEDICAL RESEARCH ARE GIVEN AT THE ANNUAL MEETING OF THE FEDERATION OF AMERICAN SOCIETIES FOR EXPERIMENTAL BIOLOGY. THIS NATIONAL CONVENTION IS THE LARGEST MEETING FOR BIOMEDICAL SCIENTISTS, AND THE WORK PRESENTED IS AN EXCELLENT CROSS-SECTION OF U.S. BIOMEDICAL RESEARCH. THE FEDERATION CONSISTS OF 6 SOCIETIES, EACH REPRESENTING A MAJOR, BASIC BIOMEDICAL DISCIPLINE--BIOCHEMISTRY, IMMUNOLOGY, NUTRITION, PATHOLOGY, PHARMACOLOGY, AND PHYSIOLOGY. ONLY MEMBERS OF THESE SOCIETIES MAY PRESENT UNSOLICITED PAPERS AT THE FEDERATION MEETING; AND THE SPEAKER MUST SUBMIT TO HIS SOCIETY A SHORT SUMMARY (225 WORDS OR LESS) OF WHAT HE PLANS TO SAY. THESE SUMMARIES ARE PUBLISHED IN A SPECIAL ISSUE OF FEDERATION PROCEEDINGS THAT APPEARS JUST BEFORE THE ANNUAL CONVENTION. ALTHOUGH THE DOCUMENT SUBMITTED IS CALLED AN "ABSTRACT", THE TERM IS A MISNOMER IN THAT THE DOCUMENT IS NOT USUALLY PRODUCED BY ABSTRACTING SOME PREEXISTING DOCUMENT. AUTHORS MOST COMMONLY PREPARE THE SUMMARY BEFORE THEY HAVE WRITTEN THE FULL TEXT OF THEIR ORAL PRESENTATION. WHEN PUBLISHED, SUCH ANTICIPATORY ABSTRACTS, THEREFORE, REPRESENT PRIMARY DOCUMENTS--CONDENSED, PRELIMINARY REPORTS THAT MAY OR MAY NOT BE FOLLOWED AT SOME LATER TIME BY THE PUBLICATION OF A MORE DETAILED REPORT. THE CORPUS FOR THIS STUDY WAS SELECTED BY TAKING EVERY 10TH DOCUMENT PUBLISHED IN THE 1962 MEETING ISSUE OF FEDERATION PROCEEDINGS, VOLUME 21, NO. 2, MARCH-APRIL (2,854 DOCUMENTS IN ALL). IN THIS SYSTEMATIC SAMPLE, EACH OF THE 6 SOCIETIES IS REPRESENTED BY 9-11% OF ALL THE DOCUMENTS SUBMITTED TO IT.

AUTHOR-INDEXING FORM

THE AUTHORS OF PAPERS GIVEN AT THIS MEETING WERE REQUIRED TO COMPLETE AN "AUTHOR-INDEXING FORM" LIKE THAT ILLUSTRATED IN FIGURE 1. THIS IS THE FORM REFERRED TO AS A VOCABULARY GUIDE THROUGHOUT THIS REPORT.

"AUTHOR INDEXING FORM"

(THE FORM CONSISTED OF 2 PAGES; HERE THE LOWER PART OF THE FIRST PAGE AND THE TOP OF THE SECOND PAGE HAVE BEEN OMITTED).*

Please study the subject-category list before marking. The list will be used primarily for the arrangement of the abstracts and for the production of the subject index to the abstracts. A secondary use will be for aid in programming.

Place the number "1" in the box at the left of the most specific category which classifies the area of your paper; the number "2" in the box at the left of the next most specific category. Do not mark more than two categories.

In the blanks at the end of the subject-category list, please supply four or more additional descriptive terms (words or short phrases) which can be used, besides the subject categories already selected, for further classifying and indexing the content of your paper. The terms you supply should preferably be nouns. Generic names of chemical compounds and drugs should be used, rather than trade names or jargon.

