

A Graphical Interface for Faceted Thesaurus Design

Uta Priss

Elin K. Jacob

School of Library and Information Science

Indiana University

Bloomington, Indiana

Abstract

This paper develops a formalization of faceted thesauri that is based on formal concept analysis. The formalization facilitates graphical displays of thesauri as line diagrams of mathematical lattices. Resulting strategies that can be employed to design a thesaurus in an alternating top-down and bottom-up approach are described and demonstrated through an example.

Introduction

In light of the tera bytes of information available both on the world wide web and in institutional and commercial databases, there is increasing demand for effective subject access systems. While there is ongoing discussion as to whether natural language or controlled vocabulary is more successful or whether full-text searching should replace automated or semi-automated indexing simply because it is cheaper, information scientists will generally agree that a well-constructed thesaurus could be usefully applied to information retrieval. Although this paper will concentrate on the design and development of thesauri within information science, additional applications for such a thesaurus are possible given its similarity to other systems of structured knowledge representation such as knowledge bases and lexical databases: each of these systems not only groups words into synonym sets but also defines the semantic relations between the synonym sets. Thus thesauri, knowledge bases and lexical databases are closely related in that each provides a means for coding semantic information into a machine-readable and computer-processable format.

If information scientists are to succeed in building common language interfaces, it seems unlikely that this can be achieved without sophisticated knowledge representation tools. Although thesauri are important tools for representation and retrieval, most existing thesauri are nowhere near "perfect" (see, for example, Fischer, 1993). Part of the reason for this may be due to a lack of flexibility in thesaurus structure and the absence of a well-defined, machine-aided method of thesaurus design and construction. Faceted thesauri are potentially flexible in that several views of the same concepts can be represented simultaneously. But faceted systems can become quite complex and the flexibility of these systems is not easily

represented in a linear display or tree hierarchy. Furthermore, without the appropriate software, it is very difficult to design a faceted thesaurus. This paper introduces strategies for the design of faceted thesauri and describes supporting software based on a mathematical theory called "formal concept analysis" (Ganter & Wille, 1996).

There is some recent research on the application of formal concept analysis to thesauri (Skorsky, 1997). Much of this work has been based on TOSCANA (Vogt & Wille, 1995), a formal concept analysis tool that allows navigation through facet-like structures called "scales" or "topics". The connection between facets and TOSCANA's scales has recently been discussed by Viehmann (1996) and Kent and Neuss (1995). Our approach adds to a TOSCANA-like initial design by incorporating such aspects as decomposition of facets into smaller facets and by providing a design environment specific to the construction of faceted thesauri. Decomposition of facets into their sub-components allows the design of facets that are domain-specific rather than universal: for example, the facet of terms for mental states of human beings would be more comprehensive than a similar facet for animals. In TOSCANA, the same facet would have to be used for humans and animals. In our approach, humans and animals would share a generic facet for mental states that would be enhanced by a further facet specific to humans. Stumme (1996) has recently addressed this problem, but his proposal is not sufficiently flexible.

Craven (1990) provides a survey of graphical displays for thesauri and discusses algorithms for arranging thesaurus entries in a display with minimal line crossovers. But current research in the area of thesaurus display does not often provide for either fully graphical or fully structured representations. Thus, for example, Johnson's (1995) hypertext thesaurus interface contains a traditional tree display and a random arrangement of related terms, or RTs. Pollitt's (1997) approach currently provides the most sophisticated interface for faceted classification systems: but his system is restricted to top-down views on the hierarchical structure and does not indicate how an individual document is related to several facets.

Formal concept analysis

Formal concept analysis (Ganter & Wille, 1996) starts with the definition of a *formal context* \mathcal{K} as a triple (G, M, I) consisting of a set of [formal] objects (denoted by G), a set of [formal] attributes (denoted by M), and a relation I between G and M (i.e., $I \subseteq G \times M$). The relationship is written as gIm or $(g, m) \in I$ and is read as "the formal object g has the formal attribute m ". A formal context can be represented by a cross table, or matrix, which has a row for each object g , a column for each attribute m and a cross, or X, in the row of g and the column of m if gIm . Figure 1 shows two examples of formal contexts. Here, "filly", "mare", etc., are the formal objects; and "female", "juvenile", etc., are the formal attributes. In a context (G, M, I) , the set of all common attributes of a set $A \subseteq G$ of objects is denoted by $\iota A := \{m \in M \mid gIm \text{ for all } g \in A\}$ and, analogously, the set of all common objects of a set $B \subseteq M$ of attributes is $\varepsilon B := \{g \in G \mid gIm \text{ for all } m \in B\}$. For example, in the left-hand formal context in Figure 1, the equations $\iota\{\text{ram}\} = \{\text{adult, male}\}$ and $\varepsilon\{\text{female}\} = \{\text{filly, mare, cow, ewe}\}$ hold.

