

Flexible Information Discovery in Decentralized Distributed Systems *

Cristina Schmidt and Manish Parashar

The Applied Software Systems Laboratory

Department of Electrical and Computer Engineering, Rutgers University

{cristins,parashar}@caip.rutgers.edu

Abstract

The ability to efficiently discover information using partial knowledge (for example keywords, attributes or ranges) is important in large, decentralized, resource sharing distributed environments such as computational Grids and Peer-to-Peer (P2P) storage and retrieval systems. This paper presents a P2P information discovery system that supports flexible queries using partial keywords and wildcards, and range queries. It guarantees that all existing data elements that match a query are found with bounded costs in terms of number of messages and number of peers involved. The key innovation is a dimension reducing indexing scheme that effectively maps the multidimensional information space to physical peers. The design, implementation and experimental evaluation of the system are presented.

1. Introduction

The recent years have seen an increasing interest in two resource sharing environments: Peer-to-Peer (P2P) computing and Grid computing. Emerging from different communities, they have a common final objective: to pool and coordinate large sets of resources [5].

In the Peer-to-Peer (P2P) computing paradigm, entities at the edges of the network can directly interact as equals (or peers) and share information, services and resources without centralized servers. Key characteristics of these systems include decentralization, self-organization, dynamism and fault-tolerance, which makes them naturally scalable and attractive solutions for applications. Similarly grid computing is rapidly emerging as the dominant paradigm for wide area distributed computing [6]. Its overall goal is to realize a service infrastructure for enabling the sharing of autonomous and geographically distributed hardware, software, and information resources (e.g. computers, data, stor-

age space, CPU, software, instruments, etc.). Both systems also present a common set of challenges: large scale, lack of global centralized authority, heterogeneity (resources, sharing policies) and dynamism [8].

A fundamental problem in these large, decentralized, distributed resource sharing environments is the efficient discovery of information, in the absence of global knowledge of naming conventions. For example a document is better described by keywords than by its filename, or a computer by a set of attributes such as CPU type, memory, operating system type than by its host name. The heterogeneous nature and large volume of data and resources, their dynamism (e.g. CPU load) and the dynamism of the sharing environment (with nodes joining and leaving) make the information discovery a challenging problem. An ideal information discovery system has to be efficient, fault-tolerant, self-organizing, has to offer guarantees and support flexible searches (using keywords, wildcards, range queries). P2P systems, by their inherent properties (self-organization, fault-tolerance, scalability), provide an attractive solution.

This paper presents a P2P information discovery system that supports complex queries containing partial keywords, wildcards, and range queries. It guarantees that all existing data elements that match a query will be found with bounded costs in terms of number of messages and number of nodes involved. The key innovation is a dimension reducing indexing scheme that effectively maps the multidimensional information space to physical peers. The system can be used to index and locate content in P2P storage and sharing systems (using keywords), as a complement for current resource discovery mechanisms in Computational Grids (to enhance them with range queries) or to query interest groups in a bulletin-board news system.

The overall architecture of the presented system is a distributed hash table (DHT), similar to typical data lookup systems [14, 17]. The key difference is in the way we map data elements¹ to the index space. In existing systems, this

*The work presented in this paper was supported in part by the National Science Foundation via grant numbers ACI 9984357 (CAREERS), EIA 0103674 (NGS) and EIA-0120934 (ITR), and by DOE ASCI/ASAP (Caltech) via grant numbers PC295251 and 1052856.

¹We will use the term 'data element' to represent a piece of information that is indexed and can be discovered. A data element can be a document, a file, an XML file describing a resource, an URI associated with a resource, etc.

is done using consistent hashing to uniformly map data element identifiers to indices. As a result, data elements are randomly distributed across peers without any notion of locality. Our approach attempts to preserve locality while mapping the data elements to the index space. In our system, all data elements are described using a sequence of keywords (common words in the case of P2P storage systems, or values of globally defined attributes - such as memory and CPU frequency - for resource discovery in computational grids). These keywords form a multidimensional keyword space where the keywords are the coordinates and the data elements are points in the space. Two data elements are “local” if their keywords are lexicographically close or they have common keywords. Thus, we map documents that are local in this multi-dimensional index space to indices that are local in the 1-dimensional index space, which are then mapped to the same node or to nodes that are close together in the overlay network. This mapping is derived from a locality-preserving mapping called Space Filling Curves (SFC) [2, 16]. In the current implementation, we use the Hilbert SFC [2, 16] for the mapping, and Chord [17] for the overlay network topology.

Note that locality is not preserved in an absolute sense in this keyword space; documents that match the same query (i.e. share a keyword) can be mapped to disjoint fragments of the index space, called clusters. These clusters may in turn be mapped to multiple nodes so a query will have to be efficiently routed to these nodes. We present an optimization based on successive refinement and pruning of queries that significantly reduces the number of nodes queried.

Unlike the consistent hashing mechanisms, SFC does not necessarily result in uniform distribution of data elements in the index space - certain keywords may be more popular and hence the associated index subspace will be more densely populated. As a result, when the index space is mapped to nodes load may not be balanced. We present a suite of relatively inexpensive load-balancing optimizations and experimentally demonstrate that they successfully reduce the amount of load imbalance.

The rest of this paper is structured as follows. Section 2 compares the presented system to related work. Section 3 describes the architecture and operation of the presented P2P information discovery system. Section 4 presents an experimental evaluation of the system. Section 5 presents our conclusions and outlines open issues and future research directions.

