

# Mercury: Supporting Scalable Multi-Attribute Range Queries\*

Ashwin R. Bharambe  
ashu+@cs.cmu.edu

Mukesh Agrawal  
mukesh+@cs.cmu.edu

Srinivasan Seshan  
srini+@cs.cmu.edu

Computer Science Department  
Carnegie Mellon University  
Pittsburgh, PA 15213

## ABSTRACT

This paper presents the design of Mercury, a scalable protocol for supporting multi-attribute range-based searches. Mercury differs from previous range-based query systems in that it supports *multiple attributes* as well as performs *explicit load balancing*. To guarantee efficient routing and load balancing, Mercury uses novel light-weight sampling mechanisms for uniformly sampling random nodes in a highly dynamic overlay network. Our evaluation shows that Mercury is able to achieve its goals of logarithmic-hop routing and near-uniform load balancing.

We also show that Mercury can be used to solve a key problem for an important class of distributed applications: distributed state maintenance for distributed games. We show that the Mercury-based solution is easy to use, and that it reduces the game's messaging overhead significantly compared to a naïve approach.

## Categories and Subject Descriptors

E.1 [Data Structures]: Distributed data structures

## General Terms

Algorithms, Measurement

## Keywords

Range queries, peer-to-peer systems, distributed hash tables, load balancing, random sampling

## 1. INTRODUCTION

Much recent work on building scalable peer-to-peer (P2P) networks has concentrated on Distributed Hash Tables or

\*This work was supported in part by a grant from the Pittsburgh Digital Greenhouse.

DHTs [22, 23, 25]. DHTs offer a number of scalability advantages over previous P2P systems (e.g., Napster, Gnutella, etc.) including load balancing and logarithmic hop routing with small local state. However, the hash table or “exact match” interface offered by DHTs, although fruitfully used by some systems [5–7], is not flexible enough for many applications. For example, it is unclear how DHTs could be modified to regain the highly desirable flexibility offered by keyword-based lookups of file-sharing applications.

The main contribution of this paper is the design and evaluation of Mercury, a scalable routing protocol for supporting *multi-attribute range queries*. In our model, each query is a *conjunction of ranges* in one or more attributes. The attributes not present in the query are assumed to be wildcards. We believe that range queries significantly enhance search flexibility in a number of scenarios. In addition to being useful for answering user queries, we find that range-based queries can also be useful in the construction of distributed applications.

A number of recent systems [11, 12, 16] have proposed distributed protocols which support range-based queries. Mercury mainly differs from these systems since it supports *multi-attribute range-based queries* and explicit load balancing.

There are two main components of Mercury's design. First, Mercury handles multi-attribute queries by creating a *routing hub* for each attribute in the application schema. Each routing hub is a logical collection of nodes in the system. Queries are passed to exactly one of the hubs corresponding to the attributes that are queried, while a new data item is sent to all hubs for which it has an associated attribute. This ensures that queries retrieve all relevant data items present in the system.

Second, for supporting range queries, Mercury organizes each routing hub into a circular overlay of nodes and places data *contiguously* on this ring, i.e., each node is responsible for a *range* of values for the particular attribute. While the notion of a circular overlay is similar in spirit to some existing DHT designs, due to our choice to support range queries by placing data contiguously, we cannot use randomizing hash functions for placing data. This requirement introduces a fundamental challenge: because Mercury cannot use hash functions, data partitioning among nodes can become non-uniform (as we explain in Section 3.2), thus, requiring an explicit load-balancing mechanism. However, the load-balancing mechanism is fundamentally incompati-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGCOMM'04, Aug. 30–Sept. 3, 2004, Portland, Oregon, USA.

Copyright 2004 ACM 1-58113-862-8/04/0008 ...\$5.00.

ble with many of the techniques that DHTs use to guarantee routing efficiency.

The solution to the above challenges forms the core contribution of this paper. Some of the interesting algorithms in Mercury include:

- A message routing algorithm that supports range-based lookups within each routing hub in  $\mathcal{O}(\log^2 n/k)$  hops when each node maintains  $k$  links to other nodes.
- A low-overhead random sampling algorithm that allows each node to create an estimate of system-wide metrics such as load distribution.
- A load-balancing algorithm (which exploits the random sampling algorithm) that ensures that routing load is uniformly distributed across all participating nodes.
- An algorithm for reducing query flooding by estimating how selective each of the predicates in a query is, based on past database insertions.