A pair (A, B) is said to be a [formal] concept of the formal context (G, M, I) if $A \subseteq G, B \subseteq M, A = \varepsilon B$, and $B = \iota A$. For a concept $c := (A, B)$, A is called the *extension* (denoted by $Ext(c)$) and B is called the *intension* (denoted by $Int(c)$) of the concept. In the right-hand example in Figure 2, $(\{cow, bull, calf\}, \{cow, animal\})$ is a concept, because $\iota\{cow, bull, calf\} = \{cow, animal\}$ and $\varepsilon\{cow, animal\} = \{cow, bull, calf\}$. The set of all concepts of (G, M, I) is denoted by $\mathcal{B}(G, M, I)$. The most important structure on $\mathcal{B}(G, M, I)$ is given by the subconcept-superconcept relation that is defined as follows: the concept c_1 is a *subconcept* of the concept c_2 (denoted by $c_1 \leq c_2$) if $Ext(c_1) \subseteq Ext(c_2)$, which is equivalent to $Int(c_2) \subseteq Int(c_1)$; c_2 is then a *superconcept* of c_1 . For example, the extension of $(\{foal, calf, lamb, filly, colt\}, \{juvenile\})$ contains $\{filly\}$ as a subset and its intension is a subset of $\{female, juvenile\}$. Therefore $(\{foal, calf, lamb, filly, colt\}, \{juvenile\})$ is a superconcept of $(\{filly\}, \{female, juvenile\})$. The relation " \leq " is a mathematical order relation called *conceptual ordering* on $\mathcal{B}(G, M, I)$. The set of all concepts with the conceptual ordering form a mathematical lattice denoted by $\underline{\mathcal{B}}(G, M, I)$.

	female	juvenile	adult	male
filly	x	x		
mare	x		x	
colt		x		x
stallion			x	x
cow	x		x	
ram			x	x
bull			x	x
ewe	x		x	
foal		x		
calf		x		
lamb		x		

	horse	cow	sheep	animal
filly	x			x
mare	x			x
colt	x			x
stallion	x			x
cow		x		x
ram			x	x
bull		x		x
ewe			x	x
foal	x			x
calf		x		x
lamb			x	x

Figure 1. Formal contexts.

Graphically, mathematical lattices can be visualized by line diagrams which represent a concept by a small circle. For each object g , the smallest concept to whose extension g belongs is denoted by γg . And, for each attribute m , the largest concept to whose intension m belongs is denoted by μm . The concepts γg and μm are called *object concept* of g and *attribute concept* of m , respectively. In a line diagram, it is not necessary to write the full extension and intension for each concept. Instead, the name of each object g is written slightly below the circle of γg and the name of each attribute m is written slightly above the circle of μm . Concepts can be given names which are written next to the concepts. To distinguish them from objects or attributes, names of concepts are surrounded by a line. It is not necessary for all concepts to have names. Figure 2 shows the line diagrams of the concept lattices for the examples in Figure 1. To read the line diagram, the extension of a concept consists of all objects which are retrieved by starting with the concept and then collecting all objects that are attached to subconcepts of that concept. Analogously, the intension is retrieved by collecting all attributes that are attached to superconcepts of the concept.

Nested line diagrams facilitate the combination of several lattices that share the same set of objects such as the two lattices in Figure 2. One lattice is chosen as an *outer structure* while the other one becomes an *inner structure*. In the example in Figure 3, the right-hand lattice

in Figure 2 has been selected as an outer structure while the left-hand lattice in Figure 2 serves as an inner structure. The concepts of the outer structure are replaced by boxes, each of which contains a copy of the inner structure. The concepts in each box are subconcepts or superconcepts of the corresponding concepts in other boxes that are connected by lines of the outer structure. For example, the object concept of “foal” is a subconcept of the attribute concept of “juvenile”. Since the bottom box does not contain any objects, it is not completely represented.

