# Retrieval from facet spaces

R. B. Allen

*Bellcore, 445 South Street, Morristown, NJ 07960, USA, rba@bellcore.com*

**SUMMARY**

**The *facet-space* approach for accessing document records organized by faceted classifications is described. The interface gives users detailed control over the facet display and it makes use of color to reduce the number of windows which need to be presented. The interface supports searching. A cluster analysis is described for organizing search return lists based on facets distances. The implementation is applied to 1381 summaries of computer science dissertations as organized by the *ACM Computing Reviews* classification system.**

## 1. INTRODUCTION

### 1.1. Digital Libraries, Classification, and Interfaces

One value of structuring document collections is for guiding browsing. Not only do the semantics of the classification system let the user identify topics of interest, but the classification system generally locates related documents near to each other. While document records have been available online for many years, there has not been much work on GUIs for OPACs (Online Public Access Catalogs) until recently. One of these is the HOPAC (hierarchical OPAC) [1,2] for the Dewey Decimal System (DDC) which introduced several novel capabilities. The HOPAC took advantage of the structure of the classification system in the interface. It also integrated a virtual-shelf view of the collection and allowed users to select attributes (e.g., year of publication and libraries which hold the document) of the document records to be displayed on the shelf. When used in conjunction with search, the interface allows the user to restrict the scope of the search and to post hits against the classification hierarchy.

### 1.2. Faceted Classifications and Thesauri

While simple hierarchical classification systems such as DDC are the basis for organizing many document collections, they are not well-suited for representing collections with multiple dimensions. Thus, many specialized collections have adopted faceted classification schemes and thesauri. Faceted thesauri and classifications
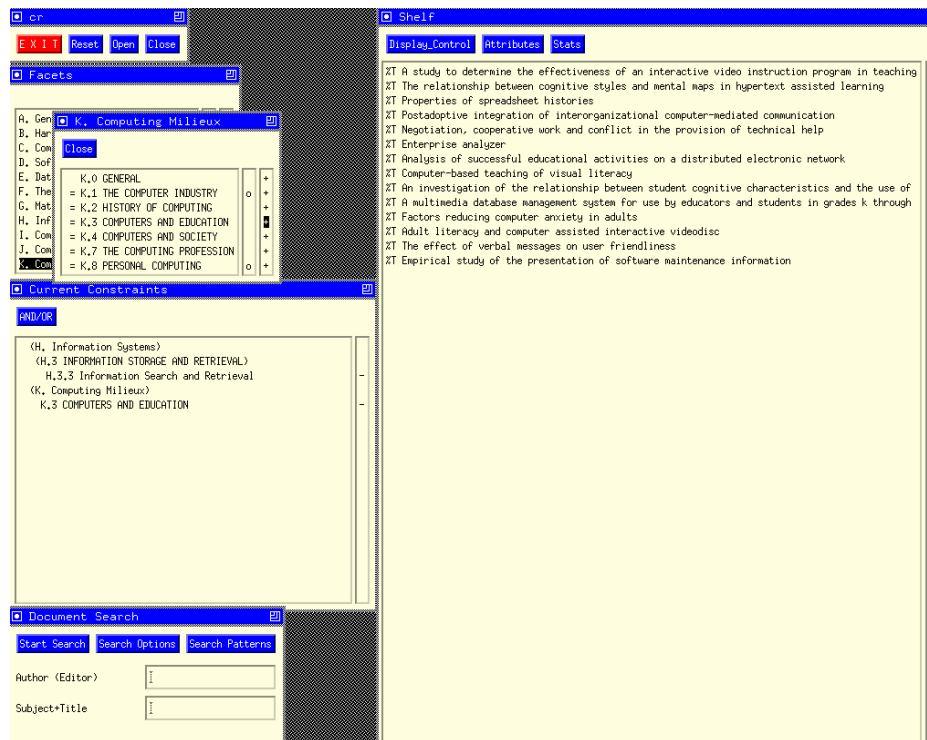
*Figure 1. Cascaded-Menu Interface.*

(e.g., [3,4]) are quite common, including: AAT, EI, ERIC, INSPEC, LISA, and MeSH.

