



Workshop on Social and Collaborative Construction of Structured Knowledge,
16th International World Wide Web Conference, Banff, Canada, May 8, 2007

Formalization, User Strategy and Interaction Design: Users' Behaviour with Discourse Tagging Semantics

Bertrand Sereno*, Simon Buckingham Shum & Enrico Motta
Knowledge Media Institute, The Open University, Milton Keynes, UK

* Now at: Centre for Advanced Learning Technologies, INSEAD, Fontainebleau, France



Licensed under Creative Commons
Attribution-ShareAlike 2.0 License

Acknowledgements:





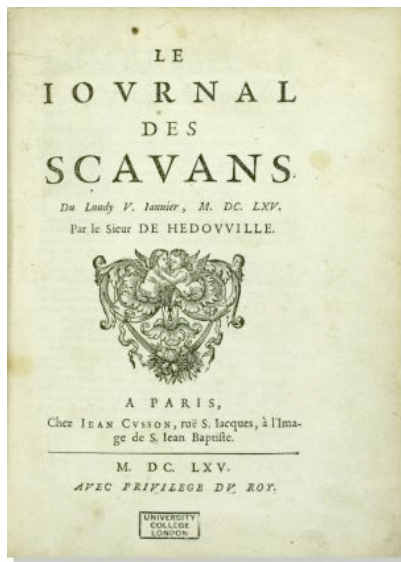
Overview

- **The problem: collaboration semantics in contested domains** — no consensus assumed; possibly not even possible
- **Previous work: Scholarly discourse as Collaborative Knowledge Structuring (CKS)**
 - *Modelling and querying Web argument structures*
- **How to help users engage in CKS?**
 - *Evaluating the ClaimSpotter tool*
- **Summary of evaluation results and design principles**
 - *Formalization / User Strategy / Interaction Design*

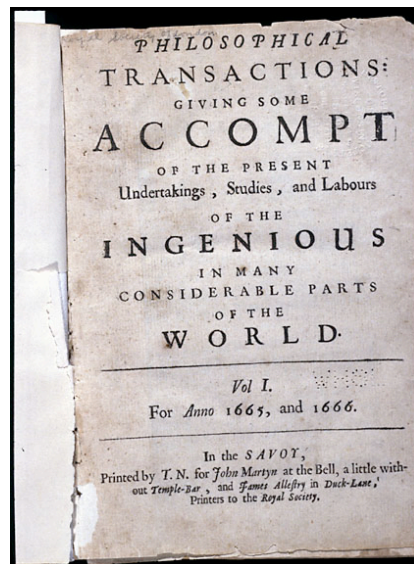
Scholarly discourse as CKS...



From:

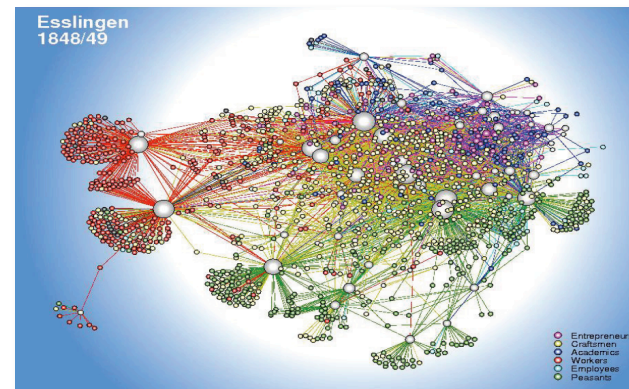


Le Journal des Sçavans
January 1665

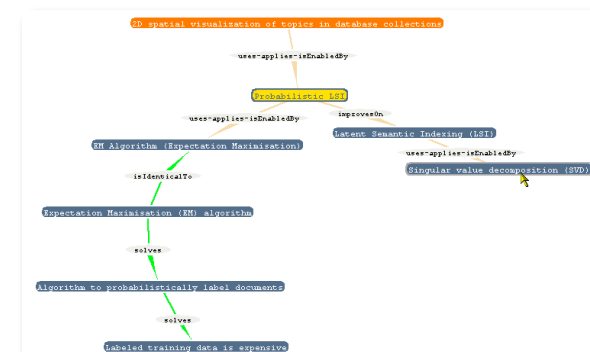


Philosophical Transactions of the Royal Society of London
March 1665

To:



Chaomei Chen, 2006: Citation analysis



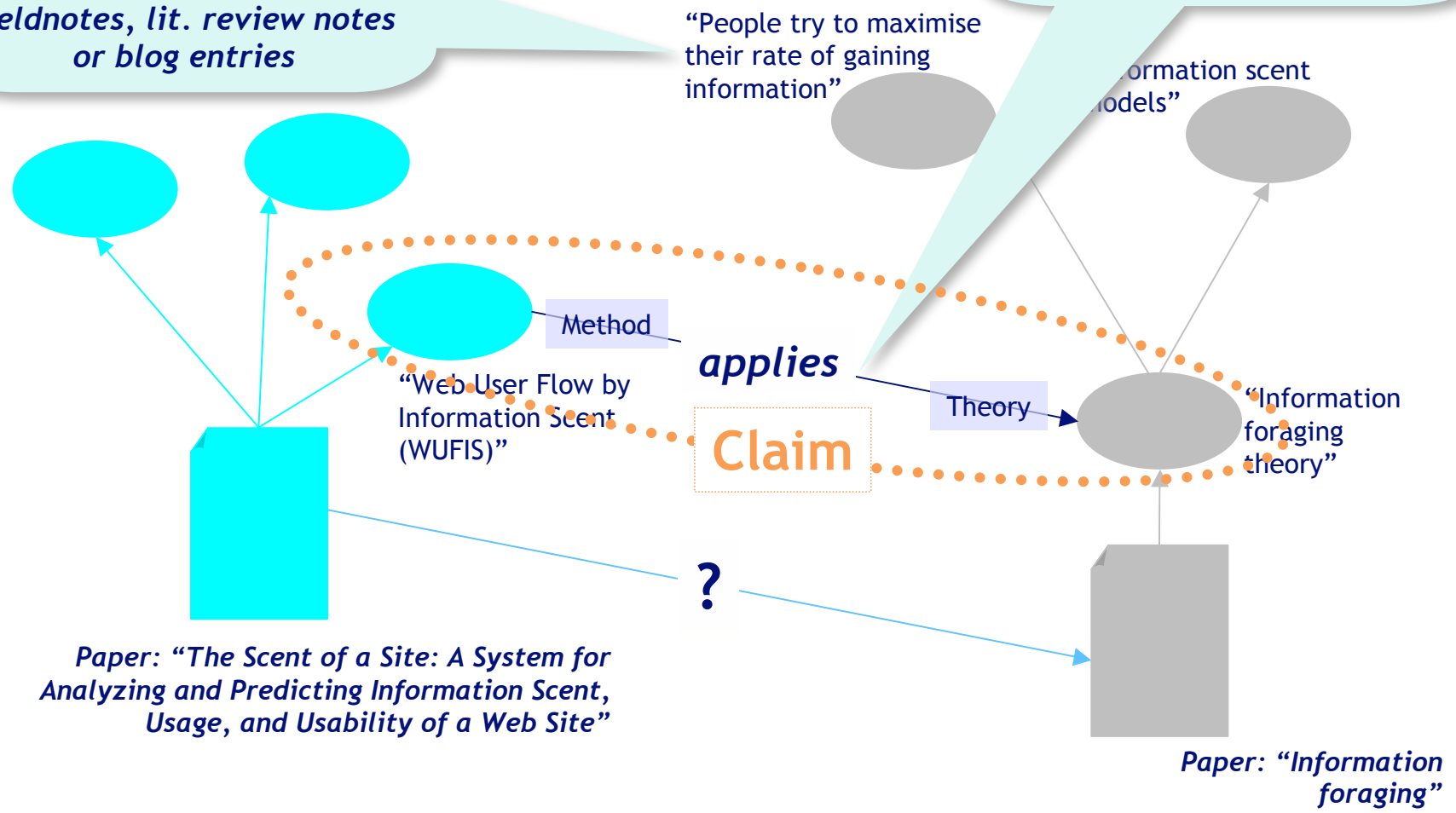
Buckingham Shum et al, 2003: lineage analysis

Scholarly discourse as CKS... Beyond document citations...



