
Jacques Maniez
Institut Universitaire de Technologie, Dijon, France

Relationships in Thesauri: Some Critical Remarks

Maniez, J.: **Relationship in thesauri: Some critical remarks.**
Int. Classif. 15 (1988) No. 3, p. 133–138, 19 refs.

After reviewing some fundamental distinctions in relationships (paradigmatic/syntagmatic, interconceptual/structural) the author proposes a functional approach for investigating the relationships in thesauri. The comparison between three closely related types of semantic fields (lexical, conceptual, thesaural) shows the specific function of relationships in all of these intellectual tools.

In information retrieval the two main functions are location of relevant concepts and search of exhaustivity. A clear distinction of these aims can contribute to solving the difficult problem of the choice of 'related terms'. It is suggested that their usefulness relies upon empirical rather than upon semantic proximity. Some practical propositions are made for the choice and display of relationships in thesauri. (Author)

1. Introduction

The usefulness of thesauri is now beyond all question, and unlike the expectation of many observers in information science, the increasing effectiveness of language processing has not stopped their proliferation.

It is striking that while a lot of thesauri have been created, and excellent manuals for thesaurus-making are published every year (1, 2), few theoretical investigations have been raised by this area in the last thirty years, as if the uncontested efficiency of the tool were in itself its justification.

Nevertheless, the focal point of the thesaurus – its network of relationships – is far from being fully clear and explicit, and it seems that lists of related terms often rely on intuition rather than on a reasoned view of their effectiveness (3). Moreover, "associative relationships" are defined in standards (4) and in reference lexica in a rather vague and negative way. For example, the fairly good *Terminology of documentation* by Wersig and Neveling (5) defines them as "relations between concepts or classes which are nothierarchical". We would like to submit some considerations on this theme, assuming that a functional approach is best suited to an understanding of intellectual tools such as thesauri.

2. Useful distinctions about relationships

Nothing is more common than the concept of relationship, as witnessed by the famous law of an unknown humourist: "Everything is in everything, and vice-versa". More earnestly, exploring the area of relations has been a fruitful practice in all sciences, as it was once again demonstrated recently in mathematics with relational algebra.

Unfortunately universal schemes, such as mathematical models, cannot be easily applied to the fields of high complexity, and if we follow the theory of integrative levels (6), information belongs to "mentefacts", which are at the highest level of complex entities. So, though it be tempting to use simple and universal keys to understand the facts of knowledge, and particularly the relations between words, concepts or subjects, it seems better to adjust our conceptual tools to the peculiar nature of our topic.

2.1 Paradigmatic versus syntagmatic relationships

Paradigmatic relationships are judgements on the relatedness of two words or two concepts, intrinsically and out of any context. For instance 'red/pink', 'love/friendship', 'illness/death' are pairs of concepts bearing some degree of similarity. In other words this type of relation is essentially comparative. Consequently it is very important in documentation since the aim of classification and indexing is to label with similar formulas documents and requests relating to similar subjects.

Syntagmatic relationships (some authors, e.g. Coates (7), call them only 'relations') are not statements but a creative process, which produces a new compound phrase or concept (a syntagm) out of the two original words or concepts. This type of relation is not permanent, but casual. Using a famous distinction from Saussure, the paradigmatic relationships belong to 'language', the syntagmatic ones belong to 'speech' (though the rules of connection take place in the grammar of the language).