SUBJECT CATEGORIES

<input type="checkbox"/> 001 Amino Acids	<input type="checkbox"/> 040 Coagulation	<input type="checkbox"/> 084 Shock	<input type="checkbox"/> 127 Site
<input type="checkbox"/> 002 Metabolism	<input type="checkbox"/> 041 Agents; factors	<input type="checkbox"/> 085 Blood Vessels	<input type="checkbox"/> 128 Drug Metabolism
<input type="checkbox"/> 003 Nutrition	<input type="checkbox"/> 042 Fibrinolysis	<input type="checkbox"/> 086 Capillary exchange	<input type="checkbox"/> 129 Endocrines
<input type="checkbox"/> 004 Synthesis	<input type="checkbox"/> 043 Platelets	<input type="checkbox"/> 087 Venous return	<input type="checkbox"/> 130 Adrenal Cortex
<input type="checkbox"/> 005 Antigen-Antibody Reactions	<input type="checkbox"/> 044 Erythrocytes	<input type="checkbox"/> 088 Wave transmission	<input type="checkbox"/> 131 Adrenal Medulla
<input type="checkbox"/> 006 Cross Reactions	<input type="checkbox"/> 045 Destruction	<input type="checkbox"/> 089 Blood Volume	<input type="checkbox"/> 132 Anterior Pituitary
<input type="checkbox"/> 007 Haptens	<input type="checkbox"/> 046 Metabolism	<input type="checkbox"/> 090 Hemorrhage	<input type="checkbox"/> 133 ACTH
<input type="checkbox"/> 008 Immunofluorescence	<input type="checkbox"/> 047 Groups	<input type="checkbox"/> 091 Transfusion	<input type="checkbox"/> 134 Control of secretion
<input type="checkbox"/> 009 In Vivo Reactions	<input type="checkbox"/> 048 Hematopoiesis	<input type="checkbox"/> 092 Cardiac Drugs	<input type="checkbox"/> 135 Gonadotropin
<input type="checkbox"/> 010 Cellular	<input type="checkbox"/> 049 Hemoglobin	<input type="checkbox"/> 093 Cardiac Muscle	<input type="checkbox"/> 136 Somatotropin
<input type="checkbox"/> 011 Pathogenetic	<input type="checkbox"/> 050 Leukocytes	<input type="checkbox"/> 094 Disorders	<input type="checkbox"/> 137 TSH
<input type="checkbox"/> 012 Non-specific Factors	<input type="checkbox"/> 051 Leukemia	<input type="checkbox"/> 095 Electrocardiography	<input type="checkbox"/> 138 Brain Hormones
<input type="checkbox"/> 013 Complement	<input type="checkbox"/> 052 Plasma Proteins	<input type="checkbox"/> 096 Cardiac Output	<input type="checkbox"/> 139 Glucagon
<input type="checkbox"/> 014 Properdin	<input type="checkbox"/> 053 Albumin	<input type="checkbox"/> 097 Control	<input type="checkbox"/> 140 Insulin
<input type="checkbox"/> 015 Precipitation	<input type="checkbox"/> 054 Globulins	<input type="checkbox"/> 098 Measurement	<input type="checkbox"/> 141 Diabetes mellitus
<input type="checkbox"/> 016 Diffusion	<input type="checkbox"/> 055 Storage	<input type="checkbox"/> 099 CV Disease	<input type="checkbox"/> 142 Mode of action
<input type="checkbox"/> 017 Immunoelectrophoresis	<input type="checkbox"/> 056 Body Water	<input type="checkbox"/> 100 Edema	<input type="checkbox"/> 143 Parathyroid
<input type="checkbox"/> 018 Quantitation	<input type="checkbox"/> 057 Bone	<input type="checkbox"/> 101 Lymph	<input type="checkbox"/> 144 Posterior Pituitary
<input type="checkbox"/> 019 Antigens; Antibodies	<input type="checkbox"/> 058 Carbohydrates	<input type="checkbox"/> 102 Cell Structure; Function	<input type="checkbox"/> 145 Diabetes insipidus
<input type="checkbox"/> 020 Antibody Formation	<input type="checkbox"/> 059 Chemistry	<input type="checkbox"/> 103 Active Transport	<input type="checkbox"/> 146 Sex Hormones
<input type="checkbox"/> 021 Determinants	<input type="checkbox"/> 060 Metabolism	<input type="checkbox"/> 104 Cell Membranes	<input type="checkbox"/> 147 Androgens
<input type="checkbox"/> 022 Microorganisms	<input type="checkbox"/> 061 Citric acid cycle	<input type="checkbox"/> 105 Cytoplasm	<input type="checkbox"/> 148 Estrogens
<input type="checkbox"/> 023 Bacteria	<input type="checkbox"/> 062 Glycolysis	<input type="checkbox"/> 106 Microsomes	<input type="checkbox"/> 149 Progestogens
<input type="checkbox"/> 024 Rickettsia	<input type="checkbox"/> 063 Hexose phosphate path	<input type="checkbox"/> 107 Mitochondria	<input type="checkbox"/> 150 Thyroid
<input type="checkbox"/> 025 Polysaccharides	<input type="checkbox"/> 064 Monosaccharide conversions	<input type="checkbox"/> 108 Nuclei	<input type="checkbox"/> 151 Iodine metabolism
<input type="checkbox"/> 026 Proteins	<input type="checkbox"/> 065 Polysaccharides	<input type="checkbox"/> 109 Cell, Tissue Culture	<input type="checkbox"/> 152 Regulation
<input type="checkbox"/> 027 Toxins	<input type="checkbox"/> 066 Small cycles	<input type="checkbox"/> 110 Cell Antigens	<input type="checkbox"/> 153 Thyroxine
<input type="checkbox"/> 028 Transplantation	<input type="checkbox"/> 067 Photosynthesis	<input type="checkbox"/> 111 Metabolism	<input type="checkbox"/> 154 Energy Metabolism
<input type="checkbox"/> 029 Autoantibodies	<input type="checkbox"/> 068 Cardiovascular System	<input type="checkbox"/> 112 Neoplasms	<input type="checkbox"/> 155 Environment
<input type="checkbox"/> 030 Tissue antibodies	<input type="checkbox"/> 069 Atherosclerosis	<input type="checkbox"/> 113 Nucleic Acids	<input type="checkbox"/> 156 Adaptation
<input type="checkbox"/> 031 Biological Oxidations	<input type="checkbox"/> 070 Experimental	<input type="checkbox"/> 114 Chemotherapy	<input type="checkbox"/> 157 Air Pollution
<input type="checkbox"/> 032 Cytochromes	<input type="checkbox"/> 071 Nutritional	<input type="checkbox"/> 115 Bacterial	<input type="checkbox"/> 158 Altitude
<input type="checkbox"/> 033 Electron Transport	<input type="checkbox"/> 072 Pathophysiology	<input type="checkbox"/> 116 Cancer	<input type="checkbox"/> 159 Hibernation
	<input type="checkbox"/> 073 Blood Flow	<input type="checkbox"/> 117 Parasitologic	<input type="checkbox"/> 160 Hyperthermia; Heat
	<input type="checkbox"/> 074 Cerebral	<input type="checkbox"/> 118 Connective Tissue Disorders	
	<input type="checkbox"/> 075 Coronary	<input type="checkbox"/> 119	
<input type="checkbox"/> 210 Synthesis	<input type="checkbox"/> 262 Neurochemistry	<input type="checkbox"/> 320 Control	<input type="checkbox"/> 366 Hepatitis
<input type="checkbox"/> 211 Transport	<input type="checkbox"/> 263 Pain	<input type="checkbox"/> 321 Disorders	<input type="checkbox"/> 367 Vitamins
<input type="checkbox"/> 212 Phospholipids	<input type="checkbox"/> 264 Peripheral Nerves	<input type="checkbox"/> 322 Diuresis; Diuretics	<input type="checkbox"/> 368 B
<input type="checkbox"/> 213 Metabolism	<input type="checkbox"/> 265 Reflexes	<input type="checkbox"/> 323 Electrolyte Excretion	<input type="checkbox"/> 369 B ₁₂
<input type="checkbox"/> 214 Synthesis	<input type="checkbox"/> 266 Axon	<input type="checkbox"/> 324 Glomerular Filtration	<input type="checkbox"/> 370 C
<input type="checkbox"/> 215 Sterols	<input type="checkbox"/> 267 Conditioned		<input type="checkbox"/> 371 Fat-Soluble
<input type="checkbox"/> 216 Metabolism	<input type="checkbox"/> 268 Spinal Cord		<input type="checkbox"/> 372 Folic Acid
<input type="checkbox"/> 217 Synthesis	<input type="checkbox"/> 269 Nitrogen Metabolism		<input type="checkbox"/> 373 Unidentified