These annotations are freeform summaries of an idea, as one would also find in researchers' journals, fieldnotes, lit. review notes or blog entries

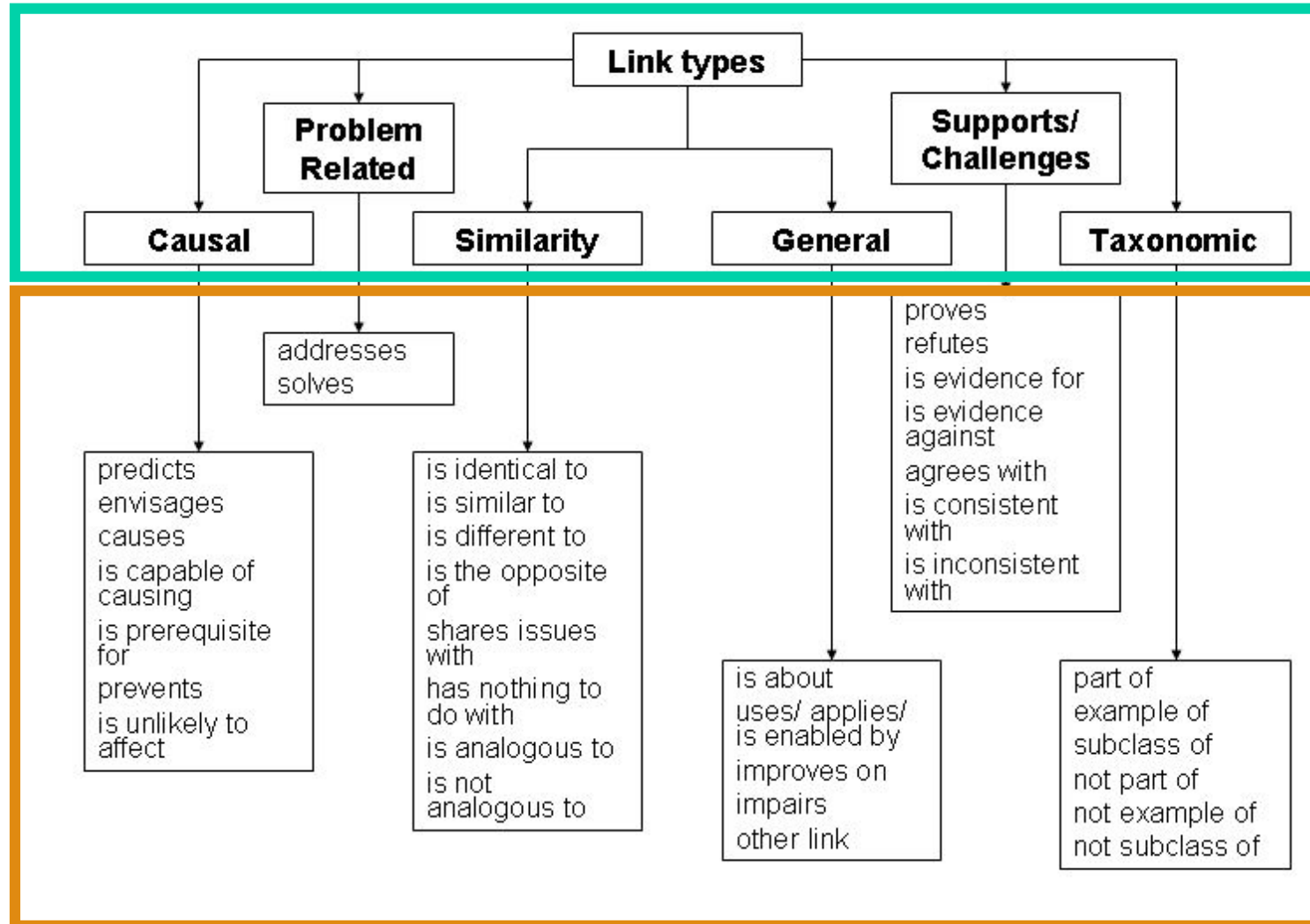
Making formal connections between ideas creates a semantic citation network → novel literature navigation, querying and visualization

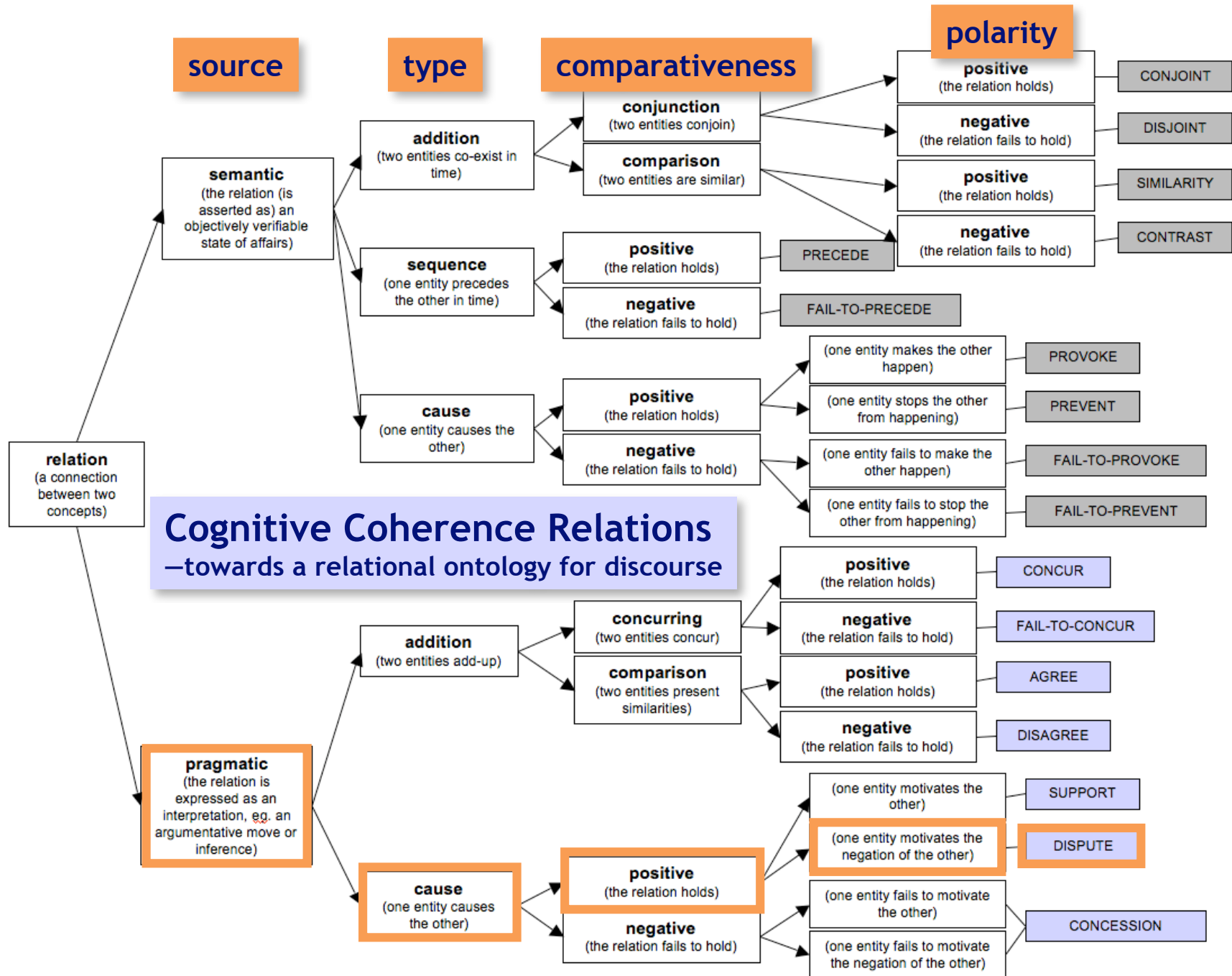




Scholarly discourse as CKS...

