

## **Maximum Coverage Probability based Query Registration and Processing in Unstructured P2P Network**

**Md Hussain Khusro<sup>1</sup>, Yasmeen Begum<sup>2</sup>**

<sup>1</sup>Pursuing M.tech (CSE) from Khaja Banda Nawaz College of Engineering, Gulbarga.  
Affiliated to VTU Belgaum, Karnataka, India.

<sup>2</sup>Prof. Yasmeen Begum, Department of Computer Science and Engineering, Khaja Banda Nawaz College of Engineering, Gulbarga.  
Affiliated to VTU Belgaum, Karnataka, India.

### **Abstract:**

Large amount of data are available in large-scale networks of autonomous data sources dispersed over a wide area. P2P is a system of acquiring data directly from the clients using a discovery process monitored by the server. As such, in such a system only information about the data and the nodes are maintained at the server and the communication is in Peer to peer manner between the clients.

If we assume a network where data is consistently being changed or new data are released and that the client is continuously generating query than, unavailability of the data at the instance of query generation leads to information loss. In order to overcome this problem, we propose a unique query processing based P2P system where query for which no data is available are stored in special nodes called Beacons. Once the data is available by some client, the beacon announces the same to the querying client.

To transfer the data by utilizing minimum bandwidth and maximum coverage, split and merge algorithm is proposed. For each downloading, a file is chunked in equal parts equivalent to number of clients. Clients start downloading the parts in parallel. Once each client has different chunks, they download the missing chunks from each other thus balancing the load at the seeder. Result show improved search time and throughput utilization for this method.

**Keywords:** P2P Network, Query Registration, Query Processing, Searching in P2P, Maximum Coverage.

### **I. Introduction:**

In recent years unstructured peer-to-peer (P2P) systems have evolved as a popular paradigm for content/resource distribution and sharing [1, 6]. Owing to the simplicity of design and flexibility towards transient node population, the real-world P2P

systems are invariably unstructured. However, most unstructured P2P content distribution systems only support a very simple model for data sharing and discovery called the ad hoc query model. A peer that is interested in discovering data items initiates a query with a set of search parameters, which is then circulated among the peers according to the specific query forwarding mechanism employed by the network. A peer receiving a query responds to the query initiator, if it has any content satisfying the search criterion. Once a query has been processed at a node, it is removed from the local buffers (some systems cache recently received queries, but for a very short duration and in an ad hoc fashion). Therefore, a query exists within the P2P network only until it is propagated to various nodes and processed by them (or for a short duration thereafter, if the network employs caching). Once a query completes its circulation, the system essentially forgets it.

While the ad hoc query model for data discovery is essential for P2P content distribution networks, it suffers from two serious limitations. First, due to its very nature, an ad hoc query is only capable of retrieving content that exists in the P2P network during the time period when it is actively propagated and processed in the network. Further, an ad hoc query can never reach a peer that joins the network after the query has completed its circulation, and hence cannot discover matching data-items on the new peer. In this scenario, the only way for a peer to discover newly added data-items would be to repeatedly issue the same query, thereby imposing unnecessary overheads on the network. Second, the ad hoc query model provides no support for peers to advertise or announce the data-items they own to other interested peers. Such capabilities are important for P2P communities where peers trade content.

These shortcomings limit the utility of the ad hoc query model for several advanced collaborative applications, such as a community of researchers sharing their recent research results or a community of amateur musicians and their patrons who are interested in buying the music produced by

the musicians. In applications such as these, participating peers would not only be interested in searching for existing content, but would also want to be pro-actively informed when content matching their interests is added to the network. Further, some communities also need a mechanism through which peers can advertise their content to other interested peers. Blind broadcast of advertisement would not only result in high overheads, but could also annoy participants who would be receiving large numbers of advertisement about data-items that they are not interested in.

An approach that can partially mitigate these limitations would be to implement a publish-subscribe (pub-sub) system on top of the unstructured overlay network. A generic pub-sub system enables its users to register subscriptions expressing their interests and to announce the occurrence of certain events by publishing them. The pub-sub system matches incoming announcements to the existing subscriptions and notifies the users that have registered the matching subscriptions. An important point to note is that the pub-sub systems attempt to provide guaranteed notification service (although it might not be possible always due to system failures).