ADDITIONAL DESCRIPTIVE TERMS

.....

.....

.....

* THIS FIGURE IS REPRODUCED FROM: SCHULTZ, CLAIRE K., WALLACE L. SCHULTZ, AND RICHARD H. ORR. COMPARATIVE INDEXING: TERMS SUPPLIED BY BIOMEDICAL AUTHORS AND DOCUMENT TITLES. AMERICAN DOCUMENTATION 16, 4, (OCT. 1965). P.299-312.

THESAURUS

THE THESAURUS EMPLOYED WHEN CRITERION AND TEST SETS WERE MANUALLY STANDARDIZED HAS BEEN PUBLISHED AS THE "GUIDE TO CURRENT TERMINOLOGY IN BIOMEDICAL RESEARCH",* WHICH REPRESENTS AN INDEXING VOCABULARY OF THE FEDERATION OF AMERICAN SOCIETIES FOR EXPERIMENTAL BIOLOGY. THE GUIDE LISTS 1,516 DIFFERENT TERMS CONSISTING OF ONE OR MORE WORDS, AND SPECIFIES THEIR CLOSEST EQUIVALENT IN THE INDEXING "LANGUAGES" USED BY THE NATIONAL LIBRARY OF MEDICINE, DEFENSE DOCUMENTATION CENTER, AND THE DIVISION OF RESEARCH GRANTS OF NATIONAL INSTITUTES OF HEALTH. THERE IS A HIGH DEGREE OF "COMPATIBILITY" AMONG THE INDEXING VOCABULARIES OF FASEB, NLM, DDC, AND NIH; THREE-QUARTERS OF ALL THE FASEB TERMS ARE READILY TRANSLATABLE INTO BOTH NLM AND NIH LANGUAGES.

