# The architecture of commKnowledge: combining link structure and user actions to support an online community

## Michael Gordon*

University of Michigan Business School, 701 Tappan Street,
Ann Arbor, Michigan 48109-1234, USA
Fax: +1-734-936-0279      E-mail: mdgordon@umich.edu
*Corresponding Author

## Weiguo Fan

Virginia Tech University, Pamplin College of Business,
3007 Pamplin Hall, Blacksburg, VA 24061, USA
Fax: +1-540-231-2511      E-mail: wfan@vt.edu

## Sheizaf Rafaeli

University of Haifa, Graduate School of Business Administration,
Room 7048 Rabin Building, Mt. Carmel, Haifa, 31905, Israel
Fax: +972-4-8249194      E-mail: sheizaf@earthlink.net

## Harris Wu and N. Farag

University of Michigan Business School, 701 Tappan Street,
Room D0263, Ann Arbor, Michigan 48109-1234, USA
Fax: +1-734-647-8133
E-mail: harriswu@umich.edu          E-mail: nfarag@umich.edu

**Abstract:** The commKnowledge project is exploring how an electronic community interested in electronic business can effectively share information. We discuss the design decisions in developing an architecture that helps automatically vet contributed items and allows access by recency, topicality, and perceived usefulness. We discuss how use graphs combine the hypertextual structure of the web with user judgments and actions to provide new means of assessing the usefulness of contributed information.

**Keywords:** Knowledge sharing; information retrieval; online community; web structure; social filtering.

**Biographical notes:** Michael Gordon is Associate Dean of Information Technology and the Arthur F. Thurnau Professor of Computer and Information Systems at the University of Michigan Business School. His research interests include: the retrieval and discovery-based uses of textual information; information-based communities; and the appropriate uses of technology to support teaching, leaning, and information sharing.

Weiguo Fan is an Assistant Professor of Information Systems and Computer Science at the Virginia Tech University. He received his PhD from the University of Michigan Business School. His research interests include: data mining and its application to textual/web data; intelligent information retrieval; knowledge management; system design and evaluation; and business intelligence.

Sheizaf Rafaeli is Director of the Center for the Study of the Information Society and a Professor at the Graduate School of Business Administration, University of Haifa, Israel. He is interested in computer-mediated communication, interaction and information sharing. He studies and builds internet-based activities such as online higher-education; publishing; political, governmental, social and commercial virtual organisations.

Harris Wu is a doctoral student in Computer and Information Systems at the University of Michigan Business School. His research interests include information retrieval and information economics.

Neveen Farag is a PhD candidate at the University of Michigan Business School in Computer Information Systems. Her research focuses on online information sharing, and specifically firm-side initiatives, which enhance online information sharing and overall online customer experience.

# 1   Introduction

No matter what your interest, there is probably information about it on the World Wide Web. The trick, of course, is to find it.

The web contains several billion pages [1] – and is growing continually. So, general purpose search engines are a necessity; but they are not the only way to negotiate the web. Specialised search engines are devoted to topics from entertainment, investment and health to parachuting, Buddhism and anthropology [2] and so ignore the vast majority of the web. Directories like Yahoo are selective, too, using paid editorial boards to sift through submitted web pages on hundreds of topics and choose and organise those that they like. Online communities are yet another means of providing access to a small corner of the web, where a group with similar interests attempts to learn from each other.

CommKnowledge is, at the same time, an online community, a website, a teaching tool and a research project. As a website, it contains a store of information about electronic *comm*erce. As a teaching tool, we seek to make this information *comm*on knowledge – to students of electronic commerce and business taking courses in school or investigating the topic on their own. A *comm*unity of interested individuals is responsible for the contents of the commKnowledge website, as well as its organisation. The commKnowledge architecture allows us as researchers to explore various ways to solicit contributions, determine which are best, and present them to the community of

commKnowledge users. Though a closed community can better ensure contributions consistent with its charter, we wanted to create an unbounded, open community accepting input from unknown sources. Though a human editorial board can help ensure pertinence and quality of contributions, we wanted to develop and test a community resource operated with minimal human intervention. Though various techniques exist for evaluating the relevance of content, we sought to devise a new means which, at the same time, pays attention to the connections among documents, which authors point out through their hyperlinks, and to the actions users take in using (and implicitly evaluating) the documents on the site. We cover these issues in this paper. We discuss both work completed and work not yet implemented.