In addition to describing the design of Mercury, we also explore how the added flexibility provided by range query lookups can simplify the state management task for distributed games. We show that, using a range query based publish-subscribe [2, 3] system built on top of Mercury, we can provide efficient and scalable dissemination of distributed state updates. We believe that a *range query* significantly enhances the application’s ability to precisely express its interests.

The remainder of the paper is organized as follows. In the next section, we compare Mercury to related work in this area. Section 3 details the basic Mercury protocol for routing data-records and queries. Section 4 presents enhancements which improve the performance of the basic protocol. In Section 5, we evaluate the scalability and load-balancing properties of the Mercury system. In Section 6, we present the design of our publish-subscribe system, and proof-of-concept distributed game. Finally, Section 7 concludes.

## 2. RELATED WORK

In this section, we compare and contrast our approach to implementing range queries with that of related systems. Our discussion focuses on two fundamental questions:

- ▷ Can we use *existing DHTs as building blocks* to implement range query predicates?
- ▷ How is our design different from other systems like SkipNet [11], etc., which also provide rich query support?

### *Using existing DHTs for range queries*

A large number of distributed hash table designs [11, 22, 23, 25] have been proposed over the past few years. They provide a hash table interface to the application, viz., **insert(key, value)** and **lookup(key)** primitives. Recent research [4, 10] has shown that, in addition to the basic scalable routing mechanism, DHTs offer much promise in terms of load balancing, proximity-based routing, static resilience, etc. Hence, it is a natural question to ask if we could implement range queries using just the **insert** and **lookup** abstractions provided by DHTs.

Our analysis, based on analyzing possible strawman designs using DHTs, indicates that the abstractions provided by a DHT are not enough for implementing range queries. Fundamental to our argument is the fact that all existing DHT designs use randomizing hash functions for inserting and looking up keys in the hash table. While hashing is crucial for DHTs in order to get good load balancing properties, it is also the main barrier in using a DHT for implementing range queries. This is because the hash of a range is not correlated to the hash of the values within a range. Hence, it is necessary to create some artificial correlation between ranges and values which is invariant under hashing.

One natural way to achieve this is to *partition the value space into buckets* and map values and ranges to one or more buckets. A bucket forms the lookup *key* for the hash table. Then, a range query can be satisfied by simply performing lookups on the corresponding bucket(s) using the underlying DHT. However, this scheme has several fundamental drawbacks. It requires the application to perform the partitioning of space *a priori*. This can be very difficult or even impossible for many applications, e.g., partitioning of file names. Moreover, load balancing and query performance is highly dependent on the way partitioning is performed. For example, if the number of buckets is too small, i.e., the partition is too coarse, queries will get mapped to a smaller set of nodes creating load imbalance. Increasing the number of buckets, on the other hand, will increase the routing hops required to answer a range query.

This indicates that while a DHT-based scheme may not be an impossibility, its implementation is likely to be *awkward* and complicated. By avoiding randomizing hash functions, Mercury seeks to remove this difficulty. At the same time, we also note that the design of Mercury is inspired from and similar in many respects to existing DHT designs. Hence, we believe that it can easily build upon recent advances in proximity-based routing and achieving resilience in DHTs [10].

### *Comparison with systems supporting rich queries*

In this section, we compare Mercury against recent systems which offer rich query support. These include SkipNet [11], PIER [12] and DIM [16].

The SkipNet DHT organizes peers and data objects according to their lexicographic addresses in the form of a variant of a probabilistic skip list. It supports logarithmic time range-based lookups and guarantees path locality. Mercury is more general than SkipNet since it supports range-based lookups on *multiple-attributes*. Our use of random sampling to estimate query selectivity constitutes a novel contribution towards implementing scalable multi-dimensional range queries. Load balancing is another important way in which Mercury differs from SkipNet. While SkipNet incorporates a constrained load-balancing mechanism, it is only useful when part of a data name is hashed, in which case the part is inaccessible for performing a range query. This implies that SkipNet supports load-balancing *or* range queries – not both.

