

2004-01-01

Using structured P2P overlay networks to build content sensitive communities

Paul Stacey

Technological University Dublin, paul.stacey@tudublin.ie

Damon Berry

Technological University Dublin, damon.berry@tudublin.ie

Eugene Coyle

Technological University Dublin, Eugene.Coyle@tudublin.ie

Follow this and additional works at: <https://arrow.tudublin.ie/engscheleart>



Part of the [Electrical and Computer Engineering Commons](#)

Recommended Citation

Stacey, Paul and Berry, Damon and Coyle, Eugene :Using structured P2P overlay networks to build content sensitive communities. Proceedings of ICPADS 2004: Tenth International Conference on Parallel and Distributed Systems, 7-9 July, 2004,p pp.281-288. doi:10.21427/d1q3-4821

This Conference Paper is brought to you for free and open access by the School of Electrical and Electronic Engineering at ARROW@TU Dublin. It has been accepted for inclusion in Conference papers by an authorized administrator of ARROW@TU Dublin. For more information, please contact arrow.admin@tudublin.ie, aisling.coyne@tudublin.ie.



This work is licensed under a [Creative Commons Attribution-Noncommercial-Share Alike 4.0 License](#)

Audio Research Group

Conference papers

Dublin Institute of Technology

Year 2004

Using structured P2P overlay networks to build content sensitive communities

Paul Stacey*

Damon Berry[†]

Eugene Coyle[‡]

*Dublin Institute

[†]Dublin Instit

[‡]Dublin Institute of Technology, Eugene.Coyle@dit.ie

This paper is posted at ARROW@DIT.

<http://arrow.dit.ie/argcon/1>

— Use Licence —

Attribution-NonCommercial-ShareAlike 1.0

You are free:

- to copy, distribute, display, and perform the work
- to make derivative works

Under the following conditions:

- Attribution.
You must give the original author credit.
- Non-Commercial.
You may not use this work for commercial purposes.
- Share Alike.
If you alter, transform, or build upon this work, you may distribute the resulting work only under a license identical to this one.

For any reuse or distribution, you must make clear to others the license terms of this work. Any of these conditions can be waived if you get permission from the author.

Your fair use and other rights are in no way affected by the above.

This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike License. To view a copy of this license, visit:

- URL (human-readable summary):
<http://creativecommons.org/licenses/by-nc-sa/1.0/>
 - URL (legal code):
<http://creativecommons.org/worldwide/uk/translated-license>
-

Using Structured P2P Overlay Networks to Build Content Sensitive Communities

Paul Stacey

Damon Berry

Eugene Coyle

*School of Control Systems and Electrical Eng.
Dublin Institute of Technology,
Kevin Street, Dublin 8, IRELAND.
{paul.stacey, damon.berry, eugene.coyle}@dit.ie*

Abstract

This paper details a proposed peer-to-peer system, which allows a user to join communities of other like-minded users in order to exchange files. Utilising the routing capabilities of Pastry, the proposed system includes an indexing service, which will facilitate the creation of virtual rendezvous points for users with similar interests (manifested by shared keywords). Users will be described by the content they store. By using Vector Space Modelling techniques users can be grouped together to form content sensitive communities. The system is built to serve as the basis for a distributed archive of research papers.

The system is designed to improve the efficiency of file searches over current p2p file sharing applications. Search requests result in a comprehensive set of relevant documents being returned as searching will be based on semantic meaning rather than literal matching.

1. Introduction

Community: a body of people having common rights, privileges, or interests [1]. Applying this notion of community to the Internet, it seems that the Internet as a whole lacks a sense of community. The World Wide Web was originally envisaged in 1986 by Tim Berners-Lee as a way for academics to share knowledge [2]. The web has become a victim of its own success in relation to that goal. While there is a huge quantity of information on the web and it is easy to find text about almost any topic, it is often difficult to distinguish “high quality” information such as peer-reviewed papers from material from less prestigious sources. The client-server paradigm may be partly to blame for this; power has been taken away from the individual and been placed in the hands of operators

of large servers. In recent times, systems such as Napster [3] and Gnutella [4] have gained huge popularity. These systems have initiated a surge of interest and research into the peer-to-peer (p2p) framework. These systems are restructuring the Internet away from the client server model to one where a client is also a server, giving individuals more freedom and control. This paper presents the core of a distributed document-sharing environment with particular focus on distributed searching and retrieval of research papers.

