

Integrating Information Seeking and Structuring: Exploring the Role of Spatial Hypertext in a Digital Library

George Buchanan
Middlesex University
Bramley Road
London, N14 4YZ
+44 (0) 20 8411 5939

g.buchanan@mdx.ac.uk

Ann Blandford
Harold Thimbleby
University College, London
Remax House, 31/32 Alfred Place
London, WC1E 7DP
+44 (0)207 679 5288/5204

a.blandford@ucl.ac.uk

h.thimbleby@ucl.ac.uk

Matt Jones
University of Waikato
Hamilton
New Zealand
+64 7 858 5174

mattj@cs.waikato.ac.nz

ABSTRACT

This paper presents Garnet, a novel spatial hypertext interface to a digital library. Garnet supports both information structuring – via spatial hypertext – and traditional information seeking – via a digital library. A user study of Garnet is reported, together with an analysis of how the organizing work done by users in a spatial hypertext workspace could support later information seeking. The use of Garnet during the study is related to both digital library and spatial hypertext research. Spatial hypertexts support the detection of implicit document groups in a user's workspace. The study also investigates the degree of similarity found in the full text of documents within such document groups.

Categories and Subject Descriptors

H.5.4 [Hypertext/Hypermedia]: User issues

General Terms

Design, Human Factors.

Keywords

Spatial hypertext, digital libraries, information retrieval.

1. INTRODUCTION

This paper introduces a combined spatial hypertext and digital library system, called Garnet. Garnet complements traditional digital library document retrieval tools with the structuring and coordinating support provided by a spatial hypertext system. Given the complementary roles of information structuring and information seeking – supported by spatial hypertext and digital libraries respectively – in physical environments, Garnet was created to explore the advantages of a combined system in an electronic environment.

Research on information seeking and use has observed that

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

HT'04, August 9–13, 2004, Santa Cruz, California, USA.
Copyright 2004 ACM 1-58113-848-2/04/0008...\$5.00.

information workers often use implicit cues that they place in their physical environment to support their own information seeking [11, 13]. Spatial hypertexts support similar activity in a digital environment, facilitating the organization of documents in both formal and informal ways [14] and patterns observed in the physical organization of written documents [13] recur in spatial hypertexts [17].

Digital libraries support information seeking, but focus upon providing their users with documents through a range of retrieval mechanisms such as indexed search and classification browsing. The retrieval of documents of potential interest is, however, only part of the wider information seeking process. Information scientists, e.g. Ellis [7] and Kuhlthau [12], emphasize the role of coordinating and selecting activities play in supporting document retrieval, and spatial hypertext facilitates a number of these tasks.

The combination of spatial hypertext and digital library raises the following questions:

- Can spatial hypertext provide an effective interface to a digital library?
- Will users demonstrate the interleaving of information seeking and structuring reported in physical environments?
- Can the information structuring performed by the user in the spatial hypertext be used to support information seeking? Particularly, can the theme of a document group be determined from the text of its members?

This paper presents initial outcomes from the study of these questions. The body of this paper proceeds in four parts. Firstly, an example scenario demonstrates Garnet in use. Secondly, we report on the design of Garnet, highlighting some issues that we encountered. Thirdly, the user study performed on Garnet is introduced and its findings discussed. Finally, the subjects' use of space and the textual analysis of the document groups that they created are reviewed.

2. GARNET IN USE

A pilot version of Garnet has been created, which is integrated with the New Zealand Digital Library Project's Greenstone software [20]. Greenstone is a comprehensive open-source Digital Library software system, supporting common actions such

as full-text and index searching, and browsing in category hierarchies. Access to the digital library system is via a remote digital library protocol. As demonstrated in our earlier work [1], the Greenstone protocol can be trivially mapped to the three other common DL protocols – Dienst, Z39.50 and SDLIP – so Garnet could readily be integrated with alternative digital library systems that employ these other protocols.

We will now demonstrate the system in use. The library material we will use is the Humanity Development Library of the United Nations, one of the widely available examples of a Greenstone library collection, which consists of several thousand pages.

2.1 Overview

In Figure 1, we see a ‘typical’ Garnet user session in progress; a ‘window’ appears inside the main browser window. This window is a ‘collection’ of materials which the user has recorded in the current, or a previous, session. Each document is represented by a rectangle containing some text, as indicated in the diagram, which we term a ‘label’ for simplicity.

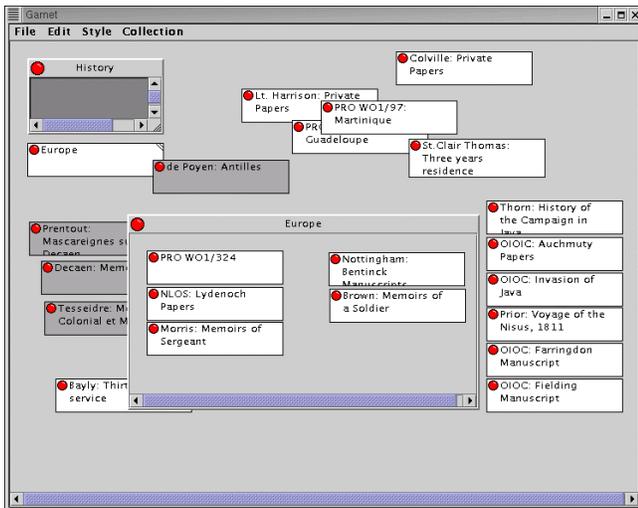


Fig. 1. A Garnet Client in Use

Within a collection, the user is free to place, size and color each document label as they see fit – the space is entirely freeform. Labels can be moved and/or copied between collections in the usual way for similar direct manipulation environments. Document labels can be added explicitly by the user or through interaction with the digital library’s search facilities.

Therefore, the user is free to use the document labels both in freeform structures of their own making inside collections, and in a more formal organization by using the explicit hierarchical forms of a set of document collections. Taking the example above, we have a collection called “Europe”, which has a column of three documents on the left-hand side and another column of two documents on the right. The column is a structure created by the user’s exploitation of space – it is not a feature enforced by the system. The column idiom can also be seen in the root collection – on the right-hand side. Some use of color can be seen here – e.g “de Poyen: Antilles” – but the relationship is not clear to us as readers of the hypertext.