Connecting freeform tags with **naturalistic connections** (“dialects”) grounded in a **formal set of relations** (from semiotics and coherence relations)





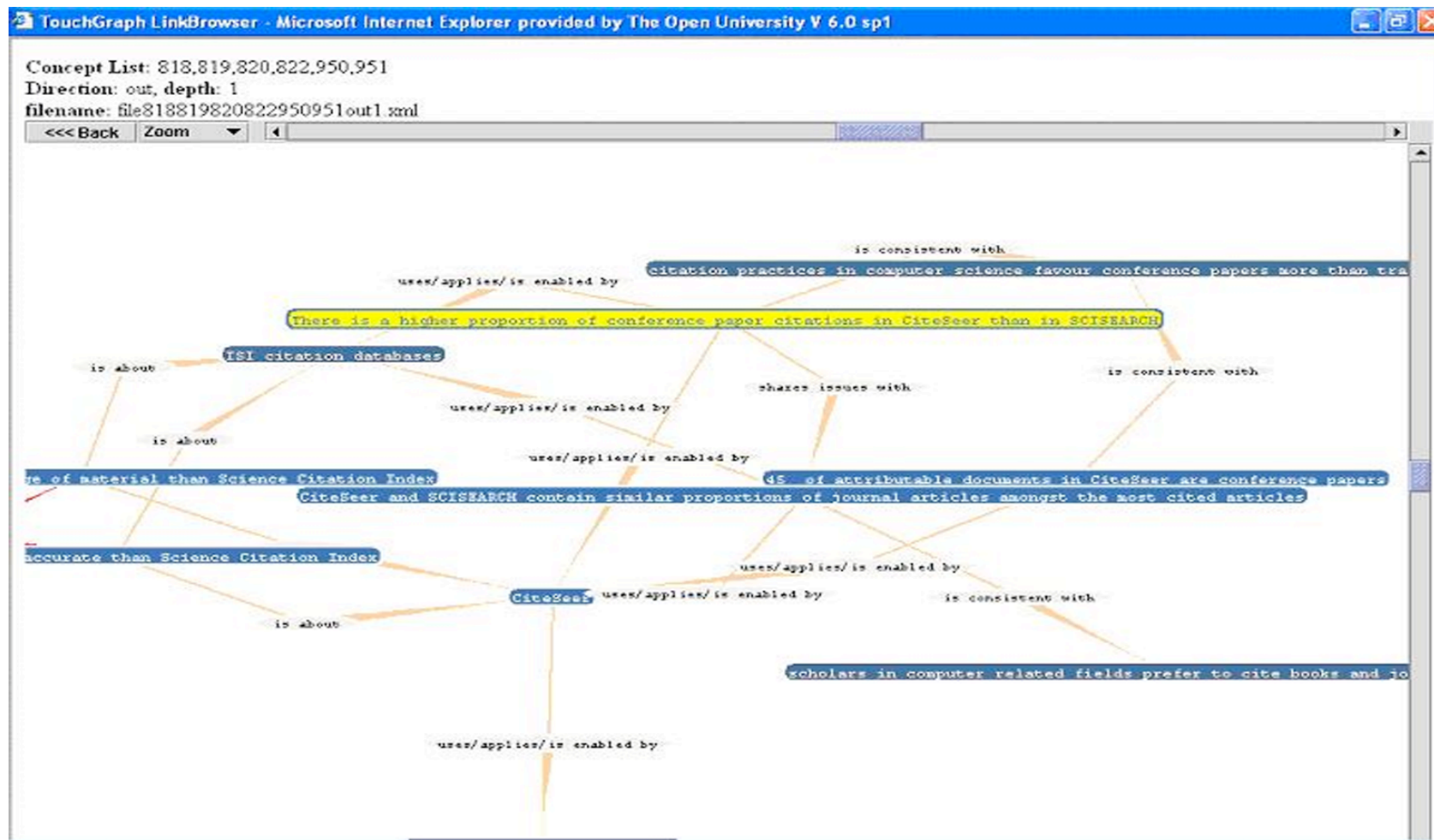


Scholarly discourse as CKS...

Visualising claims and arguments

The link-tracking service allows the user to specify structures

e.g. show tags one link out from any tag on the left hand of a claim containing the string "CiteSeer"





Scholarly discourse as CKS...

Querying on argumentation structures

find **discover** **advanced** **claiMaker**

machine learning

Perspective In contrast agree

Neural network text categorizer Depth

machine learning Depth

[About](#) - [ClaiMaker](#) - [Problems](#) - [Help](#)









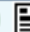


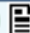




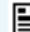







Scholarly discourse as CKS...

“What papers contrast with this paper?”

1. Extract concepts for this document
2. Trace concepts on which they build
3. Trace concepts challenging this set
4. Show root documents

The key issues you are concerned with:	
445	Decision Forest classifier   
446	Decision Forest classifier improves on C4.5 and kNN   

The related issues you may be concerned with:	
446	Decision Forest classifier improves on C4.5 and kNN   
515	Instance based learning   
511	Decision tree learning   
277	decision trees and naive Bayes perform well for text categorization   

The following claims disagree ...				
1	Optimised rules outperform Naive Bayes and decision trees   	 disagrees with 	decision trees and naive Bayes perform well for text categorization   	 3621  2



The point is... we think these kinds of structures are worth having

But can users create them?



How to help scholars engage in CKS?