* SCHULTZ, CLAIRE K., COMPILER AND EDITOR, GUIDE TO CURRENT TERMINOLOGY IN BIOMEDICAL RESEARCH. FEDERATION PROCEEDINGS VOL. 24, NO. 4, JULY-AUGUST, 1965.

APPENDIX C

SUBJECTS PARTICIPATING IN STUDY TRIALS

Criterion group

The criterion sets employed in this study were established by a group of potential users of the indexing for the document corpus -- in this case, members of the professional association from which the documents in this corpus were obtained, the Federation of American Societies of Experimental Biology, as described in Appendix A. The criterion group was a sample of the membership selected to represent the document authors' peers. Two active research workers from each of the six disciplines in the Federation were selected by Dr. Milton Lee, Executive Officer of the Federation, on the basis of their recognized standing in the research community and on the likelihood that they would be willing to participate in the study. To facilitate holding a meeting at Federation Headquarters in Bethesda at which the study could be explained and uniform instructions could be given, the original selection was limited to scientists in the Bethesda area (National Institutes of Health and Naval Medical Research Institute). One of the 12 scientists originally selected had to withdraw; he was replaced by a research worker in a pharmaceutical company, who was also well known in his discipline.

Author-indexers

As described in Appendix A, each of the authors of the documents in the corpus had supplied indexing terms when he submitted his paper. The author sets consisted of these indexing terms. The titles these authors had given the documents supplied the title sets.

Professional indexers

This group consisted of eight professional indexers, all of whom were experienced in working with biomedical documents. With one exception, they were senior personnel from indexing services or from information service departments of pharmaceutical companies. The exception was an indexer who was currently working directly with biomedical scientists in a university setting. This group represents a sample of the universe of such indexers selected largely on the basis of friendship with the present investigators.

Non-professional indexers

The non-professional indexers employed in this study consisted of two groups of second-year medical students from the same school, none of which had any experience in indexing. Ten students comprised Group A; there were nine students in Group B.

APPENDIX D

Procedures for Manual Implementation*

Editing and recording term uses

A term-use matrix, similar to that shown in Table I, was created for each document however, only two test sets - the author and title sets are shown here. In the present study, term usage by the professional indexer group and the non-professional indexer group were similarly recorded. Each term use by each of the 12 members of the criterion group and by the author was indicated by an X; the "presence" of a term in the title was similarly recorded. The two members of each of the six disciplinary pairs making up the criterion group were designated A and B.

Weighting

The criterion set of terms for a document consists of all the different terms used by members of the criterion group to describe that document; thus 285 criterion sets were established, one for each document. For the document illustrated in Table I, the criterion set contained 13 terms. Each term in a criterion set was assigned a weighting factor by one of two schemes -- in Scheme #1 the weight was equal to the number of criterion group members who had used it to describe the given document; whereas in Scheme #2, this number was squared. This weighting procedure is illustrated in Table I. Note that the weighting of a term was not affected by whether the author had or had not included it among the indicia he supplied, or by whether it was supplied by the document title. Terms in test sets that had not been used by at least one member of the criterion group we will refer to as "zero terms", since they were given a weight of 0.

Scoring

The raw score for each test set was calculated by adding the weights for all terms in the set. For example, for the author set shown in Table I, the raw score is 20 when the terms are weighted by Scheme #1, and 84 by Scheme #2. Since the number of terms, and the weighting of these terms, varies from one criterion set to another, the constraints on the raw scores also vary. To facilitate comparisons, we converted the raw scores into percentages of the highest score that could be awarded ("maximal score"), i.e., the sum of the weights for all terms in the criterion set. For the document illustrated in Table I, if the author set had contained all of the 13 terms in the criterion set, its raw score would have equaled the maximal score for this document (28 by Scheme #1, 96 by Scheme #2) the actual raw score was 71% of the maximal score when Scheme #1 was employed, and 88% with Scheme #2.

* This material is adapted from the published description of the first application of the method [American Documentation 16, 4, (October, 1965)].