2. Related Work

Existing information storage/discovery systems can be broadly classified as unstructured or structured. Unstructured systems (such as Gnutella [21]) are based on flood-

ing techniques and process queries by forwarding them to neighboring peers. While unstructured systems are relatively easy to maintain and can support complex queries, they do not guarantee that all matches to a query in a practical-sized system will be found. Furthermore, overheads of flooding can be significant, and a number of techniques have been proposed to reduce these overheads [4, 8, 11]. Structured data lookup systems (e.g. CAN [14], Chord [17]) use structured overlay networks and consistent hashing to implement Internet-scale distributed hash tables. These systems offer search guarantees and bound search costs, but at the cost of maintaining a rigid overlay structure. Data placement and overlay topology are highly structured and only data lookup using unique data identifiers is supported. Structured keyword search systems [7, 18] extend the data lookup protocol with a distributed inverted index to support keyword searches.

The system presented in this paper is structured and uses an SFC-based indexing scheme. It maps data elements to peers using all keywords, and consequently, when resolving a query only the data elements that match all the keywords in the query are retrieved. This system also supports searches with partial keywords, wildcards, and range queries.

To our knowledge, there exists one other P2P resource discovery system that uses the Hilbert SFC [1]. Unlike our approach, this system uses the inverse SFC mapping, from a 1-dimensional space to a d-dimensional space, to map a resource to peers based on a single attribute (e.g. memory). It uses CAN [14] as its overlay topology, and the range of values of the resource attribute (1-dimensional) is mapped onto CAN’s d-dimensional Cartesian space. This system was designed to enhance other resource discovery mechanisms with range queries. In contrast, we use SFC’s to encode the d-dimensional keyword space to a 1-d index space which can then be mapped to any peer overlay. For example, we can map and search a resource using multiple attributes.

3. System Architecture and Design

The architecture of the presented P2P information retrieval system is similar to data-lookup systems [14, 17], and essentially implements an Internet-scale distributed hash table. The architecture consists of the following components: (1) a locality preserving mapping that maps data elements to indices, (2) an overlay network topology, (3) a mapping from indices to nodes in the overlay network, (4) load balancing mechanisms, and (5) a query engine for routing and efficiently resolving keyword queries using successive refinements and pruning. These components are described below.

3.1. Constructing an Index Space: Locality Preserving Mapping

A key component of a data-lookup system is defining the index space and deterministically mapping data elements to this index space. To support complex keyword searches in a data lookup system, we associate each data element with a sequence of keywords and define a mapping that preserves keyword locality. The keywords are common words in the case of P2P storage systems, and values of globally defined attributes of resources in the case of resource discovery in computational grids.

These keywords form a multidimensional keyword space where data elements are points in the space and the keywords are the coordinates. The keywords can be viewed as base- n numbers, for example n can be 10 if the keywords are numbers or 26 if the keywords are words in a language with 26 alphabetic characters. Two data elements are considered “local” if they are close together in this keyword space. For example, their keywords are lexicographically close (e.g. *computer* and *computation*) or they have common keywords. Not all combinations of characters represent meaningful keywords, resulting in a sparse keyword space with non-uniformly distributed clusters of data-elements. Examples of keyword spaces are shown in Figure 1.

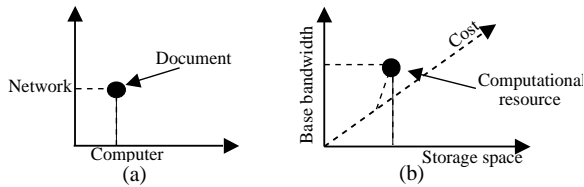


Figure 1. (a) A 2-Dimensional keyword space for a storage system. The data element “Document” is described by keywords “Computer:Network”. (b) A 3-Dimensional keyword space for storing computational resources, using the attributes: storage space, base bandwidth and cost.

To efficiently support range queries and queries using partial keywords and wildcards, the index space should preserve locality and be recursive so that these queries can be optimized using successive refinement and pruning. Such an index space is constructed using the Hilbert SFC as described below.

3.2. Hilbert Space-Filling Curve

A Space-Filling Curve (SFC) is a continuous mapping from a d -dimensional space to a 1-dimensional space $f: \mathbb{N}^d \rightarrow \mathbb{N}$. The d -dimensional space is viewed as a d -dimensional cube, which is mapped onto a line such that

the line passes once through each point in the volume of the cube, entering and exiting the cube only once. Using this mapping, a point in the cube can be described by its spatial coordinates, or by the length along the line, measured from one of its ends.

The construction of SFCs is recursive. The d -dimensional cube is first partitioned into n^d equal sub-cubes. An approximation to a space-filling curve is obtained by joining the centers of these sub-cubes with line segments such that each cell is joined with two adjacent cells. An example is presented in Figure 2 (a), (c). The same algorithm is used to fill each sub-cube. The curves traversing the sub-cubes are rotated and reflected such that they can be connected to form a single continuous curve that passes only once through each of the n^{2d} regions. The line that connects n^{kd} cells is called the k^{th} approximation of the SFC. Figures 2(b) and 2(d) show the second orders approximations for the curves in Figures 2(a) and 2(c) respectively.

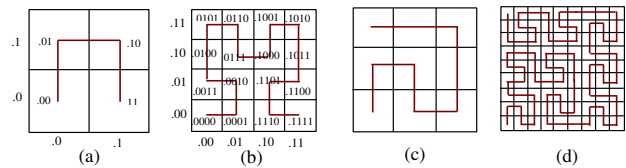


Figure 2. Space-filling curve approximations for $d = 2$: $n = 2$ (a) 1st order approximation (b) 2nd order approximation; $n = 3$ (c) 1st order approximation, (d) 2nd order approximation.