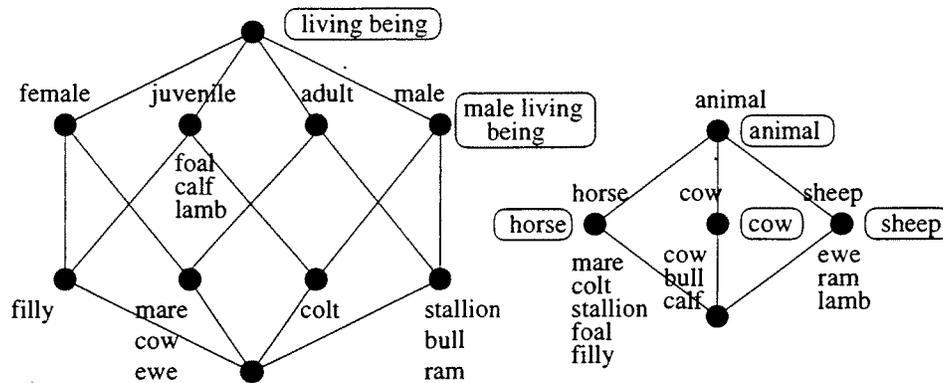


Figure 2. Line diagrams of concept lattices.

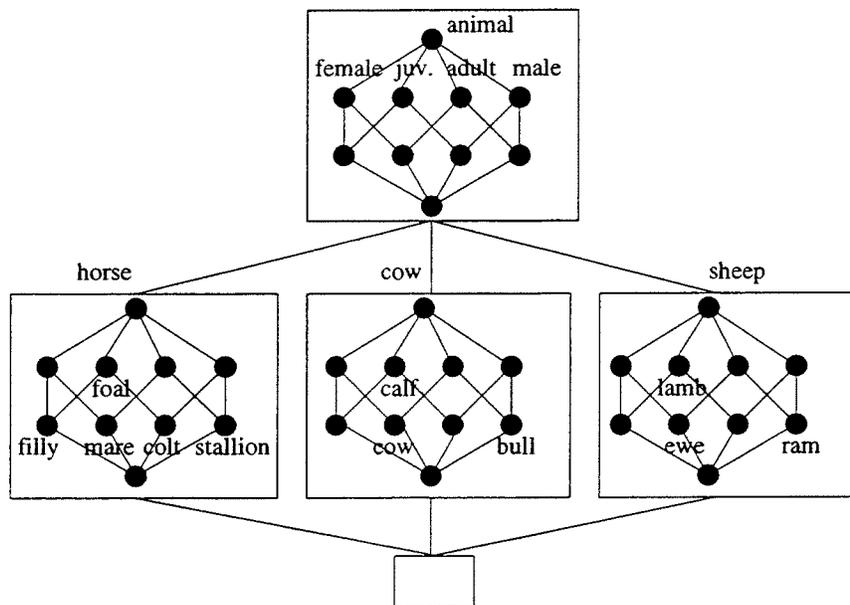


Figure 3. A nested line diagram.

The structure of a faceted thesaurus

Using formal concept analysis, the constituents and relations of a faceted thesaurus can be formally defined (see Priss (1996) for a more detailed discussion on linguistic thesauri). Each

of the following definitions consists of a structural component that describes the mathematical and formal properties and a semantic component that describes the nature of the objects, attributes and terms involved. Without a formal ontology, which would exceed the scope of this paper, the semantic components will necessarily remain somewhat vague.

Definition 1

A *thesaurus* consists of a concept lattice (or a set of concept lattices), a set of terms and a set of binary relations with the following conditions:

a) Its formal concepts are called *classes*. Terms are “names of classes”. Each class is denoted by at most one *preferred term*. Further terms can belong to each class. These are *synonyms* of the preferred term. The relation among terms that is induced by the subconcept/superconcept relation among the classes (or concepts) is called *broader term/narrower term relation* (BT/NT). In addition to that relation, further binary relations can be defined among the classes (or concepts). These binary relations introduce *semantic relations* among terms, which means that, if a relation holds between two terms, then the same relation must hold between any synonyms of the two terms. Some of these relations are given specific names, such as *part/whole relations* (BTP/NTP); other semantic relations are usually called *related term relations* (RT).

b) Formal objects, attributes and their relations should be *denotative*, which means that objects and attributes are instances or prototypes of “real” objects, processes, events, time units, qualities, etc., and the relations among them depend on what is denoted (pointed to) and not on what is connoted (implied) by the verbal expressions. Sufficient objects and attributes should be chosen so that the relations BT/NT, RT, etc. correspond to the commonly accepted relations in the domain to which the thesaurus is to be applied.