2.2 Example search

Let us now follow a simple sequence of interactions, starting with an example search. For our purposes, we are going to investigate snail farming, in an attempt to discover whether we have the appropriate resources to consider that form of agriculture. With Garnet loaded, we start a new search in the Greenstone system and we enter the simple query “snail”. In Figure 2, a simple collection window appears with a number of document labels appearing one beneath the other, similar to a typical web-based result list.

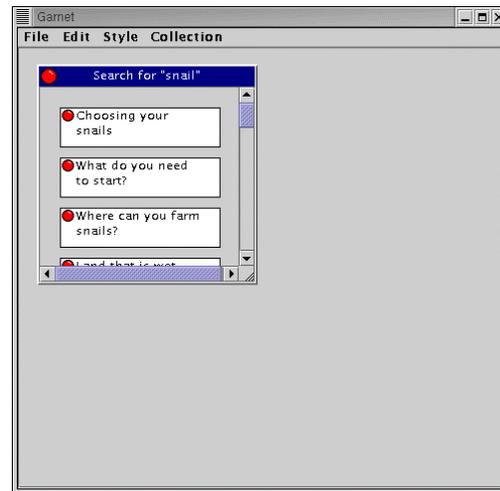


Fig. 2. A simple search

To read a document, the user double-clicks on its label. Garnet then displays the document in a separate window. On reading the first two documents, we decide that we’d like to keep the second document (“What do you need to start?”), and we move it to our root workspace window – simply dragging the document from the “Search for ‘snail’” window onto the main Garnet window.

The first document, however, seems a bit advanced, and we can delete it from the list, clicking on the small red ‘blob’ on its top left corner. As a result of this, the later documents move upwards. Should we wish to return to the search results at a later date, by default these changes would be retained.

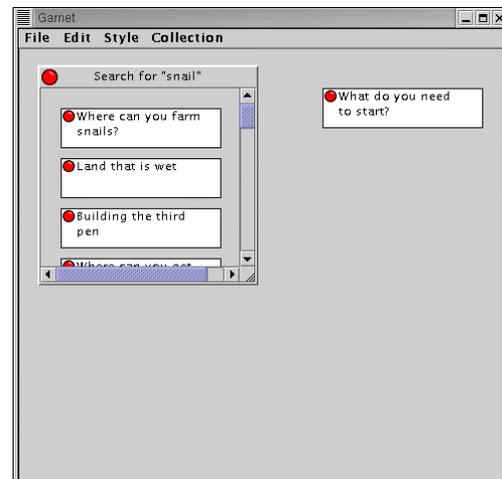


Fig. 3. the workspace after alterations

2.3 Demonstration of “scatter results”

In the previous search, we performed a plain search. Garnet, however, can exploit the organization done by the user in a novel manner. We can “scatter” a set of documents (including search results) from a selected window over the existing layout of documents in the workspace. “Scattering” places the search documents near to groups of existing documents with which they have a strong similarity.

Continuing our previous example, we have now selected a few more useful-looking documents, but let us suppose that a couple of questions remain unanswered.

Suppose we have a plentiful supply of bananas which we would like to use, but we are not sure whether this food would be appropriate. If we were to do a naïve search, on “banana”, the initial results do not well match our particular interest (Figure 4).

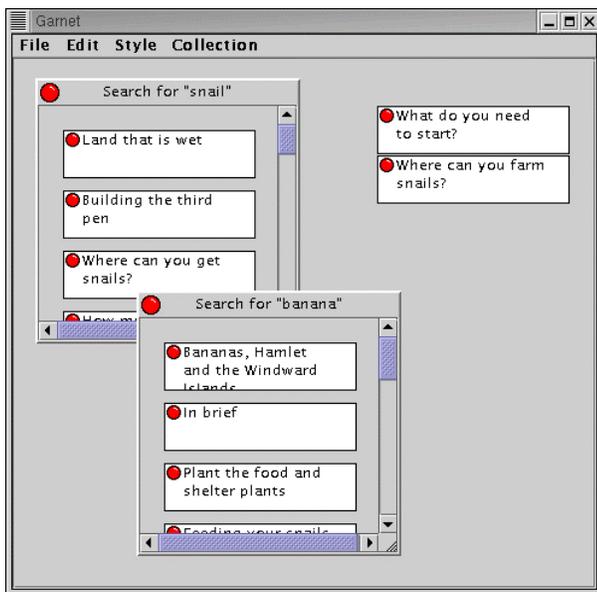


Fig. 4. Situation before “scatter”

In fact, documents which relate to our interest can be found in both the “snail” and “banana” searches. However, these documents of interest may not appear at the very top of either list. Normally, we would have to try and re-work our query manually to make it more targeted. In the case of Garnet, we could use the ‘scattering’ feature to discover any material similar to documents we have already selected for storage. Or, in other words, Garnet can generate existing search terms or filtering to represent our user’s interests, based on the workspace layout they have already created.

Viewing Figure 4 again, note the third item from the top of the new “Search for ‘banana’” list: “Plant the food and shelter plants”. This item is related to the two documents on the main workspace (for clarity we’ve chosen something that is visible in this example). If we now do a “scatter”, (Figure 5), a subset of the “banana” search results appear on the main collection. This small subset, which appears in a light gray below (Figure 5), has been found by Garnet to be a close match to the existing pair of documents, which appear in white. Suggestions are always

displayed in this gray color, and below and to the right of the group of documents which they are believed to be similar to.

We can now investigate the two suggested documents which are similar to the previously selected pair. As it happens, these documents would confirm that ripe bananas can indeed be used to feed snails. If we wanted to permanently add one or other suggestion to the workspace, we can click on the ‘blob’ which appears on the top right corner of each of the suggestions.

If we no longer wish to see the existing suggestions, or when another set of documents is scattered, the current suggestions are cleared.