TABLE D-1 Illustrative Term-Use Matrix for One Document

Term ^a	Criterion Group										Test Author	Sets Title				
	Biochemist		Physiologist		Pharmacologist		Pathologist		Immunologist				Nutritionist		Term Weight	
	A	B	A	B	A	B	A	B	A	B			A	B	# 1	# 2
Chemotherapy—Cancer (116)		x			x		x		x		x		7	49	x	
Neoplasms				x		x						x	4	16	x	x
Mast cells				x					x				3	9	x	x
Cell, tissue culture—Neoplasms (112)	x												2	4	x	
Chemotherapy (114)								x			x		2	4		
Benadryl										x			2	4		
Histamine-antihistaminics (195)				x					x				2	4		
Bacteria				x						x			2	4	x	
Cell structure; function (102)				x									1	1		
Cell, tissue culture (109)	x												1	1		
Drug action (123)	x												1	1	x	
Hosts													1	1		
Metabolism										x			1	1		
Peritoneal fluid													1	1	x	
Mice													0	0		x
													0	0		x

^aThe parenthetical numbers following certain terms refer to the code numbers of the corresponding subject categories listed on the author-indexing form used by both authors and criterion groups; in these cases, the term given here may represent an "expansion" of the term that appeared on the form (see text). Terms not followed by a number are "write-ins," i.e., words or phrases not included among the subject categories listed on the printed form.

^bThis member of the criterion group did not respond to document No. 100.

APPENDIX E

COMPUTER IMPLEMENTATION

AFTER TEST INDICIA HAVE BEEN OBTAINED IT IS POSSIBLE TO ACCOMPLISH ANY OR ALL OF THE ADDITIONAL PROCESSING STEPS BY MEANS OF COMPUTER PROGRAMS, SUCH AS THOSE CONSTRUCTED FOR CARRYING OUT THIS STUDY.* FOR MACHINE PROCESSING THE FIRST REQUIREMENT IS THAT THE DATA BE MADE MACHINE-READABLE, THAT IS, KEYPUNCHED. THE INPUT FORMAT USED FOR THIS STUDY CONSISTED OF A DOCUMENT IDENTIFICATION, AN INDEXER IDENTIFICATION AND ONE INDEXING "TERM"--AN ALPHABETIC OR NUMERIC EXPRESSION--OF ANY LENGTH, UP TO THE CAPACITY OF A PUNCHED CARD. THE SPECIFIC KEY-PUNCHING INSTRUCTIONS USED ARE GIVEN IN TABLE E-1.

IF MACHINE-EDITING IS TO BE DONE, OR EVEN IF IT IS NOT, THE NEXT PROCESSING STEP IS TO "TAG" THE KEYPUNCHED INDICIA SO THAT EVERY "WORD" CAN BE IDENTIFIED WITH ITS DOCUMENT, INDEXER, AND POSITION WITHIN THE INDEXER'S TOTAL RESPONSE TO THE DOCUMENT. THE "TAGGED" UNITS (WORDS) ARE SORTED SO THEY WILL BE PROPERLY ORGANIZED FOR MATCHING EITHER THE THESAURUS, IF MACHINE STANDARDIZING IS DONE, OR IF EDITING IS TO BE OMITTED, THE CRITERION SET.#

STANDARDIZING INDICIA CAN ACCOMPLISH ANY OF THE FOLLOWING:

(1) ELIMINATE WHAT ARE CONSIDERED "NONSUBSTANTIVE" WORDS, SUCH AS CONNECTIVES OR (2) CHANGE WORD VARIANTS SUCH AS SINGULAR AND PLURAL FORMS, INTO A "STANDARD" FORM, OR (3) CONVERT WHAT ARE CONSIDERED SYNONYMOUS EXPRESSIONS INTO A SINGLE "STANDARD" EXPRESSION, OR (4) ADD ADDITIONAL, POSSIBLY MORE GENERIC, WORDS TO INDICIA, SUCH AS "CARBOHYDRATES" IN RESPONSE TO THE TERM "GLUCOSE". TO PERFORM SUCH TRANSFORMATION ON THE RAW INDICIA THERE MUST FIRST BE A SET OF "REWRITE" RULES FOR ALL ANTICIPATED ENTRY TERMS \$ AND ALSO A PROGRAM WHICH MATCHES

* PROGRAMS WRITTEN IN FORTRAN IV AND PL/I FOR USE ON THE IBM 360/67 COMPUTER. THE INVESTIGATORS CAN MAKE THESE PROGRAMS AVAILABLE TO INTERESTED PERSONS.

THE CRITERION DATA WILL HAVE BEEN GIVEN THE SAME TREATMENT AS THE INDICIA, PRIOR TO ANY MATCHING OF THE CRITERION AND INDEXING SETS.