Denotative attributes of “dog” could, for example, be “barks” or “is a carnivore”, whereas “is a three letter word” or “is sometimes perceived as a threat” are not denotative attributes of “dog”. The formal objects and attributes are usually omitted in the final display of a thesaurus. But, if they are utilized during the design phase of a thesaurus, they should be stored for future modifications. Instead of a name, a concept can have a description composed of its objects and/or attributes. In this case, there is no term attached to the concept in the line diagram. These descriptions correspond to what ISO 2788 calls “node labels”: dummy terms that are not assigned to documents when indexing, but are inserted into the systematic display to establish relationships.

Definition 2

A *facet* consists of a concept lattice (or a set of concept lattices), a set of terms and a set of binary relations with the following conditions:

a) Its formal concepts are called *foci*. Terms are “names of foci”. Each focus is denoted by at most one [*preferred*] *term*. Synonyms and relations are defined

exactly as they are for a thesaurus. Facets can be nested such that a facet can be composed of smaller facets. A *base-line facet* is a minimal facet that consists of a set of concepts (or foci) and their common superconcept. A base-line facet must fulfill the conditions that the concepts are immediate subconcepts of the superconcept in the original facet and that their extensions are disjunct in the original facet.

b) A facet should represent a viewpoint on, or aspect of, the objects of a domain. This means that the set of terms of a facet should be selected according to a common criterion, such as “age”, “size”, “user preference” or “value”. The set of foci should represent that aspect or viewpoint as evenly and completely as possible and at the same level of specificity. The set of foci “very tiny”, “tiny”, “very little” “little”, “normal”, “huge” is incomplete because “large” is missing. It is also uneven because the facet is more detailed at the small end of the scale than at the large end.

Sometimes a facet is referred to by the name of the uppermost focus. For example, the facet on the right-hand side of Figure 2 could be called “animal facet”. On the other hand, we prefer to name a facet by a criterion for its formation and not by the uppermost focus. For example, the facets in Figure 2 should be called “gender/age” and “type of animal” and not “living being” and “animal”. It is important to clearly distinguish names of facets, foci, objects and attributes. An abstract facet could have a uppermost focus “gender/age” and subfoci “male gender”, “juvenile age”, and so on. But if abstract facets are applied as they are in Figure 2, “gender/age” and “type of animal” are neither attributes nor names of foci, but criteria for dividing the attributes of the combined contexts in Figure 1 into these two sets.

Definition 3

A *faceted thesaurus* is a thesaurus that is decomposed into a set of facets.

Depending on how the thesaurus is decomposed or on how the facets are to be combined to generate the thesaurus, several types of faceting can be distinguished. We suggest three types of faceted thesauri, but we do not intend to claim that this is an exhaustive analysis of the possible types of faceting.

First, facets can be sublattices of the faceted thesaurus that result from selecting subsets of the sets of attributes, objects or terms. For example, terms could be divided into “novice”, “intermediate” and “expert” knowledge and facets that are restricted to one of these sets could be generated. Another example would be to combine several languages or sublanguages (jargons) of one language in a faceted thesaurus and then decompose it into one facet for each language or sublanguage.

Second, facets can result from partitioning objects and/or attributes into disjoint sets (or partitions) and then creating relations between the partitions. For example, the temporal objects “1980” to “1989” could be combined into one object “the 1980s” and so on. In

this case, objects and attributes are not selected but are simply re-grouped. The resulting facets would be more general than the original facet. It is important, however, according to Definition 2, that all objects and attributes of a facet exist at the same level of specificity.

Third, all facets of a faceted thesaurus can have the same set of objects, but their sets of attributes and terms are disjoint. This is the type of faceted thesaurus that is described in the following sections of this paper and corresponds to the process of “scaling” in formal concept analysis. For example, people can be distinguished according to their gender, size, or hair color. Each of these aspects – “gender”, “size” or “hair color” – represents a different facet. (These aspects are called “multi-valued attributes” in formal concept analysis.) Different facets can be combined via their objects or attributes and can be represented graphically as nested line diagrams. The same set of objects can have different intensions and can be associated with different terms in different facets. For example, “man with blond hair” and “tall person” can refer to an identical set of objects that is the extension of different concepts in different facets. Each attribute and each term must belong to exactly one facet. If the same term should be used in two facets, it must be disambiguated by adding the facet name in parentheses behind the term.

Design of a faceted thesaurus

Traditional thesaurus construction (Aitchison & Gilchrist, 1987; Batty, 1989) consists of the following steps: 1) define scope and usage of the thesaurus; 2) identify the sources of the vocabulary; 3) collect terms; 4) form classes and the relational structure; 5) translate terms into standardized format; 6) choose a display format; 7) develop a notational structure (optional). The approach described in this paper is primarily concerned with the fourth step in this process and, to a lesser degree, with the third step.