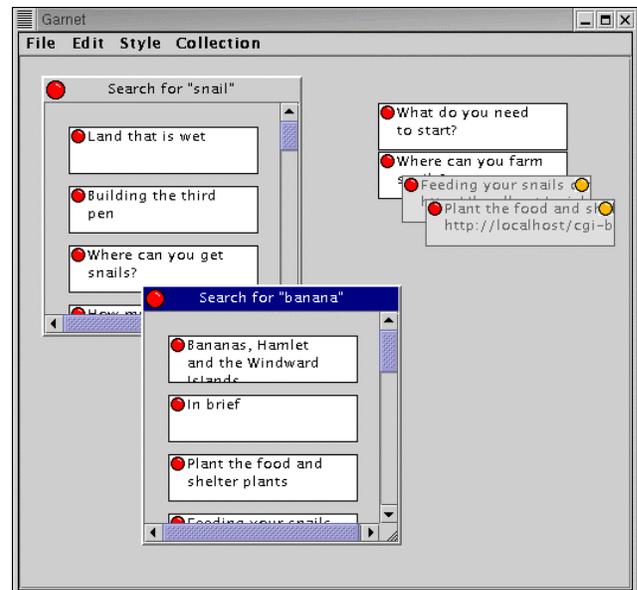


Fig. 5. After “scatter” – note the shaded document labels added in comparison to Fig. 4

In this section, the basic functions of Garnet have been introduced. We will now discuss the architecture and operation of Garnet, and follow that with a report of a user study of Garnet in use.

3. SYSTEM DESIGN

Garnet combines a spatial hypertext workspace with access to the facilities of a digital library. The spatial hypertext element is a simple spatial hypertext editor, similar to the VIKI system of Marshall and Shipman [14]. Access to the digital library system is via a remote digital library protocol. As mentioned earlier, Garnet can readily be connected to any of the four common DL protocols.

Garnet also processes documents that appear in its workspace (e.g. as an item in a search result set) so that text-matching facilities such as the “scatter” feature demonstrated above can be supported. Each document group in Garnet’s workspace has a textual representation calculated for it as its membership changes – some further details of which will be explained later. The text-matching facilities use a similarity engine built upon a common associative text-matching algorithm.

3.1 Architecture

A system schematic of Garnet is shown in Figure 6 below. The features additional to existing spatial hypertext systems are indicated in gray.

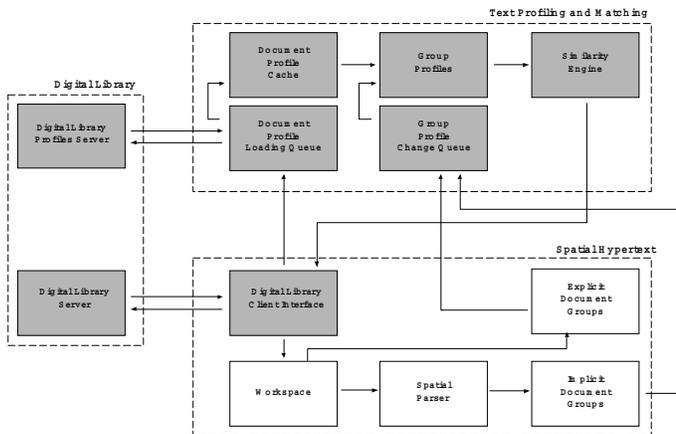


Fig. 6: Basic system architecture

The items in white are analogous to the components of existing spatial hypertext systems such as VIKI and VKB. Garnet possesses a spatial parser similar to that of VIKI [17]. The spatial parser facilitates the identification of implicit structures created by the user – for instance, the columns described in Figure 1. Explicit document groups are the ‘collections’ also discussed in Section 2.1 that the user creates directly through the interface.

Garnet’s additional items include: extensions to the core *spatial hypertext* system, an external *digital library*, and a *text profiling and matching* subsystem. We now briefly discuss these in turn.

The digital library components on the left communicate with the spatial hypertext element through the Digital Library Client element that we have added to the core spatial hypertext system. The DL Client also supports the representation of DL elements such as search result lists and browsing structures in the spatial hypertext’s workspace. The DL server components usually run on a separate machine, and are not part of Garnet per se – they are generic DL components.

The DL Profile Server is an optional element that provides pre-processed textual profiles of documents in a form that Garnet can later use for text matching. The profile server reduces computational overhead when a document is loaded, but the system will function if the profile server is not available – the profiles are then generated by Garnet’s profile matching system.

The Profile Repository stores text profiles for both individual documents and for groups of documents, be those groups explicit (i.e. ‘collections’ of documents that are akin to folders in a filing system) or implicit (e.g. a ‘column’). The Similarity Engine provides matching between profiles and can also send requests to the DL Client software – e.g. to place documents after a ‘scatter’ or to perform a separate search for documents similar to those in a chosen group.

The processing of profiles could consume considerable amounts of processor time. Therefore, both spatial parsing and the resulting identification of the representative text of document groups is evaluated lazily when the user is not interacting with the system, unless a user action requires this data to be brought up to

date. This provides for a fluid, responsive interaction with the user even in larger workspaces.

3.2 Interface Design

Garnet needs to present the facilities of the digital library within the spatial hypertext workspace. It also possesses additional features that track the user’s activity within the digital library (e.g. a search history list) that also appear in the workspace. Neither of these appears in traditional spatial hypertext systems, and the introduction of them may disrupt the fluidity of the spatial hypertext idiom for Garnet’s users. Some visual interfaces to digital libraries do use a subset of spatial hypertext features, but as we discuss in [4], major omissions are present.

The items in the workspace that represent features or documents within the digital library and those items that represent other system features have necessarily different behaviors from documents placed in the workspace by the user. For example, a user cannot reorganize a library’s classification hierarchy or the placing of documents within it, whereas documents could be added to, moved within or removed from the user’s workspace with impunity. These different behaviors could result in confusion if a user is faced with two items of similar appearance but of different types. Conversely, providing a consistent set of controls and interactions provides a more useable system. For instance, a collection of documents in the user’s workspace – akin to a folder in a typical filing system – plays a similar role to a node in a classification hierarchy, and the behavior of the two should reflect this.