Researchers have studied the problem of implementing P2P-based pub-sub systems on unstructured overlay networks [7, 17]. However, most of these systems require the underlying P2P networks to be organized according to specific architectures, and hence they cannot be used in generic overlays. Many of these systems also require the peers to maintain intricate index structures which add significant complexity to the design of the P2P network. This additional complexity can adversely affect the flexibility, efficiency, and scalability of the unstructured P2P system. Furthermore, it also makes the design, implementation, and management of P2P content distribution networks harder.

## **II. Related Works:**

The work presented in this paper is primarily related to two fields, namely P2P networks [6, 9, 14] and publish subscribe systems (event-delivery systems) [3,5], both of which have been very active areas of research in the past few years.

Pub-sub systems can be classified into two broad categories:

(1) topic-based – wherein users join specific topic groups in which all the messages related to the topic are broadcast; and (2) content-based – wherein users specify their interests through predicates. With the aim of enhancing scalability, efficiency and scalability several distributed pub-sub systems have been proposed [3, 5]. Recently, P2P computing models have been utilized for this purpose. Researchers have explored two strategies for

constructing P2P-based pub-sub systems, namely (a) adopting a structured P2P network like Chord [15] or CAN [13] as the underlying substrate, and utilizing its indexing schemes for mapping subscriptions and events to nodes of the P2P systems [10, 16]; (b) organizing the nodes of the P2P system into specialized topologies and/or embedding application specific distributed index structures within nodes of the P2P network [7, 17, 19]. The Sub-2-Sub system [17] organizes the peers into clusters using an epidemic-style algorithm such that nodes with similar subscriptions are put into the same cluster. The publisher of an event joins the corresponding cluster and disseminates the event to the cluster members.

The proposed system differs from the above works in terms of motivation, goals and approach. The goal of the above systems is to improve the various performance parameters of pub-sub systems and they use P2P-based techniques as a means towards this end. In contrast, our goal is to enhance the P2P data sharing systems, and continuous queries (that bear similarity to pub-sub model) is a means towards that end. Second, the above pub-sub systems cannot be implemented on top of generic P2P networks; they need specialized overlays (specific topologies and/or indexing mechanisms). Contrastingly, our system does not need any complex distributed indexing structures, nor does it impose any topological constraints on the overlay network. Finally, the above systems are essentially pub-sub systems, and hence guaranteed notification is one of their design goals. Our system provides best-effort notification, which is in tune with design principles of unstructured P2P networks. P2P-DIET [11] supports both ad-hoc and continuous queries, however, it assumes a super peer-based overlay.

In short, the work presented in this paper has several unique aspects, and it addresses an important problem in the area of P2P data sharing systems.

## **III. Our Approach**

### **3.1 Problem Formation**

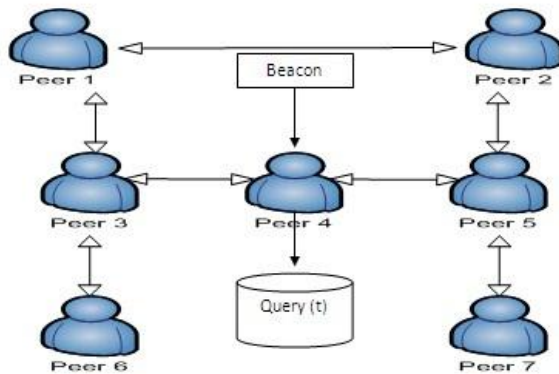
Peer to peer networks uses different computers or peers to share the files amongst themselves rather than keeping the files in a single server. Hence same file may be downloaded from different nodes in a peer. But due to non-central nature of the communication, such network loses control and finding out group of nodes relevant to data sharing is low. Therefore here we propose a system to

- 1) Index the files at the server
- 2) A search engine to respond to the queries.
- 3) Maximum neighborhood based beacon selection for registering queries



- 4) Chunking based parallel file download for load balancing
- 5) Achieve scalability for new information

### 3.2 Proposed System



**Figure 1: Pure Decentralized P2P Content-Sharing Architecture**

The above fig1 shows decentralized P2P content sharing system where a user register a query at peer4 called beacon node at time t. Once the data is available by some client, the beacon announces the same to the requesting client.

Further in present system, a client can download only the information that is available at that instance of time. But in the proposed system, a client can download the information after the information is added at a later time. Further the system proposes a chunking based technique with maximum coverage. Generally in P2P system a client downloads the information from the best possible path, so the coverage is least. Hence if any intermediate client seeks the information in between a session, then without having to re-establish a fresh session, nodes can join the existing transmission and start downloading the needed chunks.

