# Printed Material as a Knowledge Representation

Dr. Bob Jansen[*]

Dr. Bob Colomb[+]

Professor Ann Henderson-Sellers[&]

Dr. Janet Gallagher[%]

Mr. John Robertson[*]

Mr. George Bray[@]

## Abstract

This paper presents results of a research program into the utilisation of printed material as a knowledge representation in Knowledge Based Systems. Knowledge Based Systems are generally unable to explain or justify their behaviour, and we attribute this to their general lack of suitable knowledge in a suitable format. This research program is evaluating the representation of printed material using hypertext and hypermedia technologies to provide a navigable hyperspace; using explicit representations of the chunks of the domain knowledge as found in the printed material to guide the navigation process.

[*] CSIRO Division of Information Technology, Knowledge Based Systems Laboratory, PO Box 1599, North Ryde, NSW 2043. Email - jansen@syd.dit.csiro.au

[+] University of Queensland, Department of Computer Science, Brisbane Qld.
Email - colomb@cs.uq.oz.au

[&] Macquarie University, School of Earth Sciences, North Ryde NSW 2043.
Email - ann@mqclimat.mqcc.mq.oz.au

[%] Defence Science Technology Research Organisation, PO Box 1600, Salisbury SA 5105 . Email - gallagher@itd.dsto.oz.au

[@] c/o CSIRO Division of Information Technology, Knowledge Based Systems Laboratory, PO Box 1599, North Ryde, NSW 2043. Email - bray@syd.dit.csiro.au

This paper has been presented as an invited presentation at HyperOZ'92 in Adelaide, 21 February 1992.

## 1. Introduction

This paper presents results of a research program into the utilisation of printed material in Knowledge Based Systems.

Knowledge Based Systems are generally unable to explain or justify their behaviour, and we attribute this to their general lack of suitable knowledge in a suitable format. This research program is evaluating the representation of printed material using hypertext and hypermedia technologies to provide a navigable hyperspace; using explicit representations of the chunks of the domain knowledge as found in the printed material to guide the navigation process.

We describe two prototype systems that we have developed to test out our ideas, the Wool Technology Dark Fibre Risk prototype and the Greenhouse prototype. We specify a data model to support the representation of research papers, and discuss a novel knowledge representation, termed an assertion, that facilitates the mapping from computationally efficient knowledge representations to cognitively efficient knowledge representations. The data model is based on the hypothesis that there are many representations for the domain knowledge available to an expert, namely text, graphics, video, audio, computer program, etc., and each of these should be available to the Knowledge Based System if we expect the Knowledge Based System to perform as an expert. It should be recognised however that the provision of more information does not guarantee more intelligent processing, and thus the Knowledge Based System must have access to better quality information.

Section 5 discusses the use of printed material as a knowledge representation and its integration with a computational environment. We present our ideas regarding the hyper-editorial work involved in the authoring of such systems, and argue that this process requires extensive computerised support. We briefly discuss some current trends in electronic support for the authoring process.

Section 6 concludes this paper with a summation of the discussion.

## 2. Background

The focus of the current research program has its roots in two earlier research projects, the Wool Technology Dark Fibre Risk Prototype (Jansen & Robertson 89, Jansen 91) and the Greenhouse Prototype (Colomb *et al* 91). The research in both these prototype

systems was aimed at using knowledge representation techniques to facilitate the navigation through a complex *hyperspace*. Hyperspace is the term used to describe the data, anchors, and linkages that together form a hypertext/hypermedia system.

The hypothesis under consideration in these prototypes was; that the provision of simpler navigation utilising explicit structures describing both the domain knowledge and the sources of the domain knowledge would provide an end user with the link between the domain knowledge and the source of that domain knowledge facilitating the production of more suitable explanations and justifications of system behaviour. This hypothesis was to be tested by the utilisation of domain documentation as a context-sensitive knowledge representation and providing linkages with other representations of that domain knowledge.

The Wool Technology Dark Fibre Risk prototype was the result of a collaborative project with the CSIRO Division of Wool Technology. The project involved the representation of several disparate sources of domain knowledge and the explicit inter-relationships between various representations of chunks of the knowledge. The domain was that of the risk of dark fibre in shorn wool.

The Greenhouse prototype built on the results of the Wool Technology Dark Fibre Risk prototype but in the domain of carbon dioxide and the greenhouse effect. This project was a collaborative research project with the Macquarie University School of Earth Sciences led by Professor Ann Henderson-Sellers[1]. The aim of this project was to capture 100 core research papers as identified by our experts and provide similar functionality to the Wool Technology Dark Fibre Risk prototype.

## 3. The Wool Technology Dark Fibre Risk Prototype

The aim of the  Wool Technology Dark Fibre Risk prototype project was to provide a computer system whereby wool growers, buyers, and  classers could grade wool lots according to their dark fibre risk, and to gain access to the latest information, both research and commercial, regarding the dark fibre risk evaluation and the effect of various wool growing practices, eg. husbandry, on the final dark fibre risk.