## 2    The commKnowledge community

As computers transform how businesses are run, more and more information is being produced on this topic, and much of it is on the web. There is information about: businesses that rely on computers to interact with their customers or suppliers; advances in the computing and telecommunications infrastructure that supports electronic business; financing new, electronic forms of business; electronic markets and auctions; locating and pricing goods; legal issues surrounding electronic commerce; and privacy.

The sources of information on these topics – and others too numerous to mention – are just as varied. Online newspapers post stories daily. Popular press magazines devote significant attention to the topic of electronic commerce and business, as do technical journals. Businesses themselves have websites that give a first-hand glimpse of the state of practice of electronic commerce, and there are, additionally, scholars, research organisations, and consulting houses who offer their observations on what is occurring and what they believe is likely to.

This wealth of information, from countless sources, certainly makes it easy to find *some* information in the area of electronic business. But it makes it that much harder to keep up with all that is going on. The commKnowledge website was developed with this issue in mind. Whilst it is impossibly ambitious to create a site that contains everything about electronic commerce and business – broadly construed – it is possible to try to develop a central node on the web that points to a great deal of this information. That is the practical aim of the commKnowledge project.

Instead of selectively crawling the web to find relevant content or employing a panel of content specialists who seek content, we rely on individuals throughout the world to locate and contribute the items they feel belong on commKnowledge. These items will ordinarily be URLs of other websites (which we call *articles*); but they can also be content especially created for commKnowledge, such as a user's 'editorial' posted as a pdf file (which we call *features*); or they can be items like an online news source, online magazine, or online newspaper that contain many relevant articles (which we call *resources*). The machinery behind the commKnowledge website allows various policies for considering what information is included on the site. Similarly, the pattern of contributions that the commKnowledge community makes to the site, as well as the community's votes about information on the site, determine how prominently different information on the site is displayed.

The commKnowledge project originated at the University of Michigan, and has received support from individuals at the University of Haifa, Israel.

## 3   Designing and populating commKnowledge

Two considerations were foremost in designing and launching commKnowledge: how to build a community and how to organise information to make the site easy to use.

To succeed in making commKnowledge a 'go to' site for information about electronic commerce, we wanted as many people as possible to contribute information. But this is fraught with potential peril. Who are these people? Are their interests truly coincident with those of the community we hope to develop? Are they making genuine contributions, or do they seek to corrupt the site by furnishing information they know will not be of interest or contains inauthentic material?

In deciding what policies to adopt in light of these questions, we followed certain general guidelines. Firstly, we wanted to use computer-based means for making these decisions rather than rely on our own judgments and opinions. Secondly, we were willing to err on the side of being overly-inclusive rather than too-restrictive in making decisions about the appropriateness of content. Thirdly, we wanted to be flexible with our policy: we expect to learn more as the project grows, we believe that different policies might be best as the project 'scales up', and we are interested in evaluating the effects of various policies governing the screening of content.

As a bootstrapping policy, we decided that all contributions would be considered on-the-money and would become part of our site. This extremely inclusive policy was aimed at gathering content for the site quickly. With this content, we hope to make the site attract new users. With new users, we hope to increase the pace at which commKnowledge is populated with new content. To attempt to get enough information to attract the first visitors to the site, we ourselves seeded the site with the vast majority of initial contributions. Nonetheless, establishing an active, engaged community of contributors is certainly a challenge.

As the site grows and we learn more about those using it, we plan to test other policies for making include/exclude decisions. In the next section we briefly discuss the concept of a 'staging' area where contributions are vetted before they are placed on the system.