### 3.2 Concepts and Notations

Consider an unstructured P2P system comprising of peers  $(P_0, P_1, \dots, P_{N-1})$ . Let  $(L_0, L_1, \dots, L_{M-1})$  represent the logical links (connections) in the network. For simplicity, we assume that the links are bidirectional. Two peers  $P_i$  and  $P_j$  are said to be neighbors of each other if there exists a link  $L_v = (P_i, P_j)$  connecting them.

We assume that each data item  $D_r$  in the system has associated metadata (represented as  $MData(D_r)$ ) that describes it. In the current context, the metadata is a list of keywords describing the data item.

Continuous query is the means through which a peer can register its interests with the network. A continuous query, represented as  $Q = (SID; Predicate; V Time)$ , is essentially a tuple of three components, namely, source ID (SID), query predicate (Predicate) and validity time (V Time).

The source ID uniquely identifies the peer issuing the query. The query predicate is the matching condition of the query, and is used by the source peer to specify its interests. In general, the predicate can be of any form such as range predicates or even a regular expression. We assume that the predicate is a list of keywords describing the content the source peer is interested in. Validity time (V Time) represents the time until which the source node is interested in receiving notifications. Peers announce their new data items through announcements.

An announcement is represented as  $Ad = (AID; MData)$ . The announcing peer ID (AID) identifies the advertising peer and the metadata (MData) is the metadata of the content being advertised. A data item  $D_r$  (and analogously its announcement) is said to match a continuous query  $Q_m$ , if  $D_r$ 's metadata contains all the keywords in  $Q_m$ 's predicate. We use the word query and continuous query [8,12] interchangeably.

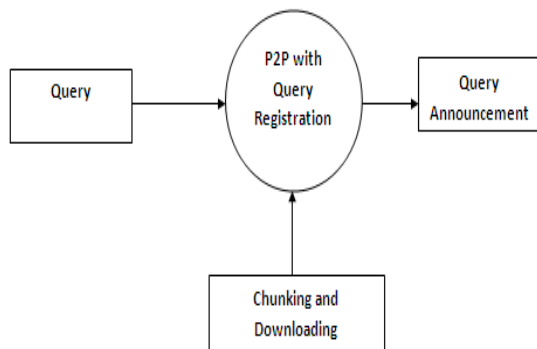
### 3.3 System Design

**Algorithm:** The proposed system is explained as bellow.

- \*Generate Random set of nodes.
- \* Index a set of information-TF-IDF (term frequency-inverse document frequency).
- \* Randomly distribute the information to the Nodes.
- \* One or more clients query for a specific information.
- \* Search and locate the clients where the information is available.
- \* Request the Seeder for downloading
- \* divide the data into equal pieces equivalent to number of lecher.
- \* Data is sent to the clients through a route from seeder to lecher. Each lecher then connects with each other and downloads the missing chunks. (Maximum Coverage Technique).
- \* If there is no data available for a query, it is registered at a node called a beacon node with maximum reachability ratio.
- \* Once the data is available, the information is announced to the lechers and the downloading begins as above.
- \* After certain time period the query is expired to maintain the integrity of the freshness of the information.
- \* For query processing, exponential time-maximum likelihood estimation query registration and forwarding is used.