1.1. Problems with Current p2p Systems

Despite their obvious popularity, systems like Napster and Gnutella suffer from many problems. Napster uses centralised indexing servers, an approach which is vulnerable to failure. Gnutella avoids Napster’s weakness by using a decentralised indexing technique, however this leaves Gnutella with the problem of locating objects within its network. Gnutella uses a flood-based search technique where each search request blindly hops across the network from one node to another searching for the requested file. As more users join, the number of nodes to be searched increases yet the number of nodes searched remains relatively small; as a result search requests do not return a true representation of the available objects stored in the system as it grows. This represents a scalability issue. In recent times, more scalable object location algorithms have emerged that are based on Distributed Hash Tables (DHT). Let us consider these useful constructs and how they can be used to facilitate searches in a distributed environment.

1.2. Existing Structured p2p Overlays

Based originally on the research of Plaxton et al. in the late 1990s [5], DHT overlay network implementations

first appeared during 2001. Projects such as Pastry [6], Tapestry [7], CAN [8] and Chord [9], all produced implementations that adhered to the principles of p2p, decentralisation, robustness and scalability. These systems may be used to form the foundation for functional p2p systems. They provide a routing substrate; a mechanism that efficiently locates objects within a certain number of routing hops. The subject of keyword searching is of particular importance if such DHT systems are to become part of the more mainstream p2p systems such as Gnutella or Kazaa [10]. Introducing keyword-searching capabilities into these systems is likely to render them a more powerful tool for file sharing than is currently available. DHT's currently provide only put and get functions. These substrates have however proved useful in building such systems as global storage facilities, including PAST [11], CFS [12] and Oceanstore [13]. PAST which as we shall see has a lot in common with the system described here, is built on top of Pastry and uses the power of Pastry to route files entered into the system to a particular point, given a file key. The file may be retrieved once the file key is known.

Before venturing further into the details of p2p technologies, it is helpful to consider a useful application of DHT systems. The application seminal to this research is that of *Content Networks*.

1.3. Content Networks

A content network is an overlay IP network that supports *content routing*. Content routing means that messages are routed based on their content rather than their IP-address. In recent years many types of content networks have been developed, including p2p networks. PAST is one such development. Content networks can be classified into many different types. In the following discussion users will be associated with identifiers that have semantic meaning. Users will also be subject to content-sensitive placement. [14] This gives taxonomy to the different types of content networks.

It was intended to build a p2p system on top of Pastry that supports the construction of communities based on the content they store. These communities will make searching for files far easier and more efficient because searches may be directed to particular communities within the network where the files are more likely to be stored. These communities will be formed using state-of-the-art Information Retrieval (IR) techniques in-order to better discover relationships between files and thus return more comprehensive search results.

The above discussion has provided a brief description of the basic technologies that are currently in use in the area of content networks and p2p. The ingredients of the proposed system will now be discussed.

2. Ingredients of the Design

The work reported here involves the construction of a p2p framework whereby users will be able to join the network and have contact with other users who are either interested in a topic or have the means to share content on that topic. It is important to describe the chosen routing substrate. Pastry will provide the p2p framework with a scalable and robust routing algorithm (described in section 2.1.). The system also incorporates a state-of-the-art Information Retrieval (IR) technique used for representing documents, known as Vector Space Modelling (described in section 2.2.). In order to group together similar users within the system, clustering techniques are employed. This enables the system to form communities of users based on a similarity index derived from Vector Space Modelling. This is the focus of section 2.3.