Pilot study: paper-based literature modelling



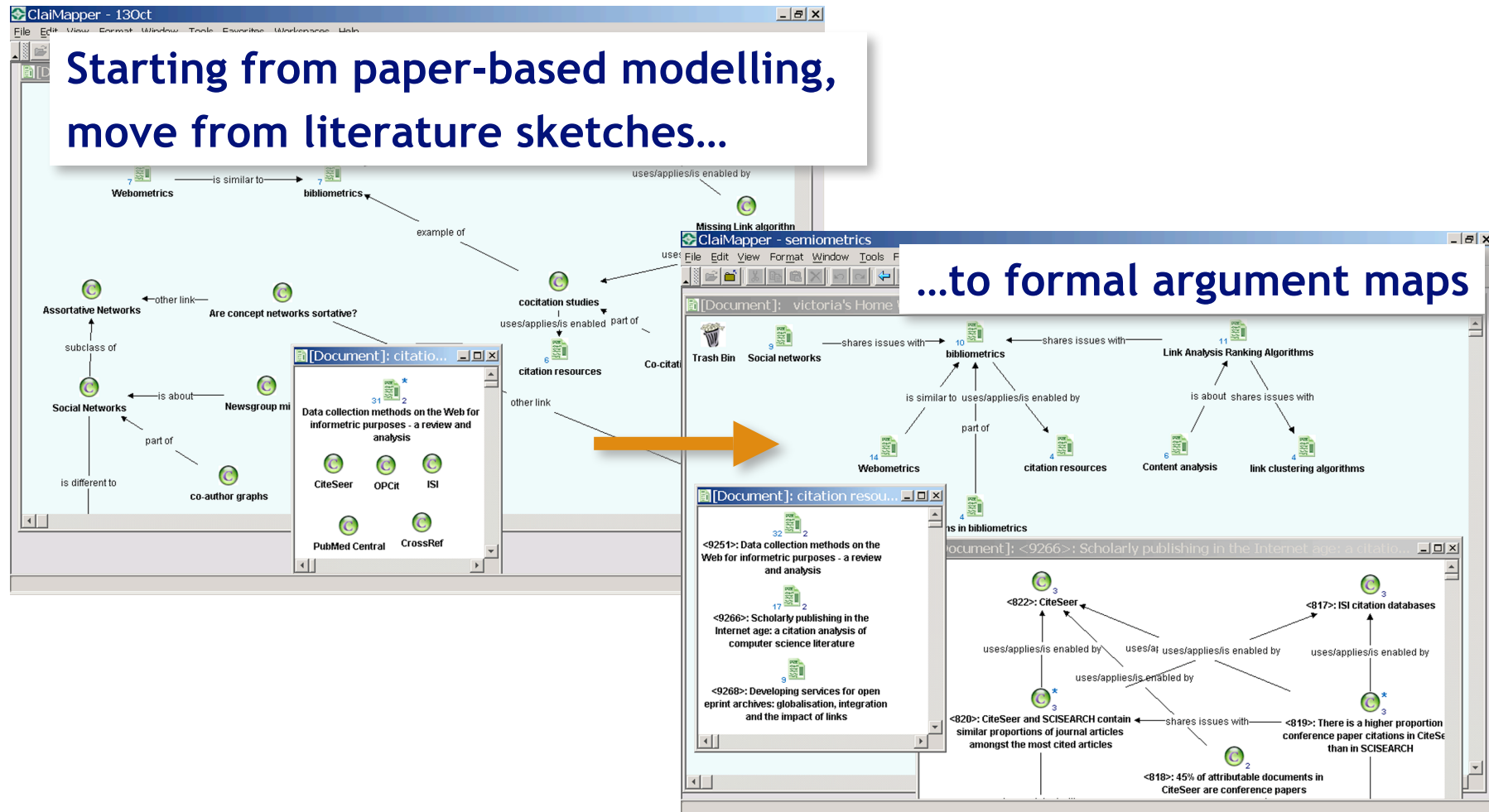
S. Buckingham Shum, V. Uren, G. Li, B. Sereno, and C. Mancini. **Computational Modelling of Naturalistic Argumentation in Research Literatures: Representation and Interaction Design Issues.** *International Journal of Intelligent Systems*, 22(1):17-47, 2006

How to help scholars engage in CKS?

From paper prototype to semiformal mapping tool



- The ClaiMapper tool





How to help scholars engage in CKS?

Pilot study: paper-based annotation

Extracting and Visualizing Semantic Structures in Retrieval Results for Browsing

Katy Börner

Indiana University, School of Library and Information Science
10th Street & Jordan Avenue, Main Library 019, Bloomington, IN. 47405 USA
E-mail: katy@indiana.edu

ABSTRACT

The paper introduces an approach that organizes retrieval results semantically and displays them spatially for browsing. Latent Semantic Analysis as well as clustering techniques are applied for semantic data analysis. A modified Boltzman algorithm is used to layout documents in a two-dimensional space for interactive exploration. The approach was implemented to visualize retrieval results from two different databases: the Science Citation Index Expanded and the Dido Image Bank.

KEYWORDS: Digital Libraries, Browsing, LSA, Conceptual Clustering, Boltzman Algorithm, Information Visualization

INTRODUCTION

The wealth of digitally stored data available today increases the demand to provide effective tools to retrieve and manage relevant data. Keyword searches over digital libraries, repositories, or the Web easily retrieve lists of several hundreds of documents. Information visualization - the process of analyzing and transforming data into an effective visual form - is believed to improve our interaction with large volumes of data. First visual interfaces to digital libraries provided full-text searching and full-content retrieval capabilities and visualized documents according to authors, time, place, or citation relationships.

A considerable body of recent research applies powerful mathematical techniques such as *Factor Analysis*, *Multidimensional Scaling*, or *Latent Semantic Analysis* to extract for example the underlying semantic structure of documents, the (evolving) specialty structure of a discipline, author co-citation patterns, changes in authors' influences in a particular field. In order to display the results of the data analysis spatially, computationally expensive techniques have to be applied to transform data analysis results to 2 or 3-dimensional coordinates. The computational expense of data analysis and visualization generation is very high. Therefore, precompiled, mostly static visualizations of fixed data sets are only displayed.

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, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
Digital Libraries, San Antonio, TX.
Copyright 2000 ACM 1-58113-231-X/00/0006...\$5.00

To our knowledge there exists no system that interactively visualizes retrieval results for browsing based on their underlying semantic structure.

DATA ANALYSIS

Latent Semantic Analysis (LSA) [4] has demonstrated improved performance over the traditional vector space techniques. It overcomes the problems of synonymy (variability in human word choice) and polysemy (same word has often different meanings) by automatically organizing documents into a semantic structure more appropriate for information retrieval. We apply LSA to extract the semantic structure of a particular database in a computationally expensive batch job.

At retrieval time, the result of a database query is hierarchically organized, based on the LSA output. Nearest neighbor-based agglomerative, hierarchical, unsupervised conceptual clustering is applied to create a hierarchy of clusters grouping documents of similar semantic structure. Clustering starts with a set of singleton clusters, each containing a single document. The two clusters most similar are merged to form a new cluster that covers both. This process is repeated for each of the remaining clusters. At termination, a uniform, binary hierarchy of document clusters is produced. The partition showing the highest within-cluster similarity and lowest between-cluster similarity is selected for data visualization.

DATA VISUALIZATION

Rather than being a static visualization of data, the interface is self-organizing and highly interactive. Data is displayed in an initially random configuration, which sorts itself out into a more-or-less acceptable display via a modified Boltzman algorithm [3]. The algorithm works by computing attraction and repulsion forces among nodes based on the result of the data analysis. Nodes may represent articles or images which are attracted to other nodes to which they have a (reference or similarity) link and repelled by nodes to which there is no link. If the algorithm does not produce a visually acceptable layout, or if the user wishes to view the results differently, nodes can be grabbed and moved.

PROTOTYPE SYSTEMS

Two systems have been implemented in Java using the data organization and visualization methods described above.

SCI-E: The first system visualizes query results from the Science Citation Index Expanded (TMI) as published by the Institute for Scientific Information®. The Citation Index

provides access to current bibliographic information and cited references in more than 5,600 journals. Querying it via the Web of Science® Interface at <http://webofscience.com/results> in an often huge number of matching documents organized in lists of ten that can be marked, saved, and downloaded for detailed study. To demonstrate a visual browser to this kind of data base we will use DAIV188, a query result data set from SCI-EXPANDED that contains 188 articles matching the topic 'data AND analysis AND information AND visualization'. The articles are represented in the usual Web of Science data output format (including author(s), article title and source, cited references, addresses, abstract, language, publisher information, ISSN, document type, keywords, times cited, etc.).

LSA was applied over keywords and abstracts of articles. As a result of conceptual clustering, the 167th partition was selected for visualization. It contains 20 clusters grouping 1 - 53 articles. Figure 1 shows the Java interface. Each book article is represented by a rectangle and each journal article by an oval. Articles are labeled by their first author. Lines between nodes visually represent co-citation links.

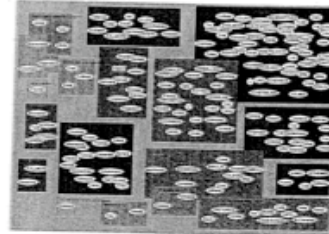


Figure 1: Java Interface to DAIV188

The 2-dimensional layout of articles corresponds to the data mining result as well as to the forces applied by the Boltzman algorithm to generate an acceptable layout. The higher the similarity of articles within a cluster the lighter its color. Each cluster is labeled by the keyword used most often.