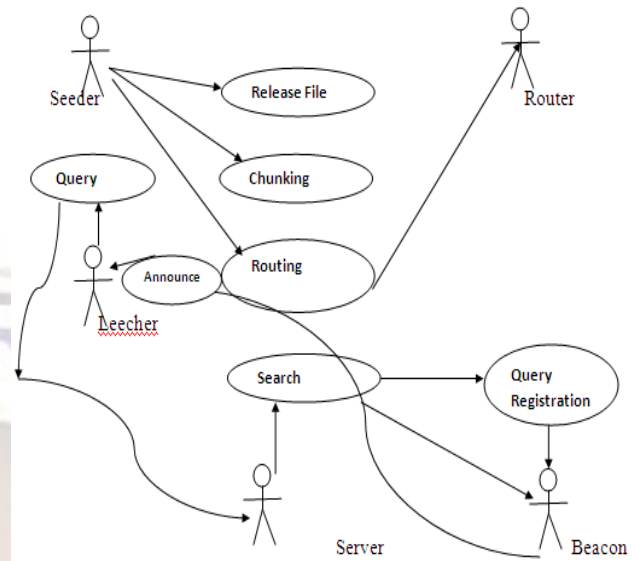
## IV. Diagrams

**First Level DFD**



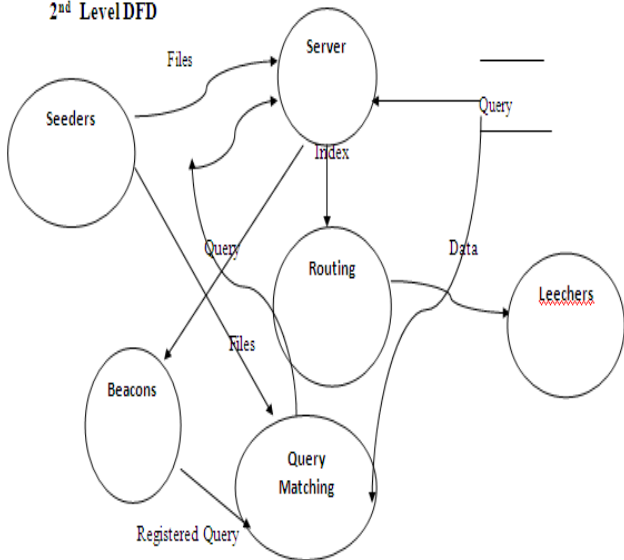
**Figure2: First Level Data Flow Diagram**

**Use Case Specification**



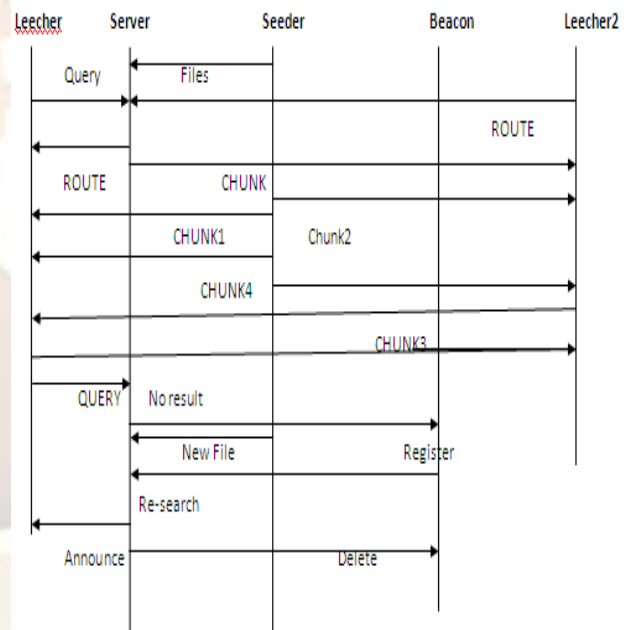
**Figure4: Use Case Specification Diagram**

**2<sup>nd</sup> Level DFD**



**Figure3: Second Level Data Flow Diagram**

**Sequence Diagram**



**Figure5: Sequence Diagram**

## V. Experiment and Results

### 5.1 Experiment:

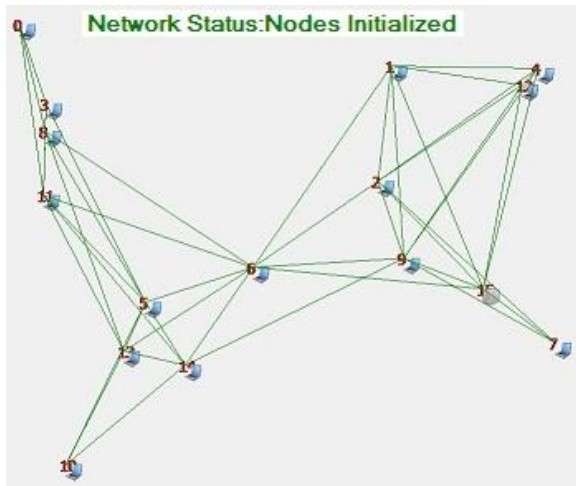


Figure6: Here we create a p2p network of 15 nodes.



Figure7: Now we perform indexing of files and announce the data randomly to peers.