2.1. Pastry

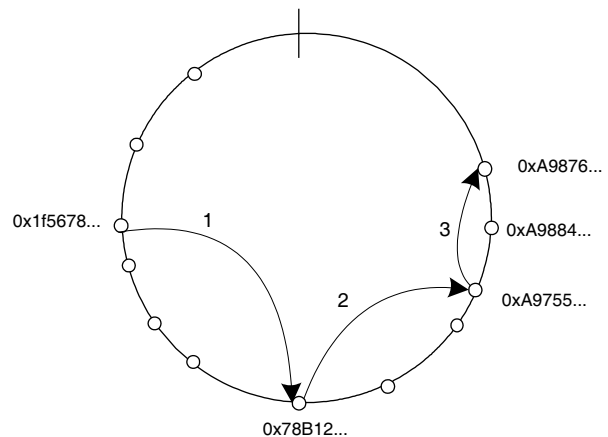


Figure 1. Pastry routes messages to nodes whose nodeIds are progressively closer to the message key.

Pastry nodes are organised around a circular id space. Each node within the Pastry network is assigned a 128-bit unique identifier that is generated typically from the cryptographic hash of, for example, its IP address and users name. E.g. Using the *Secure Hashing Algorithm* (SHA-1) [15], a string such as “computer science” will produce the hash code “0bbb843c75b8cb93ceb9d5594e208668484448ee”. Pastry has the ability to route messages between nodes when given a message key. The message is routed to the node whose nodeId is numerically closest to the key of the message. Pastry’s routing algorithm is efficient, scalable and robust.

The organisation of nodes around Pastry’s ring is random. This is due to the way in which nodes are

assigned their id's. A hashing algorithm is employed for the purpose of generating the unique codes; a commonly used hashing algorithm is SHA-1. SHA-1 is non-reversible, collision-resistant, and has a good avalanche effect. The avalanche effect of hashing algorithms means that given two very similar strings, two very different and non-numerically close hash codes will be produced. Pastry uses this feature as a way of achieving load balancing within the network by randomly placing nodes around the network

2.2. Vector space modelling

It has been shown that a document may be represented as a *feature vector*. This modelling of a document as a vector is called "Vector space modelling"[16]. In its simplest form, each document is represented by the *term-frequency* (TF) vector $d_{tf} = (tf_1, tf_2, \dots, tf_n)$, where tf is the frequency of the i th term in the document. A widely used refinement to this model is to weight each term based on its *inverse document frequency* (IDF) in the document collection. This is commonly done by multiplying the frequency of each term i by $\log(\frac{N}{df_i})$, where N is the total number of documents in the collection, and df_i is the number of documents that contain the i th term (i.e. document frequency). This leads to the $tf-idf$ representation of the document in equation 1.

$$d_{tfidf} = (tf_1 \log(\frac{N}{df_1}), tf_2 \log(\frac{N}{df_2}), \dots, tf_n \log(\frac{N}{df_n})) \quad (1)$$

In [17] it was shown experimentally, that any measure used should be normalized by the length of the document vectors. In order to account for documents of different lengths, the length of each document vector is normalized so that it is of unit length, i.e. $\|d_{tfidf}\|_2 = 1$. There are two major similarity metrics that facilitate vector space modeling of documents [18]. One of them is the *angle-based* metric that uses for example the cosine of the angle between the vectors. The cosine function is given in equation 2. Documents may then be compared and a similarity index can be established between two documents.

$$\cos(d_i, d_j) = \frac{d_i \bullet d_j}{\|d_i\|_2 * \|d_j\|_2} \quad (2)$$

Where " \bullet " denotes the "dot product" of two vectors. Since the document vectors are of unit length, the above formula simplifies to

$$\cos(d_i, d_j) = d_i \bullet d_j \quad (3)$$

The ability to represent and compare documents using vector space modelling is a very useful tool, enabling inter-document relationships to be determined more accurately than through the use of keyword matching alone.

