

Compatibility of Indexing Languages in an Online Access Environment: A Review of the Approaches¹

Marcia Lei Zeng

School of Library and Information Science

Kent State University

P.O. Box 5190

Kent, OH 44242-0001

1. INTRODUCTION

With the increase of cross-database searching and inter-system operation and the involvement of patron searching, compatibility of indexing languages is becoming an issue of great practical impact and is experiencing renewed interest. An indication of this increasing concern is the current projects being carried out by the national libraries and other national and international organizations such as the Unified Medical Language System (NLM), Universal Agricultural Thesaurus (NAL, CABI, and NAO), and Integrated Multilingual Thesaurus for Social Science (Unesco).

The controlled vocabularies used for subject indexing are as many and as varied as the materials they index and the audiences they serve. Most thesauri have been developed independently for specialized catalog, index, or information services since the early 1960s, when differences among their vocabularies did not appear to pose a problem for their users. For today's world of online access, however, the multiplicity of thesauri creates complications significantly affecting at least two types of information services. In the first instance, it affects online bibliography services such as DIALOG which offers hundreds diverse databases through a single interface. Second, it affects integrated library systems such as NOTIS which provides subject authority control for more than one logically separated subject vocabulary. In an automated environment, the absence of compatibility can add substantially to the costs of time, intelligence, and money, while it can bring unsatisfied search results to end users.

As early as in 1964, the Sixth Institute on Information Storage and Retrieval, held by The American University, had the theme "Systems Compatibility for Scientific and Technical Information". Compatibility, Convertibility, Cooperation and Standardization constituted the principal concepts around which the institute was planned (Newman 1965). In 1971, Unesco published its *UNISIST Study Report on the Feasibility of a World Science Information System*. One of the major recommendations pointed out that "The attention of scientists, learned societies and information science associations should be drawn to the need for joint efforts in developing better tools for the control and conversion of natural and indexing languages in science and technology" (Unesco 1971, 97). Several international conferences have since focused discussion on compatibility concerns since then including "Ordering Systems for Global Information Networks" (1975 The Third International Study Conference on Classification Research), "Overcoming the Language

1. This paper and the presentation are based on the article "Achieving Compatibility of Indexing Languages," which is published in A. Kent, ed., *Encyclopedia of Library and Information Science* Vol. 50, 1992, with some updating and revising.

Barrier" (1977 at Luxembourg); "Unified System of Information Retrieval Languages" (1977 at Latvia); and "Compatibility between Indexing and Retrieval Languages" (1982 at Ohio) (Svenonius 1983; Smith 1991).

Projects for the standardization of thesaurus construction led by international and national standard agencies¹ have played an important role in promoting compatibility of indexing languages. However, compatibility no longer always depends on standardized thesauri, because many more progressive organizations had developed their own standards before international and national standards were developed. Studies on achieving compatibility based on existing indexing languages have contributed significantly to positive trends that affect our online services today.

Many possible approaches to securing compatibility of indexing languages have been studied (please refer to Appendix 1. Chronology of Research Projects). In a report prepared for Library of Congress Processing Services on multiple thesauri in online library bibliographic systems, i.e., Mandel; (1987) categorized four basic approaches to providing access to databases indexed by different vocabularies: a) *segregated files* — databases using different thesauri are searched separately; b) *mixed vocabularies* — terms from all vocabularies are retrieved together in subject searches; c) *integrated vocabularies* — techniques used to relate different thesauri can be used to develop syndetic structures that would aid retrieval in a multiple database environment; and d) *front-end navigation* — features being developed in "intelligent" interfaces to online retrieval can be designed specifically to aid in searching from multiple databases. In an earlier report for Unesco's General Information Program, Lancaster and Smith; (1983) grouped approaches to securing vocabulary compatibility into five categories: mapping, intermediate lexicon (i.e., switching language), integrated vocabulary, microthesauri, and macrovocabularies.

Approaches reviewed in this article will be based on the above two important reports, however, the article would suggest that recently it has become possible to draw a distinction in terms of focus between two groups of approaches: mechanisms; that enable the establishment of compatible vocabularies (to be discussed in Section 3) and automated systems; that help to relate separately created thesauri to one another (to be discussed in Section 4).

2. DEFINITIONS AND EXPLANATIONS

The subject indexing process; involves two quite distinct intellectual steps: the "conceptual analysis" of a document and the "translation" of the conceptual analysis into a particular vocabulary or indexing language, "a controlled set of terms selected from natural language and used to represent, in summary form, the subject of documents" (ISO 5964-1985). Such a vocabulary might be a list of subject headings, a classification scheme, a thesaurus, or simply a list of "approved" key words or phrases (Lancaster and Smith 1983, 6).

In the 1971 Unesco UNISIST Study Report, compatibility was defined as "a quality of systems whose products can be used interchangeably, notwithstanding differences in notation, structure,

1. This includes the publication of ISO 2788, ISO 5964, etc. For details see reference: Lancaster 1986, pp. 29-33.

vocabularies are likely to have great differences from collection- or role-oriented ones in their vocabulary size, specificity, and subject coverage. 3. *Initial purpose of vocabulary usage.* Thesauri are mainly developed for mechanized retrieval, while a few might take manual retrieval into account. Thus leads to differences in the percentages of non-single-terms.

Vocabulary aspects. Vocabulary aspects are also very influential on the conversion among thesauri, but they are flexible. Vocabulary size, specificity, entry format control, number of entry points, and precision of expression are variable among thesauri and therefore make the results of conversion between any two of the thesauri variable.

Structure aspects. Structure aspects have an indirect influence on conversion. Most thesauri use extra structure to display terms and term relationships in addition to an alphabetical list, such as hierarchical display, subject category display, graphical display, etc. It is believed that the more a thesaurus displays its vocabulary, the more helpful to conversion.

Other aspects which influence compatibility more or less come from characteristics of language. For instance, multilingual conversion brings many special problems. Other influences may come from the frequency of term occurrence in the indexing and searching process; the policy of index term assignment; and the differences of the presentation of the same concept in a thesaurus and in the indexing and searching processes. These aspects are very changeable and are influenced by some subjective factors.

3. MECHANISMS FOR THE ESTABLISHMENT OF COMPATIBLE VOCABULARIES

When integrated retrieval systems became possible in the 1970s, there were few technical problems to merging databases which had been indexed with different thesauri. However, it was soon found that these files could not be searched by an identical strategy because a single concept might be quite differently represented in the different vocabularies. The conventional way of solving this problem is to develop mechanisms; which would enable the establishment of compatible vocabularies. Integrated vocabularies and metathesauri, intermediate lexicon, microthesauri, and macrovocabularies have been significant contributions. Other examples also exist but these at least illustrate the situation.