In Garnet, therefore, some visual cues were added to indicate whether part of the workspace represents a digital library item or the user’s own workspace. In the case of collections, user collections appear on a light gray surface whilst system collections appear on a dark gray surface. Visual cues are also used to indicate when a document is placed as a result of a scatter. Here color is also used – a shade of grey determined by textual similarity – but in addition, “scattered” documents are always placed in a regular pile on the bottom right of the document group that they match. See Figure 5 above for an example of this.

Coordination between the user and the digital library system is simplified by the fact that the DL subsystem only presents items into the workspace as a result of explicit user interaction. Though a user cannot determine where the item will appear (positioning of items follows an overlapping pile strategy similar to folders in most filing systems), the position is predictable and Garnet ensures that it is visible to the user.

Where possible, access to digital library features is found in the menu structure of Garnet. Placing further items into the main workspace would reduce the area available for the user to use for their own placement of documents.

4. USER STUDY

Garnet provides a novel interface for a digital library, providing facilities for both information structuring and traditional information seeking. Though earlier studies in physical environments noted the frequent interleaving of these activities, we are not aware of any similar study in an integrated digital environment. We undertook an initial study to identify salient issues in integrating the digital library and spatial hypertext elements of Garnet. In the initial study of Garnet, therefore, we

undertook a qualitative study to elicit design consideration and identify problems worthy of further investigation.

Our study followed a pattern of similar probing studies already established in our previous work, e.g. [3]. A panel of ten subjects was recruited, each studying a degree in psychology or computer science at final year honors level or above. These subjects were frequent users of digital documents. Given that information structuring has only been closely observed in the information seeking of skilled information workers, we believed that those with more casual information seeking skills and needs would be less realistic subjects. The subjects also had no prior exposure to spatial hypertext systems, which permitted us to capture the initial expectations of how they could benefit from an information structuring tool.

Subjects were initially screened in a pre-study questionnaire to capture their information seeking and structuring skills. Then, they were introduced to Garnet in a brief 10 minute tutorial, followed by an open ended period of self-directed exploration. The main study was then undertaken, with the subjects and their activity on the computer being recorded on video tape. At the conclusion of the main study, a post-experimental interview and questionnaire captured the subjects' impressions, views and experiences. Where users were asked to express an opinion, scoring was on a seven-point likert scale.

Each subject was given the same task for the main study – a simple information-seeking task (to find papers that would be good source material for a literature review on digital libraries). They were given a brief description of digital libraries and a list of related topics to assist the selection of their initial queries. After completing the initial digital library topic task, a further requirement for documents upon human-computer interaction as a theme in digital libraries was introduced, and subjects asked to obtain specific information on that. They first used the “scatter” tool described above to support this task, before embarking upon an independent search for this material. For this task, subjects used a digital library collection of over ten thousand computer science technical reports.

4.1 Results

Subjects were asked to compare their experiences of working with Garnet with a number of familiar systems. We also observed their pattern of work and their organization of documents during the study. We will first report the effectiveness of Garnet as a DL interface, before moving on to the patterns that we observed in our subjects' use of Garnet.

4.1.1 Accessing Digital Library Features

Subjects reported that basic digital library tasks such as searching were comparable in ease-of-use with the same features in a web-based digital library. No subject reported, or was observed, experiencing problems with these features. This strongly suggests that the Spatial Hypertext interface of Garnet does not impede access to digital library features.

Subjects were also asked whether they had problems distinguishing parts of the system that they could manipulate – e.g. documents in their own workspace – with parts where they could not – e.g. in browsing structures of the library. Given the known problems of different modes of operation in human-

computer interaction, we were concerned that this could prove a major problem. However, all subjects denied having a problem with this. There are some contributory factors that may have influenced this. Firstly, all parts of the workspace which included a view upon a digital library component – e.g. a search result set or a browsing node – was very regular in appearance, containing a column of documents or other items, and had a different color background. Compared with the more freeform organization preferred by our subjects, the contrasting regularity of system items in the workspace provided a subtle distinction to the users' own creation. The distinction between system-owned and user-owned items may also have been generally assisted by the fact that many operations could be achieved on both system- and user-owned objects of the same type, minimizing the scope for unexpected behavior.

Subjects were also asked to rate the particular visual representation of documents, search lists and other items individually. All items were rated positively: however some useful and interesting ideas were suggested, e.g.:

Firstly, six of the ten subjects independently expressed their wish to be able to alter the title of documents. We had not allowed for this, as it is at odds with the nature of a digital library where documents are not normally editable. It is, however, very much within the nature of spatial hypertexts. Explanations included opaque titles of documents, and that titles often did not fit the immediate task of the user. This suggests that even for users who, like our subjects, have not been exposed to spatial hypertexts before, some spatial hypertext features that disrupt digital library expectations may be an important contribution of integration.

Secondly, five subjects requested a more visual access to digital library features that were obtained from outside the workspace – e.g. the launch of new queries. Here, the preference could be explained both from the persistent appearance of such elements in web interfaces to digital libraries and the visual interactive style of spatial hypertext. In retrospect, this was a design error on our part, but one that could recur in other contexts.

To summarize, our subjects found no difficulties using Garnet to access DL facilities. They were able to distinguish between system- and user- owned areas with apparent ease, quickly recognizing the different behaviors of each. Our users identified areas for improvement, such as being more consistently focussed on the workspace presentation of tools, and permitting more editing of items than our digital library origins led us to believe.

4.1.2 Patterns of Behavior

We were also interested in how users followed their information seeking and information structuring tasks throughout the study. This was captured through both the video recording and post-experimental interview.

A first point of interest is that subjects closely interleaved information seeking and information structuring. Once a subject decided to keep a document, even provisionally, it was immediately moved onto the spatial hypertext workspace. Organization of the document was performed at the same time. This simple pattern was observed in every subject.