A second major design issue was how to organise our content. Three organising schemes seemed useful: time, users' opinions, and topic. Our considerations in organising the information one sees upon 'walking in the front door' show how these organising schemes interact with each other.

On a website like commKnowledge, users will be interested in returning to the site only if its content changes. So, organising the information that users see when they enter the site in reverse chronological order was a compelling option. But commKnowledge also depends on its members' efforts to evaluate content. By listing the most recent items first, we would not be taking advantage of this capability. On the other hand, listing the *best*-judged items first would mean that newly posted information would not be near the beginning of the list. In fact, being unevaluated, it might languish after all already-judged items, remaining unseen and unevaluated forever. Finally, our repository could also be broken down topically. Though all information on commKnowledge is about electronic business, it is not all the same. Some information is about the supporting hardware, for

instance, while other information is about societal issues. So, a topical organisation was compelling, too.

We settled on making our *Recent Contributions* section something akin to our 'headlines', where a frequent user of the system would spot newly posted information and could work his way back to older items. But we also decided to organise content by topic in our *Contributions* sections – using a non-hierarchical set of categories into which a user classifies his contributions. And we use a 'best first' ordering of items *within* each 'leaf' category based on *use graphs*. Use graphs combine aspects of social filtering [3] – where users' judgments are instrumental in determining the content of the site; the hyperlink structure of the web [4]– where certain documents enjoy special status because of the documents they link to or from; and users' actions within the commKnowledge site.
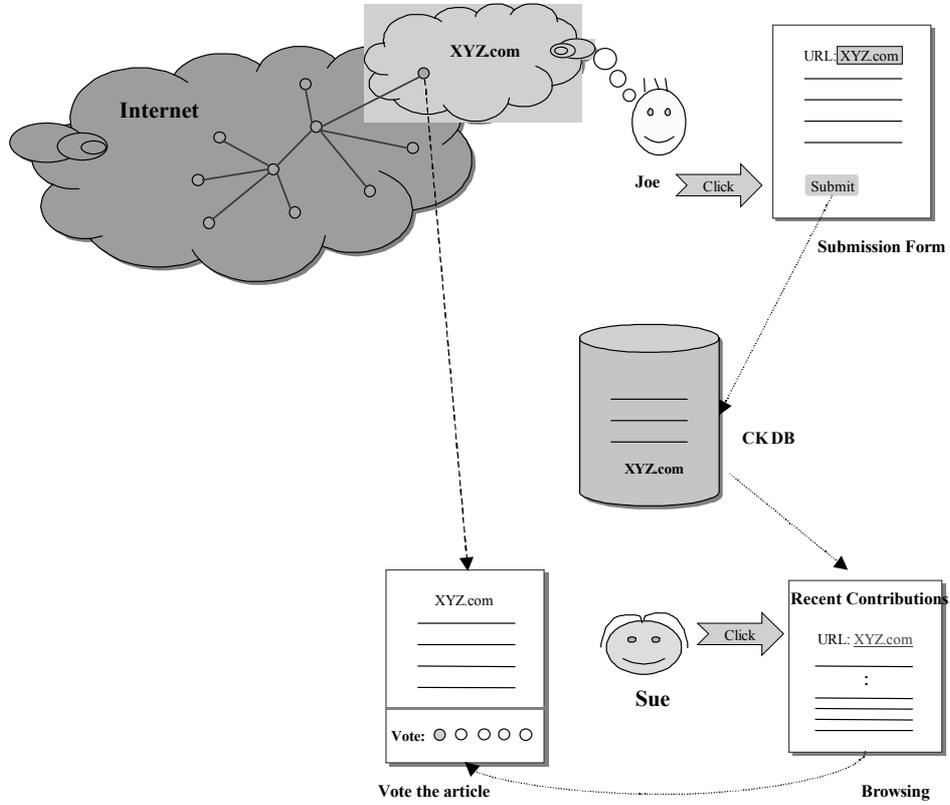
## 4 Functional description of commKnowledge