---

[1]Dr. Mervyn Jones, formerly of the School of Earth Sciences should also be acknowledge as a major partner in this project.
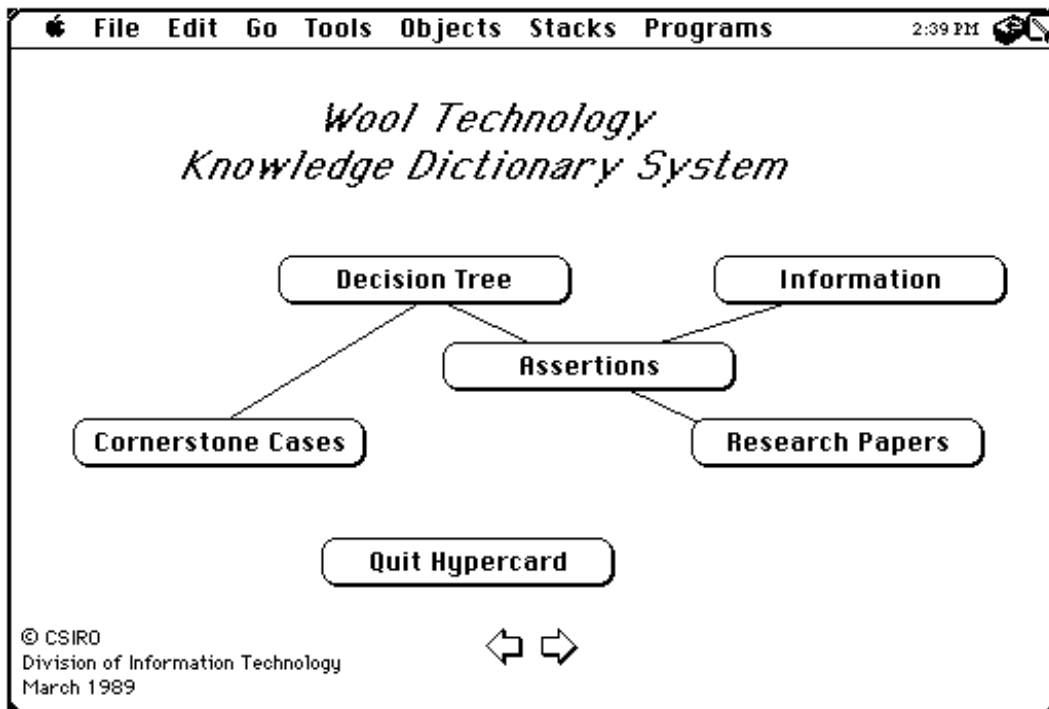
Figure 3.1 - System architecture of prototype environment

This system consisted of a set of different representations of the domain knowledge as shown in figure 3.1, each representation supporting various processing functions. The decision tree was the computational representation, the cornerstone cases represented important cases within the domain that caused the decision tree to be amended. The research papers were a static textual representation of the domain knowledge. The challenge in this environment was to link each of these representations together and use them in better ways so as to support explanation and justification of system behaviour.

The research hypothesised the existence of a novel knowledge representation which we termed an *assertion* (Jansen 91). An assertion was defined as an important statement made by an author in any section of a paper and was hypothesised to provide the link between the chunks of domain knowledge in the Knowledge Based System and the sources for that domain knowledge. The original authors were asked to read their papers[2] and identify their assertions. These were extracted and inserted into the assertions database which acted as a mapping function between nodes of the decision tree and the various sections of the research papers. A total of 150 assertions were identified by the authors, an explosion rate of 30:1.

---

[2] In this prototype, we only used five research papers.

A lesson learned from this work was that the initial definition of an assertion was too vague. In fact we began to question the ability of authors to consistently extract assertions out of their papers[3].

## 3.1 The Research Papers Model

The research papers database consisted of research papers in computer readable form. This prototype used a 'standard' and simple representation of a research paper as shown in figure 3.2.
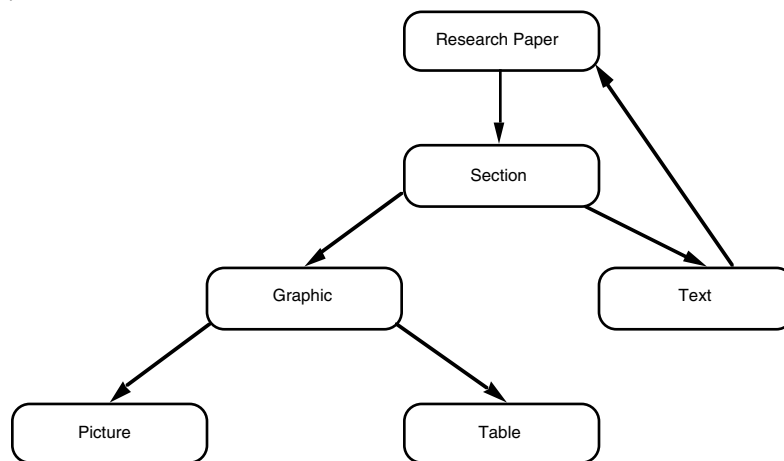


Figure 3.2 - hierarchic decomposition of a research paper. Note that the link between 'text' and 'research paper' objects depicts a reference to another paper found in a piece of text.
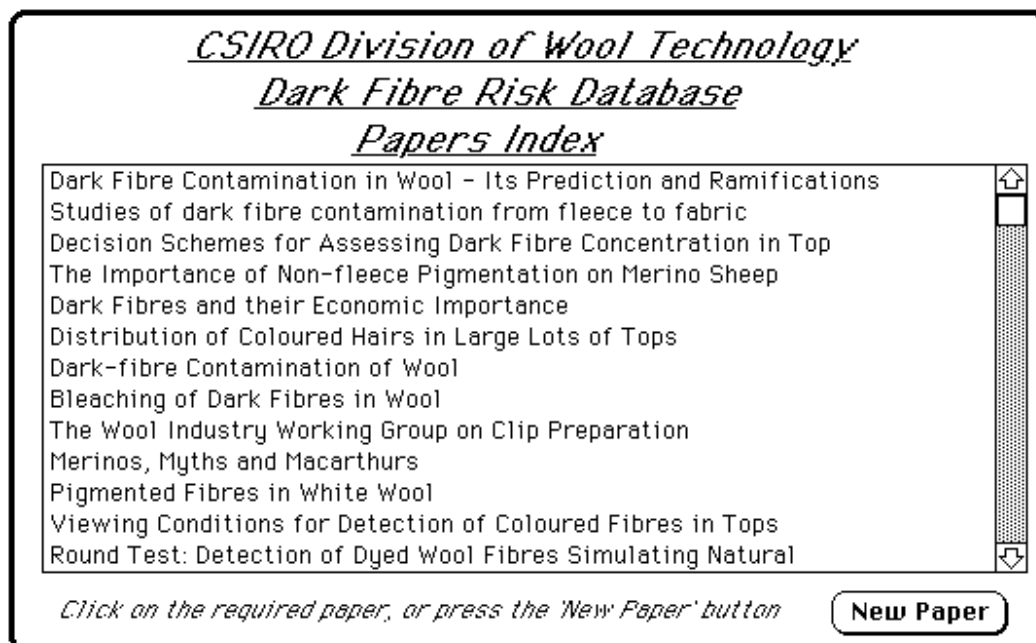


Figure 3.3 - research papers database index. This is the entry point to the research papers database. Each paper's title is shown in the scrolling field of this card. Papers are selected by clicking on the title. This will cause navigation to the paper's heading page, as shown in figure 3.