According to Batty (1989), the following bottom-up approach should be used in the formation of classes and the development of a relational structure: terms that have been collected and recorded on paper slips (or an electronic equivalent) are grouped into clusters by forming piles of slips indicating related terms. The approach advocated here differs from Batty’s in that it combines a bottom-up strategy based on an analysis of objects (extensions) of classes or concepts (term identification) with a top-down strategy that utilizes general knowledge about a generic class hierarchy.

object	storage media	recording techniques	content format	storage format
Mozart’s Flute Concert	CD	stereo	music	electronic
Map of Indianapolis	sheet of paper		map	print
Peirce Manuscripts	CDs		text	electronic
Shakespeare’s writings	book series		text	print

Figure 4. A sample of document types.

As a first step, several typical objects to which the thesaurus is to be applied are selected. This set of objects need not be complete in any sense. As an example, Figure 4 shows four

objects that represent documents in a library collection. Typical facets and the values that the objects have in these facets are then determined. While it is not necessary to select a complete set of objects, it is necessary to select a sufficient number of facets with a sufficient number of values such that the objects are distinguishable by the values and that the intensions of the classes can be adequately described by the facets. At this point, it is necessary, according to Definition 2, to review each column in the table in Figure 4 to ascertain whether the terms in the column belong to the same aspect (or facet) of the objects. This approach keeps the selection of terms independent of the construction of the conceptual ordering. Replacing a term by a synonym does not change the conceptual structure. Furthermore, "node labels" can be added at any point. This means that, if desired, attributes can be added for the sake of completeness even if there are no terms for these attributes commonly used in the domain to which the thesaurus is to be applied.

In the next step, a top-down viewpoint is adopted. The facets are abstractly modified with the help of CODA [COnceptual Design Application], a graphical computer interface which we are developing for this purpose. CODA facilitates the editing of classes and their relationships, the composition and decomposition of facets, and the determination of relationships between the different facets that can be combined in nested line diagrams. Figure 5 shows the first output that CODA generates from the table in Figure 4. In these examples, the phrase at the top of each facet is not an attribute but the name of the facet. All terms are written under their concepts and the surrounding lines are omitted. The bottom nodes in the lattices are usually omitted since no objects can have contradictory attributes.

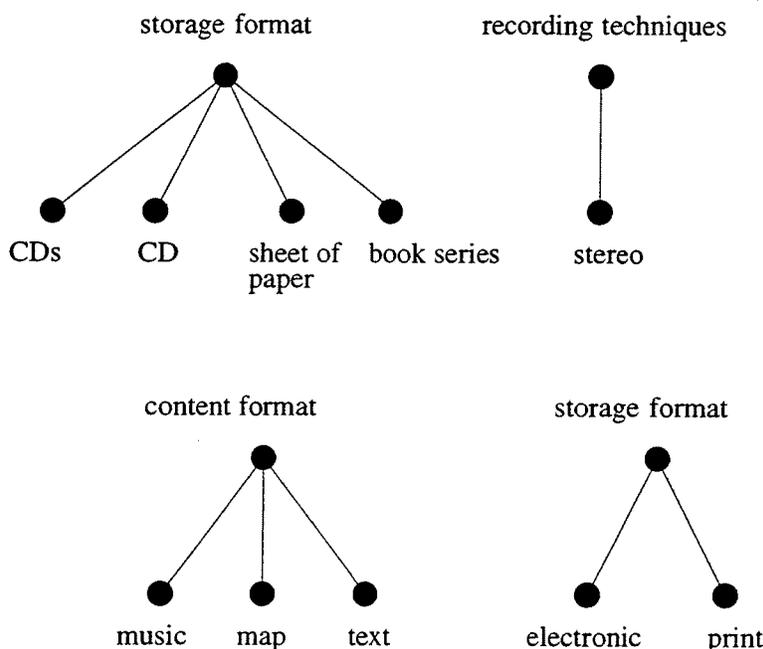


Figure 5. Automatically generated facets.