DIDO: Another instantiation of the system enables users to browse search results from the Dido Image Bank, <http://www.dlib.indiana.edu/collections/dido/> provided by the Department of the History of Art, Indiana University. Dido stores about 9,500 digitized images from the Fine Arts Slide Library collection of over 320,000 images. Each image in Dido is stored together with its thumbnail representation as well as a textual description. LSA was applied over the textual descriptions exclusively. For demonstration purposes the set of images matching the keyword descriptor 'MONET' were retrieved and displayed for browsing. It contains 21 documents inclusive two portraits of Claude Monet drawn by Edouard Manet (see Figure 2).



Figure 2: The MONET Cluster

Thumbnail representations of images have been fetched from the Dido Database showing some of Monet's favorite themes such as haystacks, cathedrals, and water lilies.

CONCLUSIONS

Initial tests show that the presented approach provides easy access to textual materials, such as articles, as well as to documents for which textual descriptions are available, such as images. Detailed user studies are in preparation. First results on using an immersive 3-dimensional CAVE environment for the interactive exploration of search results are presented in [3]. An extended version of this paper as well as colored, full-size versions of Figures 1 and 2 are accessible at <http://ella.slis.indiana.edu/~katy/DL00>.

ACKNOWLEDGMENTS

Robert Goldstone, Mark Steyvers, Helen Atkins, and Eileen Fry have been valuable discussion partners. The SVDPACK [2] by M. Berry was used for computing the singular value decomposition. The research is supported by an High Performance Network Applications grant of IU. Collaborators are Andrew Dillon and Margaret Dolinsky.

REFERENCES

- Alexander, Garcia, and Alder. Simulation of the Consistent Boltzman Equation for Hard Spheres and Its Extension to Dense Gases, *Lecture Notes in Physics*, Springer Verlag, 1995.
- Berry, M. et al. SVDPACK (Version 1.0) User's Guide, University of Tennessee Tech. Report CS-93-194, 1993 (Revised October 1996).
- Börner, K. Visible Threads: A smart VR interface to digital libraries. *Electronic Imaging 2000, Visual Data Exploration and Analysis*.
- Landauer, T. K., Foltz, P. W., & Laham, D. Introduction to Latent Semantic Analysis. *Discourse Processes*, 25, 259-284, 1998.

How to help scholars engage in CKS?



- The **ClaimSpotter** annotation tool: Web 2.0-style tagging with optional community/system tag recommendations

The screenshot displays the ClaimSpotter web application interface. At the top, there is a navigation bar with links: Login, History, Add a document, Standard, Alternate, .dot Export, Help, About. Below this is a control bar with filters: More Ideas, Concepts: My (1), Relations: (3), Argument Zones: (checkbox), Importance: None, Term(s): (input), find, clear, Reset. The main content area shows a document titled "Trusting Information Sources One Citizen at a Time" by Yolanda Gil, Varun Ratnakar. The document text is annotated with colored boxes and numbers 1 through 8. A dropdown menu is open over the text "This paper describes an approach...", showing a list of relationship types such as "supports", "proves", "refutes", "is evidence for", "is evidence against", "agrees with", "disagrees with", "is consistent with", "is inconsistent with", "taxonomic", "part of", "example of", "subclass of", "similarity", "is identical to", "is similar to", "is different to", "is the opposite of", "shares issues with", "is about", and "clear left | remove | clear right" (7). On the right side, there is a "Copy in..." section with a list of concepts and a "6" next to it. At the bottom, there are "Submit" and "Reset" buttons. The interface also includes a sidebar with a "Document" section (4) and a "TABLE OF CONTENTS" with links to Abstract, Introduction, Information Analysis in TRELIS, Source Attribution and Description, Deriving an Assessment about a Source, Helping Users, Select Sources, Related Work, Conclusions, and References.



A user-centred approach to the formative evaluation of a CKS tool

- **Research question:**
 - how do annotators approach the task of using a new Web tool to semantically annotate a document with its key contributions?
- **Focus**
 - ..is on how the tool functionality and UI affordances shape tagging behaviour
- **Quantitative and qualitative analysis**

Example *claims* (tag triples) from users



- Domain ontology *is about* A hierarchy of URIs on multiple levels
- Universal physical access *is unlikely to affect* Digital divide
- Hypertext node juxtaposition *is analogous to* Cinematic shot juxtaposition
- [Evidence] In the Bristol trial, the awareness of the presence of other players was correlated with how much our participants enjoyed the game as well as with how engaged they felt *is consistent with* Presence awareness of many other people is capable of causing, feel good factor
- Magpie moves away from hypermedia towards open service-based architectures *is evidence for* [Magpie *improves on* COHSE]

User study: selected results



- See paper for details
- and the PhD for complete account

B. Sereno. A Document-Centric Semantic Annotation Environment to Support Sense-Making. PhD Thesis, Technical Report KMI-06-13, Knowledge Media Institute, The Open University, UK, September 2005.

[<http://kmi.open.ac.uk/publications/pdf/KMI-TR-06-13.pdf>]



Tag length similar for novices and experts (64% 1-3 words)

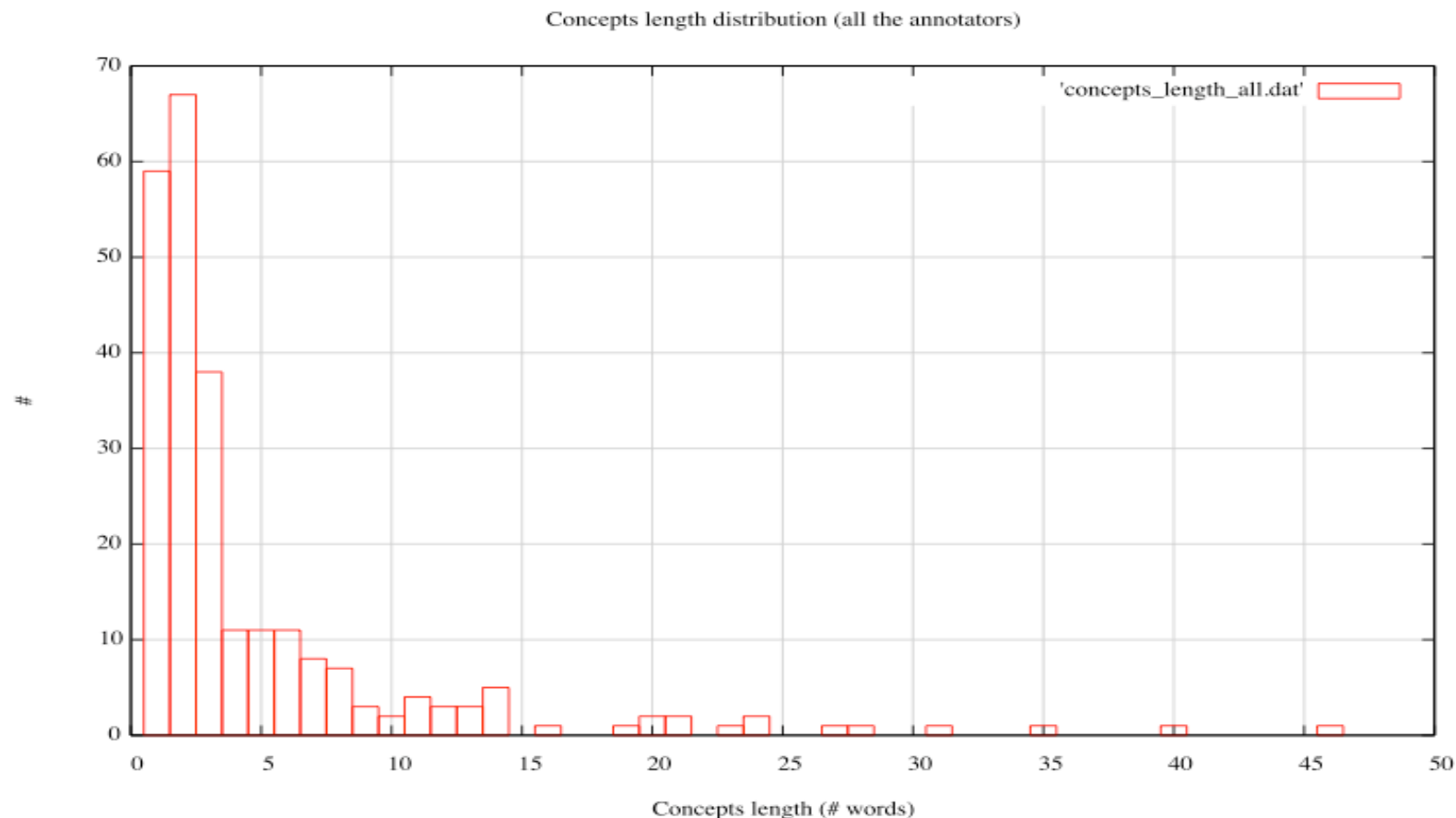


Figure 6.2: Concepts length distribution graph for all the annotators. 164 concepts out of 257 are composed of three words or less.



