

Filing and Finding Computer Files

Bonnie Nardi, Ken Anderson, Thomas Erickson

Advanced Technology Group
Apple Computer, Inc.
1 Infinite Loop
Cupertino, CA, U.S.A. 94043
nardi@apple.com, andersonkt@aol.com, thomas@apple.com

ABSTRACT

This paper describes an interview that investigated how people organize and find electronic documents on their computers. Fifteen Macintosh users were interviewed regarding their problems and approaches to filing and finding information. We found no evidence that users are having serious trouble finding files on their personal computers. We uncovered patterns of behavior that seem best described in terms of three types of information: ephemeral, working, and archived.

KEYWORDS

Filing, Finding, File Organization, Information Retrieval, Reminding, Information Types

INTRODUCTION

There are many studies of the ways paper documents are organized in offices [5, 2, 9, 11, 8, 10, 1]. However, little work has been done on the organization of computer files. In fact, we haven't been able to find any published research on this topic. Perhaps investigators who are inclined towards this sort of study—requiring field work and observation—are drawn to the seemingly richer (and more visible) area of paper-based filing. Or perhaps they have shied away from electronic filing, believing that computer-based systems account for only a small proportion of filing activity, and are still too primitive and constrained to produce phenomena of interest. Our studies reveal computer-based filing to be an interesting arena of behavior, and one which has increasing relevance for the design of future computing systems.

THE FIELD STUDY

We conducted a series of interviews focusing on how people actually organize and find the information on their computers. We interviewed fifteen Macintosh users, including managers, graphic artists, programmers,

Note: This is a draft of this paper. I believe that it is very close to the final version, but if you need to quote the paper or quibble with fine points, you had best search out the final version, which I don't have. —TDE, 3/98

administrative assistants, and librarians. We had people with as much as 1500 megabytes available (plus servers) and as many as 31,000 items on their computers. We had people with as little as 80 megabytes available and 2,400 items on their systems.

We asked our users about their problems and approaches to organizing and finding files. We had them give us a 'tour' of their machine. Finally, we asked them to find a file we had noted during the tour, and observed their activities as they tried to find the file. The interviews were videotaped in the users' offices so that we could see their work environment and their computer systems.

Most of the users we interviewed are Apple employees (though some had arrived very recently). This is a strong bias and normally we would argue against such a skewed sample. However, we generally found very conservative behaviors. While it is clearly unwarranted to take members of such a technically acculturated environment as representative of most users, we do believe that they indicate the boundary of the leading edge of electronic filing and finding practices.

FILING INFORMATION

Virtually every study on the filing of paper documents [5, 2, 9, 11, 8, 10, 1] indicates that no two people file in the same way. Although in our study it was true that no participant used the same labels, file names, tree structure, or folder structure as any other participant, we did find a universal pattern of information handling.

This pattern is characterized by three types of information: ephemeral, archived, and working. Each entails differing patterns of finding and filing, as we discuss below. It is important to remember that "information" is not a monolithic category, but requires contextualizing, which can be done by examining the three information types.

The pattern of information handling we found is similar to that found over ten years ago by Cole [2] who studied the use of paper files in workers' offices. Cole conducted structured interviews with thirty users to determine which factors most influenced information storage and retrieval

behaviors. Her goal was to outline a desirable computer based information storage and retrieval system. She concluded that any system must take into account: (1) the characteristics of the information with which users interact (2) the way users prefer to organize their information, and (3) the role of the spatial position of information found in the physical office.

What is particularly interesting for us is that Cole found that there were different information types: “action information,” “personal work files,” and “archived information.” The types of information found in paper-oriented offices are thus very similar to what we found in computer files today (ephemeral, working, archived). Cole did not elaborate further on each type of information. We shall do that now, with data from the study.