In Figure 5, CODA has generated a hierarchy template for each column of the table in Figure 4. The column headings in Figure 4 become the names of the facets in Figure 5 where they are written above the top nodes. The attribute values in Figure 4 become terms in the facets

in Figure 5. Each concept hierarchy represents one facet. Since the objects from Figure 4 are omitted, an abstract (top-down) viewpoint is taken and the developer can now modify the facets manually. Facets are stored in a database and can be reused so that common facets, such as facets for time, number, size, etc., can be selected from the database and do not have to be generated anew for each application. The same facet template can be used several times for one application as long as the terms are disambiguated by the facet names. For example, the same template can be used for “year of author’s birth” and for “year of publication” if they are disambiguated such as “1966 (birth)” and “1966 (publ)”.

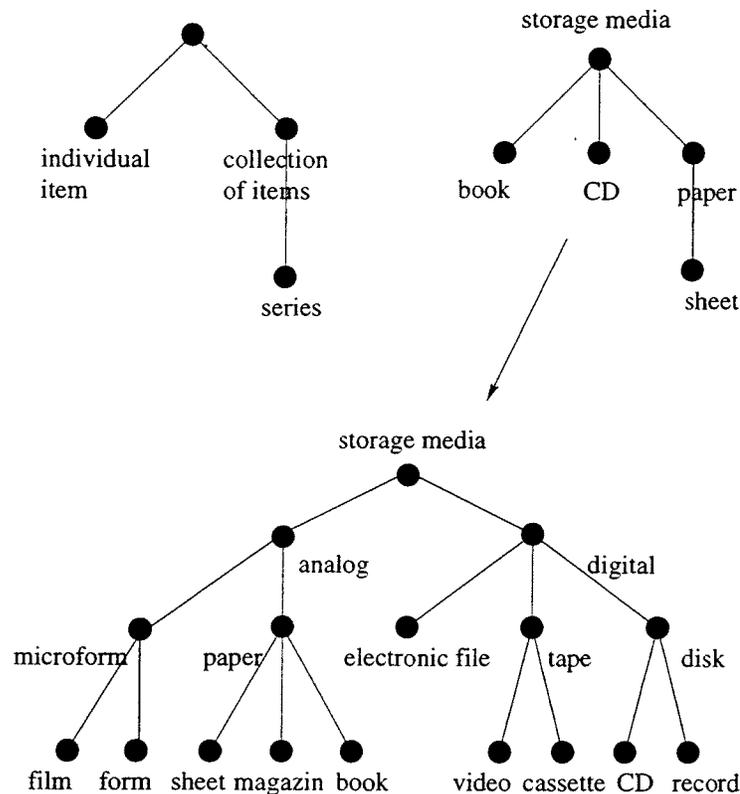


Figure 6. Manual editing of facets: first version.

Obviously, the facets in Figure 5 are not yet complete. Figure 6 shows the developer’s first attempt to improve the initial facets. The developer notices that individual items (singular forms) and collections of items (plural forms) for the same terms should probably not appear in one facet since they represent different aspects. Therefore, the original facet is split into two separate facets. The other facets are then edited and combined. Although not all of the facets are shown in Figure 6, the bottom example indicates an attempt to combine several aspects into one facet. There is, however, at least one obvious error in this display: CDs can be analog or digital. Switching to a data-driven approach that applies the facets to a set of objects, such as those identified in Figure 4, allows the developer to identify the problem. The developer must now modify the facet structure further. Figure 7 shows the result of this action: these facets do not reveal any inconsistencies when applied either to the objects from Figure 4 or to a larger set of objects. Depending on the set of objects, however, the facets may still be incomplete and thus require further additions or modifications.