Faceted classifications are hierarchical with individual nodes having "Broader Terms" (i.e., parents) and "Narrower Terms" (i.e., children). In addition, facet classifications often also include "Used For" text descriptions which describe how a term is applied and "Related-To" links which connect facets across the hierarchies. Faceted classifications are often used to organize controlled vocabulary for thesauri; that is, the facet node labels are used a keywords. The expression "facet system" is used here to describe the combination of the faceted classification and an associated document collection. In faceted systems, documents (or other information objects) may appear under several facets, but not every document is represented on every facet. The browsing interface described here provides both powerful browsing tools, but also an interrelated shelf. The variety of features in a faceted system gives it power but also adds complexity for browsing and searching document collections.

An individual facet is a tree and a faceted classification system is a collection of trees. These basic tree structures are extended by Related-To links and "poly-hierarchies". Polyhierarchices are facet nodes with multiple parents and are not explored in this interface. Documents classified with a facet classification may also have multiple parents (i.e., they may be assigned multiple facet nodes). Thus, the document browsers described here essentially may be considered to support browsing of objects with multiple inheritance.

The details of implementation of facet classifications and their application to document collections as facet systems differ greatly. These differences can be important for interface design; among the relevant ways they may differ are:

- Whether the facets show hierarchical structure as part of their labels (e.g., numbers or letters).
- Whether facet labels can be easily identified independently from their context in the facet hierarchy. If facets cannot be understood separately, then it is helpful to display clues such as the the labels of the parents.
- Differences in breadth and width of the classification. If the facet hierarchy is both broad and deep, a display of the opened facets may overfill the screen.
- Whether there are explicit "General" categories at each node. As Godert [5] has pointed out, there may be cases in which inheritance to children is not desired, as when an overview is required, rather than a detailed examination based on some subnode topic. For instance, a user may be interested only in comprehensive documents covering all the subnodes and not documents focusing on specific subnodes.
- Whether the facets are assigned consistently (e.g., professionally) to documents.
- Whether specific dimensions identify the relationships by which they are are subdivided (e.g., whole-part relationship, generic relationship, instance relationship).

For the prototype interface developed here, the *ACM Computing Reviews* (CR) classification for the computer-science literature was used [6]. The CR Classification has a few "Related-To" links, but it does not have explicit "Used For" descriptions. The document records used to test the interface were 1381 doctoral dissertations cited in the *ACM Computing Archive* [7] as published in 1992.

### 1.3. Cascaded-Menu Interface for Faceted Classification

A cascaded-menu interface for a faceted classification system was described briefly in [1,2]. Figure 1 shows this interface applied to the CR Classification and the corpus of document records. This interface had three major components: cascaded-facet menus, constraint lists, and a document shelf. Major toplevel categories are chosen from the Cascaded Menus at the upper left of Figure 1. These selections open cascaded menus which display lower-level categories. When the "+" to the right of the facet label is selected, the facet is added to the Current Constraint List (left middle in Figure 1).

The Constraint List presents the facets which are active, and the Shelf is updated with articles that match the constraints. To show the context of the selected constraint labels, the parents of the constraints are displayed in parentheses on the Constraint List. The constraints propagate to all their descendants. Constraints could be dropped from the Constraint List by clicking on the "−" on the right side of the widget.