An important property of SFCs is *digital causality*, which comes directly from its recursive nature. A unit length curve constructed at the k^{th} approximation has an equal portion of its total length contained in each sub-hypercube; it has n^{k*d} equal segments. If distances across the line are expressed as base- n numbers, then the numbers that refer to all the points that are in a sub-cube and belong to a line segment are identical in their first $k*d$ digits. This property is illustrated in Figure 2 (a), (b).

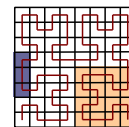


Figure 3. Clusters on a 3rd order space-filling curve ($d = 2$, $n = 2$). The colored regions represent clusters: 3-cell cluster and 16-cell cluster.

Finally, SFCs are *locality preserving*. Points that are close together in the 1-dimensional space (the curve) are mapped from points that are close together in the d -dimensional space. For example, for $k \geq 1$, $d \geq 2$, the k^{th} order approximation of a d -dimensional Hilbert space

filling curve maps the sub-cube $[0, 2^k - 1]^d$ to $[0, 2^{kd} - 1]$. The reverse property is not true, not all adjacent sub-cubes in the d -dimensional space are adjacent or even close on the curve. A group of contiguous sub-cubes in d -dimensional space will typically be mapped to a collection of segments on the SFC. These segments called clusters are shown in Figure 3.

In our system, SFCs are used to generate the 1-d index space from the multi-dimensional keyword space. Applying the Hilbert mapping to this multi-dimensional space, each data element can be mapped to a point on the SFC. Any range query or query composed of keywords, partial keywords, or wildcards, can be mapped to regions in the keyword space and the corresponding clusters in the SFC.

3.3. Mapping Indices to Peers and the Overlay Network

The next step consists of mapping the 1-dimensional index space onto an overlay network of peers. This step is similar to existing data-lookup systems. In our current implementation we use the Chord [17] overlay network topology. In Chord each node has a unique identifier ranging from 0 to $2^m - 1$. These identifiers are arranged as a circle modulo 2^m . Each node maintains information about (at most) m neighbors, called *fingers*, in a *finger table*. The i^{th} finger node is the first node that succeeds the current node by at least 2^{i-1} , where $1 \leq i \leq m$. The finger table is used for efficient routing.

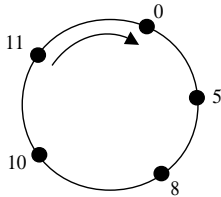


Figure 4. Example of the overlay network. Each node stores the keys that map to the segment of the curve between itself and the predecessor node.

In our implementation, node identifiers are generated randomly. Each data element is mapped, based on its SFC-based index or key, to the first node whose identifier is equal to or follows the key in the identifier space. The node is called the *successor* of the key. Consider the sample overlay network with 5 nodes and an identifier space from 0 to 16, as shown in Figure 4. In this example, data elements with keys 6, 7, and 8, will map to node 8, the successor of these keys. The management of node joins, departures, and failures is described below.

Node Joins: The joining node has to know about at least one node already in the network. It randomly chooses

an identifier from the identifier space and sends a join message with this identifier to the known node. This message is routed across the overlay network to the successor of the new node (based on the new identifier). The joining node is inserted into the overlay network at this point and takes a part of the successor node's load. The cost for joining is $O(\log_2^2 N)$ messages².

Node Departures: The finger tables of the nodes that have entries pointing to the departing node have to be updated. The cost is $O(\log_2^2 N)$ messages.

Node Failures: When a node fails, the finger tables that have entries pointing to it will be incorrect. Each node periodically runs a stabilization algorithm where it chooses a random entry in its finger table, checks for its state, and updates it if required.

Data Lookup: Nodes are efficiently located based on their content. Data lookup takes $O(\log_2 N)$ hops. In our system partial queries and range queries will typically require interrogating more than one node, as the desired information may be stored at multiple nodes in the system.

3.4. The Query Engine

The primary function of the query engine is to efficiently process user queries. As described above, data elements in the system are associated with a sequence of one or more keywords (up to d keywords, where d is the dimensionality of the keyword space). The queries can consist of a combination of keywords, partial keywords, or wildcards. The expected result of a query is the complete set of data elements that match the user's query. For example, (computer, network), (computer, net*) and (comp*, *) are all valid queries. Another query type is a range query where at least one dimension specifies a range. For example if the index encodes memory, CPU frequency and base bandwidth resources, the following query (256 - 512MB, *, 10Mbps - *) specifies a machine with memory between 256 and 512 MB, any CPU frequency and at least 10Mbps base bandwidth.

3.4.1. Query Processing

Processing a query consists of two steps: translating the keyword query to relevant clusters of the SFC-based index space, and querying the appropriate nodes in the overlay network for data-elements.

If the query consists of whole keywords (no wildcard) it will be mapped to at most one point in the index space, and the node containing the result is located using the network's look-up protocol. If the query contains partial keywords and/or wildcards or it is a range query, the query identifies a set of points (data elements) in the keyword space that correspond to a set of points (indices) in the index space. In

² N is the number of nodes in the system

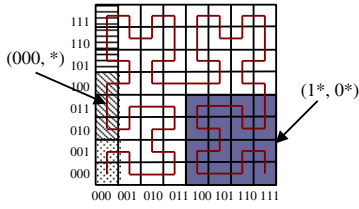


Figure 5. Regions in a 2-dimensional space defined by the queries (000, *) and (1*, 0*).

Figure 5, the query (000, *) identifies 8 data elements, the squares in the vertical region. The index (curve) enters and exits the region three times, defining three segments of the curve or clusters (marked by different patterns). Similarly the query (1*, 0*) identifies 16 data elements, defining the square region in Figure 5. The SFC enters and exits this region once, defining one cluster.