---

[3] This was the subject of a follow up research project and is described in  Rantanen 91

Author
R. A. Foulds
CSIRO Division of Wool Technology
Sydney Laboratory

Publication
Proceedings of the 7th International Wool Textiles Res
Conference, Tokyo, 1985, 65-74
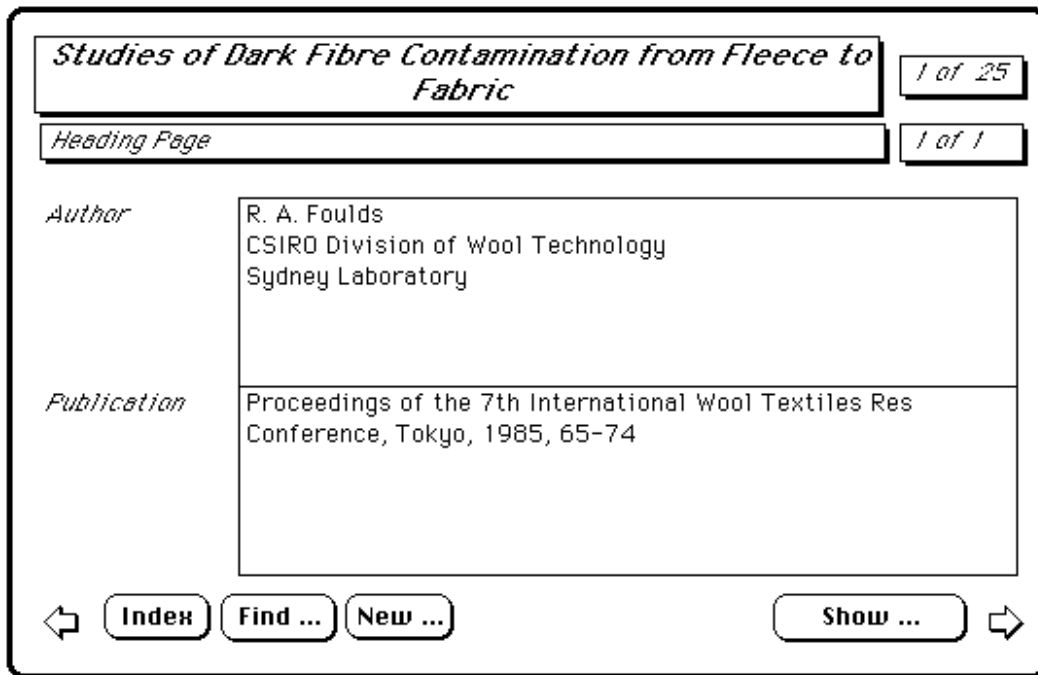
[ Index ] [ Find ... ] [ New ... ]     [ Show ... ]

Figure 3.4 - a research paper's heading page. This page provides publication details of the paper. Access to various parts of the paper is effected by the forward/backward browse arrows, or the 'Find ...' button.

It can be clearly shown that a dark fibre test for greasy wool which would relate to this contamination level is not viable as it would require an impracticabbly large sample [1] to be representative of a given lot of wool. Where contamination of merino wool is caused by urine-staining we consider a reasonable assumption to be that the stained wool will be distributed in the form of discrete clumps of size similar to wool staples, at the rate of say, four per bale. The impossibility of economic sampling of raw wool for this order of contamination intensity is in contrast to that of other fleece characteristics such as fibre diameter, yield, staple length and staple strength. It can be shown that the chance of collecting a dark staple in one coring is 1 in 1000 [1]. (Alternatively, one can say there are 999 chances out of 1000 of missing a dark staple with one core.)

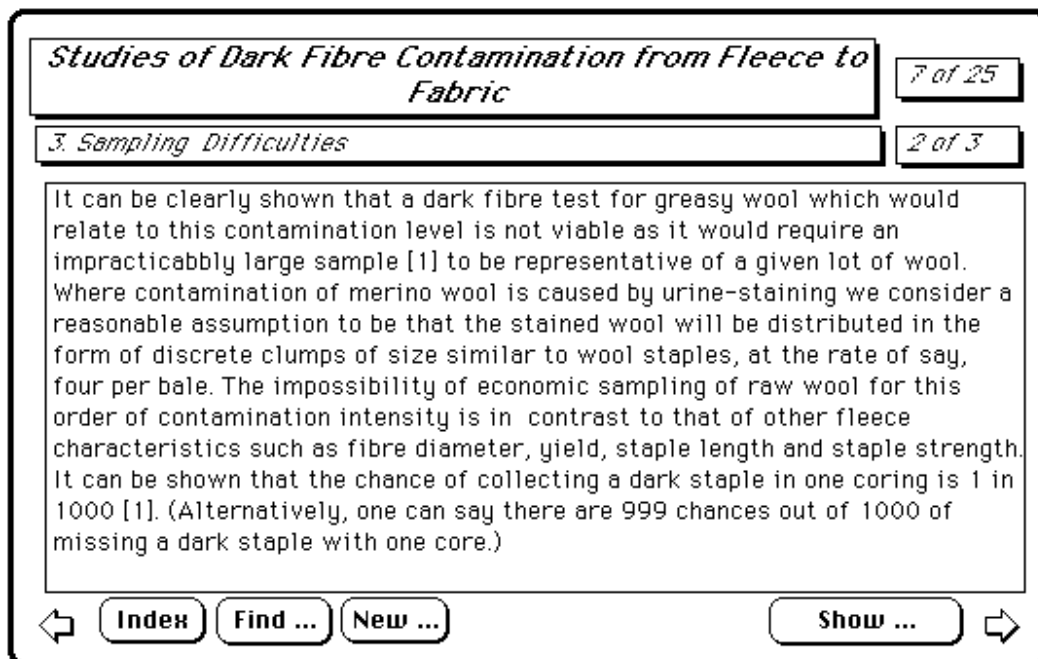[ Index ] [ Find ... ] [ New ... ]     [ Show ... ]

Figure 3.5 - a section of text from a research paper. As sections of text can be arbitrarily large, each section spans a number of cards. Hence the two 'x of y' fields in the upper right of the card. The upper one shows how many cards there are in the paper, and which card this is (card 7 of 25), and the lower shows how many cards there are in the section, and which this is (card 2 of 3). The section title is given as '3. Sampling Difficulties'.