The cascaded-menu interface had a number of difficulties. Several of these difficulties revolved around the design of the menus so that browsing was cumbersome.
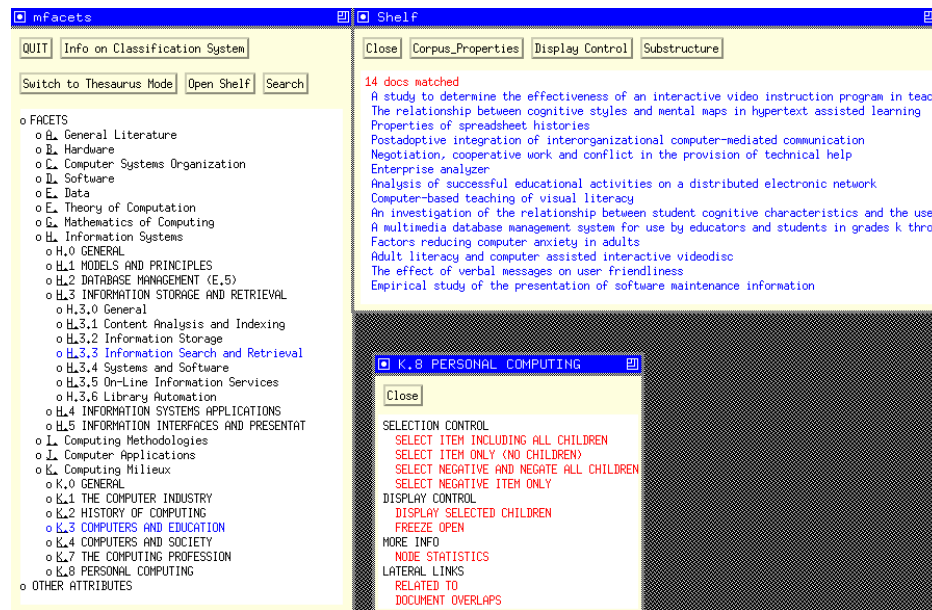
mfacets

QUIT   Info on Classification System

Switch to Thesaurus Mode   Open Shelf   Search

o FACETS
  o A. General Literature
  o B. Hardware
  o C. Computer Systems Organization
  o D. Software
  o E. Data
  o F. Theory of Computation
  o G. Mathematics of Computing
  o H. Information Systems
    o H.0 GENERAL
    o H.1 MODELS AND PRINCIPLES
    o H.2 DATABASE MANAGEMENT (E.5)
    o H.3 INFORMATION STORAGE AND RETRIEVAL
      o H.3.0 General
      o H.3.1 Content Analysis and Indexing
      o H.3.2 Information Storage
      o H.3.3 Information Search and Retrieval
      o H.3.4 Systems and Software
      o H.3.5 On-Line Information Services
      o H.3.6 Library Automation
    o H.4 INFORMATION SYSTEMS APPLICATIONS
    o H.5 INFORMATION INTERFACES AND PRESENTAT
  o I. Computing Methodologies
  o J. Computer Applications
  o K. Computing Milieux
    o K.0 GENERAL
    o K.1 THE COMPUTER INDUSTRY
    o K.2 HISTORY OF COMPUTING
    o K.3 COMPUTERS AND EDUCATION
    o K.4 COMPUTERS AND SOCIETY
    o K.7 THE COMPUTING PROFESSION
    o K.8 PERSONAL COMPUTING
o OTHER ATTRIBUTES

Shelf

Close   Corpus_Properties   Display Control   Substructure

14 docs matched
A study to determine the effectiveness of an interactive video instruction program in teac
The relationship between cognitive styles and mental maps in hypertext assisted learning
Properties of spreadsheet histories
Postadoptive integration of interorganizational computer-mediated communication
Negotiation, cooperative work and conflict in the provision of technical help
Enterprise analyzer
Analysis of successful educational activities on a distributed electronic network
Computer-based teaching of visual literacy
An investigation of the relationship between student cognitive characteristics and the use
A multimedia database management system for use by educators and students in grades k thro
Factors reducing computer anxiety in adults
Adult literacy and computer assisted interactive videodisc
The effect of verbal messages on user friendliness
Empirical study of the presentation of software maintenance information

K.8 PERSONAL COMPUTING

Close

SELECTION CONTROL
    SELECT ITEM INCLUDING ALL CHILDREN
    SELECT ITEM ONLY (NO CHILDREN)
    SELECT NEGATIVE AND NEGATE ALL CHILDREN
    SELECT NEGATIVE ITEM ONLY
DISPLAY CONTROL
    DISPLAY SELECTED CHILDREN
    FREEZE OPEN
MORE INFO
    NODE STATISTICS
LATERAL LINKS
    RELATED TO
    DOCUMENT OVERLAPS

*Figure 2. Facet-Space Interface with Two Terms Selected (H.3.3 and K.3).*

The cascaded facets did not allow the user to have two menu labels open from various parts of the hierarchy at one time. There was no way to inhibit propagation of selections to children (see [5]); nor did the system let the user select the negatives of facets such that *no* documents with those facets would be included. Having the Constraint List detached from the facet hierarchy meant the user could not easily see the relationship between the menu items and the constraints. Other problems were associated with the Shelf and search specifications. On the Shelf, document titles were not ordered. When there were too many documents, it was difficult to find anything useful. Fractional search hit counts were assigned according to the number of facets in which the hit documents were included. However, this appeared to be confusing for users.