Each cluster may contain zero, one or more data elements that match the query. Depending on its size, an index space cluster may be mapped to one or more adjacent nodes in the overlay network. A node may also store more than one cluster. Once the clusters associated with a query are identified, straightforward query processing consists of sending a query message for each cluster. A query message for a cluster is routed to the appropriate node in the overlay network as follows. The overlay network provides us with a data look-up protocol: given an identifier for a data element, the node responsible for storing it is located. The same mechanism can be used to locate the node responsible for storing a cluster using a cluster identifier. The cluster identifier is constructed using SFC's digital causality property. The digital causality property guarantees that all the cells that form a cluster have the same first i digits. These i digits are called the cluster's *prefix* and form the first i digits of the m digit identifier. The rest of the identifier is padded with zero.

The node that initiated the query can not know if a cluster is stored in the network or not, or if multiple clusters are stored at the same node, to make optimizations. The number of clusters can be very high, and sending a message for each cluster is not a scalable solution. For example, consider the query (000, *) in Figure 5, but using base-26 digits and higher order approximation of the space-filling curve. The cost of sending a message for each cluster can be prohibitive.

3.4.2. Query Optimization

Query processing can be made scalable using the observation that not all the clusters that correspond to a query represent valid keywords as the keyword space and the clusters are typically sparsely populated with data elements. This is because we are using base- n numbers as coordinates along

each dimension of the space and not all the base- n numbers are valid keywords. The number of messages sent and nodes queried can be significantly reduced by filtering out the useful clusters early. However, useful clusters cannot be identified at the node where the query is initiated. The solution is to consider the recursive nature of the SFC and its digital causality property, and to distribute the process of cluster generation at multiple nodes in the system, the ones that might be responsible for storing them.

Since the SFC generation is recursive, and clusters are segments on the curve, these clusters can also be generated recursively. This recursive process can be viewed as constructing a tree. At each level of the tree the query defines a number of clusters, which are refined, resulting in more clusters for the next level. The tree can now be embedded into the overlay network: the root performs first query refinement, and each node further refines the query, sending the resulting sub-queries to the appropriate nodes in the system.

Consider the following example. We want to process the query (011, *) in a 2-dimensional space using base-2 digits for the coordinates. Figure 6 shows the successive refinement for the query and Figure 7 shows the corresponding tree. The leaves of the tree represent all possible matches for the query.

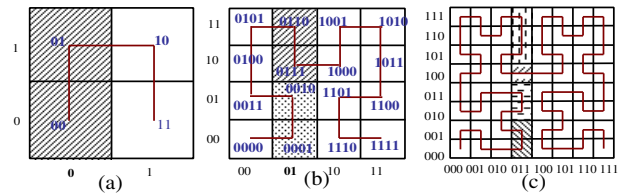


Figure 6. Recursive refinement of the query (011, *). (a) one cluster on the first order Hilbert curve, (b) two clusters on the second order Hilbert curve, (c) four clusters on the third order Hilbert curve.

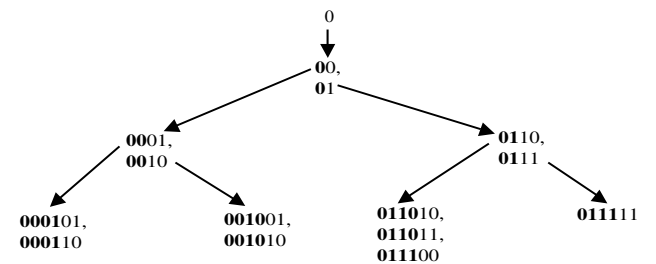


Figure 7. Recursive refinement of the query (011, *) viewed as a tree. Each node is a cluster, and the bold characters are cluster's prefixes.

The query optimization consists of pruning nodes from the tree during the construction phase. As a result of the

load-balancing steps (see Section 3.5), the nodes tend to follow the distribution of the data in the index space - i.e., a larger number of nodes are assigned to denser portions of the index space, and no nodes for the empty portions. If we embed the query tree onto the ring topology of the overlay network, we can prune away many of the nodes that do not contain valid data elements, knowing that their children do not exist in the system.

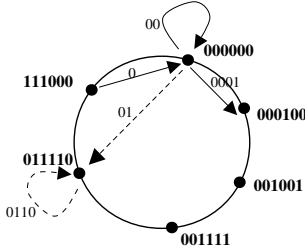


Figure 8. Embedding the leftmost tree path (solid arrows) and the rightmost path (dashed arrows) onto the overlay network topology.

Figure 8 illustrates the process, using the query in Figure 6 as an example. The leftmost path (solid arrows) and the rightmost path (dashed arrows) of the tree presented in Figure 7 are embedded into the ring network topology. The overlay network uses 6 digits for node identifiers. The arrows are labeled with the prefix of the cluster being queried. The query initiated at node 111000. The first cluster has prefix 0, so the cluster’s identifier will be 000000. The cluster is sent to node 000000. At this node the cluster is further refined, generating two sub-clusters, with prefixes 00 and 01. The cluster with prefix 00 remains at the same node. After processing, the sub-cluster 0001 is sent to node 000100. The cluster with prefix 01 and identifier 010000 is sent to node 011110 (dashed line). This cluster will not be refined because the node’s identifier is greater than the query’s identifier, and all matching data elements are stored at this node.

A second query optimization is used to reduce the number of messages involved. It is based on the observation that multiple sub-clusters of the same cluster can be mapped to the same node. To reduce the number of messages, we sort the sub-clusters in increasing order and send the first one in the network. The destination node of the sub-cluster replies with its identifier. The sending node then aggregates all the sub-clusters associated with this identifier and sends them as a single message routed to the destination.