Ephemeral Information

The basic characteristic of ephemeral information is that it has a short shelf life—there is a very limited time in which it is of value. It may be good for a day, a week, or sometimes two weeks, but it is rarely information that users want to file away. Sometimes the shelf life is so short that the information must be dealt with within a day, or it becomes irrelevant. For example, Mary¹, an administrative assistant, puts all of her incoming email and documents on her computer desktop during the day. At the end of the day she makes sure that her desktop is clear of all files. If a file is held over for more than a day, she puts it along the right hand side of the screen, so that she will be reminded to work on it the next day.

Ephemeral information is rarely created by the user; it typically comes from outside sources such as mailing list email, news wires, on-line bulletin boards and databases. Users will typically scan through these sources, selecting a few pieces of information that are relevant to their current needs and interests. Examples of ephemeral information include downloaded news articles, email announcements, email directed to large groups, and product information. Ephemeral information may also include transcribed information (e.g., phone numbers) captured on the increasing number of computer-based facilities for dealing with ephemeral information, such as electronic ‘yellow stickies,’ note pads, calendars, and To Do lists.

The problem for users is how to organize information of this nature. If the item is one that informs, but only for a limited time, where to file it and what to do with it become problematic. A common solution is to keep ephemeral information on the desktop or in a folder “loosely” filled with other folders on the top level. In essence, the information need not be filed at all. Keeping it visible allows it to easily grab the user’s attention and thus function as a reminder. For example, one user pulled items off the Internet related to virtual reality. After downloading several items, he began clustering them in the middle of his

screen, because of their similar content, and because he wanted to remind himself to send them to his father and a friend who were interested in virtual reality.

Ephemeral information is often incorporated into To Do Lists. A To Do list may be a single document, with information transcribed into it, or it may be a folder of files, each of which serves as a To Do item. The primary importance of a To Do item is for the action it indexes rather than its particular content. It acts as a trigger—a phone number that needs to be added to a message book, a message that needs a response, or a file containing instructions to be followed. People often have rituals associated with To do lists such as regularly looking them over (often at the start and end of the day), and checking items off the list. In our study, every person had a method of handling To Do items.

Users who are ‘forced’ to file ephemeral information in non-visible places often forget about, or fail to use it. Some applications—particularly electronic mail applications—automatically save email in a particular folder, or save all messages within a single file. This created a problem for many users: email would be pushed too far down the stack, and would fail as a visual trigger. None of our users had discovered a good way to deal with this problem. Non-ideal strategies for coping included responding to all mail immediately, or putting messages back into the incoming mail basket.

Working Information

Working information is information that is relevant to the user's current work needs, and tends to be frequently accessed and shared. It is created by the user, or the user’s co-workers. Its shelf life can be measured in weeks or months (although the important bits are transformed into archived information and retain their utility for much longer). Examples of working information include memos, meeting notes, working papers, and presentations.

Working information is organized by spatial location and by category. Users generally report that they have no problem finding working information. They come to know the files because of repeated, frequent, and dynamic interaction with them. Working papers are created, modified, and edited; presentations are assembled and fine-tuned; budgets are proposed, rejected, modified, cut, and submitted. In each case, the interaction involves direct and repeated reference to the materials in their spatial location, resulting in the gradual accumulation of contextual cues about where the information is located and how to find it again. As the information is accessed less frequently, the spatial awareness of location may give way to the category structure of the information. This process of moving toward more categorization is emergent as the person works more and more with the materials, and as the materials may grow. In short, spatial location plays the primary organizational role while the information is being used a lot, with categorization emerging as the work nears

¹ All names are fictitious.

completion.

Archived Information

Archived information stands in contrast to ephemeral and working information. It is only indirectly relevant to the user's current work. The shelf life of archived information is measured in months or years. It is highly structured and infrequently accessed. Most of the people we interviewed said that they access archived information somewhere between once-a-month and a few times per year. Archived materials tend to reflect a completed whole, rather than a current process. Examples of archived information include technical reports, customer files with a complete history of the relationship, and archives of all email related to a particular project.