A document on a new subject or of uncertain role would often be placed in a particular group in the workspace before being reorganized to another position later in the subject's work. This behavior mirrors the patterns of work previously observed in physical environments [11]. However, two subjects (8 and 10) focused on a single miscellaneous column, minimizing their organization work within the task. In interview, one reported that they would organize their documents more precisely at the end of their detailed reading, and before doing any final searching. The other subject stated that they would probably not organize documents within a task, though they would organize documents between separate tasks.

When subjects identified a theme in two or more documents, this would result in a new group being created in the workspace. However, the consequences did not stop there. In half of all cases, the creation of a new group would result in the user doing a new query to the digital library to attempt to obtain similar documents to add to that group.

Given these behaviors, the organizational activity of information structuring was certainly interlinked in a manner that resonates with previous information seeking and spatial hypertext research [7, 12, 15].

4.1.3 Visual Organization of Workspace

In [14], Marshall and Shipman identified common visual idioms used in Spatial Hypertext. The subjects in our study predominantly used columns in organizing their documents. Nine of the ten subjects used columns, whereas only one used a row. Piles – irregular, overlapping stacks of documents – appeared in only one workspace. The vertical form of the search result sets may have influenced our subjects' own choice of visual structure, as may the form of the available workspace and the document labels. Furthermore, six of the nine subjects who used columns added explicit labels to each column to explicitly identify its theme or topic. The post-experimental interview revealed that users exploited the user-created labels to clarify the theme of each document group for their own benefit. This was perceived as a more straightforward step than creating a collection for a topic.

4.1.4 Spatial Hypertext Facilities

Subjects embraced the ability to organize documents on their workspaces. When asked what benefits they perceived in this, answers included: Subject 4, "I can see a document on the desktop without having to go back"; Subject 7, "being able to store stuff and organize them is good...this way you can have stuff that relates between a couple of areas". Seven subjects specifically mentioned the advantages of having an overview of what they collected, and eight reported storing documents as being an important benefit over traditional Web-based DL interfaces.

Subjects also made positive comments over the tangible, drag-and-drop interaction of the interface: e.g. Subject 9, "I really like the ability to manipulate here and move them around and take them off"; Subject 6, "You just drop stuff where you want it".

Other advantages reported included supporting deciding which search to do next, remembering which searches had already been done and prioritizing documents in perceived order of importance. All these are activities previously reported in physical environments, and claimed as potential advantages of spatial hypertext.

4.1.5 Discussion

Our study clearly suggests that spatial hypertext's information structuring facilities are supportive of traditional information seeking in a digital library. Subject's patterns of workflow under observation matched the interleaved patterns observed in [7, 11, 12, 13, 16] and subjects themselves reported some of these patterns themselves in the post-experimental interview.

Our subjects also demonstrated known patterns in spatial hypertext, despite none having used any similar system before (the closest simile was that two had used 'MindMap' software). This corroborates existing hypertext research and suggests that our subjects demonstrated typical rather than exceptional behavior.

The visual, gestural interaction of spatial hypertext was particularly noted as an advantage by the participants, and suggested changes such as editing document titles and presentation of search facilities on the workspace are consistent with spatial hypertexts and visual DL interfaces like DLITE [5].

5. EXPLOITING ORGANIZATION IN A SPATIAL HYPERTEXT

The potential benefits of computation over hypertext have been demonstrated in recent years through examples such as the identification of web communities [8] and the PageRank algorithm for ranking web search results [2]. As noted by Frank Shipman [18], computation over *spatial* hypertext is a much less explored area.

As noted before, studies of information seeking in physical environments revealed that the organization of documents – be it formal or informal – played a key role in users' management of their information seeking. Users implicitly encoded their past and present interests and their future seeking intentions in their placing and grouping of documents. Also, in electronic environments, the visual organization of objects can reflect topical themes that a user tracks over time, as observed in [10, 15].

If a group of documents has a common theme, then the documents may share common words that characterize that theme. Users of Garnet can read individual documents, and can readily see the titles of documents in the Garnet workspace to extract keywords associated with a group. However, users are unlikely to systematically extract every relevant keyword from every group. On the other hand, simple textual processing would readily extract key common words automatically.

Given full-text access to one or more digital libraries Garnet can process the original text of each document in its workspace. We decided to exploit this to generate textual profiles of the common words of each document group. This textual profile is thus a product of computation over a spatial hypertext. We explored the possibility of using this profile to support either information structuring – the core activity of spatial hypertext – or information seeking – the core activity of a digital library.

5.1 Implementation

As noted in Section 3.1, Garnet creates a textual profile of each document group found in its workspace. The computation of this profile occurs in three phases. The first phase involves

identifying the implicit and explicit structures in the hypertext, assisted by Garnet's spatial parser. In the second phase, the documents that are represented in each structure are then processed textually to extract descriptive terms that, hopefully, represent the common theme of the structure. Finally, these textual representations are used to assist information seeking or structuring by performing some form of textual processing – e.g. refining a search or sub-dividing a group of documents.

The operation of spatial parsers has already been discussed in the spatial hypertext literature [17], and this paper will therefore pass over this phase of our computation over hypertext.

After the spatial parser has identified all the document groups in the workspace, Garnet creates a textual representation for each group. For each group, Garnet extracts the words that appear in every document in the group (i.e. the intersecting set of words). Garnet excludes words in a stopword list of very common terms from each group representation.

This simple representation of words that are common to a group has proven effective for text matching [21]. It also avoids a range of potential technical problems. One such problem arises if the term frequency information for each word is used. When using a range of remote DLs, each library would have a different rate of occurrence for any given word. Normalizing word frequency weights raises conceptual and practical difficulties. Furthermore, such frequency information is usually either unavailable or only obtainable via laborious and inefficient means – only partial support is available in any of the four standard DL protocols.

An alternative to using individual words would be the use of phrases common to the group, but Zamir et al have demonstrated that this is of only a small benefit at a higher processing cost [21].