3.1. Integrated Vocabularies and Metathesauri

In a multi-thesauri environment, an alternative is to provide user with a composite or integrated vocabulary. An initial and fairly cursory way of this approach is to construct a master list of all terms drawn from all of the source thesauri. Recently a word "metathesaurus" is widely used. Rada (1990) briefly stated that "a metathesaurus transcends a set of thesauri." Roulin (1991) defined "a metathesaurus is a set of reference thesaurus and various indexing languages that are attached to it either via common concepts hierarchically structured in the same way (in the case of 'sub-thesauri' and 'associated thesauri') or only via relationships established between the indexing terms that remain proper to each indexing language (in the case of 'national vocabularies')". An integrated vocabulary or a metathesaurus may exist in a form other than a traditional thesaurus.

physical carriers, etc., without any special 'conversion machinery'." Two related terms were defined as well. **Convertibility** was defined as "a quality of systems whose products can be made interchangeable through 'conversion' programmes." **Conversion** was defined as "the process of transforming information records, with regard to transcription encoding, data structure, etc., so as to make them interchangeable between two or more services or systems using different conventions and media" (Unesco 1971; 147). The implication of these terms for indexing languages is demonstrated by Dahlberg's (1983) definition that **compatibility of an ordering system** is the quality which permits the elements of one such a system to be used together interchangeably with the elements of another ordering system.

From a formal point of view, a concept, as an element of an indexing language, consists of three parts: *the referent*, *the characteristics* predicated of a referent, and the *verbal or coded expression* denoting the referent and comprising the predicated characteristics in its designation in form of a term, a descriptor, or a notation.

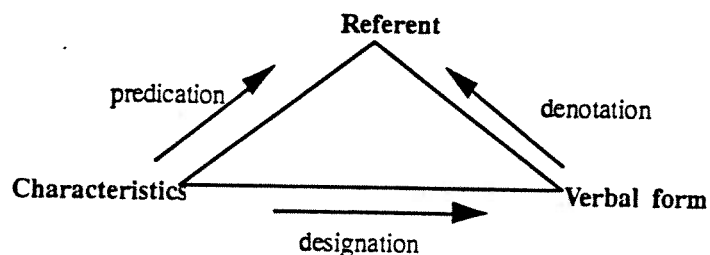


Figure 1. Dahlberg' Concept Triangle. --Source: Dahlberg 1983.

From this perspective, the **conceptual compatibility** between elements of controlled vocabularies can be compared in relation to three measures of strength of correlation. Listed from the strongest type of correlation to the weakest these include: a) *conceptual coincidence* — two concepts match in all their characteristics; they can also be called "equivalent"; b) *conceptual correspondence* — two concepts match in many of their characteristics; also called "similar concepts"; c) *conceptual correlation* — two concepts are set into "correlation" where the kind of correlation can be indicated, e.g. by mathematical symbols (Dahlberg 1983).

Many factors contribute to the compatibility of thesauri. The extent of overlap in the subject matter, specificity, and vocabulary size are commonly believed to play important roles (Lancaster 1986, Dahlberg 1981). To summarize the factors that influence the compatibility of thesauri, we may deal with the factors from the following aspects (Zeng 1990):

Principle aspects. Principles, which have the most direct influence on the compatibility, are decided at the first step of thesaurus design. They may include at least three parts, that is: 1. *Initial term-structure principle.* A term-structure based on either pre-coordination or post-coordination was chosen at the beginning. Meanwhile, the percentage of non-single-terms in a post-coordinated vocabulary was also decided. 2. *Original orientation of vocabulary design.* Vocabularies might be designed for certain collections, or for particular projects or roles. The discipline-oriented

PROCEEDINGS OF THE 3rd ASIS SIG/CR CLASSIFICATION RESEARCH WORKSHOP

An example of integrated vocabulary was the "Vocabulary Guide" constructed by the Women's Educational Equity Communications Network (WEECN). It is lessentially an alphabetic display, in a single sequence, of all terms related to WEECN's interests that appear in each of five databases (Lancaster 1986, 194-195).

Figure 2. Example from WEECN VOCABULARY GUIDE

FEAR (MeSH)	
RT	CONFLICT (ERIC)
	CONFLICT (PA)
	CONFLICT (SSIE)
	CONFLICT (TEST)
	FEAR (ERIC)
	FEAR (SSIE)
	FEAR (TEST)
	PSYCHOLOGICAL PATTERNS (ERIC)
	STRESS(PSYCHOLOGY) (TEST)
	STRESS-BEHAVIORAL ASPECTS (SSIE)

Explanation: For each entry, it shows synonyms, near-synonyms, and otherwise-related terms from all vocabularies of five databases: Mesh (Medical Subject Headings), PA (Psychological Abstracts), SSIE (Smithsonian Science Information Exchange), and TEST (Thesaurus of Engineering and Scientific Terms) -- Source: Lancaster 1986.

BRS' TERM database, though not named as an integrated thesaurus, functioned as such a tool in which terms of different thesauri are organized into concept records. It also included hierarchical information and free-text searching suggestions (Schwartz and Eisenmann 1986).

A similar approach used by integrated library systems is to use an automated system for subject authority control and for mixed vocabulary searching. This will be discussed in Section 4 of this article. The difference between integrated vocabulary and on-line mixed vocabulary identified by this author is that in the latter the equivalencies among vocabularies (synonymy and near-synonymy) are not established.

A *Unified Medical Language System (UMLS)* supported by the National Library of Medicine (NLM) was reported in 1987 (Humphreys, 1988). Experimental linkings of *Medical Subject Headings (MeSH)* with the *SNOMED (Systematized Nomenclature of Medicine)*, *CMIT (Current Medical Information and Terminology)*, and *PDQ (the National Cancer Institute's retrieval system)* thesauri were conducted. A metathesaurus has been used to store medical concepts and terms in a canonical form to which multiple existing vocabularies and classification systems would be mapped. Each Meta-1 record contains three types of information: basic facts, relationships, and usage data (Rada 1987; 1990). However, the goal of *UMLS* goes beyond a single medical vocabulary, it attempts "to make the myriad of classifications of medical knowledge invisible to the user while providing a single logical path to a broad range of biomedical information sources" (*NLM News* 1987, 4-6). In Europe, Directorate XIII-B3 of the European Economic Community has outlined a five-step procedure to create a metathesaurus which the commission calls the "Multi-function, Multilingual Thesaurus Database". A list of about 1,000 thesauri in prominent use in the world has already been compiled. Connecting the thesauri that have been cataloged by DG XIII-

Figure 3. A BRS/TERM Record
for the Concept of "Depression" in its Psychological Sense.