## 1.4. Overview of Facet-Space Approach

The facet-space interface was developed to remedy some of the problems of the cascaded-menu interface. Integration for faceted systems may be enhanced in many ways such as search, graphical views, better integration of constraints, and better statistics. As in the cascaded-menu interface, the documents on the shelf are those which match multiple constraints. In these, a mixture of classification system browsing, thesaurus-term identification, and object retrieval were integrated into one system.

Figure 2 shows the basic interface with two facets opened (**H** and **K**), two subfacets selected (**H.3.3** and **K.3**), the Shelf with the documents matching those facets, and Facet Options for another subfacet (**K.8**). The interface as shown in Figure 2 has been implemented in the X-Window System. Some peripheral features,

such as the cluster analysis shown in Figure 4, have not been integrated into the implementation.

## 2. FACET-DISPLAY WIDGET

The Facet-display widget is on the left side of the figure. Two items, shown in blue but not visible in a black-and-white figure, have been selected. Other items the user is viewing are shown in black. In this interface, the complexity was controlled by leaving out those facets that are not active and introducing fine-grain control over the display of items.

### 2.1. Facet-Display Widget Buttons

Several buttons are provided at the root widget. The *Info on Classification System* button provides details about the classification system; in the current implementation, the CR system would be described. The *Switch Mode* button switches between Thesaurus Mode and Document Mode (see Sections 4, 6). The *Open Shelf* button controls whether the Shelf (see Section 3) is displayed. The *Search* button brings up the search widget (see Section 6).

### 2.2. Facet Labels

Clicking on a facet opens its children if they are closed. To conserve space and confusion, opening children also causes all non-selected parents to close. Thus, only open items or selected items are displayed. If the children are already open and there are no selected terms, clicking on the label closes the children. The facet labels were truncated so the facet display would not cover the Shelf.

The toplevel of the facet space is indicated by the term "FACETS". This has the same controls provided to the individual facets. "OTHER ATTRIBUTES" can also be used to sub-select the documents on the shelf and provide further restrictions on the shelf display (see [2]). Attributes include the year of publication and the journal in which the article was published. Of course, at least one of the attributes must be active for each item (e.g., each document has a year of publication). Other attributes (not implemented here) could include institution of origin, collections to be included, and the type of document (e.g., journal article, masters or doctoral thesis). Of course, several of these attributes could be hierarchically organized and could be managed similarly to the facets.

The presence of child facets is shown by an underscore at the beginning of the facet label. The underscore is doubled in length if any of the children have been selected (Section 2.3.1).

### 2.3. Facet Options

Menu-element icons are displayed beside each item and clicking on one of these opens options for that item. The four groups of options are described in the following subsections. The options are tailored to the classification system and item. For instance, leaf nodes do not have options for opening and closing children. In the

bottom center of Figure 2, the Facet-Option widget for facet **K.8** has been opened.

### 2.3.1. Facet-Selection Control

The first group of options allows different ways of selecting the item as a constraint for the shelf display. There are positive (all documents which possess the attribute are included) and negative selections (no documents which possess that attribute are included). For both positive and negative selections, it is possible to have either item-only or inheritance of selections.

Although the facets are selected graphically, they essentially implement Boolean constraints on the documents. Thus if two unrelated constraints are selected, the documents they return must AND the selected facets as constraints. The constraints are more complex when facets are selected with all their children. In that case, the system ORs of facets with their children and the AND with other facets. Of course, still more complex arrangements are possible (but not supported in the current interface). such as ORing facets of "Related-To" terms.

Color gives an extra dimension of the display and it is useful to characterize state of the interface. In this interface, color is used to indicate which options are have been selected and, thus, replaces the Constraint List in the cascaded-menu interface. When selected, the item changes color (red or blue in the current implementation) to provide rapid visual indication of the selected facets.