Archived information is generated from a user's working information. All of our users archived by project; thus archiving tended to occur at the end of a project. Archiving generally involves three tasks: selecting the information to be saved; organizing the files and folders; placing the archive in a particular location. The selection step was often the most difficult. All of our users pared down the information they had after a project was over, removing files they felt would no longer be important in understanding the project. Many would have been glad to keep everything; however, they also realized that they would never actually go back to some of it. Some admitted there were times when they should have kept more of the information from a previous project, especially when it was something that could help them on a current project. However, they also noted that what was thrown out was rarely something they couldn't recapture or come close to recreating. Organizing presented other problems, but was less severe in its consequences (after all the files wouldn't be completely gone, just in a different location). Organizing allowed a more logical and personally meaningful relationship between the files and folders. Each user tended to have organizational structures that he or she repeated, making organizing a little less difficult and finding much easier. Placing simply involved moving files to a permanent location. Sometimes the location was a volume devoted to archived material, sometimes a tape backup, an external drive, or, less frequently, a diskette. Whatever the medium, there is always a designated, marked space for archived information. People reach a point of diminishing returns with respect to filing.

The overall style of organizing information varied from person to person. Some people used deep hierarchies (up to five folders deep), while others preferred only two levels of folders. Some people used labels while others did not. All users organized most or all of their files by project. In some cases, chronological organizations were used, such as all correspondence from January or from 1993. Some people interleaved the two types of organizations, filing by year and then by project within each year. In general, people reach a point of diminishing returns with respect to organizing information. Every user in our study discussed

how they had started elaborate filing schemes at one time or another but had failed to follow through with them. It just wasn't worth it.

Within a category, such as project, each user's organizational scheme had similar structures. For example, Tom, a manager, had three different projects. Under each project he had a personnel folder, a budget folder, and a working folder. This enabled him to avoid expending cognitive energy to form new areas, except as needed, and to have a general mental map of the structure of any project. Another user, Henry, claimed to be a counter example: "I've tried to get organized but it doesn't make sense. It [filing] doesn't fit the way I work . . . My filing is always emergent." However, in reality he repeated his general folder structure for each project, so that each subfolder tended to look like the other project subfolders. These personal, ritualized ways of organizing information provided the groundwork for later attempts to find the information.

FINDING INFORMATION

In our study, we asked people if they had problems or were frustrated in trying to find files and information. The study participants reported that they had no difficulty in finding files. We then asked them to tell us how they would find a file. Later in the interview, we asked them to find a specific file that we had observed while they gave us a tour of their machine. The verbal descriptions of how they would find a file and the way they actually tried to find the file matched closely. We found little variation in the way in which people found files. The predominant pattern for finding was:

1. look in a particular location
2. look in a different location
3. use the Find (or equivalent application) command
4. use a text search program

Often only step 1 was needed. Steps 2–4 were instigated as required and in the order stated.

In this section we discuss each method for finding, and then turn to a few general issues.

Location

Location seems to be of paramount importance in finding files. Often the file was in the first place the user looked. Generally, users reported that using location to find files is desirable because it is fast and reasonably accurate. They were able to quickly and comfortably get to the proper location. As Helga, a study participant, explained, "I usually know where something is . . . If it isn't where I think it is, then I'm pretty sure it will be in another spot. . . . I can find it within one of two places." The same knowledge of personal file structure was evident with almost every study participant. People had a good idea where files were on their computers; they learn the locations as they routinely access their files. They have learned the location as they have learned the content.

Archived information was most difficult to find. General

memory loss was part of the problem; memory fades over time for files which are infrequently accessed. In addition, as information was archived, it was often placed in a different location and embedded in a reorganized subset of its working context. Nevertheless, if a file wasn't in the first area the user looked, it was generally in the second area. When searching a folder, users tended to use its content to provide contextual cues such as the names, types, and number of files. If the file wasn't where they were searching, the search process helped to suggest a likely alternative.