The commKnowledge website obtains and displays its contents dynamically. Users of the site use web forms to post new items. These forms are used to collect the item itself (usually a URL) as well as additional information such as a user-supplied summary. Once this item is submitted the commKnowledge site is immediately dynamically updated. In future versions of the system, submissions will be placed in a 'staging area' before they are added to the site. Trusted users will then make decisions about their suitability for the site. The section of the paper entitled Evaluating commKnowledge Content and Users describes how these users can be determined. We use the terms *content*, *content item,* or *item* to designate something that has been contributed to and now appears on commKnowledge.

When a registered commKnowledge user views a content item, it is displayed in a large content frame. A second frame encourages the user to vote on its usefulness, read or make comments, or vote on others' comments. This information becomes part of what other users see when they look at the item, and it is also the input we will use later to screen content and determine the order to display content. See Figure 1.

Users can sort lists of content items in various ways: reverse chronologically, by submitter, by 'type' (article, resource, or feature), by course (for the use of instructors), by a document- or user-'hub' score, or by a document- or user-'authority' score. (We discuss hub and authority methods in the section entitled Evaluating commKnowledge Content and Users.) We also provide a full-text search capability for negotiating the content on the site.

**Figure 1**    Flow from item submission through dynamic presentation
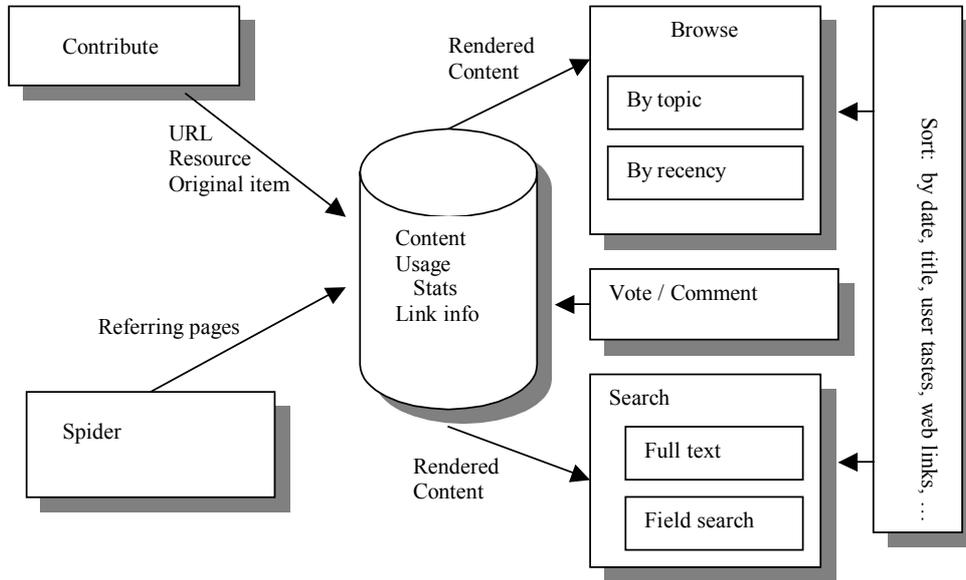


## 5    CommKnowledge architecture

CommKnowledge is organised around a relational database implemented in Oracle. When users make submissions to commKnowledge, the information they provide is stored in many relations (tables). These tables are what the system uses to draw its content from to create html pages for users to view. The database also serves as a repository for other information the system needs to run, such as information about the link structure among contributions and about voting histories.

Figure 2 provides a system diagram of commonKnowledge.

**Figure 2** CommKnowledge system diagram



The main modules of commKnowledge are:

## *Contribute*

A user may contribute an online article, resource or an original contribution. When he clicks the *contribute* button, our web server sends to his browser an online submission form. Using this form, he provides information about the article such as: its category, its author, a brief summary, and its source.