This basic distinction (Paradigmatic vs syntagmatic) is generally acknowledged in information science. However, some authors (8) state that it is merely conventional and that at a deep level they do not differ from each other. Personally we think that it is essential because it differentiates two kinds of clearly distinguishable mental operations: comparing and linking. But the confusion is easily understandable because the two kinds of relationships are closely interrelated:

a) To be combined, two concepts must be compatible, which implies that they are neither too similar nor too different. For instance it is impossible to form a good syntagm between two synonyms ('to drink a drink', 'to walk a walk') as well as between incompatible terms ('the logarithm of a fly'). In short the combination of terms relies also on paradigmatic conditions, which could be called 'syntactico-semantic' (fig. 1).

b) The creation of a new syntagm generates in itself a genus/species relationship between each of the original terms and the new expression. Example:

cancer + lung → lung cancer
cancer → specific: lung cancer
lung → specific: cancerous lung

In documentary languages, syntagmatic relationships are less obviously important than paradigmatic ones, and depend on the type of language:

- in fully precoordinated languages such as enumerative classifications (e.g. the L.C.) no syntax is needed.
- in fully postcoordinated languages such as conven-

tional thesauri the syntagmatic relations are implicit and not differentiated

- on the other hand, syntagmatic relations are explicit and of basic importance in faceted classifications and in syntactically-structured languages such as SYNTOL or PRECIS.

PARADIGMATIC	SYNTAGMATIC	EXAMPLES
-identity	no syntagm	human man
-synonymy	(or very few	to chase-pursue
-antonymy	syntagms)	good-bad
-specificity	few syntagms	rat-animal

*-ASSOCIATIVITY	MANY SYNTAGMS	
* organic part		hand-finger
* usual agent		baker's bread
* usual means		paint-brush
*		

*COMPATIBILITY	MANY SYNTAGMS	
* casual part		hand-protectio*
* casual agent		mother's bread*
* casual means		painting with *
*		fingers *

dissimilarity	No syntagms	eating a theorem

Fig. 1: Correlations between paradigmatic and syntagmatic relationships

2.2 Interconceptual versus structural relationships

Paradigmatic relationships may themselves be divided into two types of relation: interconceptual and structural. The former link any two terms of a language (see above). The latter link every term to the structure of the language by assigning it to one or several basic classes. For instance, the term 'travel' may be assigned to the category of nouns in a lexicon, to the field 'Transportation' in a thesaurus, to the facet 'Process' in a faceted classification. Some authors call it the dependence relationship (e.g. Van Slype (1), p. 54–55). This type of relationship relies on the main level of division of the language, which determines its structure (divided by fields or by facets). Using a geographical comparison, they locate every town or village in relation to one main area.

In classifications the choice between fields and facets is fundamental because it determines the way of expressing subjects: synthetic or 'analytico-synthetic' (Ranganathan). On the other hand in thesauri, which are by nature post-coordinate languages, this choice is of minor importance, except if the thesaurus is completed with syntactic rules: in such a case facets are preferable since they are in themselves an outline of a syntactic structure. But this type of categorization also has some drawbacks, as we shall see further on.

3. A functional view of paradigmatic relationships

In all languages (in the broadest meaning of the term) the paradigmatic organization is very helpful to users because it simulates the mental network of relationships which supports our understanding of the world and our intellectual activity. This usefulness is obvious in three types of intellectual tools: analogical lexica (interverbal), terminologies (interconceptual), indexing

languages (intertopic). Furthermore the worlds of language, knowledge and information retrieval are so closely interrelated that often the same terms are used in every domain (e.g. 'semantic field', 'thesaurus'). However each tool gets special features corresponding to the special aims of the operation for which it is used. A brief comparison of paradigmatic relationships in these three intellectual implements can help us to specify the case of the thesaurus.

3.1 The analogical dictionary (e.g. *Rogel's Thesaurus*)

The scope of this dictionary is clearly defined in the preface (9): "It is a vocabulary on a large scale, categorized by topics. With such an arrangement it offers a choice of words to fit any given context".