One might argue that the query-load imbalance in SkipNet can be corrected by using *virtual servers* as suggested in [21]. However, it is unlikely to help in this regard for the following reason: for effective load-balancing, the number of virtual servers needed must be proportional to the *skew* (ratio of *max* to *min*) in the load. The scheme pro-

posed in [21] assumes that the load skew results from the standard deviation of random hash function distributions, which is typically very small ( $\mathcal{O}(\log n)$ , see [1]). However, in our case, the skew results from differences in *query workload*, which can grow quite large. Hence, larger number of virtual servers would be required increasing routing hops by about  $\log(s)$  where  $s$  is the skew. Moreover, the scheme would fare even worse for range lookups since it would increase the number of distinct nodes accessed for processing the query by  $\mathcal{O}(s)$ .

The DIM data structure [16] supports multi-dimensional range queries by embedding them into a two-dimensional space and using a geographic routing algorithm. However, the routing cost scales only as  $\mathcal{O}(\sqrt{n})$ , which while reasonable in a medium-sized sensor network, is quite expensive for larger scales. Furthermore, the “volume expansion” that occurs while projecting from higher dimensions onto two-dimensions can be quite large resulting in more flooding of the query. Also, queries containing wildcards in certain attributes get flooded to all nodes. On the other hand, Mercury, like most databases, uses query selectivity mechanisms to route through the attribute hubs of highest selectivity thereby significantly reducing flooding for queries containing only a few attributes.

The PIER system is a distributed query engine based on DHTs. While PIER provides support for some types of rich queries, such as equi-joins, it is unclear how to extend PIER to support efficient resolution of queries containing range-based predicates.

All the above systems and recent work on balancing load in such systems [1, 21] treat load on a node as being proportional to the *range of identifier values* the node is responsible for. In other words, they assume a uniform data distribution which is sensible for DHTs which use cryptographic hashes. Mercury, on the other hand, defines load on a node as the number of messages routed or matched per unit time, and supports explicit and flexible load balancing.

We note that Mercury uses a *leave-join* style load balancing algorithm that is similar to [21]. Karger and Ruhl [13] as well as Ganeshan et. al. [8] have concurrently proposed similar leave-join based load balancing mechanisms. The protocol described in [13] is randomized and requires communication with  $\log n$  random nodes in the system. Ganeshan et. al. propose a deterministic protocol which requires the maintenance of skip graph overlays. In the face of skewed node range distributions (which any load balancing protocol will yield), these tasks are far from trivial, and a significant part of the Mercury protocol is aimed at addressing this difficulty. In general, many approaches to diffusion-based dynamic load balancing [9] require information about the dynamic distributed network that is hard to obtain, in order to make informed decisions. Instead, Mercury uses light-weight sampling mechanisms to track load distribution within the overlay.

### 3. MERCURY ROUTING

In this section, we provide an overview of the basic Mercury routing protocol. Although the protocol as described in this section is similar in spirit to existing structured overlay protocols (e.g., Chord, Symphony, etc.), we present it for completeness. At the same time, it is important to note that it works well in only a limited set of conditions. In Sec-

tion 4, we significantly extend this basic routing protocol to work for a wider set of operating points.

#### 3.1 Data Model

In Mercury, a data item is represented as a list of typed attribute-value pairs, very similar to a record in a relational database. Each field is a tuple of the form: (**type**, **attribute**, **value**). The following types are recognized: **int**, **char**, **float** and **string**.<sup>1</sup>

A query is a *conjunction* of predicates which are tuples of the form: (**type**, **attribute**, **operator**, **value**). A disjunction is implemented by multiple distinct queries. Mercury supports the following operators:  $<$ ,  $>$ ,  $\leq$ ,  $\geq$  and  $=$ . For the **string** type, Mercury also permits *prefix* and *post-fix* operators. Figure 1 presents an example.

<i>float</i>	<i>x-coord</i> = 50	<i>float</i>	<i>x-coord</i> < 53
<i>float</i>	<i>y-coord</i> = 100	<i>float</i>	<i>x-coord</i> > 34
<i>string</i>	<i>player</i> = "john"	<i>string</i>	<i>player</i> = "j*"
<i>string</i>	<i>team</i> = "topgunz"	<i>int</i>	<i>score</i> = "*"
<i>int</i>	<i>score</i> = 76		
Data		Query	

Figure 1: Example of a data item and a query as represented in the Mercury system.