3.5. Balancing Load

As we mentioned earlier, the original d -dimensional keyword space is sparse and data elements form clusters in this space instead of being uniformly distributed in the

space. As the Hilbert SFC-based mapping preserves keyword locality, the index space will have the same properties. Since the nodes are uniformly distributed in the node identifier space when the data elements are mapped to the nodes, the load will not be balanced. Additional load balancing has to be performed. We have defined two load-balancing steps as described below.

Load Balancing at Node Join: At join, the incoming node generates several identifiers (e.g., 5 to 10) and sends multiple join messages using these identifiers. Nodes that are logical successors of these identifiers respond reporting their load. The new node uses the identifier that will place it in the most loaded part of the network. This way, nodes will tend to follow the distribution of the data from the very beginning. With n identifiers being generated, the cost to find the best successor is $O(n \cdot \log_2 N)$, and the cost to update the tables remains the same, $O(\log_2^2 N)$. However, this step is not sufficient by itself. The runtime load-balancing algorithms presented below improve load distribution.

Load Balancing at Runtime: The runtime load-balancing step consists of periodically running a local load-balancing algorithm between few neighboring nodes. We propose two load-balancing algorithms. In the first algorithm, neighboring nodes exchange information about their loads, and the most loaded nodes give a part of their load to their neighbors. The cost of load-balancing at each node using this algorithm is $O(\log_2^2 N)$. As this is expensive, this load-balancing algorithm cannot be run very often.

The second load-balancing algorithm uses virtual nodes. In this algorithm, each physical node houses multiple virtual nodes. The load at a physical node is the sum of the load of its virtual nodes. When the load on a virtual node goes above a threshold, the virtual node is split into two or more virtual nodes. If the physical node is overloaded, one or more of its virtual nodes can migrate to less loaded physical nodes (neighbors or fingers). An evaluation of the load balancing algorithms is presented in Section 4.

4. Experimental Evaluation

The performance of our P2P information discovery system is evaluated using a simulator. The simulator implements the SFC-based mapping, the Chord-based overlay network, the load-balancing steps, and the query engine with the query optimizations described above. As the overlay network configuration and operations are based on Chord [17], its maintenance costs are of the same order as in Chord. An evaluation of the query engine and the load-balancing algorithms is presented below.

4.1. Evaluating the Query Engine

The overlay network used to evaluate the query engine consists of 1000 to 5400 nodes. 2-dimensional (2D), and a

3-dimensional (3D) keyword spaces are used in this evaluation. Finally, we use up to 10^6 keys (unique keyword combinations) in the system, each of which could be associated with one or more data elements. We measure the following:

Number of routing nodes: the nodes that route the query. Some of them also process the query.

Number of processing nodes: the nodes that actually process the query, refine it, and search for matches. The goal is to restrict processing only to those nodes that store matching data elements.

Number of data nodes: the nodes that have data elements matching the query.

Number of messages required to resolve a query. When using the query optimization each message is a sub-query that searches for a fraction of the clusters associated with the original query.

The types of queries used in the experiments are:

Q1: Queries with one keyword or partial keyword, e.g. (computer, *) for 2D, (comp*, *, *) for 3D.

Q2: Queries with two to three keywords or partial keywords (at least one partial keyword), e.g. (comp*, net*) for 2D, (computer, network, *) for 3D.

Q3: Range queries.

4.1.1. Evaluating a 2-dimensional keyword space

This experiment represents a typical P2P storage system where the number of keys and data elements in the system increases as the number of nodes increases. The system size increases from 1000 nodes to 5400 nodes, and the number of stored keys increases from $2 \cdot 10^5$ to 10^6 .

The results for experiments using six different type Q1 queries (query1 - query6) are plotted in Figure 9. Each query resulted in a different number of matches.

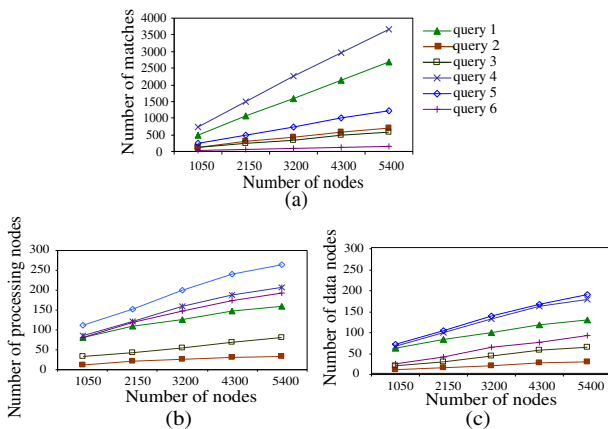


Figure 9. Results for query type Q1, 2D: (a) the number of matches for the queries, (b) the number of nodes that process the query, (c) the number of nodes that found matches for the query.

As seen in Figures 9 (b) and (c), the number of processing and data nodes is a fraction of the total nodes and increase at a slower rate than the system size. The number of processing nodes does not necessarily depend on the number of matches. For example, to solve a query with 160 matches (query6) can be more costly than solving a query with 2600 matches (query1). This is due to the recursive processing of queries and the distribution of keys in the index space. In order to optimize the query, we prune parts of the query tree based on the data stored in the system. The earlier the tree is pruned, the fewer processing nodes will be required and the better the performance will be. For example, if the query being processed is (computer, *) and the system contains a large number of data elements with keys that start with “com” (e.g., company, commerce, etc.) but do not match the query, the pruning will be less efficient and will result in a larger number of processing nodes.