When the user clicks the *submit* button, several validation checks are automatically applied to the submission by JavaScript to ensure its quality. For example, a spider/agent automatically detects an invalid URL. After it passes all validity checks, the submission form is posted to the web server via the HTTP protocol. From there, it is parsed and its contents are automatically inserted into the *item* table, which stores all the information related to this particular submission. Once the table has been updated by the insertion, the new submission is available for others. See Table 1 for an abbreviated schema for an online article or resource:

**Table 1**     Item table schema

| Field | Meaning |
| --- | --- |
| Submission-ID | Key field of submission |
| Title | User supplied title to be displayed on the screen |
| Summary | User supplied summary of the submission |
| Posting Time | When the submission was posted |
| Authors | Author(s) of the item submitted |
| Article Title | Original title of submission |
| URL | URL of submission |
| Source | Source (magazine, newspaper, etc.) of submission |
| Submitter-ID | User who made the submission |
| Category | User-supplied topical category of item |

Upon insertion, a submission launches a spider that downloads the item's contents from the internet and extracts all its hyperlinks. The hyperlinks to other sites are stored in a *connectivity_info* table. This connectivity information is used for calculating authority and hub scores. The downloaded pages are also stored for further processing late at night. At that time, pages that point to the submitted page are identified by a commercial search engine, and all of their hyperlinks are extracted and stored in the connectivity_info table. As before, hyperlinks within the same website are ignored. See Table 2.

**Table 2**     Connectivity_info table schema

| Field | Meaning |
| --- | --- |
| Source URL | URL of page containing link |
| Target URL | URL of page linked to |

### View contributions

A design goal of commKnowledge is to make sure that users can make the best use of the information being shared by the community. Two modules are designed to support ease of navigation of the knowledge base: *Recent Contributions* and *Browse by Category*. The Recent Contributions module lets users see the most recent submissions, while the Browse by Category module offers users capabilities to explore content categories.

### Recent Contribution

The *Recent Contributions* display page contains the 15 most recent submissions to the system, with links to earlier contributions. These contributions are dynamically and automatically extracted by a server side script from the *item* table based on their *submission time*.

### Browse by category

When a user clicks on the *Contributions* button on the main page, the main categories of the contributions are displayed. Currently, these are *b2b, b2c, financial, statistics, hardware & software, societal & government, legal, international, education, miscellaneous*. When he/she clicks on one of these main categories, a script on the server

side will 'pull' the items within this category from the database, reformat them as html code using a standard format template, and send them back to the client's browser in a display list.
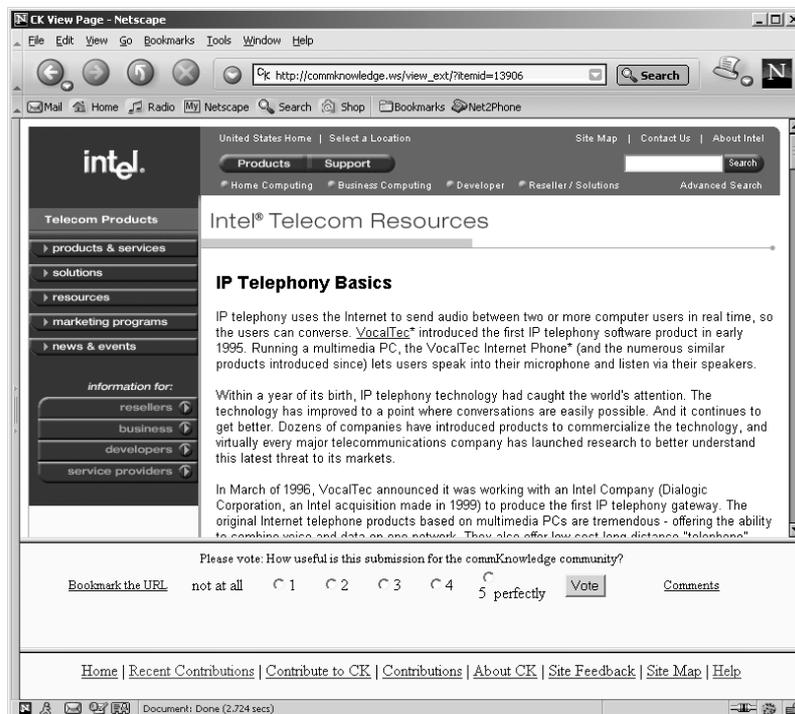
## *Sorting*