2.3. Clustering

Clustering is the unsupervised classification of patterns (observations, data items, or feature vectors) into groups (clusters) [19]. Document clustering was originally investigated as a means of improving the performance of search engines. Since then document clustering or cluster-based techniques have been used in domain identification in such areas as radio news and imaging. Cluster analysis allows the identification of groups, or clusters, of similar objects in multi-dimensional space. *Hierarchic clustering* has been put forward for its efficiency and effectiveness in information retrieval. There are numerous document clustering algorithms. *Agglomerative Hierarchical Clustering* (AHC)[20] algorithms appear to be the most commonly used. However these algorithms have proven to be slow when applied to large document sets. *Linear time clustering* algorithms have been suggested as the best candidates for large document sets, these clustering algorithms include the K-means algorithm [21] and the single pass method [22].

Given a set of feature vectors, a clustering tree can be constructed. Vectors that are deemed similar, using a similarity index such as the cosine function are placed near each other on the tree and those that are less similar are positioned further away. There are many methods for constructing clustering trees. Jain et al. [19] provide a good overview and discuss the pros and cons of each method.

3. Putting the Ingredients Together

Having considered Pastry, Vector Space Modelling and Clustering, the next step is to explain how these ingredients can be incorporated to create a p2p system that supports decentralised, scalable and robust content sensitive communities.

3.1. Building a Decentralized indexing service

One of the key elements of the system is the *Pastry routing schema*. Pastry provides the system with a

scalable, robust routing algorithm that can route to any node in $\lceil \log_{2^b} N \rceil$ steps on average, where N is the number of nodes and b is a configuration parameter. As stated previously, Pastry routes messages within the Pastry network to nodes whose nodeIds are closest to the key of the message. The system discussed here is based around a decentralised indexing mechanism where nodes that are “interested” in a certain topic may join communities of like-minded users. In order to facilitate this, a “rendezvous” point where nodes can discover others that have “similar interests” is required. As will be shown, Pastry provides a routing algorithm that enables the system to be built without any central control.

3.1.1. Creating indexing nodes

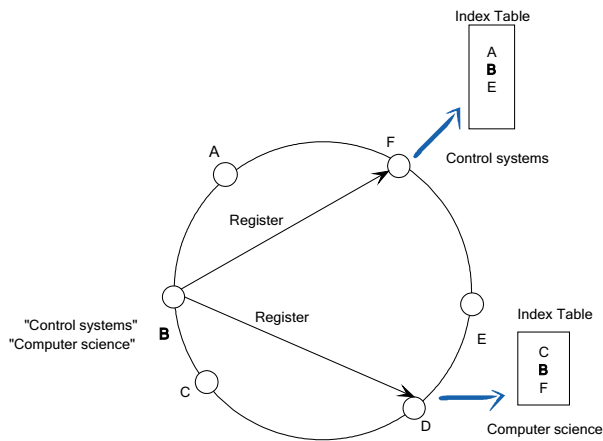


Figure 2. Node B routes a register message with 2 message keys that are the cryptographic hash of the keyphrases throughout the Pastry ring. Index tables registering all nodes sharing the same keyphrases are constructed.

Given a string of characters, the SHA_1 hashing algorithm will produce a 160-bit hash code representing the string. This property forms the basis for the indexing service.

Once a node has calculated its unique nodeId it may join the Pastry network. The joining procedure is provided by Pastry. Any node joining the network will have a set of *keyphrases* associated with it that best describes the node’s “topics of interest”. These phrases will serve as the basis for discovering nodes sharing similar content. Each of these keyphrases is hashed to get a hash code for each phrase. These hash codes will be used as message keys so that Pastry can route them around the Pastry ring to a live node whose nodeId is numerically closest to the 128 most significant bits of the 160-bit key. A registry-message is constructed; this message contains the nodes details such as its IP-address