Note that even under these conditions, the results are good. A keyword search system like Gnutella [21] would have to query the entire network using some form of flooding to guarantee that all the matches to a query are returned, while in the case of a data lookup system such as Chord [17], one would have to know all the matches a priori and look them up individually.

Figure 9 (c) plots the number of data nodes which are a subset of the processing nodes. The number of data nodes is close to the number of processing nodes, indicating that the query optimizations effectively reduce the number of nodes involved. Furthermore, comparing the number of matches to the number of data nodes demonstrates the clustering property of the Hilbert SFC index space, which can be defined as a ratio of the number of matches for a query to the number of data nodes that store these matches.

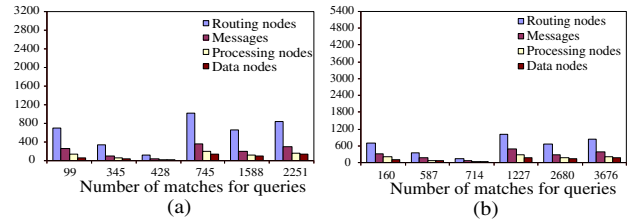


Figure 10. Results for all the metrics, 2D: (a) for a 3200 node system and $6 \cdot 10^5$ keys, (b) for a 5400 node system and 10^6 keys.

Figure 10 plots all the measurements for two system sizes and demonstrates the scalability of the presented system. Each group of bars represents the results for a particular query. The maximum value, plotted on the vertical axis, is the size of the network. As seen in the graph, the processing nodes are a small fraction of the routing nodes, and a very small fraction of the entire system. The processing node population does not increase as fast as the system. The

number of data nodes is close to the number of processing nodes. The number of messages used is almost twice the number of processing nodes, which is as expected.

Figure 11 shows the corresponding results for experiments with 5 different type Q2 queries. As expected, the results are significantly better than those for type Q1 queries. This is because query optimization and pruning are effective when both keywords are at least partially known.

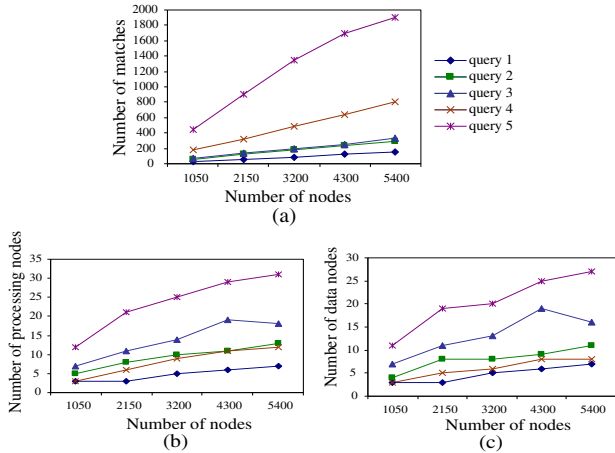


Figure 11. Results for query type Q2, 2D: (a) the number of matches for the queries, (2) the number of data nodes.

4.1.2. Evaluating a 3-dimensional keyword space

The experiment conducted for the 2D keyword space was repeated for a 3D keyword space. Overall, the results are similar to the 2D case. The growth of the system, either in the number of nodes or in the quantity of data stored, or both, affects the behavior of the queries in the same way. The only difference is the magnitude of the results. As described in Section 1, documents that share a specific keyword will typically be mapped to disjoint fragments on the curve (clusters). In the 3D case the number of such fragments is larger than in the 2D case - 3 keywords used instead of 2 and result in a "longer" curve. Consequently, the results obtained for the 3D case for all the metrics have the same pattern as the 2D case but a larger magnitude.

The results for experiments in the 3D case follow a similar pattern to the 2D cases presented in Section 4.1.1. The main difference from the 2D cases is the magnitude of the graphs; results for the 3D case may be larger by two to three times. This is expected and is due to the fact that for the same types of queries there are more clusters in the 3D case than in the 2D case.

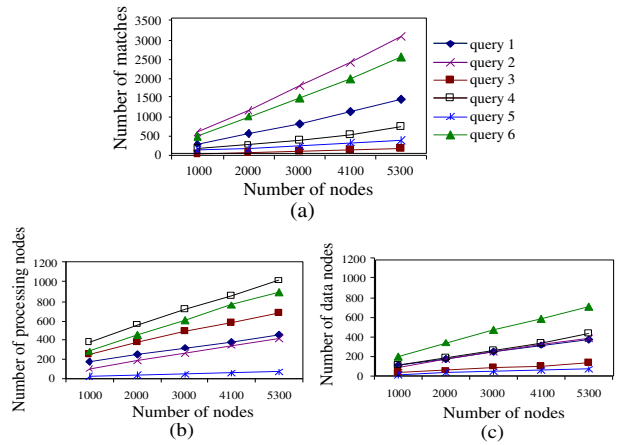


Figure 12. Results for query type Q1, 3D: (a) the number of matches for the queries, (b) the number of nodes that process the query, (c) the number of nodes that found matches for the queries.

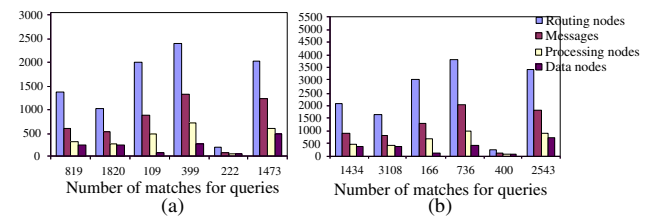


Figure 13. Results for all metrics, 3D: (a) for 3000 node system and $6 \cdot 10^5$ keys, (b) for 5300 node system and 10^6 keys.