#### 3.2 Routing Overview

Mercury supports queries over multiple attributes by partitioning the nodes in the system into groups called *attribute hubs*. This partition is only logical, i.e., a physical node can be part of multiple logical hubs. Each of the attribute hubs is responsible for a specific attribute in the overall schema. Hubs can be thought of as orthogonal dimensions of a multi-dimensional attribute space. The first routing hop determines which dimension to route through. The rest of the routing is unidimensional and is based on the values of a single attribute of the data item. We note that this mechanism does not scale very well as the number of attributes in the schema increase and is suitable only for applications with moderate-sized schemas.

To simplify the description, we will use the following notation: let  $\mathcal{A}$  denote the set of attributes in the overall schema of the application.  $\mathcal{A}_Q$  denotes the set of attributes in a query  $Q$ . Similarly, the set of attributes present in a data-record  $D$  is denoted by  $\mathcal{A}_D$ . We use the function  $\pi_a$  to denote the value (range) of a particular attribute  $a$  in a data-record (query). We will denote the attribute hub for an attribute  $a$  by  $H_a$ .

Nodes within a hub  $H_a$  are arranged into a circular overlay with each node responsible for a *contiguous range*  $r_a$  of attribute values. A node responsible for the range  $r_a$  *resolves* all queries  $Q$  for which  $\pi_a(Q) \cap r_a \neq \phi (= \{\})$ , and it *stores* all data-records  $D$  for which  $\pi_a(D) \in r_a$ . Ranges are assigned to nodes during the join process.

##### *Routing Queries and Data-Records*

Queries are passed to exactly one of the hubs corresponding to the attributes that are queried. In other words, a

<sup>1</sup>Our basic data types are sortable, enabling us to define numeric operations (addition and subtraction) on them. Care needs to be taken when handling **string** attributes.

query  $Q$  is delivered to  $H_a$ , where  $a$  is *any* attribute chosen from  $\mathcal{A}_Q$ . We will see in Sections 4.3 and 5.4 that although choosing any attribute hub suffices for matching correctness, substantial savings in network bandwidth can be achieved if the choice is done more intelligently using query selectivity. Within the chosen hub, the query is delivered and processed at all nodes that could potentially have matching values.

To guarantee that queries locate all the relevant data-records, a data-record  $D$ , when inserted, is sent to *all*  $H_b$  where  $b \in \mathcal{A}_D$ . This is necessary because the set of queries which could match  $D$  can arrive in any of these attribute hubs. Within each hub, the data-record is routed to the node responsible for the record's value for the hub's attribute.

Notice that we could have ensured correctness by sending a data-record to a single hub in  $\mathcal{A}_D$  and queries to all hubs in  $\mathcal{A}_Q$ . At first glance, this might appear to be a better choice since data-records could be much bigger in size than queries and replicating them might be more expensive. However, recall that a query can get routed to multiple locations within each hub depending on its selectivity. This, combined with the fact that many queries may be extremely non-selective in some attribute (thereby, flooding a particular hub), led us to choose a design with data-records broadcast to all hubs.

Concerning the cost of replicating data, we note that it is not necessary to replicate entire data records across hubs. In cases where doing so would be prohibitively expensive, a node within one of the hubs can hold the data record while the other hubs can hold a *pointer* to the node. This would reduce the storage requirements, at the cost of one additional hop during query resolution.

The above arguments about query workloads and replication cost notwithstanding, Mercury could be easily modified to support situations where replicating queries is a more appropriate design choice.

Within a hub  $H_a$ , routing is done as follows: for routing a data-record  $D$ , we route to the *value*  $\pi_a(D)$ . For a query  $Q$ ,  $\pi_a(Q)$  is a range. Hence, for routing queries, we route to the *first* value appearing in the range and then use the contiguity of range values to spread the query along the circle, as needed.

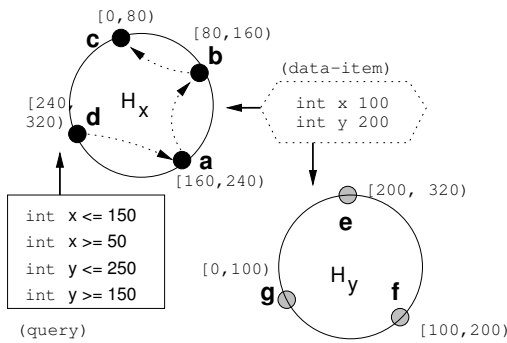


Figure 2: Routing of data-records and queries.