The macrostructure is made of 990 topics distributed between six major classes (abstract relations, space, matter, intellect, volition, affections) and about thirty sections. For instance 'book' is a topic (a 'head') pertaining to the major class 'Intellect' and to the section 'Means of communicating ideas'. This framework forms a small hierarchical classification, completed with the set of conventional 'parts of speech' (nouns, verbs . . .). But this macrostructure is more an intellectual refinement than a functional part of the system, since the reader looks up the index, and not the classification scheme, to find the place of a word. On the other hand the semantic microstructure, that is the 'lexical field' surrounding each topic is of basic interest:

- it lists the lexical environment of the heading, as comprehensively as possible, including synonyms and shades of meaning (more than 200 words for 'book'),
- the structure of the list is clear but implicit: it holds a series of sublists (from the nearest to the furthest: book writing, textbook, edition, library, author . . .), each of which bears the close environment of its head. So the only one explicit relationship is that which links a word to the number of its topic, and the cross-references are limited to words pertaining to several topics (ex: THESAURUS: see LIST, 87. DICTIONARY, 559.).

In short, a lexical field aims at listing together in a well-ordered scheme all the words of a language referring to the same topics, in order to facilitate the most accurate choice of words. Let us stress that terms are not defined and that relationships are left implicit.

3.2 Conceptual terminology (e.g. *Wersig's "Terminology of documentation"* (5))

For some decades terminology has become an important discipline, with terminologies having become efficient tools.

Here the scope is not to list comprehensively the words of a language pertaining to the same domain, but to index all the concepts of a field and to link every concept to one word of a language (monolingual terminologies) or to n words of n languages (multilingual terminologies).

6.1 Explication of requests

It is a well-known fact that every message tells more than its literal content. For instance the statement "I was in London last week" tells implicitly "I was in England last week". This fact is not a drawback in the context of human communication, where an utterance is immediately decoded and interpreted. But in the context of documentary mediation and double coding everything must be explicated. Equivalence and genus/species relationships can explicate the implicit face of each term.

Equivalence. The relationship of equivalence has so often been explored as to deter further comments. Generally in thesauri it is processed in a restrictive way with 'USE' and 'USED FOR' expressions, and this practice is reasonable for it limits economically the number of descriptors. But from a logical viewpoint it would be just as efficient to automatically bind equivalent terms, as is the case in multilingual automated thesauri.

This relationship affords essential protection against silence, since it generates this implied statement: 'If you are interested in "A" then you are *intrinsically* interested in "B"'.

But the difficult point about it is to state the degree of dissimilarity beyond which two terms are no longer equivalent. It has often been demonstrated that strict equivalence of terms cannot exist (see (15) and (12), p.239–242). Let us only point out that the requirements of information retrieval are opposite to those of linguistic utterance. While nuances of meaning contribute to the richness of speech, the presence of too closely related terms in a thesaurus can create silence and should be avoided.

Hierarchical relationships. For pragmatic reasons, this rather vague attribute applies altogether to the genus/species and to the whole/part relationships, which are often mixed in thesauri. In principle, these relations are clearly distinct: a thumb is a finger (genus/species) but *is not* a hand (whole/part). It follows that a statement which is true for a genus is necessarily true for a species, while the truth of the whole does not necessarily imply the truth of the part. "If you are interested in fingers then you are *certainly* interested in thumbs – though to a lesser extent –" but: "If you are interested in hands then you are *probably* interested in fingers".

The genus/species relationship is the basis of classification and of syllogism. It is asymmetrical and transitive, so that it creates a perfect hierarchy.

In an indexing language, this relationship brings the adequate remedy to the implicitness of a query, chiefly when it is implemented automatically: if so, a query about limbs will also return documents about legs or arms. (Let us stress, by the way, the semantic and linguistic proximity of 'implicit' and 'implication').

Unfortunately the difference between genus/species and whole/part relationships is less clear in practice than in theory. In some situations the difference between the two relations becomes almost imperceptible. This is the case, in particular, when an abstract concept

summarizes different sciences (mathematics, physics, computer science . . .).

Geometry, for instance, is a part of mathematics, but it is also a species of mathematical sciences, i.e. of mathematics. Moreover, when a 'whole' is a kind of cybernetic system, in which every part is intimately related to all the others (for instance a living organism), it is difficult to be interested in the overall system and not in its parts: that is why in medical thesauri the whole/part relation is generally mixed with the genus/species relation.