\$ INCLUDING THE TRIVIAL REWRITE RULE THAT RETAINS SOME TERMS IN THE SAME FORM AS WHEN ENCOUNTERED IN THE INPUT. WORDS NOT OF INTEREST CAN BE OMITTED FROM THE THESAURUS AND AUTOMATICALLY DELETED FOR REASON OF NON-MATCH, BUT SINCE THIS PRACTICE DELETES "NEW" WORDS OF INTEREST, IT IS A BETTER PRACTICE TO ACTIVELY DELETE UNWANTED TERMS AND "SAVE" NONMATCHING WORDS ON A SEPARATE LIST THAT IS PUT OUT FOR HUMAN REACTION. WITH THIS APPROACH IT IS ALSO POSSIBLE TO DECLINE TO MAKE THESAURAL REWRITE RULES FOR CERTAIN SEMANTICALLY AMBIGUOUS TERMS AND WAIT UNTIL THEIR FULL CONTEXTS ARE KNOWN (AT PROCESSING TIME) TO INSTRUCT THE COMPUTER HOW TO DEAL WITH SPECIFIC OCCURRENCES.

ENTRY TERMS WITH THE THESAURUS, AND THEN CARRIES OUT THE THESAURUS REWRITE RULES AS THEY ARE ENCOUNTERED. THE THESAURUS USED IN THIS STUDY EXAMINES ONE INPUT WORD AT A TIME, BUT IT ALSO EXAMINES CONTEXTS, SO THAT EXPRESSIONS SUCH AS "AMINO ACIDS" OR "CITRIC ACID CYCLE" CAN BE RETAINED AS UNITS. MORE INFORMATION IS GIVEN ABOUT THE RECENTLY DEVELOPED CONTEXT-DEPENDENT STANDARDIZING TECHNIQUE USED IN THIS STUDY IN A SEPARATE PAPER. *

THE FOLLOWING EXAMPLE WILL SERVE TO ILLUSTRATE AN EXAMPLE OF STANDARDIZATION, THE FOLLOWING PHRASE, "AMINO ACIDS IN RUMINANT NUTRITION" WOULD BE EXAMINED ONE WORD AT A TIME, AS EACH WAS ENCOUNTERED IN ITS ALPHABETIC ORDER WHEN PROCESSING A LIST OF ALL INDIVIDUAL WORDS MAKING UP THE INDICIA OF A TEST CORPUS. "ACIDS", ENCOUNTERED FIRST, WOULD BE HELD FOR POTENTIAL "PARTNER WORDS" TO BE ENCOUNTERED IN THE LATER PORTION OF THE ALPHABETIC LIST. "AMINO", ENCOUNTERED NEXT, WOULD, BY MEANS OF ITS SEQUENTIAL "TAG" BE IDENTIFIED AS A "PARTNER" OF ACIDS AND THE REWRITE RULE WOULD CAUSE "AMINO ACIDS" TO BECOME THE STANDARDIZED INDEX TERM. "IN" WOULD BE DELETED FROM THE LIST AS SOON AS ENCOUNTERED; "NUTRITION" WOULD BE HELD FOR POTENTIAL MATCH WITH "PARTNER WORDS". WHEN "RUMINANT" WAS PROCESSED IT WOULD BE IDENTIFIED AS A PARTNER WORD OF "NUTRITION" AND THE REWRITE RULE WOULD TRANSFORM "RUMINANT NUTRITION" INTO THE STANDARDIZED TERM "ANIMAL NUTRITION".

IF STANDARDIZATION IS NOT DONE, NON-SUBSTANTIVE WORDS SUCH AS "AND", WORD VARIANTS SUCH AS "ENZYME", "ENZYMES", "ENZYMAL", "ENZYMATIC", AND SYNONYMS SUCH AS "HEART" AND "CARDIAC" ARE CANDIDATES FOR MATCH. EVERY SUCH WORD IS TREATED AS UNRELATED TO THE OTHER WHEN THE CRITERION AND INDEXING SETS ARE MATCHED FOR SCORING. AS A RESULT, IF THE CRITERION SET CONTAINS, FOR EXAMPLE, "CARDIAC ARREST" AND THE INDEXING SET CONTAINS "HEART ARREST", THE INDEX SET WILL NOT GAIN ANY SCORE, BECAUSE OF MISMATCH-- BUT THE SAME INDEXING SET COULD GAIN SCORE BECAUSE A TRIVIAL WORD SUCH AS "OF" DID MATCH IN THE TWO SETS.

THE PURPOSE OF THE SCORING PROGRAM IS TO PERFORM THE MATCHING OPERATION, FOR ONE DOCUMENT AT A TIME, BETWEEN THE INDEXING SET(S) AND THE CRITERION SET. THE PROGRAM CAN HAVE VARIOUS OPTIONS, AS IS TRUE FOR THE SCORING PROGRAM USED IN THIS STUDY, WHICH INSTRUCT THE PROGRAM TO CALCULATE SCORES FOR A SINGLE TEST SET OR GROUP OF SETS, FOR SINGLE DOCUMENTS OR FOR GROUPS OF DOCUMENTS. ANOTHER KIND OF OPTION INSTRUCTS THE PROGRAM TO CALCULATE STANDARD DEVIATIONS, MAKE "T" TESTS, OR PERFORM OTHER STATISTICAL COMPUTATIONS, AS THE REQUIRED DATA BECOME AVAILABLE DURING PROCESSING. THE PROGRAM CAN BE USED TO PRINTOUT DETAIL ABOUT THE MATCHING PROCESSES.