Figure 3.6 - an example of an 'active' table, where the table is stored in a spreadsheet program instead of in the body of the text. Access to the table is via action points in the text.

This structure was implemented in the 'Research Papers' database, using the layout as shown in figures 3.3, 3.4, 3.5, and 3.6.

Figure 3.5 also shows the method of linking citations and graphics into the text. This figure has citation points (eg. [1]) imbedded into the text in a similar way to printed papers. In this database, the citation is defined to be an action point, so that to go to the reference, the user merely clicks on the citation, and a hidden process automatically takes them to the heading page of the referenced paper. Graphics are treated in an identical way. References to figures and tables in the text are defined as action points with an associated procedure to go to the required graphic. The prototype phase has experimented with the concept of '*active tables*' where, rather than storing the table in a passive way as columns of numbers as on a printed page, the table has been implemented using a spreadsheet program, and thus the numbers can be manipulated by the user. This would increase the communications bandwidth between the user and the knowledge source. Manipulations that were supported included the graphing of the table, or a subset thereof, in a variety of ways as determined by the user, and statistical analysis as provided by the spreadsheet tool. Obviously in using active tables, mechanisms must exist to prevent the user from changing the stored form of the table,

but in the loaded form, the user may make any changes, supplying a limited decision support/what-if capability.

## 3.2. Assertions

It was assumed that, from a data modelling point of view, an assertion, as stated previously, could be used to map a particular node of the inference tree to a particular area of a research paper. The situation in reality is more complex than this. An assertion may refer to many nodes or paths through the inference tree as well as to many sections of the papers. This was, however, deemed beyond the scope of the prototype which aimed at showing the usefulness of the assertion as a knowledge representation. As a further simplifying measure, the system did not recognise duplicate assertions. Instead each assertion was manually edited to ensure uniqueness.
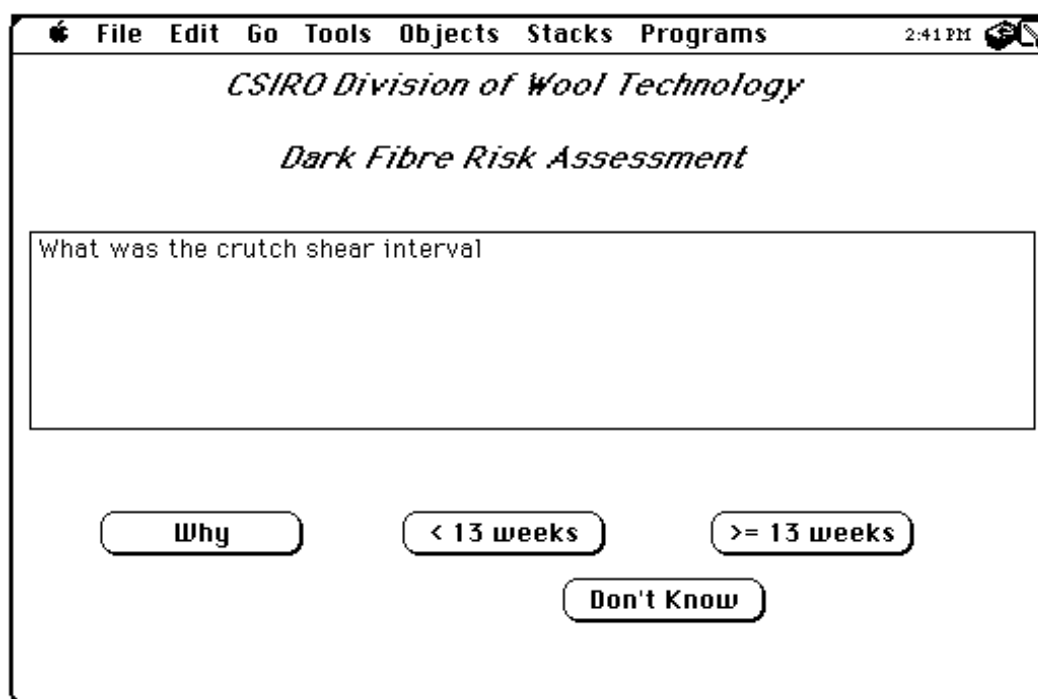


Figure 3.7 - question posed by expert system with the availability of a "Why" facility

The production of the assertion representation of the papers enabled a more intelligent, context dependent explanation and justification facility to be implemented in this environment. Functionally, the user, when asked for the value of a data item, might elect to request an explanation. At this point, there are interesting philosophical questions regarding the semantics of the 'why' facility, namely what does 'why' actually mean in this context? Can it be answered by showing the current rule? Is the user asking what the data item represents, for example what does the term "crutch-shear interval" in figure 3.7 actually mean? What is the importance of the value 13 for crutch-shear interval, why

not 42? Is the user interested in the relationship between the data item, the value to be supplied, and the result of the Knowledge Based System (ie. the dark fibre risk rating supplied)? Each of these possibilities can be supported by the assertion representation, assuming of course that all the assertions have been represented.