and *node-vector* (described in section 3.2.1). The same registry-message is routed several times throughout the network for each keyphrase. Each registry-message uses the 160-bit codes generated from the hashing of the keyphrases as keys. When the messages have arrived, each destination node is required to register the new node (see Figure 2). The effect of this is that every other node stating “control systems” as a keyphrase will send a register message with the same key to be routed to the same node, as in Figure 2, this results in all these nodes being registered in the same index table. In the case of Figure 2 the registering node is node F. If an index does not exist at node F, one will be created. The node whose nodeId corresponds to the 128 most significant bits of the hashed keyphrase will now serve as the rendezvous point or indexing node for all other nodes using the same keyphrase. Nodes will register at the same point for two reasons:

- A hashing algorithm given the same input string will always produce the same output key. Therefore a registry-message will always get the same key for the same keyphrase.
- Pastry’s routing algorithm routes messages to the live node whose nodeId is numerically closest to the message key. Therefore two nodes will always end up registering at the same point once they share a keyphrase.

The above only holds true of course when the indexing node remains live on the network. To deal with node failures and hence loss of indices, it is proposed to replicate the index table among the indexing nodes k nearest neighbours (k is a configuration parameter that determines the number of neighbouring nodes where the index will be replicated). There are a number of other systems (e.g. PAST) that use the properties of DHT systems for similar purposes. When a file is inserted into PAST, Pastry routes the file to the k nodes whose node identifiers are numerically closest to the 128 most significant bits of the file identifier (fileId). These nodes then store the file. Other global storage systems built on top of DHT’s include OceanStore [13] that is built on top of Tapestry, and CFS [12], which uses Chord.

3.2. Building Content Sensitive Communities

In this section we describe how nodes are compared based on stored content, the organising of the p2p network into communities and how this can be used to form a two-layer network.

3.2.1. Comparing Nodes Based on Content Stored

Consider a node storing a set of documents that share the same subject content. It can be said that a node’s set

of documents can be classified under one general heading. This heading will have a relation to the subject content of each of the documents stored. Another way to look at this general heading would be to describe it as the “average subject” of the documents. Consider again the situation where each document has an associated vector representation, derived from the *td-idf* representational model. It is now possible to generate an *average vector* of these document-vectors. This average vector is representative of the average subject content of a particular node’s document set. It therefore tells something about the subject content the node “is interested in”. A way of deriving such a vector comes from centroid-based classification [23]. Given a set of *S* documents and their vector representations, a centroid vector *C* can be defined, which is the average vector of the set of documents. This is given by equation 4.

$$C = \frac{1}{|S|} \sum_{d \in S} d \quad (4)$$

These average vectors are called *node-vectors*. Nodes may now be succinctly represented based on the content they store. It is also possible to compare a pair of node-vectors to assess a similarity index between the corresponding pair of nodes by using the cosine measure as described in section 2.2. These tools provide a way of comparing and hence grouping nodes that are similar within the network. This is described in the following section.

3.2.2. Organising Network into Communities

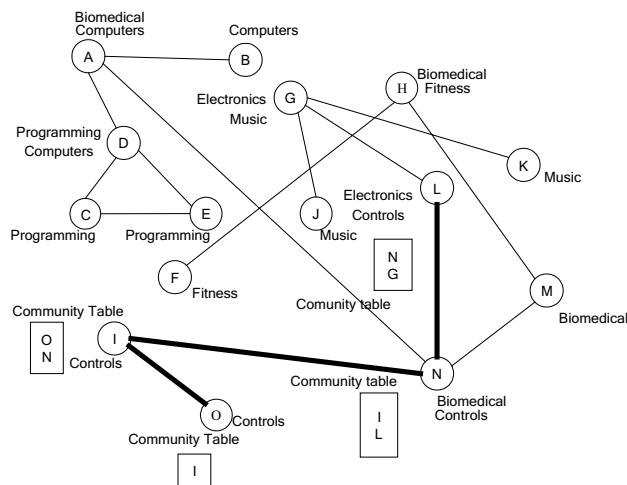


Figure 3. Shows the effect of the community layer formed by community tables.