TI	DEPRESSION (PSYCHOLOGY)
ER	DEPRESSION-PSYCHOLOGY (1978+).
ME	DEPRESSION (OF MODERATE INTENSITY). CONSIDER ALSO: ANTIDEPRESSIVE-AGENTS. DEPRESSIVE-DISORDER+ (FOR PROMINENT OR PERSISTENT DEPRESSION).
NM	DEPRESSION.
PS	DEPRESSION-EMOTION+.
SO	CONSIDER:DEPRESSIVE (129010).
MN	NEURASTHENIA. DEPRESSION-INVOLUTIONAL.
PN	ANACLITIC-DEPRESSION. ENDOGENOUS-DEPRESSION. NEUROTIC-DEPRESSIVE-REACTION. POSTPARTUM-DEPRESSIVE. PSYCHOTIC-DEPRESSIVE-REACTION. REACTIVE-DEPRESSION.
FT	AGITATED DEPRESSION. DYSTHYMIA. MELANCHOLIA. DEPRESSION. DEPRESSIVE. DEPRESSED. DESPONDENT. WEEPINESS. DESPONDENCY. DEJECTION. GLOOMY. SADNESS. DESPAIR. DISCOURAGED. LOWERED MOOD TONE. DISCOURAGEMENT. DEPRESSIVE SYMPTOMS. MODERATE DEPRESSION. ANTIDEPRESSANTS. DEPRESSANTS. DEPRESSING. DEFEATIST. HOPELESS. WITHOUT FAITH. WITHOUT HOPE. CONSIDER ALSO: SELF DENIAL. SELF ABNEGATION. SELF NEGLECT SELF ABUSE. WITHDRAWAL. SELF DEPRECIATION. FEELINGS OF WORTHLESSNESS. LEARNED HELPLESSNESS. LONELINESS. GUILT. LISTLESSNESS. UNHAPPINESS. SELF DEPRECIATION. SELF ACCUSATION. HYPOCHONDRIA. SELF DESTRUCTIVE TENDENCIES. STUPOR.

Explanation: Controlled terms are shown from the thesauri used by the ERIC (ER), MEDLINE (ME), NIMH (NM), PSYCINFO (PS), AND SOCIOLOGICAL ABSTRACTS (SO) databases. Narrower terms are also shown for MEDLINE (MN) and PSYCINFO (PS). Terms or phrases are suggested for free-text searching of this concept (FT). -- Source: Piternich 1990, 414.

PROCEEDINGS OF THE 3rd ASIS SIG/CR CLASSIFICATION RESEARCH WORKSHOP

B3 and allowing easy access by users could be an enormous benefit to those who want to share terminology (Rada 1990). In China, a *Chinese Classified-Thesaurus* is nearing its completion. This project, sponsored by the Peking Library, the national library of China, aims at an integrated vocabulary based on the national standard classification *Chinese Library Classification (CLC)* and the recommended national standard thesaurus *Chinese Thesaurus (CT)*. The new vocabulary provides for each of the classes or divisions from CLC corresponding descriptor(s) from CT, and vice versa (Zeng 1991).

One of the obvious problems of this kind of vocabulary is the updating of the vocabulary according to any changes made by any of the source vocabularies. In addition, in a bibliographic retrieval system such as BRS, the file carrying the integrated vocabulary is separate from the document files, a user have to go back to the desired file and re-type all terms found in the integrated vocabulary. Nevertheless, the creation of an integrated vocabulary needs a lot of intelligent work on mapping terms of various thesauri to one another. As pointed by Lancaster (1986), mapping two vocabularies on to each other, requires only one mapping operation (as shown in Situation I in Fig.4). However, when more and more vocabularies are to be mapped, the situation becomes increasingly complex. For example, a four-vocabulary-mapping requires twelve separate mapping operations (see Situation II in Fig.4). An alternative is to construct an intermediate lexicon, or, a "switching" language that can be used to convert from any one vocabulary to another (see situation III in Fig.4, where A-F are vocabularies to be converted, and X is an intermediate lexicon).

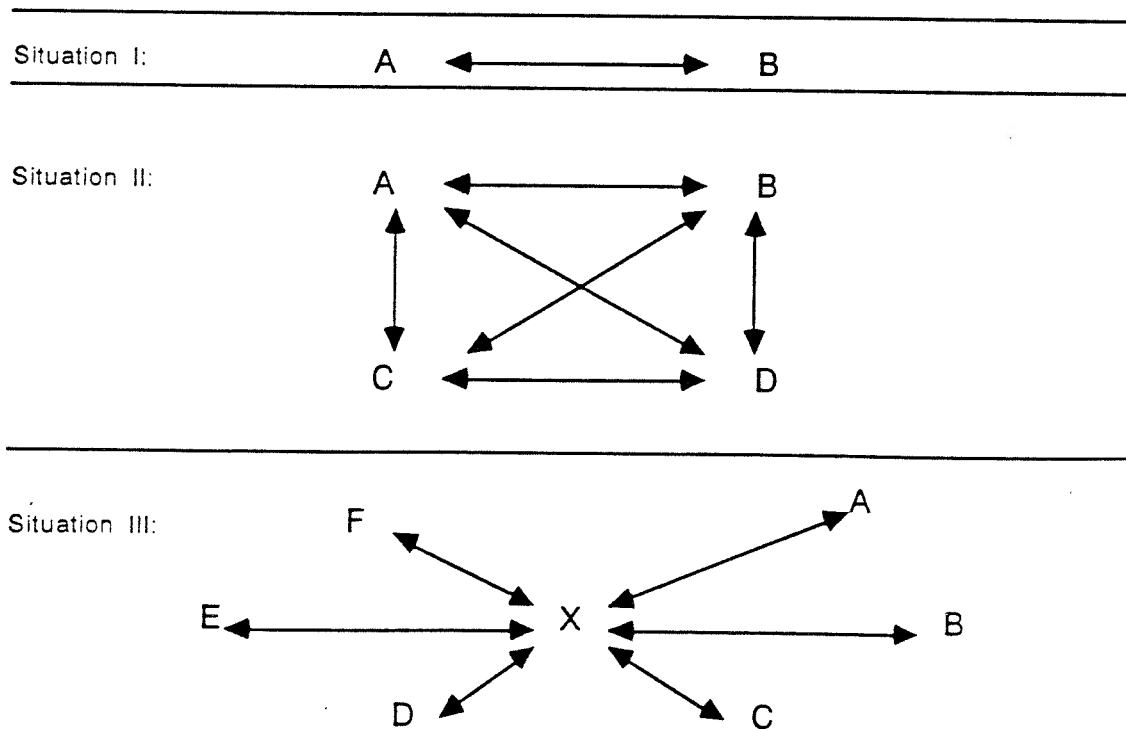


Fig 4. Possible Mapping Operation (adopted from Lancaster's *Vocabulary Control for Information Retrieval* pp. 182, 183 and 189).