The Find Command

The Find command allows users to search for files by name, or by other attributes such as size, creation date, date of last modification, and so on. No one routinely finds files using Find. Find was generally used in two cases (1) after failing to find a file by location, or (2) searching file servers for known applications or files. In both cases, the location of the file was unknown. A common example of this is the automatically installed application support file containing fonts, preferences settings, etc. Since the users neither placed the files in their location, nor interacted directly with the files, they don't know where they're located. The user needs to search through the extensions, preferences, and control panel subfolders, as well as the top level of the system folder. In cases like these, Find was the tool of choice.

The problems with Find are fairly well known and documented at Apple. Principal among them are the difficulty of recalling the exact name of a file. Study participants would like to have a "smart" Find that accurately finds a file based on a the best guess of a file's name. People attempt to name their files in a way that will help them remember the name. Susan, who had recently begun to use a Macintosh, and Mary, who had switched over from her last job where she had Windows, both commented on how important it was for them to be able to name a file that made sense. In fact, several of the people who considered themselves advanced Mac users said the same thing. They liked having the ability to name files in ways that made sense to them. Could they recall the names of the files when they wanted to find them months or weeks later? Not always. What did the ability to name their files

really get people if they could not recall the name?

Two things about naming are important. The first is name recognition. Although the file name might not be recalled verbatim, when searching for the file the name would be familiar when seen. The names of files are mini-indexes into the file. The other advantage of descriptive names is that when using Find, users could often remember at least part of the file name. Having part of the name enabled users to either hit the right file or find a file similar to or nearby the file they were trying to find.

Text search

"When I'm really desperate I use 'Retrieve It'," confessed a study participant. The use of terms such as "desperate," "frustrated," "frantic" expressed the feeling of being thwarted by the computer when searching for a file. These terms came up in the interviews in connection with the use of text search. Text searching was not satisfying to our users because it was slow, it produced too many wrong hits, and it was difficult to decide on the exact text string needed to bring up a document. On the other hand, these same users almost universally agreed that searching by text on the computer was better than any search method they used for paper documents.

With very fast, intelligent text search, this method of finding might well be more acceptable to users. However, intelligence, or the appearance of intelligence, comes along with its own set of problems. Systems that present documents in the order of their 'relevance,' for example, may prove very disconcerting if the system's definition of "relevance" (typically determined by a statistical algorithm), does not match the user's definition [4]. It remains an open question as to whether users will feel as comfortable with "logical" search as with "physical," location-based search.

A Note About Views

Since files are frequently found by inspection, as we have described, the "Views" feature on the Macintosh is important. Views allows users to display the contents of a folder in one of seven organizations: by name, date, size, kind, label (a user defined ordering with eight values), or icon or small icon where the user defines the organization via direct manipulation (see figure 1).

A common way to use viewing was to have the top level of the hard drive(s) viewed in icon mode, to that the user can have complete control over the location of the file relative to its cohorts. Beyond the top level, the kind of information that was present determined how it was likely to be viewed. Folders that had disparate kinds of information frequently used the icon views because the kind of information can be discerned from the form and appearance of the icon. In folders with similar information, alphabetic views were used. An exception to this rule was Adobe PhotoShop; PhotoShop has special icons that present thumbnail pictures of what is in the file. These mini-pictures do not provide very much information about the file, but it is enough information for someone who knows the files to recognize their content.

Users' familiarity with the Macintosh played a role in their use of views. In general, the greater the familiarity, the more usage of varying views. Most people selected folder views to maximize their ability to quickly scan the contents. Several users reported that the medium size icon combined with alphabetic view provided the best ration of visible information to screen space. A couple of people had customized their views so that "kind" was not showing—"kind" could be determined by the icon in the medium or large views). Most of these users also dropped the label marker—if labels were used, the user would be able to tell what the label was by the color of the medium icon. With "kind" and "label" fields visually suppressed, users were able to scan the same amount and types of information in less time than pure text, and just as importantly, in about half the screen space.

Is Location Really Helpful in Finding Computer-based Information?