Assuming that the node has a specific topic of interest, all documents stored will be related in some way to this main topic. A node-vector can therefore be used to compare nodes that store similar content. If each node stores a *community table* containing a list of nodes that are similar (based on the similarity metric from the vector space model described above), document collections of a similar content will then be implicitly linked or grouped together (Figure 3). This could be seen as organising the network into domains whose boundaries are not strictly defined but have a “fading” effect as we jump from one community table to the next. By moving along the path *O→I→N→L* (highlighted in Figure 3), the document collection moves from “Controls” to “Biomedical” and then to “Electronics”. Node *A* can be seen as the node within the network that stores documents relating both to “Biomedical” and “Computers” and thus is the point within the network where the two domains overlap. Any linked group of nodes can be seen as a domain. Each domain can be categorised based on the nodes that have created connections between each other. As they all store similar content, these “communities” of nodes are *content sensitive* in nature because only nodes that store content that is similar to the content stored in the community as a whole will become part of it.

3.2.3. Two Layer Network

The network is maintained and organised by employing two layers, the Pastry layer and the Community Layer. Pastry maintains routing tables and *leaf sets* (leaf sets are tables containing neighbouring nodes within a certain number of hops that each node “knows” about [6]). This layer is a means for nodes to find indexing nodes of certain subject areas. The nodes are organised randomly due to the nature of the hashing algorithm employed to create their unique nodeId. The second layer, the *community layer*, is more organised. The ability of the second layer to organise itself is a direct result of the indexing service built on top of Pastry. When a node is added to an indexing node’s index table, the indexing node is provided with the new node’s node-vector and IP-address. Comparing its node-vector to that of the already indexed nodes, the indexing node places the details of the new node in the appropriate place within a clustering tree. The cluster tree is formed by grouping “similar” nodes “close” together, nodes are determined similar if their node-vectors are determined close by a similarity metric such as the cosine measure. The new node then uses similar nodes within the cluster tree to populate its community table (see Figure 4). All nodes added to the community table are then contacted and asked to add the

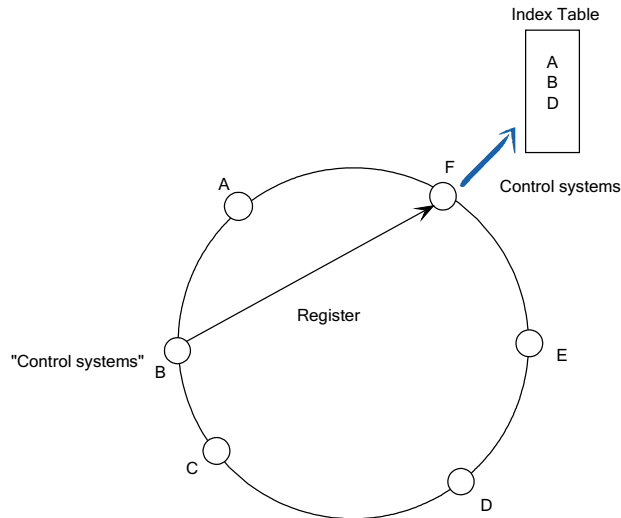


Figure 4A. The new node B populates its community table with the other nodes whose node-vectors are most similar to its own.

new node to their community table. This is done for each keyphrase. After this procedure has been carried out, the new node has links to other nodes that have similar node-vectors and hence have a good probability of being interested and sharing similar content. This means that nodes that “have similar interests” know about each other and can share content directly. This type of organisation makes searching within p2p networks much more efficient as all the content is grouped together into communities and so searches can be directed to a specific area of the network instead of being flooded blindly throughout the network.

3.3. Searching for Documents

Pure flood-based searches have become an essential feature of unstructured p2p networks such as Gnutella [4] and LimeWire [24]. These systems rely on flooding of search messages throughout the network in order to locate files stored by nodes. Within these systems pure flood search requests are given a Time to Live stamp (TTL), this TTL sets the number of hops a search message is allowed to execute before “dying”. Pure flood-based search methods have proved inefficient and non-scalable and result in bottle necking within the Internet. However, Pure flooding has been shown to scale well within the Gnutella network up to 10,000 nodes [21].