In spite of this slight logical confusion due to practical reasons the hierarchical relationships in thesauri are a powerful means of avoiding silence, allowing one to extend a search to *implicitly relevant documents*. With automated thesauri they can be systematically used on request.

6.2 Vagueness of the notion of subject

This point needs further explanation for it is generally obliterated. The notion of subject (= topic, aboutness) is difficult to grasp (see (12), p.253–263). In Wersig's Terminology (5) it is defined as "any concept or concept combination looked at as an entity". If a subject is nothing else than a concept, the degree of vagueness of a subject depends on the nature of the concepts it is composed of, either 'hard' and clearcut (e.g. aliveness, sex, SO4H2 . . .) or 'soft' and vague (youth, redness . . .). Personally, we think that a subject is more than a concept: it is a point of interest, of concern or research, and this personal feature gives the notion a subjective connotation that is not included in the notion of concept. In other words, if we are interested in one subject, our interest cannot be exactly limited by its formulation into terms of concepts.

For instance the interest in "the measurement of earthquakes" (hard concepts) does not stop exactly at the frontier of the concept itself. So it is the notion of subject per se that is vague (= non-discrete), and between the strict extension of a subject (area of formal relevance) and that of non-subject (area of irrelevance) there lies a penumbral area, that of *probable relevance*. This explains the fundamental difference between looking for a concrete object (an apple), or a fact (Newton's birthday) and the search for documents dealing with a subject.

6.3 Effectiveness of associate relationships

While equivalence and genus/species relations give a simple and sure device – based on the logic of classes – for covering the zone of implicit relevance, it is far more difficult to delimit the area of potential relevance of a document from the semantic environment of the terms which express its topic. The more so with a post-coordinated language which does not list synthetic subjects, but rather concepts, the combination of which can generate an infinite number of topics. Then a term closely related to another one can be of no use in the context of a compound topic. For instance "dictation" towards "spelling". If I am interested in "acquisition of spelling" the adjunction of "dictation" can be helpful,

but if I want documents about "the spelling of foreign terms" it is useless.

The use of listing so-called related terms in thesauri has been so far firmly preserved in spite of the uncertainty of its usefulness, and though no rule can fix the limits of sound 'relatedness'. For instance the French standard AFNOR Z47-100 (16) states wisely that related terms should be mentioned in a thesaurus 'only if they are necessary for information retrieval's sake' but no rule is given to fix the criteria of this necessity. The result is that every thesaurus designer must intuitively establish for himself the limits of relatedness, and sometimes the category of associative relationships looks like a set of heterogeneous terms. Admittedly many inventories of associated relationships have been carried out (3), but these accurate examinations establish a list of categories without any appraisal of their usefulness. For instance, the IRTET - Information Retrieval Thesaurus of Education Terms - (see (17) p. 80-84) records thirteen species of relatedness.

It seems to us that these cases could be reduced to three categories, which could each be processed in a particular way:

- A - semantic near-synonymy relatedness
- B - semantic overlapping
- C - extra-semantic relatedness

A) Near - synonymy relatedness

This case is a slight extension of equivalence relationship. It happens that several nearly coextensive terms exist in a thesaurus in order to allow shades of meaning in indexing and searching, for instance "automatic translation" and "computer aided translation". In our opinion it would be better to place this type of terms into the category of 'implicitly relevant terms'. If the searcher wished to extend the result of a query, these terms would be automatically connected on request with the basic term as well as the specific terms. For instance the basic term "life-long education" could be completed with "adult education" (near synonym) and with "retraining" (narrower term). Of course, this kind of relation is reciprocal.

B) Overlapping terms

This is the typical case of associated relationships, that is of pairs in which each term is involved in the definition of the other. For instance "book" towards "edition, library, reader, book-seller . . .". In our opinion these pairs are rarely efficient as an extension of the basic term because the portion of semantic discrepancy inside the pairs is important. So the best way would be not to mention them systematically inside the hierarchical list of descriptors and to display them only in the tables of semantic fields.