Dumais and Jones [3, 6] have suggested, on the basis of

controlled experimental work, that location information is of limited utility in computer-based filing and finding. How can this be reconciled with our finding that users like location-based filing, and even seem to be able to effectively find files using it?

There are a variety of answers to this question. First, it is well established that whether through inability, lack of experience, or sheer perversity, people do not always choose the 'best' method for achieving their ends. The literature has many examples of the use of non-optimal methods, particularly in problem solving and decision making. However, we are not convinced that this is yet another example of non-optimal human behavior.

It should be noted that the relationship between the experimental work and the situation we studied is not straightforward. The experiments tried to separate spatial and symbolic components of filing: that is, they independently manipulated whether files were named, and whether they could be located on a two dimensional area. The condition with the worst performance was the location-only condition, in which filed and retrieved objects were unnamed—only their location could be manipulated. But the situation we studied was much closer to the "name + location" condition—since on the Macintosh all files have names. And, indeed, in one experiment, the "name + location" condition showed the best performance, although it was not significantly better than "name only" condition [3]. Thus, our findings are not actually at variance with the results of this work.

However, our interpretations do differ. Dumais and Jones [3, 6] concluded that location information appears to add little to performance, since removing spatial cues did not hurt filing and finding behavior, while removing symbolic

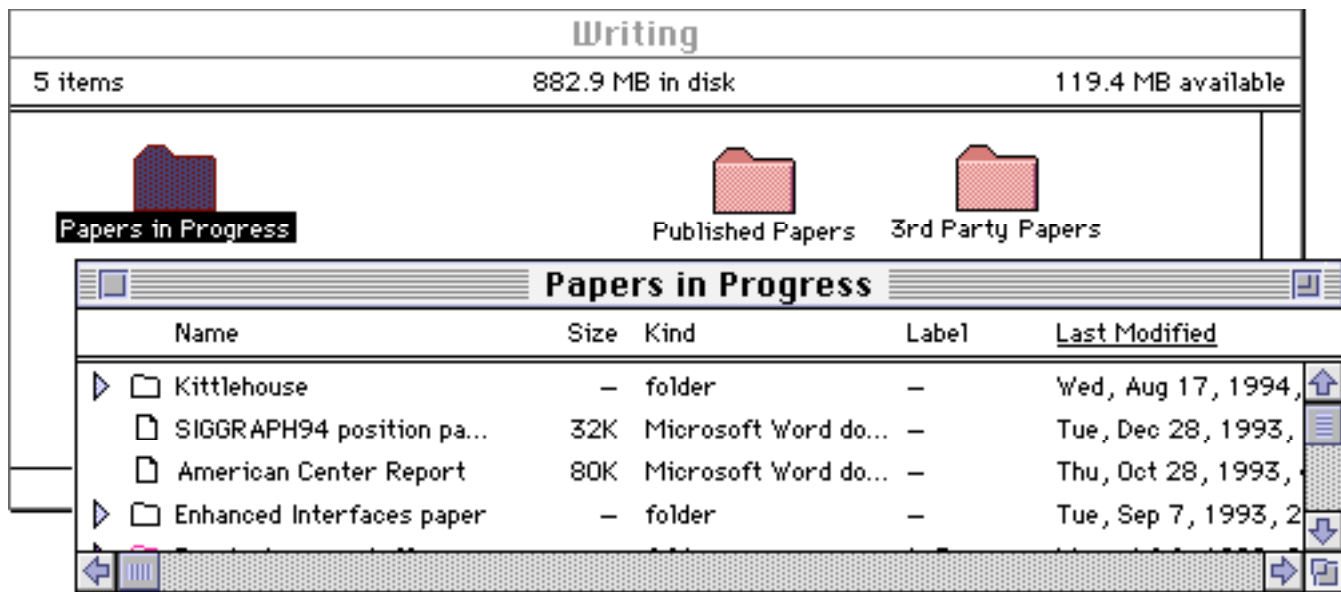


Figure 1. An example of different views on the Macintosh: view by date in the foreground, icon view in the background.

cues (i.e. names) did. On the other hand, our subjects often reported knowing where a file was, even when they didn't recall its name; and they certainly thought and spoke in terms of file locations. So, we return to the question of whether location makes a difference.