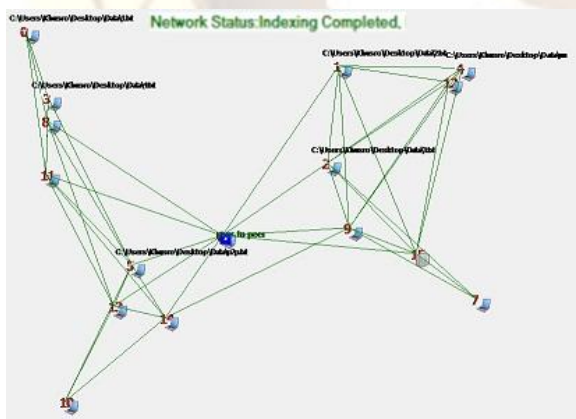


Figure8: Here we show announcement of data files at peers and perform search that generate a query for which data is not available in the network and query is registered at beacon node indicated with blue rectangle in above figure8.

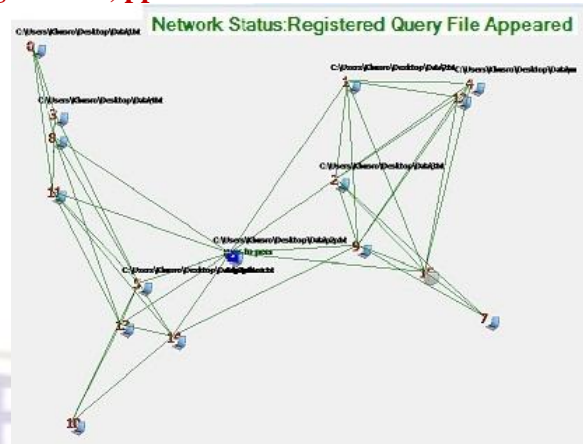


Figure9: Now a matching data for registered query is appeared in the network, our system automatically notify a registered query file appeared and its location.

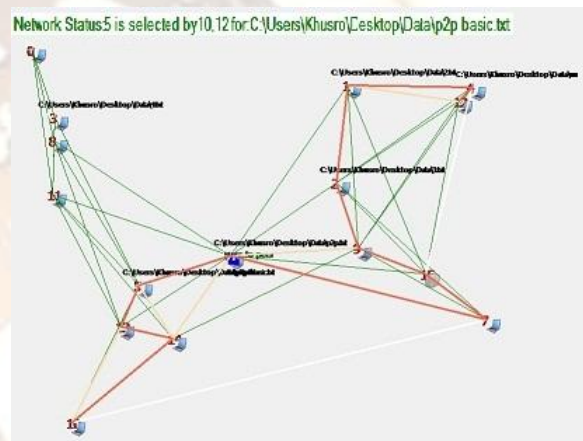


Figure10: Now we create a route from seeder to leecher(s)(requesting clients) with maximum coverage and divide the file into chunks equivalent to number of leechers and perform transmission.

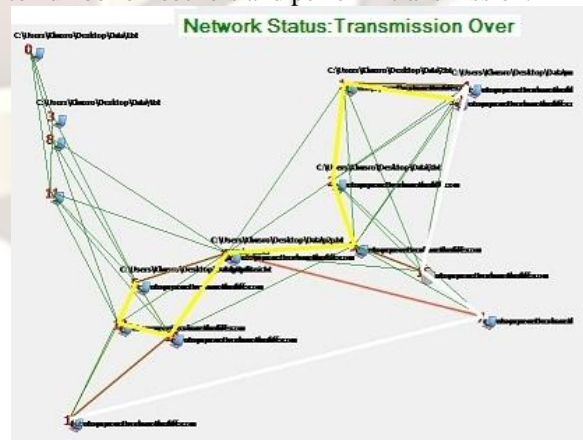


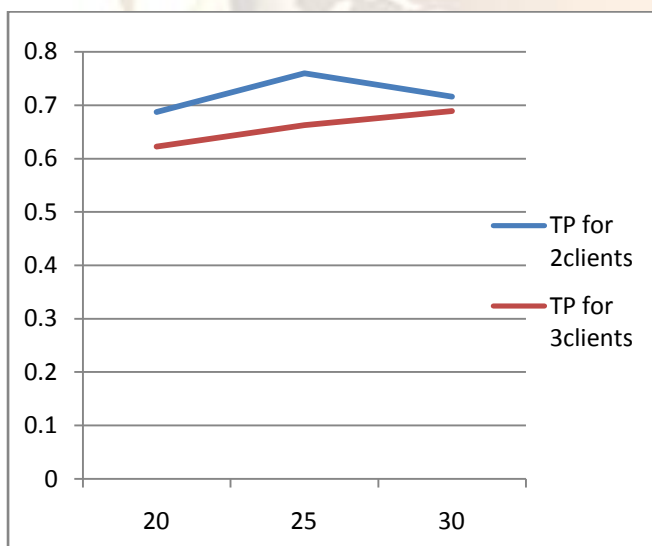
Figure11: Finally we get the data for registered query and measure its latency and chunks and determine the throughput.