Tag reuse

most of them used only twice in this study (1 hour)

Reused concepts
<i>a community-based project that wired four computing centres (hubs) in a lower socio-economic urban area a research project aiming to explore the potential of spontaneous social behaviour and playful group interaction in public spaces A set of recommendations to make the process as painless as possible a tool that assists users with interpreting the web resources a wireless location based multiplayer game access Accessing information efficiently ACE ACM Digital Library Adding formalised knowledge to a document Adding information to help sense-making analysis AquaLog awareness CitiTag ClaiMaker ClaimSpotter cognitive overhead Cognitive overload in ClaimSpotter cognitive strategies coherence COHSE Collective sensemaking Data-Flow data-flow model Digital divide Discourse ontology Document annotation domain hierarchy don't-want-tos Eprint archives ePrint services espotter explores Formalization overhead GATE have nots Holding an internal model is troublesome How people approach documents impact of the social context Information environments organised via digital libraries Information-driven reading interpretation and information gathering Linking Magpie mobile technologies navigation of web resources non-users OpCit Point-driven reading Presence awareness reasons why some people choose not to compute. Recognising entities likes names and organisations in a document robust services required for large-scale information environments ScholOnto Semantic services Semantic Web Sensemaking social experiences and group play START stories Story-driven reading subscribes to survey The Compendium approach The Fujitsu hub wiring experiment The information in there does not exist in the document The Internet This paper universal physical access use of semantic information User studies VIPERS</i>

Table C.3: Concepts reused by the annotators.

Transcript analysis

to derive themes, sub-categories and codes

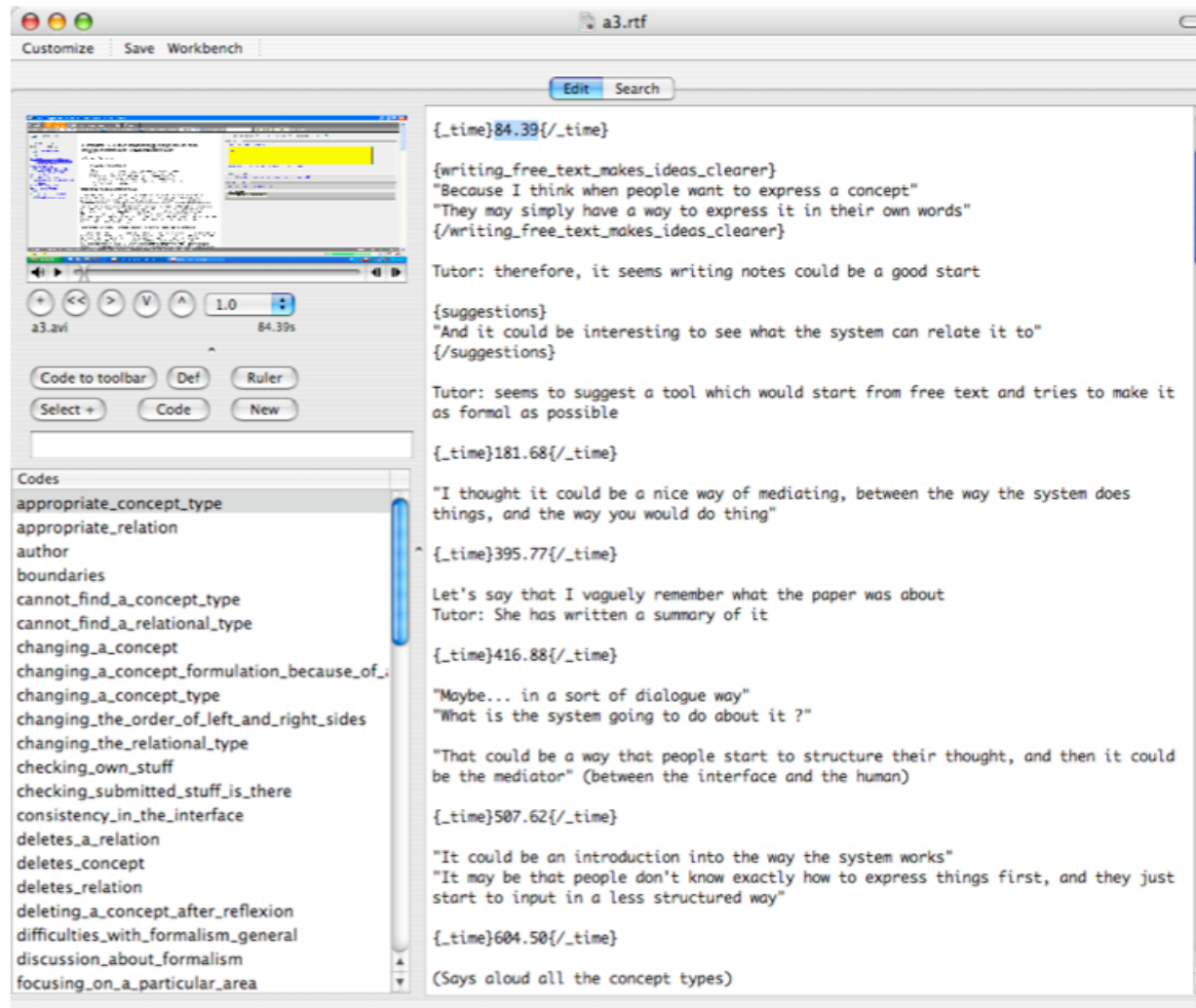


Figure 6.7: A coding session in TAMS: text chunks (main window, right side) are selected and assigned a code (selected from the bottom-left window, or created from scratch.)



Themes from the user study

- **Formalization**
 - the cognitive task of structured tagging
- **User Strategy**
 - how users approach the semantic annotation task
- **Interaction Design**
 - how behaviour is shaped by the tool's affordances



Formalization

the cognitive task of structured tagging

- Looking for the right tag type...

ClaimSpotter 0.3.9.3 | Annotate - Mozilla Firefox

http://137.108.25.237/claimspotter/0.3.9.3/index.php/section-H-1

Login History More Ideas Concepts: Patterns: Arg. Zones: Importance: >5 Term(s): search clear Reset Help About

Document

TABLE OF CONTENTS:

- Abstract
- INTRODUCTION
- THE CONCEPT
- DESIGN
- USER EXPERIENCE

STUDIES

- Spontaneous behaviours
- Presence awareness
- Experience and everyday life
- FURTHER WORK
- References

of communication technologies, as we see everyday on school playgrounds. Presence enabled technologies create new prospects for play, for adults as well. **In our research these are the boundaries we explore: what kind of engaging social experiences can emerge in the real world based on the awareness of individuals participating in a parallel virtual experience?** Does virtual presence penetrate physical presence in any way?

THE CONCEPT

The 'CitiTag' project is focused on social experiences and group play in public spaces, based on the awareness of other peoples' presence, through the use of mobile technology.