While we haven't the data to answer this question, we would at least like to leave the question open by noting some of the limitations of the experimental line of investigation.

Any experimental approach to investigating human behavior must sacrifice some of the complexity and reality that characterizes daily life to gain control over the experimental variables. Dumais and Jones' investigations differ from our observations and the reports of our users in several ways. The task in the experiment was to file forty news articles over a short time (probably less than an hour), and to find some of them again — some within the same time period, and some about two weeks later. Obviously, this pattern is not characteristic of most real world situations; although information may arrive in bursts, both the arrival and retrieval of information will generally occur over longer periods of time. Also, the news articles used in the experiment are what we would call ephemeral information — something that is not typically filed. However it may not even be fair to call it ephemeral information, in that as far as we know, the news articles had no particular relevance to the subjects' jobs or occupations (they were described as "homemakers"). Another difference is that in the retrieval task, when the subjects indicated

icons of the trash can, hard disks, and other computational objects along the right hand side of the screen, and the menu bar along the top. This begins to differentiate the space, to give it behavioral contours. Many users tend to keep open windows toward the left side of the screen, so that they can access the icons to the right without having to move windows. Users may accumulate items to throw away near the trash can before they are ready to commit to destroying them. What users actually do is not this issue; rather the issue is that even the most minimal computer screen has behavioral contours that arise from its interaction characteristics—it is not an anonymous blank space. And certainly, as the user develops a hierarchy of folders and files containing his or her personal information, the space of the computer screen becomes more complex, and more meaningful. All this is to suggest that trying to separate spatial and symbolic components of filing may not be a very good match with the real world. To us, space and symbol, location and naming, are deeply entwined.

A NOTE ON CULTURAL ISSUES

We have provided many details about the way users find and file. To step back for a moment, we can think about this information in a broader cultural light. The cultural theme that emerges from the find and file study is that users regard their computers as tools over which they want to maintain control.

Users described a sense of loss of control when they could not find files and had to resort to techniques such as text search. This was reflected in what they said and how they

Table 1. An Information Typology

where they thought an article was located, they received no feedback on whether they were correct, whereas such feedback is a natural and continual aspect of real world file finding.

Finally, the flat area on which filing was done had no personal 'meaning' to the subjects. Depending on the particular experiment, the space was either a blank area, an area with pictures of office objects such as desks and tables, or a mockup of a real office. The subjects had never interacted with the space before; it had no information of personal relevance; it was not a space in which they ever had, or ever would, carry out real tasks. The space had no behavioral characteristics. In contrast, the Macintosh places

said it, using words such as "frantic," "desperate," "frustrated." Tools such as Find, On-Location, and Retrieve-It are last resorts for finding. They do not give people using computers the feeling of power. The feeling of mastery over the computer comes with the knowledge of what is on one's computer and where it is located. We should keep in mind that filing and finding mechanisms should impart this sense of mastery that seems to derive from direct knowledge of where important documents are, and direct ways of making them visible, rather than indirect abstract methods that return an often disappointing list of irrelevant files.

CONCLUSION

We suggest that understanding the organization and retrieval of computer files can be aided by thinking of information as being of three types: ephemeral, working, and archived. Each type is defined by the characteristics of the information, the relation of that information to the user's activities, and the way in which the user interacts with the information (see Table 1 for a summary).

Researchers take it for granted that users want better ways to organize and retrieve files. We hope we have shown that ordinary users are relatively happy with their computer-based filing systems. The reason that researchers see the world differently is that they have so much more archived personal material than everyone else. For researchers, most information has an incredibly long shelf life. One doesn't like to get rid of things because they might come in handy, someday. It's impossible to know which things one will need in the distant future, so one builds as extensive a library as possible.