5.2 Exploiting Document Group Profiles

Garnet can use the text profile of a document group as follows:

- Documents in a collection or search result set can be 'scattered' (see 2.1 above) to place them next to a group of documents in the main workspace to which they bear a strong textual resemblance. This supports information structuring and also uses the user's existing organization to further filter the search result set. Documents are matched using the group-matching algorithm in [21].
- Given a select group of documents, the user can ask Garnet to find similar documents – this will invoke a search in the digital library using the representative keywords of the group, supporting information seeking.
- If a user selects a group of documents, these can be subdivided automatically using a proven clustering technique [21] that uses the textual representation used by Garnet. This provides further support to information structuring in the workspace.

Further examples of how the extracted text could be used includes: automatically generating descriptive labels for document groups; matching against document groups in other workspaces; matching against classifications within a digital library; suggesting document label color; etc.

However, existing research does not supply any evidence on whether the textual profiling of document groups is useful – for instance, if a group is topically heterogeneous, then use of it to match other texts may result in meaningless or unhelpful results.

Though some use of the text of informal spatial hypertext groupings seems to exist [19], there is no clear account of the implementation or these functions, nor a study of the textual properties of them. Furthermore, there also seem to be a lack of studies on the textual properties of user-generated document groups in the digital library and information retrieval communities.

5.3 Textual Analysis

Answering the question of how useful the textual profiles of the user's document groups could be is not straightforward. We elected to approach this through a number of more pointed subsidiary questions:

- TQ1. Are the groups of documents textually consistent?
- TQ2. Do 'clearly' miscellaneous groups exist, do they have distinct properties?
- TQ3. Can users make sense of / do users approve of outputs when group profiles are used for searching?

Starting from the expectation of secondary notation, we were interested in the degree to which documents in the same group shared the same words. In the field of information retrieval, such a property would be assumed to indicate topical consistency.

Conversely, subjects in previous spatial hypertext studies demonstrated patterns of chaotic organization – for instance the creation of 'miscellaneous' piles of documents which were unsorted. Could textual analysis identify such miscellaneous piles, as they were understood by the users that created them?

This final issue of miscellaneous piles shares, in common with TQ3, the requirement to elicit user response to the system. This aspect of the usefulness issue was addressed through part of the User Study already described in Section 4.

The membership of implicit groups in a spatial hypertext is recognized by the heuristics of the spatial parser. The spatial parser's heuristics are prone to some degree of error, resulting in a lower precision of group membership. We decided to focus upon this problematic area to derive a 'worst scenario' base-case.

5.4 Method

A key problem in analysis was the choice of how to evaluate the quality of topical consistency in the organization of each subject's workspace. Many techniques for textual analysis exist. Following from our experience with search engines and clustering techniques [3], we decided to use measures for textual consistency found in document clustering. Clustering algorithms are created to automatically organize groups of texts into topical themes, where topical similarity is taken to correlate to textual similarity.

This gave the advantages that: we could compare against the organization done by a computer; the comparison was quantifiable and not prone to human subjectivity and could readily be reproduced by others. The nominal clustering results

allowed some insight into the overall consistency of each user in terms of the documents that they selected, and a means through which we could compare users who chose differing documents.

We used two well-accepted clustering algorithms: Word Intersection Clustering and Scatter/Gather. The first identifies the words that appear in every document in a group, without giving extra significance to rare words. This apparently simple technique has given rise to one of the most effective clustering algorithms available [21]. The second algorithm uses the established Term Frequency/Inverted Document Frequency (TFIDF) weighting of words common in text search engines. The Scatter/Gather clustering algorithm uses this weighting for calculating the quality of its clusters. Scatter/Gather has been evaluated technically and in user studies as a support for human information work [6, 9]. This algorithm required generating extra term frequency information not normally held by Garnet.

In the case of each algorithm, the actual organization chosen by a subject can be compared against the groups that the clustering algorithm would have selected. Furthermore, each document group could be judged in quality using the internal scoring system for group quality found in each clustering algorithm. Finally, a range of random organizations was generated and scored to obtain a ‘noise’ baseline. If users scored close to this level, then the resulting group profiles would be meaningless.

The two algorithms vary in their method for evaluating textual consistency within a group of documents. If the conclusions reached by the two algorithms correspond, then we can have a higher confidence in our findings.

One difficulty arises in assessing the quality of organization done by each subject. As each subject was free to choose their own documents in the course of a partially defined search task, the documents selected by each will vary. Therefore, straightforward comparison between subjects is problematic. However, the overall score for the optimal organization obtained by each clustering algorithm can be used to provide a basis for comparing the relative performance of each subject.

As noted above, we asked each subject in our user study to comment upon and rate the suggestions provided by Garnet when they used the “scatter” facility in the course of our experiment. Each subject was also asked to identify groups that they used to store miscellaneous documents in.

5.5 Results

Due to lack of space, we will present only the results from the Word Intersection clustering algorithm [21] in detail. The results from the Scatter/Gather algorithm give a very similar picture.

When viewed as a whole, the organizational patterns of the subjects prove to be surprisingly consistent. Subjects created a small number (typically 3-5; mean 3.6, SD=1.5) document groups. The mean number of documents per group was 2.7 (SD=1.8). Subjects also had an average 1.3 singleton documents in their workspace. As noted above, two subjects demonstrated a very different strategy, using a single long list for most of the documents in their workspace. A third had five singleton documents – no other subject had more than two.

In comparison, the clustering algorithms created two to six groups (m=4.8; SD=1.3). Each group generally contained between two and five documents. The largest group created was of eight documents, with the mean being 2.85 (SD=1.15). Only three workspaces resulted in one singleton.

Therefore, the overall pattern of organization was similar both between subjects and across the human subjects and the two clustering algorithms. The clustering algorithms produced fewer singletons, and were more consistent than the human subjects. No statistically significant difference was obtained between the overall organization patterns of the subjects and the clustering algorithms.

Individual subjects performed variably. When compared against a random organization, four subjects scored little better or worse than a random organization: Subjects 5, 7, 8 and 10. Subjects 8 and 10 had not performed much internal organization on their workspaces, so the poor result here is unsurprising. Subject 7’s low score is readily explained by their high number of singleton documents – heavily penalized by the clustering algorithms. In the case of subject 5, little consistency in their document selection is suggested by the low nominal score for their workspace.