Various sorting criteria are provided to allow users to view display lists in different ways. Items may be sorted by submission time, submission title, submitter, submission type, or course number. In addition, items may be sorted in ways that depend on either or both documents' link structure, and users' actions.

## *Rendering dynamic pages*

Searching, viewing Recent Contributions, and viewing content by Contributions (i.e., category) all create summary display pages containing lists of content items, briefly described. When a user clicks on one of these items, a request is sent to the server to get the URL corresponding to the item's *submission id* from the *item* table. A server script also dynamically generates source code for a three-frame display page, and sends the html source code back to the client's browser. The source code for the top frame of the display page contains the URL of the submission, which the client's browser uses for fetching content from the web to display. The middle frame is a voting frame, and the bottom frame is a navigation frame, both of which the commKnowledge server script has generated. See Figure 3.

**Figure 3** CommKnowledge display page

*Vote/comment on a contribution*

In addition to sharing content, we strive to elicit members' opinions and comments to make this content more useful. This feedback serves as input to a variety of our algorithms for ranking content items and users.

When the user *submits* his vote, it is inserted in an *item_vote* table. Java scripts prevent null-votes, ensure no user votes twice on the same item, and update the average score for the submission in the *item* table. This updated information is available immediately to other users. See Table 3.

**Table 3**      Item_vote table schema

| Field | Meaning |
|---|---|
| User ID | Identifier of user |
| Item ID | Identifier of content item |
| Vote | User's vote for this item |

Users can also click on a *comment* hyperlink in the voting window to make a comment about or read other people's comments on this submission. Comments related to this submission are pulled from an *item_comment* table via a server script and are sent back to the client's browser for display in a popup window. An input text form is also provided to allow the user to make a comment on this submission. The comment pop-up window also allows a user to reply to others' comments or vote on other people's comments. All of these actions are handled by a single server script. Comments, replies, and votes on comments are stored on the server and are handled using the method of HTTP POSTing. See Tables 4 and 5.

**Table 4**      Item_comment table schema

| Field | Meaning |
|---|---|
| Comment_ID | Identifier of comment |
| User ID | Identifier of user |
| Item ID | Identifier of content item |
| Comment | User's comment on item |

**Table 5**      Item_comment_voting table schema

| Field | Meaning |
|---|---|
| Comment_ID | Comment being voted on |
| User ID | User voting on comment |
| Vote | User's vote on comment |

*Hub and authority computations*

Each night, several actions occur to update the commKnowledge database. First, 50 pages containing links *to* any item posted that day are fetched by a call to a search engine. For each of these pages, all outbound links are extracted and stored in the *connectivity_info* table. With this new link information, together with the *item_vote* and *item_comment* vote tables, we calculate new document-document, user-document, and user-user hub and authority scores.  See the next section for details about computations.
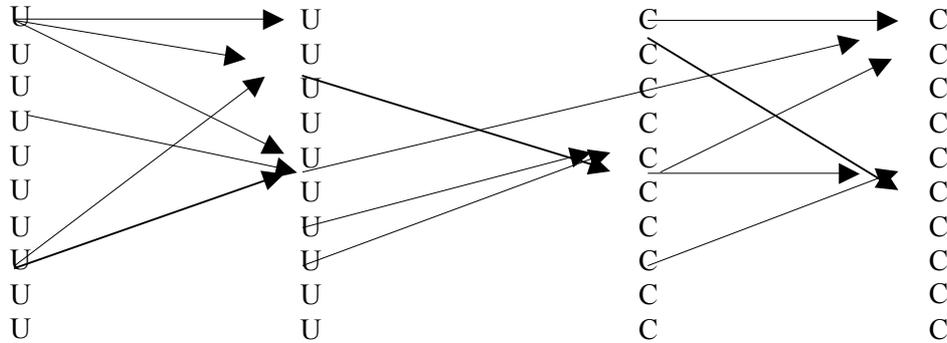
## 6   Evaluating commKnowledge content and users