In this prototype, the request for an explanation, expressed by pressing the "Why" button, led to a search of the assertion list using a fixed keywords-out-of-context (KWOC) facility, retrieving assertions containing predefined keywords. The retrieved list was displayed to the user (figure 3.8).
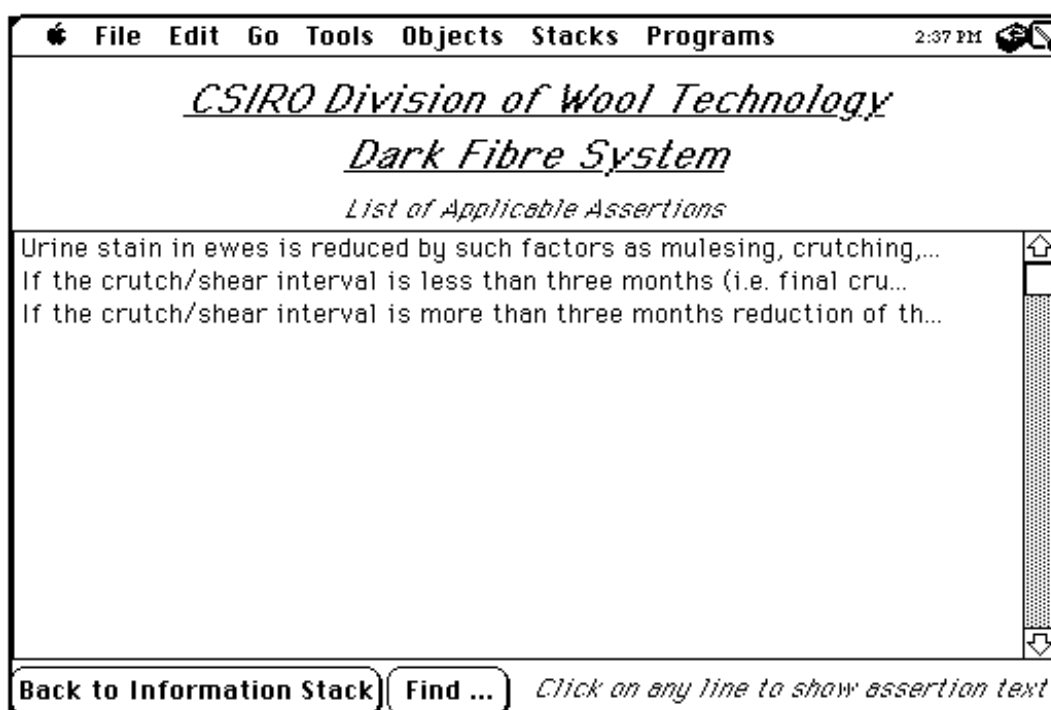


Figure 3.8 - List of assertions relevant to the question

If required, the user could select a particular assertion from the list and request its full display, instead of the first 70 characters as in the assertion list (figure 3.9). In this case, the user was presented with the full text of the assertion and its unique reference number.

Once the full text of the assertion was displayed, the user could request a justification of the assertion, leading to a display of that section of text from a research paper that caused the assertion to be represented (figure 3.5). If the text of the assertion was an exact match of any part of the text in the section of the paper being displayed, then this had been edited to ensure readability in its out-of-context representation. Note that if the assertion text matched the source text, the source chunk was highlighted.
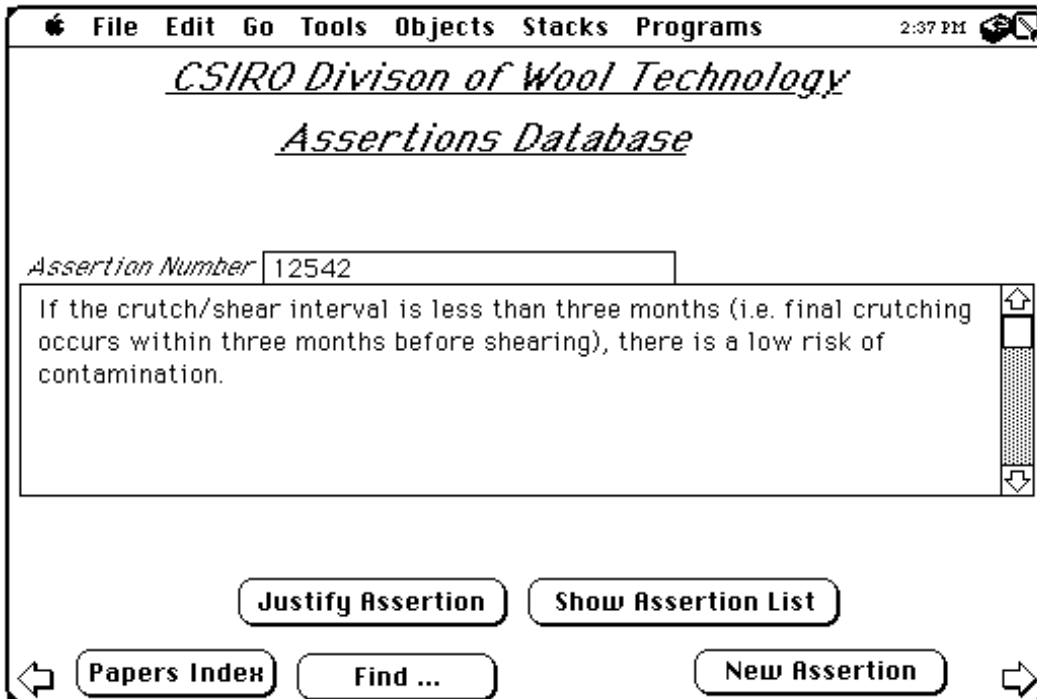
Figure 3.9 - full text of one of the relevant assertions, ie. number 2 in figure 3.8
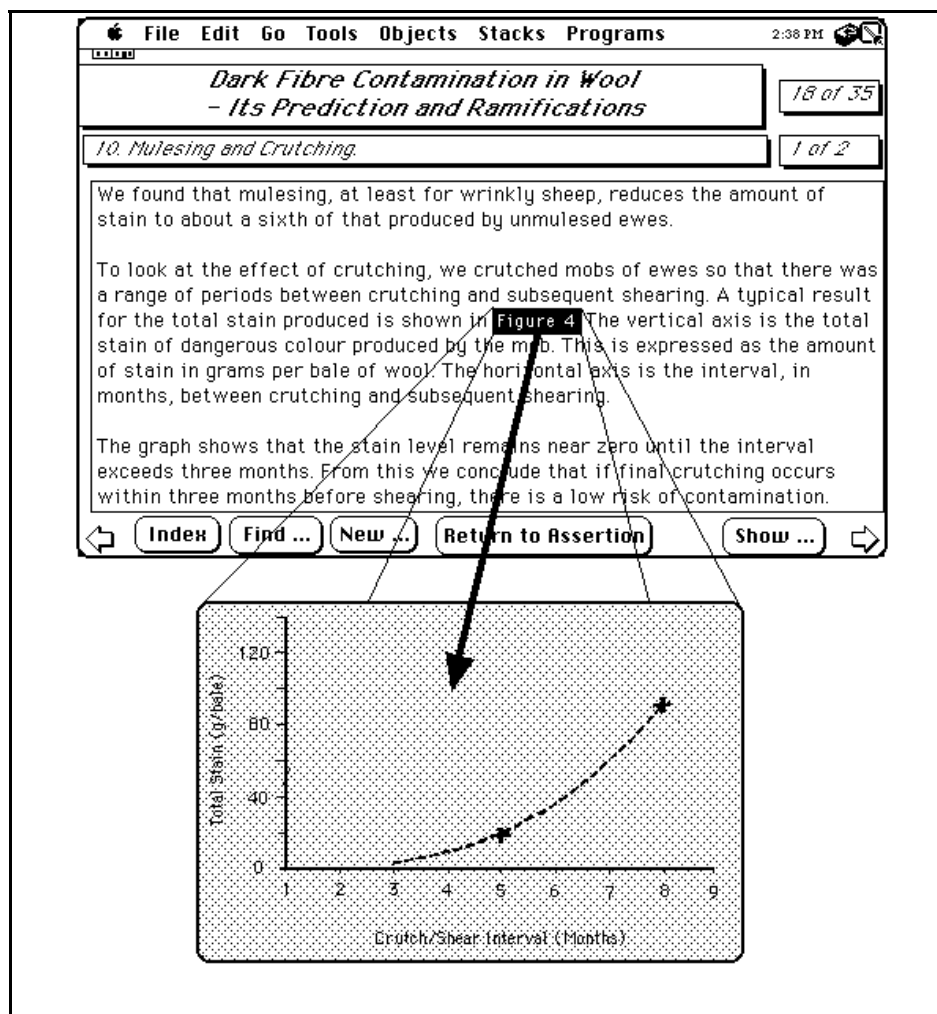
Figure 3.10 - section of text that produced the assertion and hypertext link to accompanying graph