Our CitiTag game has been inspired by the simplicity, spontaneity and instant fun of 'playground tag' [1]. We have further developed the 'tag' concept to encourage emergent social behaviours in an urban context. **City space is used as a playground and passers-by can become the usual or unusual suspects in a novel experience.**

DESIGN

CitiTag is a multiplayer, wireless location-based game, played using GPS (Global Positioning System) and handheld, IPaq PocketPCs connected to a wireless network. The game has been designed for many people, potentially as an everyday experience one could have in the future with a mobile phone while walking about in a city centre. As a player of CitiTag, you belong to either of two teams (Reds or Greens) and you roam the city, trying to find players from the opposite team to 'tag'.

You get the opportunity to 'tag' someone when you get close to them. You can also get 'tagged' if someone from the opposite team gets close to you. If this happens, you need to try and find a team member in the vicinity to set you free, to 'untag' you.

Each game event (e.g. someone is close and you can tag/untag them) appears as an alert on the IPaq screen with a sound (Figure 1). The player can then tap the screen with his or her thumb to respond to the event.

The project is motivated by the hypothesis that very simple game rules based on presence states (e.g. I am Green and 'tagged') can result in an enjoyable social experience, stimulated by real world interaction among players. Another hypothesis is that certain interactions develop once a

notes

add | remove last

concepts

import | add | remove last

type	name
definition	ness (n) of many other people remove [X] [L] [X]
analysis	Feel good' factor remove [X] [L] [X]
approach	experiences and group play remove [X] [L] [X]

analysis

approach

assumption

data

definition

evidence

hypothesis

language

methodology

model

opinion

phenomenon

problem

solution

theory

Submit Reset

User Yanna, annotating Urban space as a playground for large scale group interaction: experiences with CitiTag ClaimSpotter 0.3.9.3

00:00:42 3:58 PM



Formalization

the cognitive task of structured tagging

- Looking for the right link type...



User Strategy

how users approach the semantic annotation task

- What granularity and type of claims? When to stop?

The screenshot displays the ClaimSpotter 0.3.9.3 web application. The browser window shows the URL <http://137.108.25.237/claimspotter/0.3.9.3/index.php>. The page title is "Magpie: Supporting Browsing and Navigation on the Semantic Web". The main content area displays the abstract and introduction of the paper. The right sidebar contains a list of semantic relations with dropdown menus for selecting concepts and terms. The bottom of the screen shows a Windows taskbar with the Start button, system tray, and a media player control bar.

Document

TABLE OF CONTENTS:

- [Abstract](#)
- [INTRODUCTION](#)
- [MAGPIE USAGE SCENARIO](#)
- [MAGPIE ARCHITECTURE](#)
- [Magpie browser extension IE plug-in](#)
- [SEMANTIC SERVICES IN MAGPIE](#)
- [On-demand semantic services](#)
- [Trigger semantic services](#)
- [Semantic "bookmarking"](#)
- [OVERVIEW OF SIMILAR WORK](#)
- [CONCLUSIONS](#)
- [References](#)

Magpie: Supporting Browsing and Navigation on the Semantic Web

John Domingue, Martin Diczor.

ABSTRACT

We describe several advanced functionalities of Magpie a tool that assists users with interpreting the web resources. Magpie is an extension to the Internet Explorer that automatically creates a semantic layer for web pages using a user-selected ontology. Semantic layers are annotations of a web page, with a set of applicable semantic services attached to the annotated items. We argue that the ability to generate different semantic layers for a web resource is vital to support the interpretation of web pages. Moreover, the assignment of semantic web services to the entities allows users to browse their neighbourhood semantically. At the same time, the Magpie suite offers trigger functionality based on the patterns of an automatically updated semantic log. The benefits of such an approach are illustrated by a semantically enriched browsing history management.

INTRODUCTION

A lot of research has gone into supporting the task of finding web resources by means of "standard" information retrieval mechanisms or by means of semantically enhanced search [6], [13]. Less attention has been paid to the task of supporting the interpretation of web pages. Annotation technologies [6], [14] allow users to associate meta-data with web resources, which can then be used to facilitate their interpretation. The annotation technologies provide a useful way to support shared interpretation, but they are very limited, mainly because the annotation is carried out manually. Hence, the quality of meta-data depends on the authors or librarians annotating the web page.

The majority of web pages are not semantically annotated. This is a great obstacle in a move towards the Semantic Web [1]. Magpie is a tool supporting the interpretation of web pages and acting as a complementary knowledge source, which a user can call upon to gain instantaneous access to the background knowledge relevant to a web resource. Magpie follows a different approach from that used by most other annotation techniques: it automatically associates a semantic layer to a web resource, rather than relying on a manual annotation.

semantic web browser Concept

Magpie Concept

addresses flip

interpretation of web resources Concept

Magpie Concept

is about flip

interpretation and information paths Concept

160 Link

improves on flip

161 Link

use of semantic information Concept

addresses flip

navigation of web resources Concept

Magpie Concept

example of flip

use of semantic information Concept

inability to use existing semantic on Concept

example of flip

problem with magpie Concept

tomatically generate semantic layer Concept

example of flip

feature of magpie Concept

submit

Submit Reset

User: Enrico, annotating: Magpie: Supporting Browsing and Navigation on the Semantic Web ClaimSpotter 0.3.9.3



User Strategy

how users approach the semantic annotation task

- No initial use of tagging aids – focus is on own tags

ClaimSpotter 0.3.9.3 Annotate - Mozilla Firefox

http://137.109.25.237/dainspotter/0.3.9.3/index.php

Login History More Ideas Concepts: Relations: Arg. Zones: Importance: None Term(s): search clear Reset Help About

Document

TABLE OF CONTENTS:

- Abstract
- Two men were driving in a car. Tellers, listeners, and points
- Point-driven listening
- Point-driven reading
- Information-driven reading
- Story-driven reading
- Strategies in point-driven reading: an example
- Empirical questions
- Experimental poetics
- References

Collectively, the above pressures can be considered a 'pragmatic frame' which places limits on the nature of the point which may be constructed by a listener, and in some sense determines whether or not a point will be successfully constructed at all.

It is worth stressing here that we refer to the listener's 'construction' of point in order to underline our view that points aren't 'in' stories, waiting to be identified by perceptive listeners, but instead are constructed by listeners on the basis of various sources of information, only one of which is the text.

It is also important to note that the listening could fail for any number of reasons. For instance, the narrator may obviously intend to make a point, but the text may fail to support the point he or she intends to make. In that case narrator and listeners might still rescue the situation by 'negotiating' the point of the story, that is, by debating what the story may properly be said to be about (Polanyi 1979). A second cause of failure may be that the listener's and the narrator's sense of cultural or generic appropriateness differ too widely. For example, a 'story' may be told which according to its narrator is pointed but which the audience finds to be pointless. The audience then responds overtly or otherwise, with the question Labov (1972: 366) said "every good narrator is continually warding off": "So what?" Narrators who tell pointless or irrelevant stories are seen as boring and inept and suffer a loss of face. And although listeners who miss points are less severely penalized in social terms, there are parallel consequences, particularly if there are other members of the audience who do 'get the point.'