The items on the commKnowledge site will not all be of equal value to the commKnowledge community. Some of the factors that may influence an item's value include: its main subject, its source or author, its length, or its style. Knowing that some items are perceived as being better than others, we want to be able to make the 'best' items the first ones that the commKnowledge community sees. Naturally, our success in doing this hinges on being able to determine which items truly are the best. The hypertextual nature of the web and users' actions within the commKnowledge repository can both be useful in making this determination.

Structurally, an item may have many incoming links from other commKnowledge content items. It may also have many outgoing links to other content items within commKnowledge. Users' actions within commKnowledge also provide information. Users submit content items. Some are favourably voted on; others are not. A user reads certain items, some multiple times, and votes or comments on content items and others' comments. These actions may or may not be consonant with the community's opinions and preferences.

We lump together these various user actions with the link structure of documents in the system by creating *use graphs*. Use graphs combine different types of *endorsement*: endorsement of one content item by another, endorsement of a content item by a user, and endorsement of one user by another. Use graphs show endorsement by means of a directed edge. Users endorse other users or content items by their actions (such as a favourable vote). One content item endorses another if it links to it.

Use graphs determine the default order in which documents are displayed within a category and provide a way for users to sort the results of a search. By pointing out the 'highest quality' users, use graphs also indicate which individuals may best be entrusted with determining which documents should be allowed to become part of the commKnowledge system. A use graph is an extension of the type of graph that Kleinberg *et al*. [4] studied in the CLEVER retrieval system. Whereas their graphs consider content-content hyperlinks only, in our system we use those links as well as user-user links and user-content links. Each of these kinds of links provides a different kind of endorsement. See Figure 4 for a picture of a use graph.

**Figure 4**   Use graph showing influential users and content



Note:   As with Clever, some content items are still hubs or authorities. But now so are
        users, especially authoritative users

In considering content-content links, we use the hyperlink structure within commKnowledge to 'score' web pages. Using methods similar to those used in the Clever system, we look for densely interconnected nodes. More precisely, we look for 'authority' nodes and 'hub' nodes — where *authority nodes* are nodes with many incoming links emanating from 'high quality' hub nodes; and *hub nodes* are nodes with many outgoing links terminating at 'high quality' authority nodes. Although these definitions are circular, there are standard mathematical algorithms [5] to determine a node's 'hub value' – how good it is at pointing out useful web pages; and a node's 'authority value' – how useful it is thought to be. Both of these measures rely on the fact that a web page tends to 'endore' another web page if it links to it. We take advantage of this structure to allow users to sort the content of commKnowledge using either authority or hub scores based on link structure.

We use the user-content links contained within use graphs in a similar fashion. 'Hub users' are those who make many recommendations to high quality nodes; high quality nodes are those that these 'high quality' users endorse. Resnick *et al*. [6] explain the built in peril of recommender systems: useful judgments may be intermixed with misleading ones, with no obvious means of corroborating judgments or ferreting out those that are a disservice. In our system, however, users rise in stature only when they recommend high quality items; failing to make a recommendation or making recommendations of items judged to be of poor quality will demote you. And commKnowledge can take seriously only the recommendations of high-stature users. In other words, by knowing who these 'hub users' and 'authoritative users' are, other users can be led to the content items that they are recommending. This, then, becomes yet another way we allow the commKnowledge community to gain access to shared content.