3.2 Intermediate Lexicon/Switching Language

An intermediate lexicon is intended to transfer the contents of documents expressed in terms of any index language to another without loss of information. It entails the mapping of two or more vocabularies to an intermediate or neutral language, such as a classification or a coding scheme. This idea was developed by J. C. Gardin and his group from 1967 to 1968 in France. A matrix had been worked out previously and the elements of six thesauri available to the field of information science were correlated by the matrix (Dahlberg 1981).

Though a number of projects for such switching systems were reported worldwide, there is no such tool has been fully implemented. The reason might be that the development of such a switching tool (usually a coding system) requires a lot of intelligent work, while its effectiveness is limited in the environment that all participating vocabularies should have similar or closely related subject coverage.

It is worth noting that, instead of working on an intermediate lexicon, the British Library investigated the switching of PREserved Context Index System (PRECIS) input strings from one language to another and used PRECIS in the British National Bibliography for several years. The PRECIS system expresses concepts in terms selected from natural languages encountered in the literature. A major feature is that each term is placed in a context-dependent order in such a way that each term sets the next term into its obvious context. Syndetic control is maintained through use of a thesaurus (DeHart & Glazier 1984). PRECIS, as an indexing system, was suggested as a switching mechanism because it could be adapted for various European languages (Austin 1977).

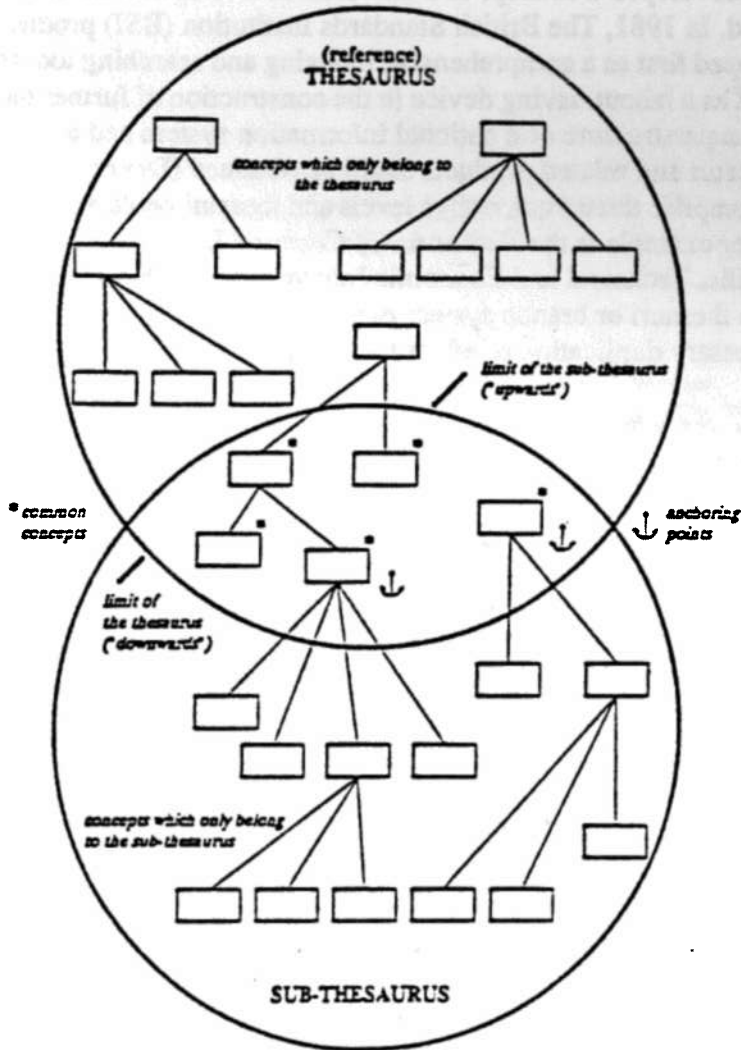
It is noted by this researcher that a visualized thesaurus could perform function of an intermediate lexicon/switching language of subject terms in a number of domains, especially in multilingual thesauri environment. The weakness of such a tool is that it will have limitation in terms of numbers of concepts and terms which can be represented in a visual thesaurus. The advantage of it is that it will not need the construction of an additional coding system. Meanwhile, it is possible to use existing visual thesauri as the basis of such a tool so that even less work will be needed.

3.3 Microthesauri

This approach treats specialized thesauri as the satellites of a superstructure. In its original application, a microthesaurus was a specialized subset of terms extracted from, and therefore compatible with, a larger thesaurus. A microthesaurus can therefore be defined as a specialized vocabulary that maps onto a broader thesaurus and is entirely included within the hierarchical structure of that thesaurus (Lancaster 1986, 198-199). This approach ensures the compatibility between an existing thesaurus or a superstructure and its satellites, and meets the need for deeper degree of specificity for refined indexing in certain subject fields. As described by Roulin (1990), the superstructure (reference thesaurus) and its satellites (sub-thesaurus) will be two non-disjointed, overlapping sets of concepts: there will be concepts that only belong either to the former or to the latter as well as concepts that are part of the intersection of both sets, since the most specific concepts in the reference thesaurus are the anchoring points of parts of hierarchical chains peculiar to the sub-thesaurus.

PROCEEDINGS OF THE 3rd ASIS SIG/CR CLASSIFICATION RESEARCH WORKSHOP

Figure 5. Thesaurus and sub-thesaurus



Example

(The asterisked descriptors belong to the common portion, the descriptors marked with an anchor are "anchoring points" and the descriptors in bold characters belong to the sub-thesaurus only.)

- COMPUTER SCIENCE *
- ARTIFICIAL INTELLIGENCE*
- COMPUTER APPLICATION* |
- COMPUTER ASSISTED DESIGN
- CONTROL TECHNOLOGY
- ROBOTICS
- COMPUTER ENGINEERING*
- SOFTWARE* |
- AUTHORING SYSTEM
- COMPUTER GAME*
- COMPUTER GRAPHICS
- DATA BASE MANAGEMENT SYSTEM
- EDUCATIONAL SOFTWARE*
- EXPERT SYSTEM
- OPERATING SYSTEM
- SPREAD SHEET

--Source: Roulin 1990, 34-35.