POINT-DRIVEN READING

Listening to an oral conversational story, then, is normally 'point-driven' recipients listen to the story in anticipation that they will be able to reconstruct it as a 'pragmatic gesture,' a 'global speech act.' Similarly, there is a type of reading which may be called 'point-driven,' because in it, too, the understander reads with the expectation that the text will enable the construction of a valid, pragmatic point. Point-driven reading is both similar to and different from point-driven listening. It is similar because in both types the construction of point is a function of the text, the comprehender's cultural, and the comprehender's generic expectations. Point-driven reading is different from listening though, for several reasons. First, oral stories are generally immediately connectable to previous discourse, whereas written stories are relatively autonomous: the reader has less

notes

add | remove last

concepts:

import | add | remove last

type	name	
n/a	point-seeking strategies	remove [X] [L] [R]
n/a	narrative surface	remove [X] [L] [R]
n/a	negotiation	remove [X] [L] [R]
n/a	coherence	remove [X] [L] [R]

relations

import | add | remove last

left concept (type), relation, right concept(type)

n/a	narrative surface	Concept
is evidence for		
n/a	point-seeking strategies	Concept
flip remove		
n/a	negotiation	Concept
is evidence for		
n/a	point-seeking strategies	Concept
flip remove		
n/a	coherence	Concept
is evidence for		
n/a	point-seeking strategies	Concept
flip remove		
n/a	negotiation	Concept
is evidence against		
n/a	information seeking strategies	Concept
flip remove		

User: Joanna, annotating: POINT-DRIVEN UNDERSTANDING: PRAGMATIC AND COGNITIVE DIMENSIONS OF LITERARY READING ClaimSpotter 0.3.9.3

00:00:00

9:57 AM



Interaction Design

how behaviour is shaped by the tool's affordances

- 'Flip' left/right tags to match the link type

The screenshot displays the ClaimSpotter 0.3.9.3 web application interface. The browser window shows the URL <http://137.108.25.237/claimspotter/0.3.9.3/index.php>. The interface includes a navigation menu on the left with sections like 'Document', 'TABLE OF CONTENTS', and 'Abstract'. The main content area displays a document titled 'ANALYSIS AND CONCLUSION' with highlighted terms like 'access', 'shows', 'survey', 'awareness', and 'jobs'. On the right, there are panels for 'notes', 'concepts', and 'claims'. The 'concepts' panel shows a table with columns for 'type', 'concept', and 'reuse'. The 'claims' panel shows a form for creating claims with dropdown menus for 'left concept', 'relation', and 'right concept'. The bottom of the screen shows a Windows taskbar with the Start button, several open applications, and a system tray with the time 10:28 AM.



Interaction Design

how behaviour is shaped by the tool's affordances

- Skimming highlighted text

ClaimSpotter 0.3.9.3 | Annotate - Mozilla Firefox

Http://137.108.25.237/claimspotter/0.3.9.3/index.php

Login | **History** | **More Ideas** | Concepts: | Patterns: | Arg. Zones: | Importance: None | Term(s): | search | clear | Reset | Help | About

Document

power of social cohesiveness that can be brought about by knowledge of the presence and location of others in both real and virtual spaces. We also know that wirelessly internetworked groups of humans can exhibit emergent prediction capabilities [2] and thus demonstrate self-organizing dynamics.

Our work (th) is motivated by the idea that the **presence awareness** (th) of many other people can enhance the 'feel good' factor of being part of a large group and thereby afford spontaneous interactions. Previous online studies [3] have shown that spontaneous social behaviours can 'emerge' among groups present in multi-user environments, even without explicit and verbal communication. The recent Flash Mobs phenomenon illustrated that people do not hesitate to perform certain acts in public together with many others, which otherwise would have been quite embarrassing. In fact, people participating in those events seemed very engaged and amused.

These acts of spontaneous play have been thought-provoking within the context of our research. Play has been inherently social, before the advent of communication technologies, as we see everyday on school playgrounds. Presence enabled technologies create new prospects for play, for adults as well. In our research these are the boundaries we explore: what kind of engaging social experiences can emerge in the real world based on the **awareness** (th) of individuals **participating in a parallel virtual experience** (th)? Does virtual presence penetrate physical presence in any way?

THE CONCEPT

The **CitiTag** (th) project is focused on **social experiences and group play** (th) in public spaces, based on the **awareness** (th) of other peoples' presence, through the use of **mobile technology** (th).

Our **CitiTag** (th) game has been inspired by the simplicity, spontaneity and instant fun of **playground tag** (th) [1]. We have further developed the 'tag' concept to **encourage emergent social behaviours in an urban context**. City space is used as a playground and passers-by can become the usual or unusual suspects in a novel experience.

DESIGN

CitiTag (th) is a multiplayer, wireless location-based game, played

notes
add | remove last

concepts
import | add | remove last

type	name	
definition	presence awareness (th) of mo	remove [x] [y] [z]
n/a	feel good factor	remove [x] [y] [z]
n/a	social experiences and group	remove [x] [y] [z]
n/a	mobile technology	remove [x] [y] [z]
approach	playground tag	remove [x] [y] [z]
hypothes	spontaneous social behaviour	remove [x] [y] [z]
approach	participating in a parallel virtus	remove [x] [y] [z]
n/a	CitiTag	remove [x] [y] [z]

relations
import | add | remove last

submit
Submit | Reset

User: Yanna, annotating: Urban space as a playground for large scale group interaction: experiences with CitiTag

ClaimSpotter 0.3.9.3

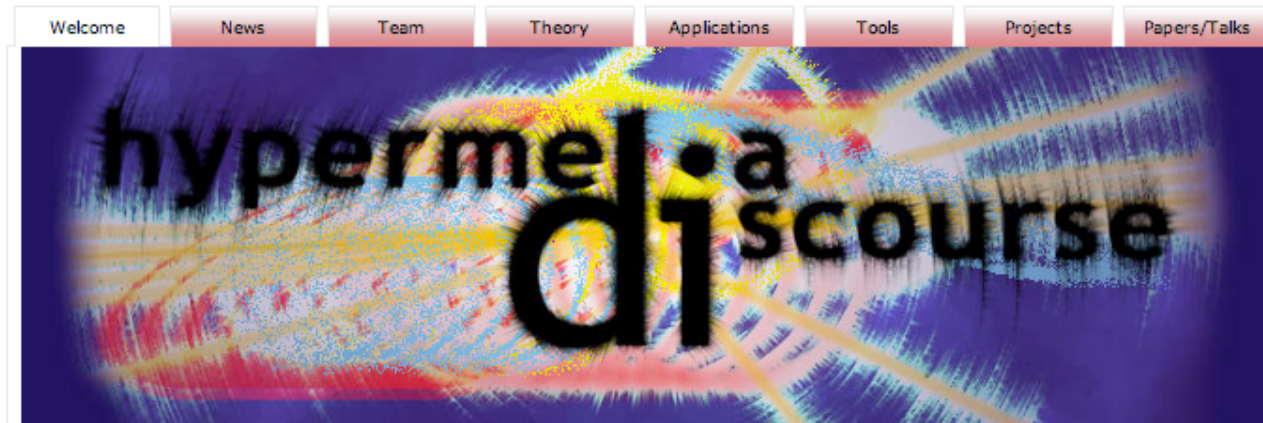
00:00:18

4:09 PM

Lessons Learnt & Design Principles



- **Untrained users can do it:** in their first hour they created coherent claims. UI design validated to this degree.
—future work: longitudinal evaluation at scale
- New users **attend to what is highlighted** for them (matching tags; primary doct.), and generally don't click down a level
—next version combines visualizations and document-centric features
- Support **incremental formalization**
—cf. use of *is-about* as a placeholder link; provide an *Other...* category and try to map automatically to the ontology
- **Users' strategies vary** — don't assume a strong workflow
a paper-based pilot study can provide insights into this
- **Web 2.0 UI simplicity:** good design needed to provide high functionality, walk-up-and-use CKS tools
—we overwhelmed some users with overlaid suggestions for tags



ClaimSpotter:

papers and demos

<http://kmi.open.ac.uk/projects/hyperdiscourse/tools/claimspotter>

Hypermedia Discourse project:

theories / tools / case studies / user studies: face-face and asynch. interaction

<http://kmi.open.ac.uk/projects/hyperdiscourse>



collaboration / semantics / usability / community informatics / argumentation

<http://www.PragmaticWeb.info>

Short/full paper submission deadline: 14 May