This situation seemed rare, as in most cases the assertion hypertext browsing capability formed the link to another representation of the knowledge, a graph, from the reference point "figure 4" (figure 3.10).

## 3.3. Results of this Prototype

This prototype showed that in any domain there will be a number of disparate representations of the domain knowledge. A similar finding has been previously reported in Kornell 87.

Each knowledge representation can be utilised to support different functions, functions not easily supported by a single representation (eg. the assertion enabled the retrieval of contextual information from the document sources). The explicit representation of all the different forms of the domain knowledge in their respective media was a suitable technique to support flexible browsing, explanation and justification of system behaviour provided that the media were accessible from the computational environment. This requirement necessitates the use of hypertext/hypermedia techniques both as a storage/access mechanism as well as supporting the user interface.

## 4. The Greenhouse Prototype

As stated previously, the Greenhouse Prototype extended the findings of the Wool Technology Dark Fibre Risk prototype but in the domain of the role of carbon dioxide in the greenhouse effect. The knowledge consisted of 10 research papers and a model of the production of $CO_2$ from the global economy. It was conceived of as a tool to make greenhouse knowledge accessible to policy planners. In particular, the model had a hierarchical hypertext front-end for entry of parameters and had its output directed to a Wingz™ spreadsheet so that it could be manipulated and graphed by the user. Important parameters and outputs were linked to research reports which describe their meaning and significance.

A second version of this prototype largely addressed issues of scale. First, the Hypercard™-style substrate was not entirely suitable for hundreds of documents. Second, considerable editorial work on the part of domain experts was needed to add the necessary structure to the knowledge. The first version had made it much easier to actually implement the links, but extracting the intermediate structures (assertions,

concepts) was quite laborious. The second version was developed for a knowledge base of about 100 documents, but its architecture can be scaled to systems containing at least 1000 documents in a coherent area of knowledge.

The second version was never fully completed but has served as a springboard for a third version currently being developed. This will be the subject of a later report. The third version is investigating how an expert uses extensive domain knowledge in supporting tasks, as well as extending the research on assertions using linguistic expertise. The central component of this version is the use of a conceptual space representing and describing the domain knowledge in its various representations as a mechanism for linking each chunk of domain knowledge. One benefit of this approach is the availability of qualitative modelling of the causal representations within the conceptual space

## 5. Printed Material as a Knowledge Representation.

Each of the prototypes described above had in common the use of printed material as one representation of domain knowledge. The research that forms the focus of the current research program[4] aims at discovering the nature, integration of, and use of printed material in an electronic environment supporting a Knowledge Based System. The support can cover areas such as: knowledge acquisition, where the printed material acts as the knowledge source for the knowledge engineer; education, where the printed material acts as the source media used to convey knowledge to an end user as a direct result of interacting with this media; or query answering, where the printed material acts as the substrate through which an end user may navigate to discover the answer to some problem or support some task.

We will confine our interest to that printed material associated with technical or reference publications.

### 5.1. Technical Documentation

Technical documentation exists in a large part as journal articles and conference proceedings, but also as monographs. In monograph form it is also called a *book*.

---

[4] The Knowledge Based Systems Program of the CSIRO Division of Information Technology of which the first author is the program manager.