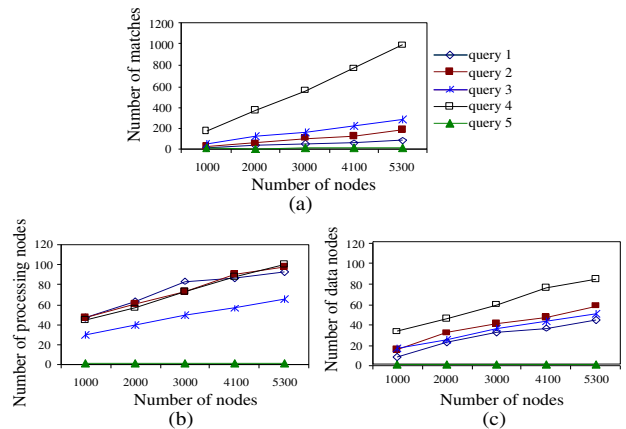


Figure 14. Results for query type Q2, 3D: (a) the number of matches for queries, (b) the number of processing nodes, (c) the number of data nodes.

4.1.3 Evaluating the range queries in a 3-dimensional keyword space

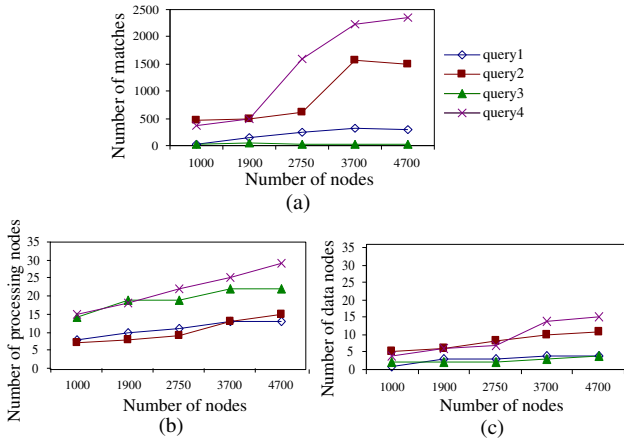


Figure 15. Results for query type Q3 (range query), of the form: (keyword, range, *). (a) the number of matches for queries, (b) the number of processing nodes, (c) the number of data nodes.

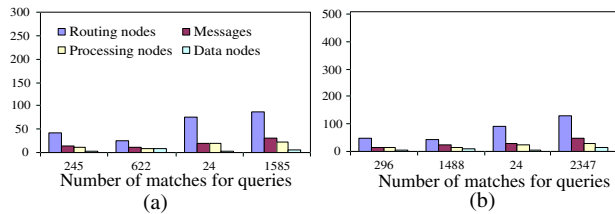


Figure 16. Results for all metrics for range queries: (a) for 2750 node system and 6×10^5 keys, (b) for 4700 node system and 10^6 keys.

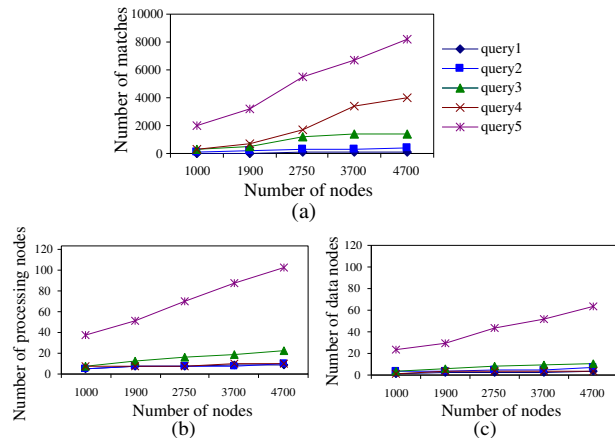


Figure 17. Results for query type Q3 (range query), of the form: (range, range, range). (a) the number of matches for queries, (b) the number of processing nodes, (c) the number of data nodes.

The following types of range queries were evaluated: (keyword, range, *) and (range, range, range). The experiments show that the results do not depend on the size of the range (because the index space is not uniformly populated), but more on the number of matches found and the distribution of the data.

4.2. Evaluating the Load Balancing Algorithms

The quality of the load balance achieved by the two load balancing operations is evaluated for the initial distribution shown in Figure 18. As expected, the original distribution is not uniform. The load-balance step at node join helps to match the distribution of nodes with the distribution of data. The resulting load balance is plotted in Figure 19 (a). While the resulting load distribution is better than the original distribution in Figure 18, this step by itself does not guarantee a very good load balance. However, when it is used in conjunction with the runtime load-balancing step, the resulting load balance improves significantly as seen in Figure 19 (b). The load is almost evenly distributed in this case.

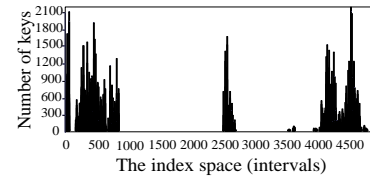


Figure 18. The distribution of the keys in the index space. The index space was partitioned into 5000 intervals. The Y-axis represents the number of keys per interval.

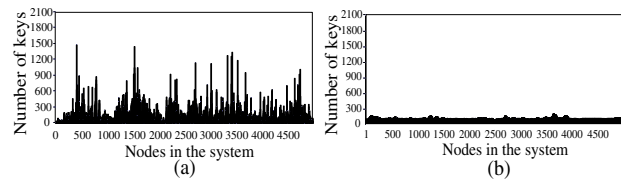


Figure 19. The distribution of the keys at nodes (a) when using only the load balancing at node join technique, (b) when using both the load balancing at node join technique, and the local load balancing.