The basic element of microthesauri approach is a superstructure. In the early 1970s, Soergel (1972, 1974) developed a concept of a universal source thesaurus, from which other thesauri could be derived. In 1981, The British Standards Institution (BSI) produced its *ROOT Thesaurus*. It was developed first as a comprehensive indexing and searching tool for technological applications and second as a labour-saving device in the construction of further thesauri. It conceivably could serve as the superstructure of a national information system and as the basis from which a whole series of thesauri and related products could be obtained (Dextre and Clarke 1981). The entire system may comprise three, four, or five levels and thesauri could exist at each level (Lancaster 1986, 201). Another example is the *Thesaurus of Common Topics (TCT)*, developed at the Institute of Scientific, Technical and Economic Information in Warsaw. It was intended to provide particular branch thesauri or branch system of thesauri with some ready sets of terms, and thus eliminate unnecessary duplication of effort (Scibor and Jabrzemska 1989).

In some respects, large comprehensive vocabularies developed by national libraries and information centers are likely to be treated as superstructures. The Chinese Documentation Standardization Committee suggested in 1979 that all special thesauri developed after 1980 should consider to be the microthesauri of the *Chinese Thesaurus*, a universal vocabulary developed by the China Science and Technology Information Institute and the National Library of China. Most of the Chinese special thesauri developed since then have used the *Chinese Thesauri* as their basic term source and structure model (Zeng 1990). In America, LCSH has been acting as superstructure as well. The *Legislative Indexing Vocabulary (LIV)*, *LC Thesaurus of Graphic Materials*, *Subject Headings for Children's Literature*, and the newly developed *Art and Architecture Thesaurus (AAT)* are examples of microthesauri which were developed as satellites of *LCSH* in order to meet the specialized needs of particular subject fields, materials, and patrons. *LCSH* compatibility codes have been used in these microthesauri (Mandel 1987, 4; 80-92). For instance, each *LIV* term is assigned an *LCSH* compatibility code as follows:

- LC — term is the same in *LCSH*
- LCX — term is a "see" reference in *LCSH*
- LCC — term is similar to one in *LCSH*
- LCD — characters match an *LCSH* term, but the meaning is different
- LCO — no match in *LCSH*.

The problem is that most existing indexing languages were developed completely independently. Although these thesauri may conform to common standards or guidelines in their structure, they have not been integrated or interconnected with a common superstructure. On the other hand, few thesauri could be expected to perform such a complex function of a superstructure as the *LCSH* does.

3.4. Macrovocabulary

The macrovocabulary approach, as described by Lancaster, is conceptually similar to that of the microthesaurus, but the method of implementation is virtually reversed. The idea is simply to create a kind of generic superstructure of terms that will subsume a group of existing thesauri, or other types of vocabularies, in diverse subject fields so as to link terms in specialized vocabularies (Lancaster 1986, 202-203).

PROCEEDINGS OF THE 3rd ASIS SIG/CR CLASSIFICATION RESEARCH WORKSHOP

Figure 6. CORRELATION OF BSO SOCIAL SCIENCE FIELDS WITH EQUIVALENTS IN OTHER SYSTEMS

BSO	DDC	BBC
530/588 <i>Social Sciences</i>	300 <i>Social Sciences</i>	K <i>Social Sciences</i>
533 Cultural Anthropology	306 Cultural & institutions	KC Cultural, social anthrop.
535 Sociology	30 1/7 Sociology	KA Sociology
537 Demography	304.6 Population (Demogr.)	-
540 <u>Political sci., polit.</u>	320 Political science	R Political science
542 Political inst. & org.	306.2 Political institutions	-
543 Pol. org. patterns &	321 Kinds of governments	RJ Political systems
544 Political history	320.9 Hist. of pol. inst.	RJS His. by subj., political
545 Pol. of part.groupings	-	-
546 Pol. of part.states &	-	RK Pol. by place
550 <u>Public administration</u>	350-359 Public administr.	RO Public administration
554 Central admin. & gov.	351 Central governments	RR " " , Central
556 Devolved admin.	352 Local governments	RS " , Regional & Local
560 <u>Law. Juridical sciences</u>	340 Law	S Law
562 Civil Law	346 Private Law	SBF Private & public law
563 Public Law, Const. etc.	343 Public l., Const., Crim.	SBF Private & public law
565 International Law	341 International Law	SDD International Law
567 Systems of Law	-	SE Common Law systems
568 Law of partic. countr.	349 Law of indiv. states &	-
570/575 <u>Social welfare</u>	361 Soc.probl. & soc.welf.	Q Social welfare
580 <u>Economics</u>	330 Economics	T Economics
581, 89 Microeconomics	338.5 Gen.Product.econ.	-
582 Macroeconomics	339 Macroecon. & rel. t.	-
584 Economic organization	330.1 Econ.,syst. & theor.	-
586 Sectorial economics	-	TT Economic syst., sect.
588 Management of enterpr.	658.2 Management of plants	TX Management of enterpr.

Source: Dahlberg 1980.

Historically, there are a number of examples of such efforts. Started in the late 1950s, the British Classification Research Group (CRG) pursued the idea of a sort of "Ur-Classification," a generalized scheme of terms not necessarily related to any one practical application in a library or documentation center. The origins of all these ideas were first announced to the world at the First International Study Conference on Classification for Information Retrieval in 1957. CRG's experience from the 1990 revision of Class J Education of the *Bliss Bibliographic Classification (BC)* suggested that a good general scheme could be compiled by integrating specialist schemes. A general, or "Ur-Classification" will provide a reservoir of terms for specialist schemes, while the special schemes provide detailed analysis and enumeration by experts in each field (Foskett 1991).

In 1978, the UNISIST program developed an element of *The Broad System of Ordering (BSO)*. It is in the form of a classification scheme, with notations, containing 4,000 terms, and is conceived as a tool to allow the interconnection of information systems, services, and centers using diverse index languages. In fact, it exists as a general superstructure rather than the more specific switching approach embodied in the intermediate lexicon (Dahlberg 1980; Lancaster 1986, 203-205).

Sager et al. (1982) and Whitelock (1982) described developments in another macrovocabulary project — a Unesco-funded project to create an *Integrated Multilingual Thesaurus* for social science. The Descriptor Bank, which lay at the heart of this project, converted a number of thesauri to a common form, included terminological and syndetic information, and resolved differences in preferred term choice, definition, hierarchy, precoordination and compounding, and language. In the preliminary study for this macrothesaurus, Meyriat (1980) analyzed 60 existing information languages in order to identify for each the scope and depth of coverage of 41 subject fields.