An exception should be made, however, with related terms which occur frequently together in titles and/or indexing formulas, for the degree of efficiency of relatedness for information retrieval cannot be fixed in abstracto but depends on actual documents. So statistics about the frequency of cooccurrence of pairs are of great value for this purpose and could greatly improve the design of thesauri.

C) Extrasemantic relations

If we agree that semantics is the study of the relations between signs and meanings, we may call extrasemantic the type of relationships between terms that rely on facts rather than on words or terms.

For instance 'braking' and 'side-slipping' are often connected in actual life, though the semantic contents of the terms are not related. Among the thirteen categories recorded in IRTET several belong to this type:

- between cause and effect (example above)
- between two terms referring to objects which bear some kind of similarity (in structure, in physical or social environment . . .)

Ex: *Eskimo/Turkish language (structural)*
Britanny/Auvergne (demographical)

The relatedness based on structural or casual similarity is not yet of great use in thesauri, but it seems of great interest and could be developed by the observations of experts in the field, and by the results of automatic clustering of terms (18).

These relationships are remarkable in that the connection can be found neither in linguistic nor in conceptual analysis of terms, but relies on hidden relatedness of real things disclosed by experience.

It is interesting to notice the parallelism between this kind of association in thesauri and the approach of expert systems: here also the source of efficient retrieval does not lie in an algorithm but in some 'expertise' which must be put into the 'knowledge base'. And this kind of relationships plays the same role as the "inference engine" since they are the abridged equivalent of an implicative judgement.

And more generally, relationships in thesauri could be considered as an archetype of the recent concept of 'hypertext' (19), in that they involve every document within a virtual set of related documents.

7. Conclusion

This critical view of relationships in thesauri leads us to the following remarks and proposals:

The traditional classification of relationships relies on logico-semantic categories. They are useful, but only as means subordinate to the aims of information indexing and retrieval.

It would be better to start from a functional view and to distinguish clearly:

- the choice of the most relevant terms of the thesaurus for the utterance of indexing formulas and request formulas (precision requirement)
- the enrichment of the request formula in order to reach a comprehensive set of relevant documents (exhaustivity requirement).

These aims could be reached thanks to a division of thesauri into three parts:

- a) a comprehensive list of descriptors and non-descriptors (with USE reference)
- b) a semantic display reflecting the macrostructure and showing in each subfield a comprehensive list of de-

scriptors with their relationships. In that respect the division by fields is better than the division by facets for it gives a more global view. It is good to complement the lists with graphic displays (arrowgraphs, semantic roadmaps . . .) which will provide a complete table of all relationships at a glance (see (1), p.96–98).

These two parts would be enough to reach the first aim: an accurate choice of the best descriptors to express the basic subject of a document or request.

c) *a restricted list of useful relationships* reserved for searching and aimed at allowing for an exhaustive retrieval. This part would eliminate equivalence relationships (as useless) and generic ones (as productive of noise). Each descriptor would be surrounded only with terms likely to extend the results of a search, either manually or automatically, without creating a high rate of noise. It would be convenient to replace conventional categories (hierarchical, associative) with functional categories:

- *implicit relevance* joining specific terms and near-synonyms.
- *probable relevance* joining two classes of ‘associated terms’:
 - those related semantically (overlapping) and by a high degree of cooccurrence
 - those related by extraseantic factors (causation/effect, structural similarity . . .)

Of course, these two levels of extension should remain under the control of the searcher, even in an automated system.

This set of relationships, based on the concept of implication, would play the role of an expert system in automated databases.

Acknowledgement

I wish to thank Mrs. Martha Melter and Mr. François Melter for helping me to revise this text.