The journal article is essentially a small object: generally consisting of about ten to fifteen pages. There are instances of articles that are much larger, but the usual constraint on journal articles is the physical size of the journal. Books are generally larger and usually result from an amalgamation of several smaller publications that together form a solid body of knowledge.

The structure of both forms is essentially identical, at least as far as the formal data storage is concerned. The main difference is found in the various indexes available to access the various parts of the publication. The model shown in figure 3.2 forms the basis for our work in storing and navigating through electronic forms of printed material. A more complete model is shown in figure 5.1 on which the Greenhouse prototype versions and the Hyperbook development were based
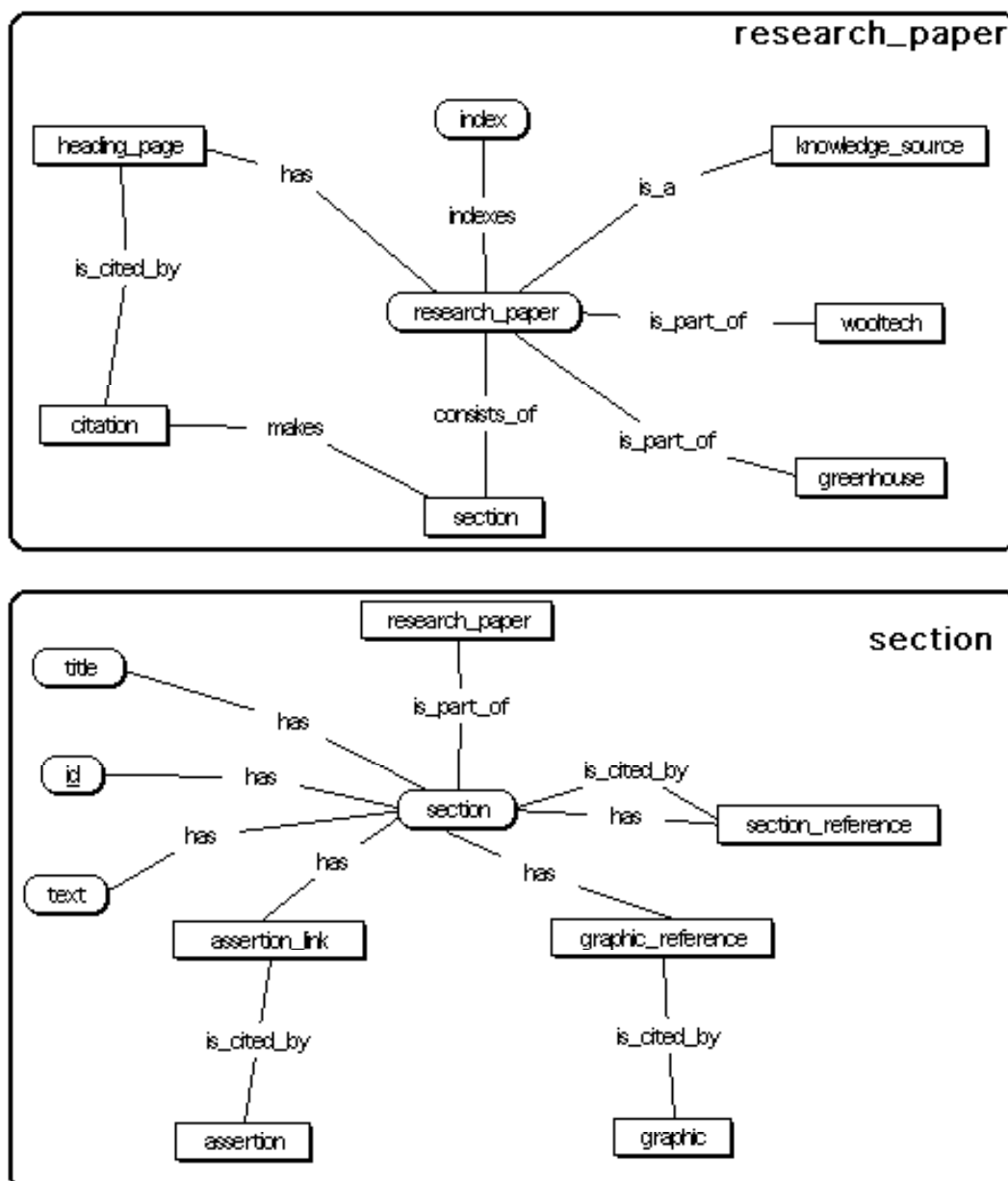
Figure 5.1 - more complete data model of printed material. The top diagram shows the structure of the research paper, whilst the bottom diagram shows the structure of a section of text.

As shown in figure 5.1, printed material has many semantically rich components, each of which can serve as a useful source for domain knowledge either in knowledge acquisition, navigation or education modes.

The traditional static representations include text, tables, graphics, and citations or references to object outside the bounds of the publication. Each of these representations serves a singular purpose: to provide the reader with chunks of knowledge.

## 5.2. Authoring Active Printed Material

The electronic representation of printed material presents the developer with the opportunity of utilising these knowledge representations in ways not available to authors and readers of traditional printed media. The static representations can be activated, or become *hyperactive*, through the use of hypertext/hypermedia technology. Activation of these representations increases the bandwidth of the communication channel between the author and the reader. The realisation of increased communication bandwidth is a major area of interest in the Man-Machine Interface (MMI) community.

Candlin & Saedi 83 discusses the authoring and reading processes as two separate discourses each having their own characteristics. The successful production of electronic material depends largely on supporting this discourse, both for the author(s) and the readers. It should be recognised that these processes of discourse are complex; so complex that current technology is only able to pay 'lip service' to them. There is no hypertext environment that has not suffered from the 'lost in hyperspace' problem, or that supports the authoring process to such an extent that the authors confidently input directly into the environment. The input of the text and graphics is but one process that faces authors. A more complex and time-consuming task is the creation of the hyper-environment, the *hyper-editorial* work; the anchors and links that enable the reader to perform associative browsing.

Our experience has shown that the engineering of the static representations into active representations is an art rather than a craft. The traditional static representations have dynamic analogues that are more than just their repetition in a dynamic environment. The process of transforming the static representations into their active analogues

requires extensive design and planning to ensure that the author's and reader's expectations can be met.