A feasibility study for developing a *Universal Agricultural Thesaurus* unifying CABI, NAL and FAO thesauri was initiated in 1989. C.A.B. International (CABI), National Agricultural Library (NAL) and Food and Agriculture Organization of the United Nations (FAO) have been involved in this effort which will improve access to agricultural information in a more cost effective way. It was stated that "without such a tool, individual institutions would continue the costly and duplicative activities associated with local subject thesaurus development and validation of taxonomic names, and researchers would continue to be uncertain as to the completeness and accuracy of their research documentation" (Andre 1989). Additionally, a number of projects for macrovocabularies were conducted in the United States, Soviet Union, and France. Their purposes included using the superstructure to suggest useful search terms to retrieve information on particular subject from a variety of databases, and linking various specialized information centers in order to achieve some measure of compatibility among the vocabularies of these centers (please refer to Appendix 1).

4. DEVELOPING AUTOMATED SYSTEMS FOR THE CREATION OF LINKAGES AMONG VOCABULARIES

Many automated systems; have been designed to aid the process of integrated subject retrieval. Most of these systems reflect an effort to map separately maintained vocabularies. These approaches tend to emphasize the phase of subject searching instead of the phase of thesaurus construction, and tend to be more dependent upon the state-of-the-art of computer support.

PROCEEDINGS OF THE 3rd ASIS SIG/CR CLASSIFICATION RESEARCH WORKSHOP

4.1. Automatic Mapping

Mapping entails the direct translation of terms in one vocabulary into corresponding terms in another. Machine support is particularly valuable for mapping terms and merging term lists, with algorithms developed for matching references, component words, stems, Boolean combinations, etc. Concepts can be mapped between two thesauri with a variety of tools. First, direct lexical matching between concept main terms can be performed; second, knowledge about the syntax or the morpho-semantics of main terms can be employed; third, the knowledge in the relationships within the thesauri themselves can be the basis for sophisticated mapping of terms from one thesaurus to another (Rada 1987). When vocabularies exist in machine-readable form, part of the conversion such as direct lexical matching can be done with common computer software tools. Automatic identification includes exact match, some variant spellings, some variations in word forms and inversion, and also, under certain conditions, it involves mapping on the basis of the cross-references and hierarchical structure of the vocabularies.

It has been convincingly demonstrated that the more alike the two vocabularies are in structure, and the higher level of subject overlap, the more conversion can be performed automatically. Wall and Barnes; (1969) found that they could automatically map 76% of a sample of *Medical Subject Headings (MeSH)* to the *Agricultural/Biological Vocabulary* because both of the vocabularies were pre-coordinated subject heading lists and had high degree of subject overlap. Another study was an experiment of convertibility between *ASTIA Subject Headings* and *AEC Subject Headings* by Hammond and Rosenberg (1962). Although the two vocabularies differed with regard to their term structure principles, they were highly compatible because of their high degree of agreement in subject coverage. In a project of the European Economic Community to connect patient records with the medical literature, *International Classification of Diseases (ICD)* of the World Health Organization, and the *Systematized Nomenclature of Medicine (SNOMED)*, a thesaurus for patient records, were mapped into *EMTREE*, the thesaurus used for indexing *Excerpta Medica*. Over 50% of *SNOMED* terms were directly mapped to equivalent *EMTREE* terms with simple matching rules. Mapping between *ICD* and *EMTREE* was performed through facets (Rada 1990).

A related possibility is that some terms which are identified by a machine as the "equivalent" might have different meaning when they are used in the different subject fields, such as "pressure" in psychology and in mechanical engineering. As a result, mapping between two vocabularies which are composed of different subject topics may lead to a more or less confusing situation. Therefore, automatic mapping is encouraged for use only when two vocabularies have a similar subject coverage.

4.2. Segregated File Searching Aids

Online vendors have made various attempts to provide users with term search or database selection aids. In a system maintaining logically separate subject indexes and subject authority files for each vocabulary, instructions usually are provided so that patrons can use different commands when searching files which were indexed by different vocabularies. As early as in the 1970s, the University of Illinois began to design a more user-oriented "transparent system," i.e., a system "containing the necessary converters or translators to help the user circumvent the need for understanding all the specific differences of databases, systems, command languages, vocabularies and access protocols" (Williams 1978, 361). In such a system, the user is not distracted by many

Figure 7. Sample Output Listing for VSS Subject Switching, for "Heavy Water" --Source: Chamis 1991, p.60

```
WELCOME TO "VSS" - VOCABULARY SWITCHING SYSTEM USING
BATTELLE'S DATA MANAGEMENT SYSTEM, BASIS
VSS CONTAINS FOUR VOCABULARY SETS
1- BUSINESS          A. ABI      B. MANAGEMENT CONTENTS
2- SOCIAL SCIENCE   A. ERIC    B. PSYCH ABSTRACTS
3- LIFE SCIENCE     A. BIOSIS  B. CA      C. MESH
4- PHYSICAL SCIENCE A. DOE     B. CA      C. EI
                   D. INSPEC  E. IRON    F. IRON
PLEASE SELECT 1 OF THE 4 SETS BY ENTERING EITHER 1,2,3,4 4
VSS PROVIDES FOR 6 SWITCHING OPTIONS:
1- SYNONYMS          2- BROWSE          3- NARROWER TERMS
4- BROADER TERMS    5- NARROWER/BROADER TERMS 6- OTHER (USER-DEFINED)
PLEASE SELECT 1 OF THE 6 OPTIONS BY ENTERING # 2
SPECIFY MAXIMUM # OF TERMS TO BE DISPLAYED PER VOCABULARY. 10
PLEASE ENTER A SINGLE SEARCH TERM OR COMMAND HEAVY WATER
SWITCH SUCCESSFUL
TERM TYPE   VOCAB          TERM
YOUR TERM   DOE            HEAVY WATER
YOUR TERM   EI             HEAVY WATER
YOUR TERM   INSPEC        HEAVY WATER
YOUR TERM   NASA          HEAVY WATER
RELATED     DOE NASA      COOLANTS
RELATED     DOE EI INSPEC NASA  DEUTERIUM COMPOUNDS
RELATED     DOE CA INSPEC NASA  MODERATORS
RELATED     DOE INSPEC      TRITIUM COMPOUNDS
RELATED     DOE            DUAL TEMPERATURE PROCESS
RELATED     DOE EI INSPEC NASA  DEUTERIUM
RELATED     DOE EI INSPEC NAS  TRITIUM
RELATED     INSPEC        FISSION REACTOR MATERIALS
WD MATCH    DOE            HEAVY WATER PLANTS
```


PROCEEDINGS OF THE 3rd ASIS SIG/CR CLASSIFICATION RESEARCH WORKSHOP

of the differences that exist among systems and within host files. Williams discussed later the wider issues involved in the design of transparent systems and described a number of approaches, which she categorized as gateways, front ends, intermediaries, and interfaces (Lancaster et al. 1989, 51).