With a more structured overlay network, flood-based searches can be used while maintaining scalability and cutting down on unnecessary query messages. This is possible by performing a focused flood search where the

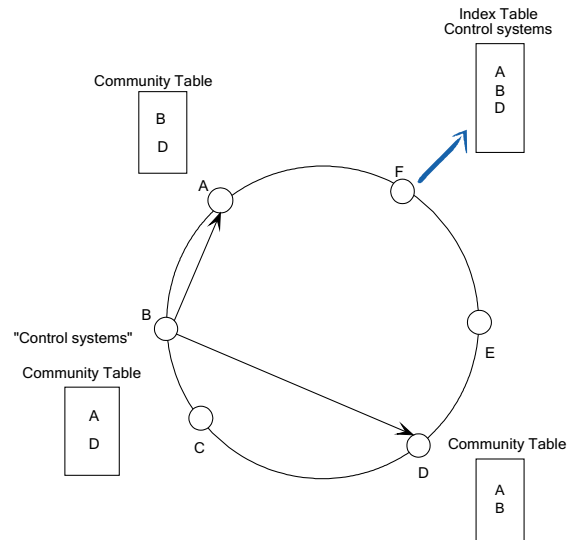


Figure 4B. Node B then informs nodes A and D that they must add B to their community tables.

search is focused on a specific area of the network. It is possible to target a particular part of the system and perform an exhaustive search on those areas that are more likely to contain the type of files being requested. This is achieved by first of all entering a list of keyphrases. A *search-vector* can be produced in the same way a file vector was calculated. In order to direct the search to an area of the network storing files relating to the keywords, the search-vector is compared locally to nodes within the community table by again using a similarity metric. The search request is then forwarded to those nodes whose node-vectors are the most similar to the search vector. The contacted node performs a flood-search on its community table using the original keywords as the search parameters. The proxy searcher then compiles a list of “hits” and returns them directly to the requesting node. The requesting node may choose to download directly any of the files that it finds or perform another search on a different part of the network.

4. Status of work

A Java package that implements the vector space modelling algorithm has been developed. The package includes the Porter Stemming algorithm [25] and the Van Rijsbergen Stop List [26], which are common refinements to vector space modelling (Figure 5).

point that needs more attention is the use of node-vectors in classifying a user's document set. This use of the node-vector could prove to be naive and a more accurate implementation of this idea may need to be investigated.