* GOPNIK, MYRNA AND CLAIRE K. SCHULTZ. METHODS FOR THESAURUS PROCESSING OF CONTEXT-DEPENDENT SEGMENTS IN LANGUAGE. SUBMITTED FOR PUBLICATION.

AND INTERMEDIATE CALCULATIONS IT PERFORMS, OR ONLY SPECIFIED RESULTS, SUCH AS PERCENT OF MAXIMAL SCORE OR POINTS PER TERM.

IF ALL OF THE COMPUTER PROGRAMS JUST DESCRIBED ARE CONSIDERED AS A SYSTEM, WITH THE STANDARDIZATION PROCEDURE OPTIONAL, IT CAN BE SEEN THAT KEYPUNCHED RAW INDICIA CAN BE FED INTO THE COMPUTER, AND THE SCORED RESULTS OBTAINED, WITHOUT ANY MANUAL PROCESSING REQUIRED.

TABLE E-1
KEYPUNCHING INSTRUCTIONS FOR INDEXING DATA

WRITE THE DOCUMENT NUMBER IN THE FIRST THREE COLUMNS. IT IS ALWAYS A THREE DIGIT NUMBER (IT WILL RANGE FROM 001 TO 285) AND IS FOUND IN THE TOP RIGHT HAND CORNER OF THE INDEXING SHEET. DO NOT SKIP A SPACE. IN THE NEXT TWO COLUMNS WRITE THE NUMBER CODE OF THE PARTICULAR INDEXER. THIS WILL ALWAYS BE A TWO DIGIT NUMBER (IT WILL RANGE FROM 00 TO 99), AND IS FOUND IN THE TOP LEFT HAND CORNER OF THE INDEXING SHEET. DO NOT SKIP A SPACE. IN THE NEXT COLUMN WRITE THE NUMBER OF THE DISCIPLINE OF THE AUTHOR OF THE DOCUMENT. THE DISCIPLINE OF THE AUTHOR IS GIVEN ON THE EXTREME RIGHT OF THE FIRST LINE OF THE DOCUMENT IN THE FASEB BOOK. USE THE FOLLOWING CODE TO NUMBER THE DISCIPLINE:

- | | |
|-----------------|---------------|
| 1. PHYSIOLOGY | 4. IMMUNOLOGY |
| 2. BIOCHEMISTRY | 5. NUTRITION |
| 3. PHARMACOLOGY | 6. PATHOLOGY |

DO NOT SKIP A SPACE. IN THE NEXT COLUMNS ENTER ONE INDEXING TERM USED BY THAT INDEXER FOR THAT DOCUMENT. THIS WILL BE EITHER A THREE DIGIT NUMBER WHICH HAS BEEN CHECKED BY THE INDEXER ON THE SHEET OR A TERM WRITTEN IN BY THE INDEXER ON THE SPACE PROVIDED AT THE END OF THE SHEET. IF THE INDEXING TERM IS A NUMBER THE FINISHED CARD WILL CONTAIN NINE DIGITS WITH NO SPACES BETWEEN THEM. IF THE INDEXING TERM IS WRITTEN, THEN YOU MAY USE AS MUCH OF THE CARD AS NECESSARY TO RECORD IT. IF THE TERM IS FROM THOSE WRITTEN IN AND YOU HAVE TROUBLE READING THE HANDWRITING OR ARE NOT SURE OF THE SPELLING, READ THE DOCUMENT AND MOST OFTEN THE TERM IN QUESTION WILL APPEAR THERE. IF YOU CANNOT FIND THE TERM IN THE DOCUMENT, AND CANNOT DECIPHER THE HANDWRITING, KEYPUNCH YOUR BEST GUESS AND SET THE CARD ASIDE TO BE CHECKED. IN SOME CASES A WRITTEN-IN TERM WILL CONSIST OF MORE THAN ONE WORD. IF THIS IS THE CASE, LEAVE ONE SPACE BETWEEN WORDS. DO NOT INCLUDE ANY PUNCTUATION, E.G., COMMAS, PARENTHESES. DO INCLUDE HYPHENS. IF THE TERM IS A CHEMICAL FORMULA WITH SUBSCRIPTS, THEN COPY IT AS IF IT WERE ALL ON ONE LINE, E.G., CO₂ BECOMES CO2. REPEAT THE SIX DIGIT DOCUMENT-INDEXER-AUTHOR NUMBER AT THE BEGINNING OF THE NEXT CARD AND, FOLLOWING THE ABOVE FORMAT, RECORD THE NEXT INDEX TERM. GREEK LETTERS SUCH AS α , CAN BE PUNCHED A FOLLOWED BY A HYPHEN. BETA CAN BE PUNCHED B FOLLOWED BY A HYPHEN. DELTA CAN BE PUNCHED AS A D FOLLOWED BY A HYPHEN. GAMMA IS WRITTEN OUT (GAMMA) FOLLOWED BY A HYPHEN.