### 2.3.2. Display-Control Options

The second group of options controls the details of the display. The "Ellipsis" label controls whether all siblings of selected widgets are shown. Ellipsis would be necessary if too many items were displayed and they scrolled down off the screen. The "Close Children" label blocks display of the lower-level selected items. The "Freeze Open" label overrides the default closing of the labels when lower levels are open.

### 2.3.3. More-Information Options

The third group of options allows the user to access more information. The "Node Stats" option provides information about how many child nodes a given node has and how many of them are selected. Node Statistics also include measures of the impact of the given node on the constraints of the document display. If a node is selected the statistics include the ratio of the number of items excluded from the display to the total because of that item. "Used For" information could also have been included as a choice, if it had been available for this facet classification system.

### 2.3.4. Lateral-Links Display Control

The fourth group of options allows the user to access lateral links such as "Related-To" items and "Overlaps" (see Section 5). It is possible to specify the ways in which the terms are related. Figure 3 shows the display for terms that overlap with **H.3.3**, which is indicated with a <-. The nodes with which there is at least one overlap to
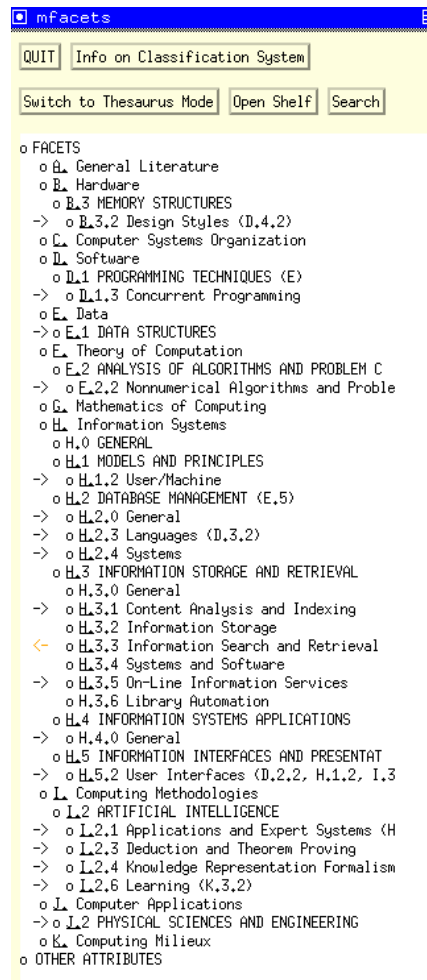
*Figure 3. Facets with Overlaps to H.3.3.*

**H.3.3** are shown with ->. In addition, the parents of the overlapping facets have also been opened. In this implementation, facets that had even one overlap were marked. An extension of the current interface might add a slider for controlling the granularity of overlaps to be displayed (see the HOPAC interface in [2]). Because there were not too many overlapping facets, it was possible to display all of them in Figure 3. An alternative to using arrows to indicate overlaps would be to post "Lateral-Link" hits on the facets, analogous to the posting of search hits in the SuperBook$^{TM}$ and cascaded-menu interfaces.

*Figure 4. Results of a Search on the Word "Language" in Document Mode.*

## 3. SHELF WIDGET

### 3.1. Document List

The books satisfying the constraints of the selected facets are displayed on the Shelf with the total number of filtered documents presented on the first line.

There are often so many documents matching the constraints that the user might become lost easily. Thus some type of structure could be generated for organizing the documents on the shelf. The HOPAC interface had books arranged in a linear order following the DDC. However, for the facet systems, there is no analogous default layout [8]. One solution would be to provide additional structure to the shelf. For instance, the shelves could be organized by some attribute such as chronological order. Another possibility would be to synthesize an ordering by subfacets (see Section 6).

### 3.2. Shelf Buttons

The *Corpus Properties* button provides details about the document collection. In this interface, a description of the *Computing Archive* and, the subset of selected documents would be described.

The *Display Control* button allows the user to set which fields of the document records are displayed. The title is always displayed, but the user can decide whether other fields such as the author and the publisher are shown.