5. Conclusions and Future Work

In this paper, we presented the design and evaluation of a P2P information discovery system that supports complex searches using keyword, partial keywords, wildcards

and range queries, while guaranteeing that all data elements in the system that match the query are found with bounded costs in terms of the number of messages and nodes involved in the query. The presented information discovery system essentially implements an Internet-scale distributed hash table. It can be used to complement existing resource discovery mechanisms in computational grids by enhancing them with range queries, by storage and archival systems to support keyword searches, and by bulletin-board systems to support discovery by interest profiles. A key component is a locality-preserving mapping from data element keyword space to the index space that is used to assign the data elements to peers. This mapping is based on recursive, self-similar SFCs. The recursive structure of the index space is used to optimize query processing and to reduce the number of nodes queried. As the mapping preserves keyword locality, data-elements may not be uniformly distributed in the index space. Dynamic load balancing schemes (applicable at node join and at runtime) were presented.

An experimental evaluation of the presented P2P storage system using a simulator was presented. The experiments demonstrated the scalability of the system, and showed that only a fraction of the total nodes in the system typically process a query and this fraction is almost the same as the nodes that store data elements matching the query. The results also demonstrated the ability of the mapping to preserve keyword locality and the effectiveness of the load balancing algorithms. We are extending this research to evaluate other network topologies and mappings, and to address issues such as hot-spots, fault-tolerance, security and resistance to attacks, and maintenance of geographical locality in the overlay network.

References

- [1] A. Andrzejak and Z. Xu. Scalable, efficient range queries for grid information services. In *Proceedings of the Second IEEE International Conference on Peer-to-Peer Computing (P2P2002)*, Sweden, Sept. 2002.
- [2] T. Bially. *A class of dimension changing mapping and its application to bandwidth compression*. PhD thesis, Polytechnic Institute of Brooklyn, June 1967.
- [3] I. Clarke, O. Sandberg, B. Wiley, and T. Hong. Freenet: A distributed anonymous information storage and retrieval system. In *Proceedings of the ICSI Workshop on Design Issues in Anonymity and Unobservability*, Berkeley, California, June 2000.
- [4] A. Crespo and H. Garcia-Molina. Routing indices for peer-to-peer systems. In *Proceedings of the 22nd International Conference on Distributed Computing Systems (ICDCS)*, Vienna, Austria, July 2002.
- [5] I. Foster and A. Iamnitchi. On death, taxes, and the convergence of peer-to-peer and grid computing. In *2nd International Workshop on Peer-to-Peer Systems (IPTPS'03)*, Berkeley, CA, Feb. 2003.
- [6] I. Foster and C. Kesselman. *The GRID - Blueprint for a New Computing Infrastructure*. Morgan Kaufmann Publishers, Inc., San Francisco, CA, 1999.
- [7] O. D. Gnawali. A keyword-set system for peer-to-peer networks. Master's thesis, MIT, June 2002.
- [8] A. Iamnitchi, I. Foster, and D. Nurmi. A peer-to-peer approach to resource discovery in grid environments. Technical Report TR-2002-06, University of Chicago, 2002.
- [9] D. Karger, E. Lehman, T. Leighton, M. Levine, D. Lewin, and R. Panigrahy. Consistent hashing and random trees: Distributed caching protocols for relieving hot spots on the world wide web, 1997.
- [10] W. Li, Z. Xu, F. Dong, and J. Zhang. Grid resource discovery based on a routing-transferring model. In *Proceedings of the 3rd International Workshop on Grid Computing*, pages 145–156, Baltimore, MD, 2002.
- [11] Q. Lv, P. Cao, E. Cohen, K. Li, and S. Shenker. Search and replication in unstructured peer-to-peer networks. In *Proceedings of the 16th international conference on Supercomputing*, June 2002.
- [12] B. Moon, H. Jagadish, C. Faloutsos, and J. Saltz. Analysis of the clustering properties of hilbert space-filling curve. Submitted to *IEEE Transactions on Knowledge and Data Engineering*, Mar. 1996.
- [13] C. Plaxton, R. Rajaraman, and A. W. Richa. Accessing nearby copies of replicated objects in a distributed environment. In *Proceedings of the ACM SPAA*, pages 311–320, Newport, Rhode Island, June 1997.
- [14] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A scalable content-addressable network. In *Proceedings of ACM SIGCOMM*, 2001.
- [15] A. Rowstron and P. Druschel. Pastry: Scalable, distributed object location and routing for largescale peer-to-peer systems. *IFIP/ACM International Conference on Distributed Systems Platforms (Middleware)*, Nov. 2001.
- [16] H. Sagan. *Space-Filling Curves*. Springer-Verlag, 1994.
- [17] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup service for internet applications. In *Proceedings of ACM SIGCOMM*, 2001.
- [18] C. Tang, Z. Xu, and M. Mahalingam. Peersearch: Efficient information retrieval in peer-to-peer networks. Technical Report HPL-2002-198, HP Labs, 2002.
- [19] B. Yang and H. Garcia-Molina. Improving search in peer-to-peer systems. In *In Proceedings of the 22nd International Conference on Distributed Computing Systems (ICDCS)*, Vienna, Austria, July 2002.
- [20] B. Y. Zhao, J. Kubiawicz, and A. D. Joseph. Tapestry: An infrastructure for fault-tolerant widearea location and routing. Technical Report UCB/CSD-01-1141, Computer Science Division, University of California at Berkeley, Apr. 2001.
- [21] Gnutella webpage: <http://gnutella.wego.com/>.
- [22] Napster webpage: <http://napster.com/>.