When compared against the nominal score for the clustering algorithms, humans must necessarily perform worse. On average, the score was only 60% of the score of the clustering algorithms on each workspace.

However, this overall picture is complicated by the particularly poor scores from particular strategies. When individual group scores are compared, humans score much closer to the clustering algorithms. Average document group scores are 21.3 for clustering algorithms, and 21.5 for humans. The final scores are affected by higher number of singletons and one further factor: miscellaneous groups.

Table 1: Sample Results of Textual Analysis

Subject	1	2	3	4	5	6	7	8	9	10
Human Created Groups										
Groups	4	3	5	5	2	5	2	1	6	1
Single docs	1	2	0	1	0	1	5	0	1	2
Score	32	35	41	29	4	39	24	26	43	12
Cluster Groups – Zamir & Etzioni										
Groups	6	4	6	6	2	5	4	5	5	5
Single Docs	1	0	0	0	0	0	0	1	0	1
Nominal	63	56	57	32	17	46	43	40	55	39
Random	12	25	20	11	14	30	29	31	16	11

Within the human-organized workspaces each workspace had one (or in one case two) low-scoring groups. These individual groups scored less than 10 – i.e. less than half the average. Closer inspection revealed a bipolar distribution – eleven groups with average scores of 7.3 against a second group with an average of 27.8. Standard deviations were 3.5 and 6.7 respectively.

Each subject had identified their low-scoring groups as being a ‘miscellaneous’ or ‘unsure’ group in the course of the experiment or during the post-experimental interview. This suggests that such ‘miscellaneous’ document groups can be readily identified through textual processing. Topically themed groups (as identified by our subjects) corresponded with the higher group scores. These topical groups achieve scores that are comparable or superior to those of cluster-algorithm generated groups.

Overall, low scores for the organization of the workspace as a whole seem to be the product of having one very low-scoring miscellaneous pile (particularly for Subjects 8 and 10) or a high number of singleton documents (Subject 7). Subject 5 seems to have performed poorly by all available measures.

However, even this picture is not as simple as it seems. Subject 7 had a consistent but atypical strategy that was elicited in our post-experimental interview: they organized by author alone. This was a consistent approach which our text profiles were able to represent. However, the ‘full-text’ analysis we performed unduly penalized a systematic organization that an analysis using metadata would have rated highly.

User satisfaction with the “scatter” facility of Garnet from the user study was generally positive. Seven of our ten subjects rated the facility as “Useful” or “Very useful”, whilst two more rated it as “Somewhat helpful”. Two of the three dissenting subjects were those who had used a single, long list in their workspace and had scored poorly in our textual analysis of the subject workspaces. All three scored notably low individual document group scores. Qualitative feedback was also supportive of the ‘scatter’ facility: Subject 3: “It was helpful and a couple of times in brought up a document that I had seen – if I’d missed a document, that was helpful”, Subject 9: “I got a couple of new categories out if it...it extends what you’ve already got”.

The “scatter” facility relies upon the individual quality of document groups and includes a quality cut-off that avoids suggestions of low certainty. Higher scoring (i.e. topical rather than miscellaneous) groups are much more likely to result in the presentation of a ‘suggestion’. The topical groups of our participants were of similar size and score to the themed groups of clustering algorithms.

5.6 Discussion

Subjects generally performed simple organizations similar in form to that which a clustering algorithm would have done. Furthermore, the quality rating of the topical groups within those organizations compared favorably with the nominal performance of two well-accepted clustering algorithms. This performance suggests a high level of textual consistency, and may be considered a positive indicator that topic extraction from document groups in a spatial hypertext may prove an effective form of exploiting computation over spatial hypertext.

Given the particularly low score of some miscellaneous document groups, spatial hypertexts could readily identify such groups and positively support the use of these when appropriate – e.g. if a document does not match the well-formed groups in a workspace, it could be suggested as an addition to the miscellaneous pile.

However, when, as occurred in the case of two of our subjects, a user embarks upon a simple strategy where most documents are in

a large miscellaneous structure, another approach may prove better. It is possible that in this case the user may be well supported by providing the option of automatic organizing tools. Clearly, such support needs independent evaluation and assessment.

User’s patterns of work naturally vary, and it would appear that the local quality of individual groups – as opposed to entire workspaces – is a better indicator of the usefulness of computation over the spatial hypertext.

6. Spatial Hypertext - Review

From the perspective of Spatial Hypertext, Garnet serves as an interesting example system that can address a number of the “Seven directions for Spatial Hypertext Research” identified by Frank Shipman [18].

Firstly, its integrated access to a digital library through a remote protocol serves both as an example of “Integrating Spatial Hypertext into the Information Environment” – particularly in regard to how spatial hypertext systems can support the wider information seeking task. It builds upon the spatial parsing found in VIKI and others [17] as a means of exporting information to other systems. As HCI researchers, we feel that the representation of the wider environment within spatial hypertexts is also worthy of study, and we have found that though representational issues may not be a problem, conflicts may emerge between the affordances of spatial hypertext and other parts of the information environment (e.g. renaming documents in the case of Garnet).

Secondly, the study of the textual profile of the document groups created by subjects using Garnet is a contribution to the issue of “Computation In and Over Spatial Hypertexts”. Our study provides some indicative evidence that the implicit structures in spatial hypertexts may prove to be exploitable supports for information seeking and structuring. Much further work remains to be done here.

Thirdly, the textual analysis of our subjects’ work is a contribution to the development of “Evaluation Methods and Practices”. Experience with our approach currently suggests that local quality – of individual document groups – can readily be assessed and is less sensitive to differing user strategies.

Spatial hypertext has had a considerable influence on the development of visual DL interfaces such as SketchTrieve [10]. However, we feel that the influence of spatial hypertext, and information structuring tools in general, on digital library research is underdeveloped. Conversely, Garnet indicates that digital library research can contribute to the development of spatial hypertext research.