## 5.2 Results



### Result1: Number of Node v/s Latency

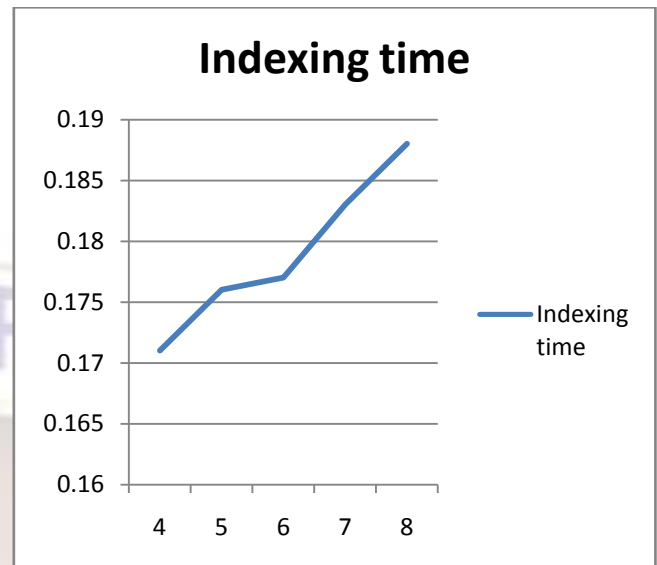
Latency is defined as the end to end delay for a node to acquire an entire file through collecting of chunks from seeder and through transformed seeders. The graph shows that the latency of the system depends upon the leechers rather than the network size. For limited leechers, the transmission time is low and for the higher leechers, the same is increased.



### Result 2: Number of Nodes v/s Throughput.

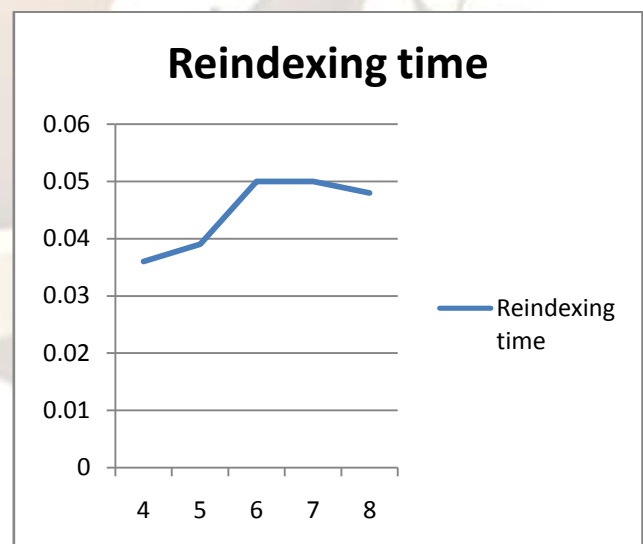
The performance graph shows that the chunking and parallel downloading and maximum coverage routing helps in enhancing the throughput. Generally throughput decreases significantly in P2P system. But in the current process, throughput is increased with increase of nodes which suggest that as the network

grows, availability of the data also increases which is most desirable property of the P2P system.



### Result 3: Number of Data Items v/s the indexing time.

This experiment was conducted by considering different length of text documents. The indexing is considered as TF-IDF score of the documents. Indexing time increases with increase of number of files and is independent of the file size. Therefore as the network grows, indexing the entire set consumes time. Hence the concept of query processing is used which eliminates the indexing for every query. Once a query is unanswered the query is stored. Hence only re indexing is needed once a new data is made available which is relevant to the query.



### Result 4: Number of Files v/s Reindexing time

The performance graph clearly explains the utility of the technique. Once a new data is matched with query it is reindexed. Numbers of reindexing



iterations are limited. Hence the system takes lesser time for announcement of new files.