There are as yet very few tools available to aid in the hyper-editorial process. Traditional authoring is widely supported by a myriad of word-processing environments, each of which is more than capable of taking words typed in by the author and storing them in electronic form. In addition, the better word-processors supply functions that enable the author to process the text in a variety of ways. For example, the word-processor I am using, enables me to plan a tome using an outlining process, build and maintain contents pages and hierarchical indexes, etc. It will not however enable me to identify and define anchors and links that can be subsequently input directly into my hyper-environment. For this process, the author must rely in the large part on traditional hand methods.

There is research to support this hyper-editorial work by the use of *mark-up languages*. A mark-up language is a language for identifying and defining the structure of printed media. Subsequent processing of the 'marked-up' text makes it possible to automatically create the anchors and links in the hyper-environment.

Another approach addresses this problem by parsing the written text looking for syntactic clues as to document structure or 'important' phrases. *Aida* , a commercially available tool, can parse printed text and produce a summary of that text. The success of Aida is directly related to the effort the authors have expended in writing their text; Aida assumes for example that the first sentence of each paragraph identifies the subject of that paragraph. Although this is seen as 'good writing' practice, this is unusual to find in most scholarly texts.

We have conducted a simple experiment to assess Aida's usefulness[5] for a particular task. As described earlier, the Greenhouse Project identified assertions from ten research papers by asking domain experts, in lieu of the authors, to identify the important points in the papers. The results indicated that even on one paper where <u>two of the experts were joint authors</u>, the amount of agreement in 'the important points' was less than 20%. We ran these papers against Aida and extracted three levels of summary, 5, 10, and 20%. Simple eye-ball comparison shows that Aida's technique of using syntactic clues has almost no overlap with the 'important points' identified by the domain experts.

---

[5] The detailed results of this experiment will be presented in a later paper.

The above result is not surprising if one recognises that if the authors of a joint paper can not agree as to what is important, how can we expect a syntactic parser to sort it out. Rantanen 91 reports on a study of the assertion recognition undertaken by the Greenhouse experts. This study concludes that each reader uses different criteria for identifying the assertions. This conclusion, taken on its own, would imply that an automatic identification process is impossible, since every reader would need to run it for themselves in relation to the current task, etc. (in fact the *context* of their current process of discourse).

The study represented in Rantanen 91 used the KSS0 tool (Shaw & Gaines 1989) developed by Brian Gaines at the University of Calgary. This tool enables a cluster analysis to be performed on data, and this was performed on the results of the interviews of the individual experts. This analysis shows a clustering of assertions into sets, and our interpretation is that although each expert uses their own criteria in identifying the assertions, they do appear to identify related areas of the conceptual space as shown by the clustering effect. Thus automatic support for this process is still feasible. The development of such support is in our opinion dependent on the use of **contextually sensitive** semantic and linguistic techniques to identify the assertions.

Another approach is the use of cooperative authoring environments. *Cooperative*, or *collaborative*, authoring (Kennington *et al* 88, Seeley & Leadbetter 88, Begeman & Conklin 88) uses advanced technology in supporting many authors in constructing a document. One technology used is hypertext, and this approach places many more problems on the authoring process. Collaborative authoring does not address the problem of anchor and linkage identification and representation *per se*, but if used in conjunction with some of the techniques outlined in this paper, may provide valuable insight into the authoring process.

## 6. Conclusion

This paper has presented results of a research program into the utilisation of printed material in Knowledge Based Systems.

We discussed two prototypes that resulted in our present understanding of the authoring and representation of electronic documents. We discussed the issues involved with the authoring process and speculated that complex computerised support is required if we are to take-on existing printed material. For new material, advanced collaborative authoring systems based on today's word-processing systems will have to

be developed. The use of collaborative authoring techniques may provide valuable insight into the authoring process. Similar technology will have to be developed to unravel the browsing and navigation process especially if the system is to support associative access.

The issue of context was briefly raised as an issue that must be addressed to provide intelligent tools and environments. The study of context, although not described in any great detail in this paper, underpins the current research program.

## Acknowledgments

## References

Begeman ML & Conklin J, *The Right Tool for the Job*, Byte, October 1988, pp255-268

Candlin CN & Saedi KL, *Processes of Discourse*, Journal of Applied Language Study, Vol. 1, No. 2, August 1983, pp103-133

Colomb RM, Robertson J & Jansen B, *CSIRO Hypertext Research Project*, in Proceedings of the Australian Database-Information Systems Conference 1991

Jansen B & Robertson J, *Management of Wool Dark Fibre Risk Knowledge Using Hypertext*, CSIRO Division of Information Technology, Technical Report TR-FD-89-05, May 1989

Jansen B, *Two Expert System Applications: Implications for Knowledge Representation for Explanations and Justifications*, in Proceedings of the World Congress on Expert Systems (WCES'91), Florida USA, December 1991

Kennington RW, Seeley DA & Morales JR, *A Collaborative Document Editing System*, Info Tech: Towards 2000, Proceedings of the ACC'88, Sydney, 89-110

Kornell J, *Formal Thought and Narrative Thought in Knowledge Acquisition*, International Journal of Man-Machine Studies (1987) 26, pp203-212

Rantanen J, *Knowledge Acquisition for Assertion-Based Knowledge Representation*, CSIRO Division of Information technology Technical Report TR-FD-91-01, 1991

Rantanen J, *Knowledge Acquisition for Assertion-Based Knowledge Representation*, CSIRO Division of Information technology Technical Report TR-FD-91-01, 1991

Seeley DA & Leadbeter M, *Supporting Group Innovation with a Hypertext Authoring System*, Info Tech: Towards 2000, Proceedings of the ACC'88, Sydney, 1988, pp 133-157

Shaw MLG & Gaines BR, *Comparing Conceptual Structures: Consensus, Conflict, Correspondence and Contrast*, Knowledge Acquisition, Vol. 1, No. 4,1989,  pp341-363