Finally, use graphs contain user-user links. The same circular interrelationship we've been describing holds among users. By commenting on others' contributions or on their comments, users are implicitly making endorsements. Some people act as 'hubs' by becoming especially good at endorsing high quality users. 'Authoritative' (high quality) users are endorsed by many of these people. By knowing who these well-functioning individuals are, we gain the opportunity to screen new contributions to the system. Since content items may be submitted by *any*one, it is useful to vet them before placing them

on the site. Authoritative users – as well as hub users who can reliably identify them – are perfect candidates for this job.

Two points should be emphasised about Figure 4. First, all the arrows in a use graph are considered *positive endorsements*. This means that all content-content links are considered positive – even though a link could actually be to an illustration of a bad example of some idea. Similarly for user-content links: these are considered positive endorsements.

Second, graph 'links' are generated and 'weighted' in different ways, depending on the type of link. Content-to-content links come about by the first content item referring to the second with a hyperlink. A user-to-content link arises by combining a user's voting on (endorsing) a content item's usefulness with information about the user's viewing behaviour. A user's vote on an item's usefulness (1 to 5) is rescaled to a numerical weight between 0.2 and 1.0. The number of times the user reads an item is converted using a nearly logarithmic transformation so that reading an item once creates a weight of 0.5 whereas a large number of reads creates a weight of 1.0.

In short, use graphs

1    can determine the order with which documents in the same category are displayed

2    can provide a new way for users to sort the results of a search

3    can even serve as a mechanism to determine which documents are allowed to
     become part of the commKnowledge system.

In our research, we are exploring how best to 'blend' the various hub- and authority-weights that we obtain from each of the link types contained in use graphs: document-document, user-document, and user-user. Our goal is to use this information in the most effective way to identify hub and authority users ('high performance' users whose judgments can be relied upon) so as to improve the identification of hub and authority documents. If we succeed, we believe we can make effective decisions about what content to include on commKnowledge and can display it in ways that best serve the commKnowledge user community.

## 7   Closing remarks

Community knowledge-sharing efforts have paid practical dividends. Open source efforts have led to the development of sendmail and Linux [7]. The Slashdot.org community and site is a vital source of technical information for software developers. Other sites serve as clearinghouses for information on various topics of professional and personal interest.

Obtaining these benefits is not always easy, however. One of the co-founders of slashdot revealed that getting your site to the point that it 'takes off' is difficult; but after that, momentum carries it forward on the weight of fresh content and a critical mass of participants [8].

The upshot of this is to recognise that a site like commKnowledge can offer significant advantages by pooling quality information from a variety of sources and providing different means of accessing that information that take advantage of the community's collective wisdom. But, the fullest benefits of the site come from 'network effects' – the notion that the value of the community increases markedly as more people

participate in it. These arise only from the investments of those who help build the community before it makes good on its promise.

Even before it reaches critical mass, commKnowledge can serve both as a personal filing system for gathering together and sharing with others information on electronic business, and as a reference to many materials devoted to this topic. By actively contributing to and using the site and encouraging others to do the same, commKnowledge can become a definitive source for a community wishing to stay abreast of developments in electronic business.

## References

1   Search Engine Watch,
    http://searchenginewatch.com/reports/sizes.html

2   Searchability,
    http://www.searchability.com/atoz.html

3   Resnick, P. and Varian, H. (1997) 'Recommender systems', *Communications of the ACM*, Vol. 40, No.3, pp.56–58.

4   Kleinberg, J.M. (1999) 'Authoritative sources in a hyperlinked environment', *JACM*, Vol. 46, No. 5, pp.604–632.

5   Golub, G. and Van Loan, C. F. (1983) *Matrix Computations,* Johns Hopkins University Press, Baltimore.

6   Resnick, P., Zeckhauser, R., Friedman, E. and Kuwabara, K. (2000) 'Reputation systems', *Communications of the ACM*, Vol. 43, No. 12, pp.45–48.

7   Raymond, E.S. (1999) *The Cathedral and the Bazaar: Musings on Linux and Open Source by an Accidental Revolutionary*, O'Reilly & Associates, Inc., Sebastopol, CA.

8   DeMaagd, K. (personal communication, 2002).