The *Substructure* button controls the presentation of the categories of the items displayed. It shows other nodes among selected terms. For single facet selections, this is the same as the overlap measure (Section 5).

### 3.3. Document-Record Display

Additional information about the documents can be obtained by clicking on their titles. This causes a summary of the document characteristics to be displayed.

## 4. SEARCHING ON FACETS IN THESAURUS MODE

The interface supports access to facets as would be used in a typical thesaurus. For example, a user might use the interface to find thesaurus terms [9]. Because there are no documents in this mode, no shelf is present. Moreover, there are different Facet Options than for browsing documents. For instance, there is no negative selection and those options do not appear in the Facet-Options list in Thesaurus Mode. There are no attributes of the thesaurus, so the Attributes label does not appear on the facet list. For keyword matches, the outputs of the search hits are easily displayed in Thesaurus Mode. The relevant nodes of the facet space are opened and the hits are highlighted.

Beyond explicit "Related-To" links, for browsing, similarity between facets could be determined by tree distance. Tree distance is not a simple indicator of similarity. First, there is no clear way to measure the distance between facets. Moreover, tree-distance is complex for systems in which several facets are active. Alternatively, distance could be derived from semantic distance for facet terms or on the words in the "Scope Notes." Basing the metric on documents in the corpus would provide still better distance measures.

## 5. OVERLAP MEASURE OF DOCUMENT SIMILARITY

[2] briefly described a measure of similarity between facets. Because most of the documents are assigned to several categories, the overlapping categories are presented in the cascaded-menu interface by selecting the "o" from the first vector on the right side of the cascaded-menu widget (Figure 1). The categories that had two or more overlapping documents with **H.3.3 Information Storage and Retrieval** were **H.2.4 Systems**, **H.2.0 General**, **D.3.2 Design Styles**, **H.5.2 User Interfaces**, and **I.2.6 Learning**. A slightly different view is obtained from the graphical presentation of the overlaps in Figure 3 which makes apparent that the topics **H.2 Database Management** and **I.2 Artificial Intelligence** are closely related to **H.3.3**. The same overlaps mechanism may be extended to the case in which several facet nodes have been selected.

## 6. SEARCHING IN DOCUMENT MODE

In [2], searches were conducted on book titles in the Dewey hierarchy. While search is not usually effective on individual titles because they are too short, using all the titles collected under a node gave reasonable results. For this interface, keyword and term-weighted searches are currently implemented.

The SuperBook text browser introduced an effective technique of displaying search hits. Search hits were posted against the hierarchy; thus, the hierarchy provided structure to organize the return list. Analogously, search hits were posted
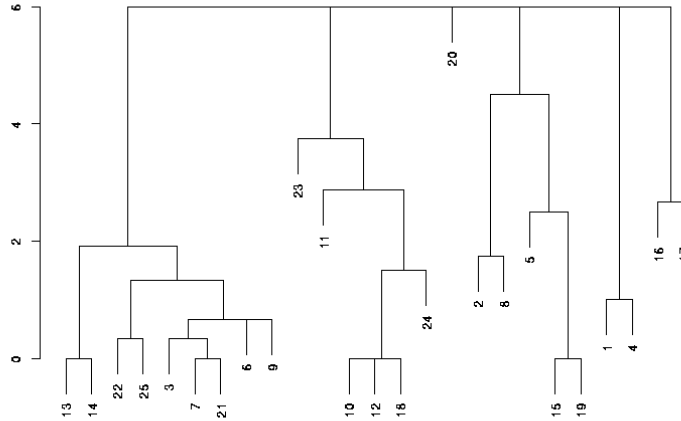
*Figure 5. Cluster Analysis of Search Return-List Documents Based on their Facet Distances.*

against the classification hierarchy for the HOPAC [2].

Displaying hits for a faceted classification is difficult because of assignment of documents to multiple facets. In fact, it may not be very meaningful to combine hits across an facets because the facets are meant to be independent dimensions. The cascaded-menu interface [2] used fractional category memberships when the hits are spread across categories.