IF A TERM WILL EXCEED COLUMN 78 LOOK FOR A CONNECTIVE EARLIER IN THE TERM WHERE IT COULD BE BROKEN INTO TWO TERMS.

Ex: AMINO ACID METABOLISM NUTRITION/OF IMMATURE
EMBRYONIC CHICKS FED METHIONINE

(TOP)

ERIC REPORT RESUME

ERIC ACCESSION NO.

CLEARINGHOUSE
ACCESSION NUMBER

RESUME DATE

P.A.

T.A.

IS DOCUMENT COPYRIGHTED?

YES ☐NO ☐

ERIC REPRODUCTION RELEASE?

YES ☒NO ☐

TITLE

EVALUATION OF INDEXING BY GROUP CONSENSUS

FINAL REPORT

PERSONAL AUTHOR(S)

SCHULTZ, CLAIRE K. AND OTHERS

INSTITUTION (SOURCE)

INSTITUTE FOR THE ADVANCEMENT OF MEDICAL COMMUNICATION, PHILA., PA.

SOURCE CODE

REPORT/SERIES NO.

OTHER SOURCE

SOURCE CODE

OTHER REPORT NO.

OTHER SOURCE

SOURCE CODE

OTHER REPORT NO.

PUB'L. DATE AUGUST -30, 1968 CONTRACT/GRANT NUMBER OEC 1-7-070622-3890

PAGINATION, ETC.

33 PAGES

RETRIEVAL TERMS

INDEXING

METHODS

INDEXERS

RESEARCH

INDEXES

CRITERIA

EVALUATION

IDENTIFIERS

CRITERION GROUP METHOD

ABSTRACT

THIS METHOD TESTS THE EFFECTIVENESS OF TEST INDEXING SETS, USING A CRITERION GROUP TO SET THE STANDARD FOR "IDEAL" INDEXING. THE CRITERION GROUP FOR A PARTICULAR APPLICATION IS CHOSEN BY THE TEST ADMINISTRATOR, CONSISTENT WITH HIS OWN CONCEPT OF WHO REPRESENTS THIS "IDEAL". MATCHING TEST SETS OF INDEXING TERMS WITH THE CRITERION SET YIELDS AS MANY DEGREES OF MATCH AS THERE ARE MEMBERS OF THE CRITERION GROUP (REFERRED TO AS CONSENSUS NUMBER). IMPORTANT VARIABLES FOR THE METHOD ARE: SIZE OF DOCUMENT SAMPLE, SIZE OF CRITERION GROUP, INDEXERS' INSTRUCTIONS, METHOD OF EDITING RAW INDICIA TO MAKE THEM COMPARABLE, AND METHOD OF WEIGHTING TERM SETS FOR SCORING.

RESULTS OF TESTING THE METHODOLOGIC VARIABLES FOR THEIR EFFECTS ON RELIABILITY, SENSITIVITY, FLEXIBILITY AND PRACTICALITY OF THE METHOD SHOW THAT "INDICATIVE TESTS CAN BE MADE AT THE 80% LEVEL OF CONFIDENCE WITH DOCUMENT SAMPLES AS SMALL AS 10 AND CRITERION GROUPS AS SMALL AS 4; 95% CONFIDENCE REQUIRED COMPARABLE VALUES AS LARGE AS 20 DOCUMENTS AND 9 CRITERION GROUP MEMBERS. THE 3 EDITING METHODS TESTED: "NONE", MANUAL, AND COMPUTER, YIELDED DIFFERENT SCORES, BUT EACH PRESERVED DIFFERENCES BETWEEN TEST SETS, SO WAS NOT IMPORTANT TO SENSITIVITY OR RELIABILITY, ONLY PRACTICALITY. SCORES CHANGED WITH DIFFERENCES IN INDEXER INSTRUCTION OR ADDITION OF A VOCABULARY GUIDE, SO THE METHOD IS SENSITIVE TO SUCH DIFFERENCE BETWEEN TESTS BUT CAN BE CARRIED OUT SUCCESSFULLY WITH ESSENTIALLY NO INDEXER INSTRUCTION AND NO GUIDE. CONSENSUS NUMBER WAS ALMOST AS USEFUL AS CONSENSUS NUMBER SQUARED FOR DETECTING DIFFERENCE IN TEST SETS; BUT THE LATTER WEIGHTING EMPHASIZED "RECALL" VALUE TO SOME EXTENT.