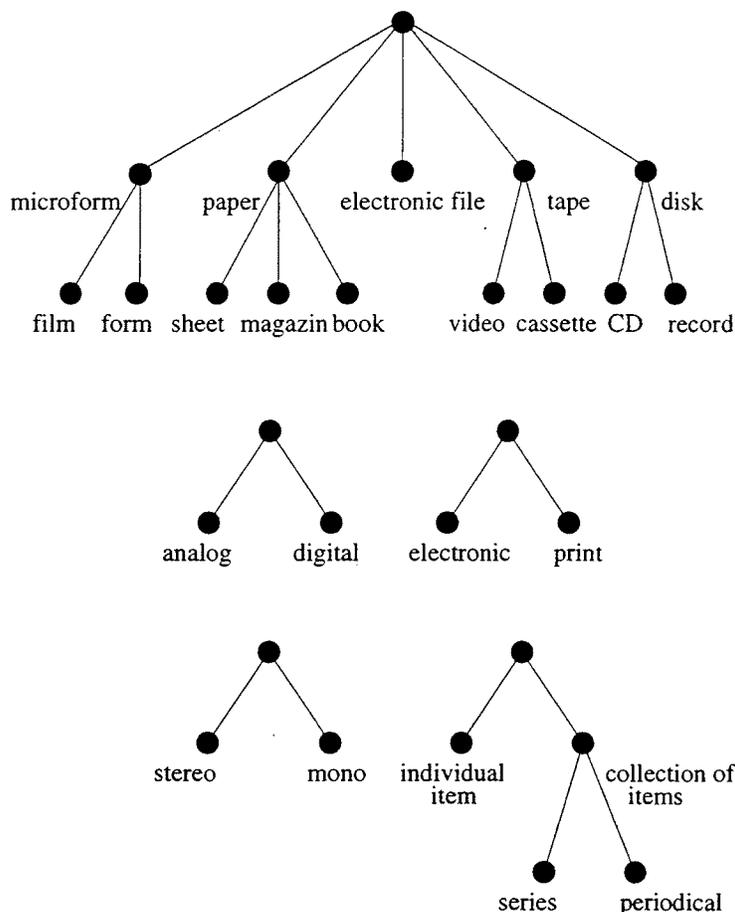


Figure 7. Manual editing of facets - final version.

The combination of two facets in a nested line diagram results in the “direct product” of the concepts represented in the individual facets. For example, the combination of the two facets “analog/digital” and “electronic/print” from Figure 7 results in four times four concepts. This is demonstrated in the left-hand diagram in Figure 8. The right-hand diagram in Figure 8 shows the same lattice as a nested line diagram. Here the bottom nodes have been omitted because their extensions are empty. Other concepts in the direct product of facets can also have empty extensions. For example, can an object be “digital” and “print” at the same time? While it would be possible to print out on paper the 0’s and 1’s of an electronic file, it is not likely that such a document would be part of a library collection. Therefore, since such a document does not exist as an object in the domain of the thesaurus, a concept that has the terms “digital” and “print” would have an empty extension. The software CODA allows the developer to omit such “dummy nodes” by restricting a concept in one facet to another concept in another facet such that, if these two facets are combined, the first concept and all its subconcepts apply only to the second concept and its subconcepts and not to any other concept in the second facet. In Figure 9, “digital” is restricted to “electronic”, which means that all objects that are “digital” must be “electronic”. The bottom half of Figure 9 shows the normal line diagram (left-hand side) and the nested line diagram (right-hand side) for the two facets given these restrictions. Compared to Figure 8 the diagrams are considerably simplified.

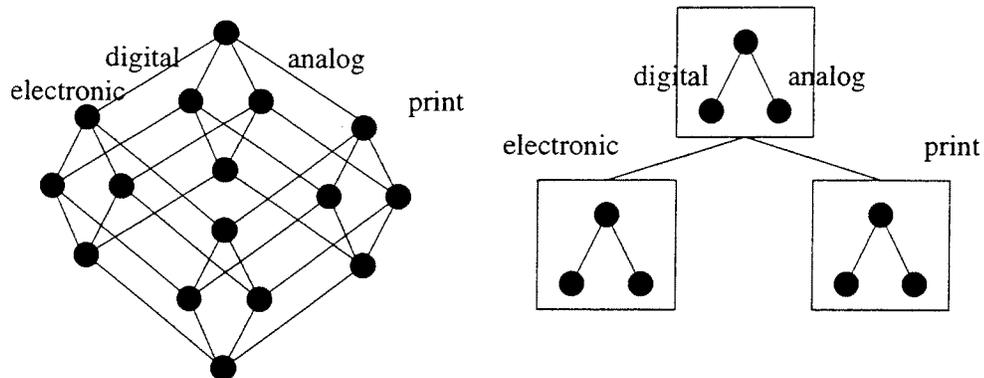


Figure 8. A direct product of two facets (left: normal line diagram, right: nested line diagram)