References

- (1) Van Slype, G.: Les langages d'indexation: conception, construction et utilisation dans les systèmes documentaires. Paris: Ed. d'organisé 1987. 277 p.
- (2) Aitchison, J.; Gilchrist, A.: Thesaurus construction. A practical manual. London: ASLIB 1987. 175 p.
- (3) Willets, M.: An investigation of the nature of the relation between terms in thesauri. In: J. Doc. 31 (1975) No. 3, p. 158–184.
- (4) Somers, H.L.: Observations on standards and guidelines concerning thesaurus construction. In: Int. Classif. 8 (1981) No. 2, p. 69–74.
- (5) Wersig, G.; Neveling, U.: Terminology of documentation. Paris: Unesco Press 1976. 274 p.
- (6) Foskett, D.J.: Classification for a general index language. London: The Library Assoc. 1970. 48 p. (Here p. 23–31).
- (7) Coates, E.J.: Some properties of relationships in the structure of indexing languages. In: J. Doc. 29 (1973) No. 4, p. 390–404 (Here p. 392).
- (8) Gardin, J.C.: Les analyses de discours. Neuchâtel: Delachaux et Niestlé 1974. p. 136–137.
- (9) Dutch, R.A.: Roget's thesaurus. Harmondsworth, Middlesex: Penguin books 1975. 712 p.
- (10) Dahlberg, I.: Conceptual Definitions for INTERCONCEPT. In: Int. Classif. 8 (1981) No. 1, p. 16–22.
- (11) Dahlberg, I.: Conceptual compatibility of ordering systems. In: Int. Classif. 10 (1983) No. 1, p. 5–8.
- (12) Maniez, J.: Les langages documentaires et classificatoires – Conception, construction et utilisation dans les systèmes documentaires. Paris: Ed. d'organisé 1987. 293 p.
- (13) Hacndler, H.: Synthetic description systems for accurate data identification and selection. Principles and methods of nutritional data banks. (In: Int. Classif. 15 (1988) No. 2, p. 64–68.
- (14) DeHart, F.E.: Topic relevance and retrieval effectiveness. In: Int. Classif. 10 (1983) No. 1, p. 9–14.
- (15) Ungváry, R.: Intensional splitting. An empirical examination of conceptual duality. In: Int. Classif. 13 (1986) No. 1, p. 9–17.
- (16) AFNOR: Règles d'établissement des thésaurus monolingues Z47–100. Paris: AFNOR 1981. 20 p.
- (17) Lancaster, F.W.: Vocabulary control for information retrieval. Washington D.C.: Information Resources Press. 1972. 233 p.
- (18) Courtial, J.P.; Pomian, J.: Un système à base de logique des associations. In: Documentaliste, 24 (1978) No. 1, janvier-février 1987, p. 11–16.
- (19) Macleod I.A.: Hypertext as a Document retrieval application. In: L'informatique documentaire No. 31, Sept. 1988, p. 21–32.

USA: Forest-Press Sponsored Brainstorming Session

Directly preceding the Conference “Classification Theory in the Computer Age: Conversations Across the Disciplines”, Albany, Nov.18-19, 1988 (see below) a one-day meeting was held on Nov.17, 1988 in Albany with 17 persons thinking about and discussing the following questions:

- (1) What are the most important trends in the areas of information retrieval and how will they affect general classification systems over the next 5-10 years?
- (2) What will be/should be the future operating environment (at the working library level) for such classification systems? Assuming that the current structure of general classification systems does not change, what approaches can be used to modify the user environment to improve the use of classification systems in such environments?
- (3) What changes in present classification systems will improve their use for subject retrieval in a manual or online environment?
- (4) What needs, other than arranging books or surrogates, can classification systems meet? How can/should current classification systems be adapted to serve these needs, or what new systems or services should be developed? Consider classification systems both for organizing materials (e.g., journal articles) and/or information (e.g., in thesauri and AI systems).

This note on the meeting is not meant to report on its results. It is hoped that Forest Press will release in some way the Conclusions drawn and the priorities for research and development identified. Meanwhile anybody interested in answering some or all of these questions is invited to write down his or her comments and to send them to the Editor of this journal (IC).