## VI. Conclusion

Peer-to-peer systems have become a popular media for sharing large amount of information among millions of users. While previous research efforts are focusing on supporting search in P2P systems, obtaining hidden and valuable knowledge from these data through data mining techniques is essential for scientific findings and many other applications. In this work, we investigate searching and query registration for fast information announcement to the nodes. We provide complexity analysis on the transmission incurred by the system. The analytic result indicates that proposed system can efficiently mitigate data to the nodes seeking the information. The system can be further improved by adopting passive replication technique.

## VII. References

- [1] Gnutella P2P Network. [www.gnutella.com](http://www.gnutella.com).
- [2] Kazaa P2P Network. [www.kazaa.com](http://www.kazaa.com).
- [3] G. Banavar, T. Chandra, B. Mukherjee, J. Nagarajarao, R. E. Strom, and D. C. Sturman. An Efficient Multicast Protocol for Content-Based Publish-Subscribe Systems. In Proceedings of ICDCS 1999, 1999.
- [4] N. Bisnik and A. Abouzeid. Modeling and analysis of random walk search algorithms in P2P networks. In Proceedings of HOT-P2P, 2005.
- [5] A. Carzaniga, D. S. Rosenblum, and A. L. Wolf. Design and evaluation of a wide-area event notification service. ACM Transactions on Computer Systems, 19(3):332–383, 2001.
- [6] Y. Chawathe, S. Ratnasamy, L. Breslau, N. Lanham, and S. Shenker. Making Gnutella-like P2P Systems Scalable. In Proceedings of ACM SIGCOMM 2003, 2003.
- [7] P. Chirita, S. Idreos, M. Koubarakis, and W. Nejdl. Publish/Subscribe for RDF-based P2P Networks. In Proceedings of the 1st European Semantic Web Symposium, May 2004.
- [8] Lakshminish Ramaswamy, Member, IEEE, and Jianxia Chen, Student Member, IEEE. The CoQUOS Approach to Continuous Queries in Unstructured Overlays. IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL.23, NO. 4, April 2011
- [9] S. Androutsellis-Theotokis and D. Spinellis. A Survey of Peer-to-Peer Content Distribution Technologies. ACM Comput. Surv., 2004.
- [10] A. Gupta, O. D. Sahin, D. Agrawal, and A. E. Abbadi. Meghdoot: content-based publish/subscribe over P2P networks. In Proceedings of Middleware 2004, 2004.
- [11] S. Idreos, M. Koubarakis, and C. Tryfonopoulos. P2P-DIET: One-Time and Continuous Queries in Super-Peer Networks. In Proceedings of EDBT, 2004.
- [12] L. Ramaswamy, J. Chen, P. Parate, and A. Meka. Lightweight Support for Continuous Queries in Unstructured Overlays. Technical report, The University of Georgia, 2006.
- [13] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A Scalable Content-Addressable Network. In Proceedings of ACM SIGCOMM 2001, Aug 2001.
- [14] P. Reynolds and A. Vahdat. Efficient peer-to-peer keyword searching. In Proceedings of Middleware 2003.
- [15] I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup service for internet applications. In Proceedings of ACM SIGCOMM 2001, Aug 2001.
- [16] P. Triantallou and I. Aekaterinidis. Content-Based Publish-Subscribe Over Structured P2P Networks. In Proceedings of the International Workshop on Distributed Event-Based Systems (DEBS), 2004.
- [17] S. Voulgaris, E. Riviere, A.-M. Kermarrec, and M. van Steen. Sub-2-Sub: Self-Organizing Content-Based Publish Subscribe for Dynamic Large Scale Collaborative Networks. In Proceedings of the 5th international workshop on peer-to-peer systems, Feb 2006.
- [18] B. Yang and H. Garcia-Molina. Improving search in peer to peer systems. In Proceedings of ICDCS 2002.
- [19] C. Zhang, A. Krishnamurthy, and R. Wang. Combining flexibility and scalability in a peer-to-peer publish/subscribe system. In Middleware 2005.

### Author's Profile:



**Mr. Md Hussain Khusro** pursuing M.Tech in Computer Science and Engineering from Khaja Banda Nawaz(K.B.N) College of Engineering Gulbarga. Affiliated to V. T. U.,

Belgaum, Karnataka, India. My research areas of interest are data mining and data warehousing.

**Mrs. Yasmeen Begum**, Professor, Department of Computer Science and Engineering, Khaja Banda Nawaz College of Engineering Gulbarga. Affiliated to V. T. U., Belgaum, Karnataka., India.