7. CONCLUSIONS

Garnet is an integrated spatial hypertext and digital library system that supports information seeking and information structuring.

Users have been observed interleaving information seeking and structuring in physical environments. We found that users demonstrate the sane interconnected, interleaved interaction in electronic environments. Our subjects appreciated this integrated approach and evaluated the spatial hypertext facilities positively whilst finding no hindrance to their use of traditional DL facilities. Other visual DL interfaces have been influenced by

spatial hypertext [5, 10] without fully embracing its idioms. We have demonstrated that a combined system supports workflows that rely upon the presence of the information structuring facilities true spatial hypertext provides.

We discovered that even the implicit document groups in spatial hypertexts can demonstrate textual consistency. Users positively rated our simple exploitation of their own organization of documents. This may prove to be one route to exploiting computation over spatial hypertext, and suggests that computation over spatial hypertext is worth further study.

Much further work remains to be done. A larger and more detailed study of user organization patterns should be done to gain a higher measure of confidence in our finding. Alternative text representations and matching techniques are yet to be investigated. The relative performance of user- and computer-generated groupings of documents as perceived by actual users should be studied in much greater depth.

8. ACKNOWLEDGMENTS

Our thanks to Middlesex University and University College London for their support of this work.

9. REFERENCES

- [1] Bainbridge, D., Buchanan, G., McPherson, J.R., Jones, S., Mahoui, A., Witten, I. H.: "Greenstone: A Platform for Distributed Digital Library Applications". Proceedings of the European Conference on Digital Libraries, Springer-Verlag, pp. 137-148, 2001.
- [2] Brin, S., Page, L.: "The anatomy of a large-scale hypertextual (Web) search engine", Proceedings of the 7th World Wide Web Conference, (WWW7) also as Computer Networks and ISDN Systems, Vol. 30:1-7, pp. 107-117, 1998.
- [3] Buchanan, G., Jones, M., Marsden, G.: "Exploring Small Screen Digital Library Access with the Greenstone Digital Library". Proceedings of the European Conference on Digital Libraries, Springer-Verlag: pp. 583-596, 2002.
- [4] Buchanan, G., Blandford, A., Thimbleby, H., Jones, M., "Integrating information seeking and structuring: exploring the role of spatial hypertexts in a Digital Library", Proceedings of the European Conference on Digital Libraries, 2004, in press.
- [5] Cousins, S.B., Paepcke, A., Winograd, T., Bier, E.A., Pier, K.A.: "The Digital Library Integrated Task Environment (DLITE)". ACM Conference on Digital Libraries, pp.142-151, 1997.
- [6] Cutting D., Karger D., Pedersen J., Tukey, J. W. : "Scatter/Gather: A Cluster-based Approach to Browsing Large Document Collections", Proceedings of the 15th Annual International ACM/SIGIR Conference, Copenhagen,, pp. 318-329, 1992.
- [7] Ellis, D., "Modelling the information seeking patterns of engineers and research scientists in an industrial environment", Journal of Documentation 53(4):pp.384-403, 1997.
- [8] Gibson, D., Kleinberg, J. , and Raghavan, P. "Inferring Web communities from link topology". Proceedings of the 9th ACM Conference on Hypertext and Hypermedia, pp. 225-234, 1998.
- [9] Hearst , M.A. ,Pedersen, J. O.. Reexamining the cluster hypothesis: Scatter/Gather on retrieval results. In Proceedings of the 19th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 76-84, 1996.
- [10] Hendry, D. G., Harper, D. J., "An informal information-seeking environment". Journal of the American Society for Information Science, 48(11), pp. 1036-1048, 1997.
- [11] Kidd, A., "The Marks are on the Knowledge Worker", Proceedings of the ACM CHI Conference, Boston, MA, pp. 186-191, 1994.
- [12] Kuhlthau, C. C., "Seeking Meaning: a process approach to library and information services" Ablex Publishing, Norwood, New Jersey, 1992.
- [13] Malone, T. W.: "How do People Organise their Desks? Implications for the Design of Office Information Systems", ACM Transactions on Information Systems, v. 1 (1), pp. 99-112, January 1983.
- [14] Marshall, C., Shipman, F. and Coombs, J.: "VIKI: Spatial Hypertext supporting emergent structure". Procs. of the ACM European Conference on Hypermedia Technology ACM Press, NY, pp. 13-23, 1994.
- [15] Marshall, C. and Shipman, F.: "Spatial Hypertext and the practice of information triage". Proceedings of the Eighth ACM Conference on Hypertext, ACM Press, New York, NY, pp. 124-133, 1997.
- [16] O'Day, V, Jeffries, R.,: "Orienteering in an Information Landscape: How Information Seekers Get From Here to There", Proceedings of INTERCHI, ACM, pp. 438-445, 1993.
- [17] Shipman, F., Marshall, C., and Moran T.: "Finding and Using Implicit Structure in Human-Organized Spatial Layouts of Information", Proceedings of Human Factors in Computing Systems (CHI '95), pp. 346-353, 1995.
- [18] Shipman, F.: "Seven Directions for Spatial Hypertext Research", First International Workshop on Spatial Hypertext, ACM Hypertext Conference 2001, Aarhus, Denmark, 2001. Online at: <http://www.csd.tamu.edu/~shipman/SpatialHypertext/SH1/shipman.pdf>
- [19] Shipman, F., Moore, J. M., Maloor, P., Hsieh, H.W., Akkapeddi, R.: "Semantics happen: knowledge building in spatial hypertext". Proceedings of ACM Conference on Hypertext, pp. 25-34, 2002.
- [20] Witten, I., McNab, R., Boddie, S., Bainbridge, D. : "Greenstone: A Comprehensive Open-Source Digital Library Software System". Proceedings of the Fifth ACM Conference on Digital Libraries, ACM Press, pp.113-121, June 2000.
- [21] Zamir, O., Etzioni, O., Mandani, O. and Karp, R. M. "Fast and Intuitive Clustering of Web Documents". Third International Conference on Knowledge Discovery and Data Mining, AAAI Press, Menlo Park, California. pp. 287-290, 1997.