8. References

- [1] <http://www.dictionary.com>.
- [2] <http://www.w3.org>.
- [3] Napster. www.napster.com.
- [4] Marius Portmann, Pipat Sookavata, Sebastien Ardon, Aruna Seneviratne. The Cost of Peer Discovery and Searching in the Gnutella peer-to-peer File Sharing Protocol. Proceedings, *Ninth IEEE international conference on Networks*, 2001.
- [5] C. Greg Plaxton, Rajmohan Rajaraman, Andrea W. Richa. Accessing Nearby copies of Replicated Objects in a Distributed Environment. *Theory of computing systems* 1999.
- [6] A. Rowstron and P. Druschel, "Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems". *IFIP/ACM International Conference on Distributed Systems Platforms (Middleware)*, Heidelberg, Germany, November, 2001.
- [7] B. Y. Zhao, J. D. Kubiatowicz, and A. D. Joseph. Tapestry: An infrastructure for fault-resilient wide-area location and routing. *Technical Report UCB//CSD-01-1141*, U. C. Berkeley, April 2001.
- [8] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A scalable content- addressable network. In Proc. *ACM SIGCOMM'01*, San Diego, CA, August. 2001.
- [9] I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup service for Internet applications. In Proc. *ACM SIGCOMM'01*, San Diego, CA, August. 2001.
- [10] Kazaa. <http://www.kazaa.com>.
- [11] A. Rowstron and P. Druschel. PAST: A Large-Scale, Persistent Peer-to-Peer Storage Utility. *HotOS VIII*, Schloss Elmau, Germany, May 2001.
- [12] Dabek, F., Kaashoek, M. F., Karger, D., Morris, R., AND Stoica, I. Wide-area cooperative storage with CFS. In Proceedings of the *18th ACM Symposium on Operating Systems Principles (SOSP '01)* Canada, October. 2001.
- [13] J. Kubiatowicz, et al. OceanStore: An Architecture for Global-Scale Persistent Storage. *ASPLOS*, Cambridge, MA, December 2000.
- [14] H. T. Kung, C. H. Wu. Content Networks: Taxonomy and New Approaches. *The Internet as a Large-Scale Complex System*, 2002.
- [15] D. Eastlake, 3rd and P. Jone. RFC 3174: US secure hashing Algorithm 1, Sept. 2001.
- [16] Salton, G., Wong, A., and Yang, C. S., A Vector Space Model for Automatic Indexing. *Communications of the ACM*, November 1975.
- [17] Willet P., Similarity coefficients and weighting functions for automatic document classification an empirical comparison. *International Classification*, 1983.
- [18] Jones, W., and Furnas, G., Pictures of Relevance: A Geometric Analysis of Similarity Measures. *Journal of the American Society for Information Science*, November 1987.
- [19] Jain, A.K., Murty M.N., and Flynn P.J. Data Clustering: A Review. *ACM Computing Surveys*, 1999.
- [20] P. Willet. Recent trends in hierarchical document clustering: a critical review. *Information Processing and Management*, 1988.
- [21] J. J. Rocchio, Document retrieval systems - optimization and evaluation. *Ph.D. Thesis*, Harvard University, 1966.
- [22] D. R. Hill. A vector clustering technique. In Samuelson (ed.), *Mechanised Information Storage, Retrieval and Dissemination*, North-Holland, Amsterdam, 1968.
- [23] Han, E-H. and Karypis, G. Centroid-Based Document Classification: Analysis and Experimental results In Proc. Of the *4th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD)* September 2000.
- [24] LimeWire. <http://www.limewire.com>.
- [25] The Porter Stemmer Algorithm. <http://www.tartarus.org/~martin/PorterStemmer>.
- [26] Van Rijsbergen, C.J. *Information Retrieval*, Butterworths, 1975.
- [27] JAMA: <http://math.nist.gov/javanumerics/jama/>.
- [28] PJX: <http://www.etymon.com/epub.html>
- [29] R. Cilibrasi, R. de Wolf, P. Vitanyi, Algorithmic clustering of music, Submitted. <http://arxiv.org/archive/cs/03030>.
- [30] Jane Grimson, William Grimson, Damon Berry, Gaye Stephens, Eoghan Felton Dipak Kaltra, Pieter Toussaint, Onno W. Weier: A CORBA-based integration of distributed electronic healthcare records using the Synapses approach. *IEEE Transactions on Information Technology in Biomedicine* 2 1998.
- [31] Jane Grimson, Eoghan Felton, Gaye Stephens, William Grimson and Damon Berry. Interoperability issues in sharing electronic healthcare records - the Synapses approach. Proceedings of *Third IEEE International Conference on Engineering of Complex Computer Systems*, IEEE 1997.
- [32] CORBAMED document: formal/01-04-04 (Person Identification Service Specification, v1.1).
- [33] C. Tang, Z. Xu, and M. Mahalingam. Peersearch: Efficient information retrieval in peer-to-peer network". In Proceedings of *HotNets-I, ACM SIGCOMM*, 2002.
- [34] Z. Xu, M. Mahalingam, and M. Karlsson. Turning heterogeneity to an advantage in overlay routing. *Technical Report HPL-2002-126*, HP Laboratories Palo Alto, 2002.