One solution to the search-hits and the shelf-organization problem is to cluster return lists such as the one in Figure 4. Potentially, this will yield a coherent 2D view of the complex facet space. Figure 5 shows a cluster analysis of minimum mean tree distances for the documents returned in the query. The minimum mean tree distance was calculated as the minimum tree distance for each pair of facets of the two documents. Some of the documents in the collection had no facets assigned. Thus, of the 31 documents in the return list only 25 were able to be included in the cluster analysis.

The cluster dendrogram effectively provide "hints" to users about important subdivisions. Thus after the search, the tree shows there are six main headings. From left to right, these correspond to the categories Natural Language Processing (NLP), Computer Uses in Education, Logic, Requirements, Parallel Languages, and Query Languages. Moreover, the fine structure of the clusters also, generally, follows the facet structures. For instance, Documents 13 and 14 are both about NLP and "Arts and Humanities". Thus, the cluster results could be used as a structure for organizing documents on the Shelf. A graphical interface similar to [10] could be developed for browsing the clustered facets.

## 7. DISCUSSION

Some features of the faceted classifications and of the interface are complex and usability data would be helpful. Some features of the interface could probably be improved. For instance rather than clicking for menu selection, it would be easier to use a sliding action. It should also be useful to keep a history of the facet combinations selected so the user could return to previous configurations.

While the emphasis in this paper has been on browsing facets spaces for documents, the interface could be useful for accessing other types of hierarchical information (e.g., class hierarchies). It could also be used for creating and assigning facets. Versioning support would be necessary for accessing previous computer-science literature with this interface.

A number of enhancements are possible by introducing still more graphical elements. For instance in the facet-display interface (e.g., Figure 3), the pointsize of the labels could be manipulated so that less important labels were smaller [11]. There could also be purely graphical views of the hierarchy. For instance, color could be used to indicate hit-density [2]. A related visualization problem is occurs for browsing intersecting trees in hierarchical Web-browser hotlists [12].

Facet systems incorporate features between keyword systems and simple hierarchies. Thus, they have the advantages of both structure and of a rich semantics. The facet-display interface for exploring facet spaces has been described. Several new interface features have been developed for the Facet-Display Widget. It appears to be much easier to use and more powerful than the earlier cascaded-menu interface.

## ACKNOWLEDGEMENTS

## REFERENCES

1.  R. B. Allen, 'Navigating and searching in digital library catalogs', in *Digital Libraries '94*, pp. 95–100, (1994). College Station, TX.
2.  R. B. Allen, 'Two digital library interfaces which exploit hierarchical structure', in *Proceedings of Electronic Publishing and the Information Superhighway*, pp. 134–141, (May 1995). Boston.
3.  J. Aitchison and A. Gilchrist, *Thesaurus Construction: A Practical Manual*, Aslib, 1987. London.
4.  B.C. Vickery, *Faceted Classification*, Rutgers University Press, 1965. New Brunswick.
5.  W. Godert, 'Facet classification in online retrieval', *International Classification*, **18**, 98–109, (1991).
6.  ACM, 'ACM Computing Reviews classification system', *ACM Computing Reviews*, **35**, 4–44, (1994).
7.  ACM, *ACM Computing Archive*, ACM, 1994. New York.
8.  R.J. Hyman, *Shelf Access in Libraries*, ALA, 1982. Chicago.
9.  Schatz, B.R., Johnson, E.H., Cochrane, P.A., and Chen, H., 'Interactive term suggestion for users of digital libraries: Using subject thesauri and co-occurrence lists for information retreival', in *Proceedings ACM Digital Libraries*, (1996). Bethesda, MD.
10. Allen, R. B., Obry, P., and Littman, M., 'An interface for navigating clustered document sets returned by queries', in *Proceedings ACM SIGOIS: Conference on Organizational Computing Systems (COOCS)*, pp. 166–171, (Nov. 1993). Milpitas, CA.

11.  H. Koike, 'Fractal views: A fractal-based method for controlling information display',
     *ACM Transactions on Information Systems*, **13**, 305–324, (1995).
12.  Wittenburg, K., et al., 'Group asynchronous browsing on the world-wide web', in
     *World-Wide Web Conference 4*, (Dec. 1995). Boston.