In our study, we found that ordinary users either archive formally in databases or sparsely in personal systems. Most users, as we saw in our study, manage this personal information quite well. They interact with information that has a much shorter shelf life. Information of relevance to a salesperson, marketing expert, or secretary is usually ephemeral or working information. Large amounts of archived information such as customer orders and legal documents is kept in carefully managed databases, not in personal filing systems. The findings of our study are echoed in that of Kidd [7], who found that many

“knowledge workers” do not rely heavily on information once it has been filed; hence retrieval is not a key task for them.

But perhaps this will change. The increasing integration of communications technologies into computing system could result in a large increase in the amount of ephemeral information available to the ordinary user. And it's difficult to see how user's could cope with their current filing and finding practices. There seems to be no alternative to text-based search in this case, although it will have to be extremely “intelligent” to match the feelings of ease and control imparted by location-based search. Intelligent agents, and the question of how to make one intelligent enough to produce good results while remaining understandable enough to be controlled by the user, remains as an interesting research challenge.

Regardless of the successes or failures of particular technologies, it is useful to consider the interaction between information characteristics, the user's tasks, and how the user interacts with the information. Thus we make the point again that information is not a monolithic category; it's useful to look at a user's mix of ephemeral, working and archived information to understand their filing and finding needs.

ACKNOWLEDGMENTS

We would like to thank Dan Rose for comments on an earlier version of this paper.

Ephemeral Information	
examples	down-loaded news articles; to do items
source	bulletin boards, mailing lists, transcribed information
relevance	relevant to current or near-future tasks
shelf life	minutes to days
principal uses	reminders; activity triggers
how accessed	not filed or loosely filed; kept visible to increase likelihood of triggering action
Working Information	
examples	working papers; project schedules
source	created by user or colleagues
relevance	relevant to current tasks
shelf life	days to weeks to months
principal use	codification of project related information
how accessed	spatially for recently used info, categorically as familiarity decreases
Archived Information	
examples	technical reports; project email archive
source	former working information that has been winnowed and reorganized
relevance	indirect
shelf life	months or years
principal use	occasional reference
how accessed	categorically and spatially (and using search tools)

REFERENCES

1. Blomberg, J., Suchman, L., Trigg, R. (Forthcoming, 1994.) Reflections on a Work-Oriented Design Project. Proceedings of the Participatory Design Conference (PDC 94), Chapel Hill, NC, October 27-28, 1994.
2. Cole, I. (1982). Human Aspects of Office Filing: Implications for the Electronic Office. Proceedings of the Human Factors Society, Seattle WA
3. Dumais, S. T. & Jones, W. A. (1985) A Comparison of Symbolic and Spatial Filing. CHI '85 Proceedings. ACM Press
4. Erickson, T. & Salomon, G. (1991) Designing a Desktop Information System: Observations and Issues. CHI '91 Proceedings. ACM Press.
5. Harris, J.E. (1980). Memory aids people use: Two interview studies. *Memory and Cognition*, Vol. 8 (1), 31-38.
6. Jones, W. A. & Dumais, S. T. (1986) The Spatial Metaphor for User Interfaces: Experimental Tests of Reference by Location versus Name. *Transactions on Office Information Systems*. ACM Press.
7. Kidd, A. (1994). The marks are on the knowledge worker. In *Proceedings CHI '94*. 24–28 April, Boston.
8. Lansdale, M. (1988). The psychology of personal information management. *Applied Ergonomics* 19, 55–66.
9. Malone, T. (1983) How do people organize their desks? Implications for the design of office information systems. *ACM Transactions on Office Information Systems* , 1, 99–112.
10. Mander, R., Salomon, G. and Wong, Y. (1992). A Pile Metaphor for Supporting Casual Organization of Information. CHI '92 Proceedings. ACM Press.
11. Suchman, L. and Wynn, E. (1984). Procedures and Problems in the Office. *Office Technology and People*, Vol 2, pp. 134-154.