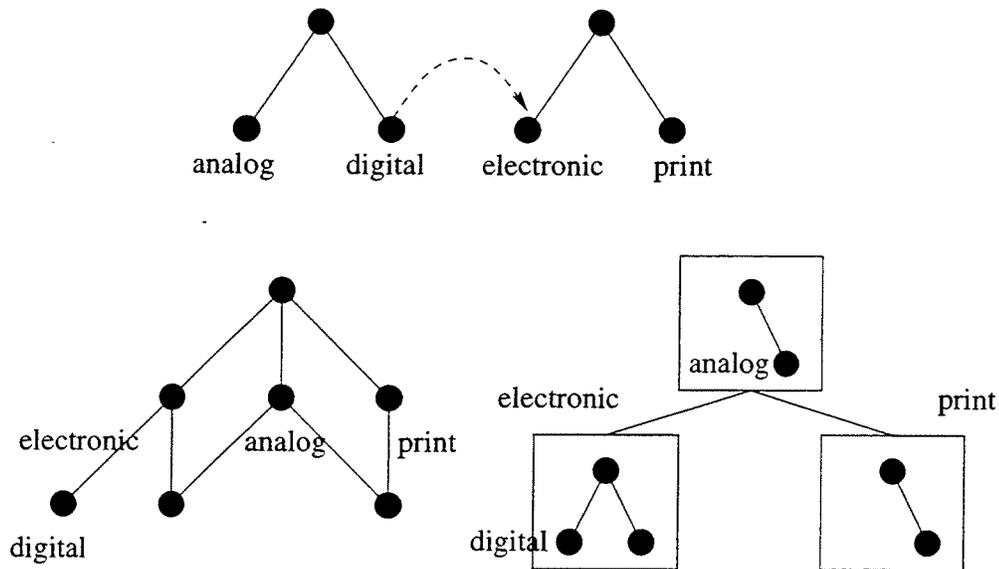


Figure 9. A combination of facets with respect to restrictions.

It should be pointed out that there are still many questions to address in an attempt to formalize faceted thesauri. This paper presents a first attempt at a formal definition of thesauri and facets in the context of formal concept analysis. The approach presented in this paper is directly applicable to thesaurus design and construction. The next step would be to perform usability testing of the system to determine how much training is required for developers and end-users before they can interpret and apply the diagrams appropriately. The system would then be compared to other interfaces for faceted thesauri.

As set out in this paper, there are at least three different types of faceting that can be applied to the terms in a domain. The examples in this paper represent only one type. Future theoretical research must also be concerned with other types of faceted thesauri. We believe that the first type of faceting, which would select terms according to different languages or user groups is especially relevant in the present interdisciplinary environment and should have numerous applications.

References

- Batty, David (1989). *Thesaurus Construction and Maintenance: A Survival Kit*. Database, Vol. 12 (1), 13-20.
- Craven, Timothy (1990). *Automatic Structure Modification in the Graphic Display of Thesauri*. Advances in Knowledge Organization, 1, 146-153.
- Fischer, Dietrich (1993). *Consistency Rules and Triggers for Multilingual Terminology*. TKE'93: Terminology and Knowledge Engineering. INDEKS-Verlag. 333-342.
- Ganter, B.; Wille, R. (1996). *Formale Begriffsanalyse: Mathematische Grundlagen*. Springer-Verlag.
- ISO 2788 (1986). International Standard 2788: *Documentation - Guidelines for the establishment and development of monolingual thesauri*. International Organization for Standardization.
- Johnson, Eric H. (1995). *A Hypertext interface for a Searcher's Thesaurus*.
<http://csdl.tamu.edu/DL95/>
- Kent, R.; Neuss, C. (1995). *Conceptual Analysis of Resource Meta-Information*. Computer Networks and ISDN Systems , 27, 973-984.
- Pollitt, Steven (1997). *Interactive Information Retrieval based on Faceted Classification using Views*. Proc. of the 6th Int. Study Conf. on Class., London, June 1997, FID, 51-56.
- Priss, Uta (1996). *Relational Concept Analysis: Semantic Relations in Dictionaries and Lexical Databases*. Dissertation, Technische Universitaet Darmstadt.
- Skorsky, Martin (1997). *Graphische Darstellung eines Thesaurus*. Deutscher Dokumentartag, Regensburg.
- Stumme, Gerd (1996). *Local scaling in conceptual data systems*. In: Eklund, Ellis & Mann (eds.). Conceptual structures: Knowledge representation as interlingua, Vol. 1115 of LNAI. Springer-Verlag, Berlin-Heidelberg, 121-131.
- Viehmann, Viola (1996). *Formale Begriffsanalyse in der bibliothekarischen Sacherschliessung*. Master's thesis, TH-Darmstadt, 1996.
- Vogt, F; Wille, R. (1995). *TOSCANA - a graphical tool for analyzing and exploring data*. In: Tamassia; Tollis (eds.). Graph Drawing. Springer-Verlag, Heidelberg, 226-233.
- Wille, Rudolf. (1997). *Conceptual Landscapes of Knowledge: A Pragmatic Paradigm of Knowledge Processing*. Proc. KRUSE Conference, Vancouver.