A well-known experiment is the Vocabulary Switching System (VSS), an experimental automated subject switching mechanism for searching multiple databases in a single natural language, developed by Niehoff et al. at the Battelle Columbus Laboratories. Though VSS exists in a form of an integrated or merged vocabulary of over ten existing vocabularies, it is considered by this author as a switching interface interposing between a user and an online system. It is a stand-alone, online database containing the subject descriptors and all the syndetic relationships found in the vocabularies included (Niehoff 1985). It transforms a user's search terms into the correct search terms for the databases to be searched. All vocabularies included in VSS were reformatted into a common input format. The unique records from each vocabulary ended up in one of five files: Concept file, Term file, Phrase file, Word file, and Stem file. As of 1985, VSS had incorporated 15 different database vocabularies in four modules (organized by discipline) and could offer various cross-database switching options and browsing capabilities. In later studies the scope of the object was expanded to service other subject areas and offered twenty-one options for "vocabulary switching", including exact matches and synonyms or alternative terms, browse, narrower terms, broader terms, etc. The system could compare a term entered by the user with terms in the vocabulary of a target database, and indicate the degree of match such as *No match*, *An exact match*, *A hierarchical match*, and *An associative match* (Piternich 1990, Chamis 1991, 56-58). Analysis of the results illuminated two points. First, the approach used was as relevant in a multilingual setting as it was in a monolingual setting. Second, different subject areas had varying degrees of vocabulary compatibility. For example, social science vocabularies were more compatible than those in other areas (Niehoff et al. 1979, 1980, and 1985).

Subject authority control is vigorously developed in the integrated library systems such as ORION, WLN, NLC/DOBIS, UTLAS, Geac BPS, NOTIS, and Carlyle (Mandel 1987). These systems currently can provide authority control for more than one logically separate subject vocabulary. Generally, they have two major functions. One function is to identify a subject source list, while the other is to validate subject headings against selected vocabularies. Validation ranges from a match of headings in order to capture cross-references (e.g., Carlyle) to complex linking mechanisms that replace headings in bibliographic records with authority record I.D. numbers (e.g., WLN, UTLAS). These systems are interesting not only for their support of multiple thesauri, but for their power to maintain control over headings and subdivisions and for their support of linking relationships among terms. For example, NLC/DOBIS and UTLAS's CATSS create a special translation to link LCSH terms, the *Canadian List of English Subject Headings*, and their French counter-parts in *Repertoire de vedettes-matiere* (Svenonius 1983; Mandel 1987, iv; 11-68).

4.3. Mixed Vocabulary Searching

Mixed vocabulary searching is a commonly used approach and has been implemented by the major online service centers in the United States. Terms from all vocabularies are retrieved together in subject searches. An example is one of the new search features of DIALOG, OneSearch. When used with DIALINDEX, it facilitates searching in many files concurrently with one search strategy (Dialog Database Catalog 1990). DIALINDEX and CROSS (developed by BRS) are files of all the

searchable terms appearing in the many databases of bibliographical records accessible through these services. In response to a term entered by a user online, such a tool will show in which databases the term occurs and how frequently it occurs in each (Lancaster 1986, 194). A similar approach has been used by online integrated library systems. For example, in ORION and DOBIS, the subject authority file index is not partitioned according to subject source. Terms from different thesauri are retrieved in the same alphabetical list, while the hits, terms and subject sources are indicated on the screen (Mandel 1987 11-15; 31-42).

Figure 8. Examples of Mixed Vocabulary Searching in Integrated Library Systems

ORION: (Terms from different thesauri are retrieved in the same alphabetical list.)

<u>hits</u>	<u>term</u>	<u>source</u>
25	Clinical psychology	(LCSH)
	Clinical psychology <i>SEE</i> Psychology Clinical	(Medical)
12	Psychology, Clinical	(Medical)
	Psychology, Clinical <i>SEE</i> Clinical Psychology	(LCSH)

DOBIS: (Subjects from all sources are displayed in an integrated alphabetical list, with the subject source clearly indicated. The number of bibliographic records linked to each heading is also shown.)

Canada - History - War of 1812	LCSH	10
Canada - History - War of 1812	CSH	50
Canada - History - War of 1812	CUT	2

Source: Mandel 1987, 13; 33.

It seems that this approach is acceptable because it is easy to implement and needs little artificial work on the existing vocabularies. But one problem this approach risks is obvious vocabulary clashes (e.g., the same term is postable in one vocabulary and non-postable in another). The seriousness of the problem depends upon the subject searching design features and the particular mix of collections and vocabularies included (Mandel 1987, v; 71-72).

5. CONCLUSION

The importance of compatibility of indexing languages in information systems is being given increasing emphasis in today's online multiple database access environment, in light of the many

PROCEEDINGS OF THE 3rd ASIS SIG/CR CLASSIFICATION RESEARCH WORKSHOP

examples of experimental extent and practical use. In order to solve problems which are beyond the scope of gradually developed international and national standards, many approaches to achieving compatibility among indexing languages have been investigated. Research projects that set new trends or reflect new approaches are reviewed in this article.

In general, efforts to develop mechanisms; that enable the establishment of compatible vocabularies can lead to more compatible information system performance, because the attention is being given to the earlier stages of the information transfer cycle. In the long run, the earlier that standardization and compatibility take place, the more cost-effective it is likely to be. But some problems remain. Intermediate lexicon and integrated vocabulary methods require creating additional tools, while microthesaurus and macrovocabulary approaches may introduce approaches some inflexibility and inability to tailor a system to local needs. Therefore, much intelligent work must be involved in the establishment of any of these compatible vocabularies. Sometimes, a measurement of the cost of, and effectiveness achieved from, the process of establishing a compatible thesaurus will call into question the value of such an effort and the extent to which it may leave the original objectives of earlier approaches to the problem of compatibility.

New developments in technology have made automated systems an alternative means to achieve compatibility. These have been especially useful for relating separately created thesauri to one another in an online multi-database access environment. Where no standards or compatible vocabularies exist, or where compatible vocabularies are felt to be too difficult to establish or too restrictive to satisfy local needs, automated systems can suggest themselves as, with some qualifications, an efficient means to support unique authority files. Undoubtedly, automated search aids that link vocabularies have emerged as feasible alternatives to dealing with the incompatibility of indexing languages. However, it is just at its embryonic stage. Most studies have concentrated on how to provide access to multi-thesauri and multi-databases. Few have addressed users' satisfaction with the search results. This is actually a very important issue because users of online bibliographic databases or online catalogs are no longer limited to professional searchers.

Notwithstanding questions which have yet to be given optimal treatment, the present work toward the establishment of compatibility among indexing languages should be seen as becoming increasingly extensive. Further improvement in compatibility will depend largely on two kinds of support. The first is the development of science and technologies such as artificial intelligence, as well as the theoretical exploration in many fields such as linguistics, psychology, and information science. The second is the continuation of efforts for standardization led by international agencies such as Unesco, IFLA, ISO as well as related national agencies of all countries.

BIBLIOGRAPHY

- Andre, P.Q.J. 1989. Universal agricultural thesaurus discussed. *Quarterly Bulletin of the International Association of Agricultural Librarians and Documentalists* 34(3): 150-151.
- Austin, D., and Digger, J.A. 1977. PRECIS: the preserved context index system. *Library Resources and Technical Services*. 21(1): 13-30.
- Chamis, A. Y. 1991. *Vocabulary Control and Search Strategies in Online Searching*. New York: Greenwood Press.

- Dahlberg, I. 1980. The Broad System for Ordering (BSO) as a basis for an integrated social sciences thesaurus? *International Classification* 7(2): 66-72.
- Dahlberg, I. 1981. Towards establishment of compatibility between indexing languages. *International Classification* 8(2): 86-91.
- Dahlberg, I. 1983. Conceptual compatibility of ordering systems. *International Classification* 10(1): 5-8.
- DeHart, F.E., and Glazier, J. 1984. Computer searching on PRECIS: an exploration of measuring comparative retrieval effectiveness. *International Classification* 11(1): 3-8.
- Dextre, S.G., and Clarke, T.M. 1981. A system for machine-aided thesaurus construction. *ASLIB Proceedings* 33(3): 102-112.
- DIALOG Database Catalog* 1990.
- Foskett, D.J. 1991. Concerning general and special classifications. *International Classification* 18(2): 87-91.
- Hammond, W., and Rosenborg, S. 1962. *Experimental Study of Convertibility Between Large Technical Indexing Vocabularies*. Silver Spring, Md.: Datatrol Corp.
- Humphreys, B.L. 1988. Unified Medical Language System: progress report. *International Classification* 15(2): 85-86.
- ISO. 1985. *Guidelines for the Establishment and Development of Multilingual Thesauri*. Geneva: ISO, 1985. (ISO 5964-1985)
- Lancaster, F.W., and Smith, Linda C. 1983. *Compatibility Issues Affecting Information Systems and Services*. Paris: Unesco General Information Program.
- Lancaster, F.W. 1986. Compatibility and convertibility. Chapter 14 of *Vocabulary Control for Information Retrieval*. 2nd ed. Arlington, Va.: Information Resources Press.
- Lancaster, F.W.; Elliker, C.; and Connell, T.H. 1989. Subject Analysis. In *Annual Review of Information Science and Technology*, edited by Martha E. Williams, 24: 35-84. Washington, D.C.: American Society for Information Science.
- Mandel, C.A. 1987. *Multiple Thesauri in Online Library Bibliographic Systems: A Report Prepared for Library of Congress Processing Services*. Washington, D.C.: Library of Congress.
- Meyriat, J. 1980. Social science information languages: a comparative analysis. *International Classification* 7(2): 60-65.
- National Library of Medicine News* Oct. 1987: 4-6.
- Newman, Simon M., ed. 1965. *Information Systems Compatibility*. Washington: Spartan Books, 1965.
- Niehoff, R., and Kwasny, S. 1979. The role of automated subject switching in a distributed information network. *Online Review* 3(2):181-192.
- Niehoff, R. et al. 1980. *The Design and Evaluation of a Vocabulary Switching System for Use in Multi-base Search Environments*. Columbus, Ohio: Battelle Columbus Laboratories.
- Niehoff, R. and Mack, G. 1985. The Vocabulary Switching System: description of evaluation studies. *International Classification*. 12(1): 2-6.
- Piternich, Anne B. 1990. Vocabularies for online subject searching. In *Encyclopedia of Library and Information Science*, edited by Allen Kent, 45: 399-420. New York: Marcel Dekker, Inc.
- Rada, R. 1987. Connecting and evaluating thesauri: issues and cases. *International Classification* 14(2): 63-69.
- Rada, R. 1990. Maintaining thesauri and metathesauri. *International Classification* 17(3/4): 158-164.

PROCEEDINGS OF THE 3rd ASIS SIG/CR CLASSIFICATION RESEARCH WORKSHOP

- Roulin, C. 1990. *Towards the European: Education Metathesaurus*. Translated by Barbe and Brett. Brussels: EURYDICE European Unit, 1990.
- Roulin, C. 1991. Sub-thesauri as part of a metathesaurus. paper presented at the 5th International Study Conference on Classification Research, Toronto, Canada, June 24-28, 1991.
- Sager, J.C. et al. 1982. Thesaurus Integration in the Social Sciences. Part II: Stages towards integration; Part III: Guidelines for the integration of thesauri. *International Classification* 9(1): 19-26; 9(2): 64-70.
- Schwartz, C., and Eisenmann, L.M. 1986. Subject analysis. In *Annual Review of Information Science and Technology*, edited by Martha E. Williams, 1: 37-61. Washington, D.C.: American Society for Information Science.
- Scibor, E. and Jabrzemska, E. 1989. Thesaurus of Common Topics. In *Encyclopedia of Library and Information Science*, edited by Allen Kent, 44: 388-395. New York: Marcel Dekker, Inc.
- Soergel, D. 1972. A general model for indexing languages: the basis for compatibility and integration. In *Subject Retrieval in the Seventies: new Directions*, edited by H. Wellisch and T.D. Wilson, 36-61. Westport, Conn.: Greenwood.
- Soergel, D. 1974. *Indexing Languages and Thesauri: Construction and Maintenance*. Los Angeles, Ca.: Melville Pub. Co.
- Smith, Linda C. 1991. UNISIST revisited: compatibility in the context of collaboratories. paper presented at the 5th International Study Conference on Classification Research, Toronto, Canada, June 24-28, 1991.
- Svenonius, E. 1983. Compatibility of retrieval languages: introduction to a forum. *International Classification* 10(1): 2-4.
- Unesco. 1971. *UNISIST Study Report on the feasibility of a world science information system*. Paris: Unesco.
- Wall, E., and Barnes, J.M. 1969. *Intersystem Compatibility and Convertibility of Subject Vocabularies*. Philadelphia, Pa.: Auerbach.
- Whitelock, P.J. 1982. A descriptor bank of social science terms. *International Classification* 9(3): 145-151.
- Williams, M.E. 1978. Online retrieval - today and tomorrow. *Online Review* 2: 353-366.
- Zeng, Lei. 1990. Establishing a compatible general vocabulary in China: the capability. *International Classification* 17(2): 91-98.
- Zeng, Lei. 1991. Research and development of classification and thesauri in China. paper presented at the 5th International Study Conference on Classification Research, Toronto, Canada, June 24-28, 1991.