Fig 2 illustrates the routing of queries and data-records. It depicts two hubs  $H_x$  and  $H_y$  which may correspond to, for example, X and Y coordinates of objects. The minimum and maximum values for the  $x$  and  $y$  attributes are 0 and 320 respectively. Accordingly, the ranges are distributed to various nodes. The data-record is sent to both  $H_x$  and  $H_y$ , where it is stored at nodes  $b$  and  $e$ , respectively. The query

enters  $H_x$  at node  $d$  and is routed (and processed) at nodes  $b$  and  $c$ .<sup>2</sup>

This routing places one additional requirement on the connectivity of each node. In addition to having a link to the predecessor and successor within its own hub, each node must also maintain a link to each of the other hubs. We expect the number of hubs for a particular system to remain low, and, therefore, do not expect this to be a significant burden. We discuss the maintenance of these links later in Section 3.4.

### Design Rationale

In this section, we discuss some of the promising alternative designs for implementing a distributed multi-attribute range-based search and comment qualitatively on the trade-offs involved.

Many DHTs [18, 25] use a cryptographic hash or random value to give IDs to nodes and data stored in the DHT. However, Mercury does not use any such cryptographic hashes or random values. This simpler mapping of data and nodes in the system allows the lookup of range predicates in queries to a collection of contiguous nodes in a hub. We note that one of the main purposes of using a cryptographic hash in existing DHTs is to assign data to nodes uniformly and randomly.<sup>3</sup> The elimination of this randomness makes load-balancing in Mercury a concern. Since there are likely to be particular ranges of an attribute that are more popular for queries and data-records, nodes responsible for these ranges from will be unfairly overloaded with both routing and computation tasks. Mercury performs explicit load balancing (see Section 4.4) by moving around nodes and changing their responsibilities according to the loads. This enables the combination of good load-balancing with support for range predicates. However, one important side effect is that the distribution of range sizes is no longer guaranteed to be uniform.

With the removal of cryptographic hashes, we could have used a variety of different DHTs as the basis for our design. Our design treats the different attributes in an application schema *independently*, i.e., routing a data item  $D$  within a hub for attribute  $a$  is accomplished using only  $\pi_a(D)$ . An alternate design would be to route using the values of all attributes present in  $D$ , e.g., treating each attribute as a CAN dimension [22]. Since each node in such a design is responsible for a value-range of *every* attribute, a query that contains a *wild-card* attribute can get flooded to all nodes. We could have merged dimensions like in the DIM data structure [16] but this would still have had similar problems for queries covering large areas. By making the attributes independent, we restrict such flooding to at most one attribute hub. Furthermore, it is quite likely that some other attribute of the query is more selective and by routing the query to that hub, we can eliminate flooding altogether.

### 3.3 Constructing Efficient Routes

Recall that most of the routing in Mercury occurs within an attribute hub (only the first hop crosses hubs.) Thus, it

<sup>2</sup>This example uses flooding to route from  $d$  to  $b$ . Sec. 3.3 introduces long-distance links, which provide a more efficient routing mechanism.

<sup>3</sup>Self-certifying names/security, and robustness to correlated failures are additional valuable properties.

is essential that the overlay structure for each attribute hub be scalable and efficient.

Simply using successor or predecessor pointers can result in  $\theta(n)$  routing delays for routing data-records and queries. Like Symphony [18], the key to Mercury’s route optimization is the selection of  $k$  *long-distance* links that are maintained in addition to the successor and predecessor links. As a result, each node has a routing table of size  $k + 2$  including its neighbors along the circle.  $k$  is a configurable parameter here and could be different for different nodes.

The routing algorithm is simple: Let neighbor  $n_i$  be in charge of the range  $[l_i, r_i)$ , and let  $d$  denote the *clockwise distance* or *value-distance* between two nodes. When a node is asked to route a value  $v$ , it chooses the neighbor  $n_i$  which minimizes  $d(l_i, v)$ . Let  $m_a$  and  $M_a$  be the minimum and maximum values for attribute  $a$ , respectively.

Then,

$$d(a, b) = \begin{cases} b - a & \text{if } a \leq b, \\ (M_a - m_a) + (b - a) & \text{if } a > b \end{cases}$$

A node  $n$  whose value range is  $[l, r)$  constructs its long-distance links in the following fashion: Let  $I$  denote the unit interval  $[0, 1]$ . For each link, a node draws a number  $x \in I$  using the *harmonic* probability distribution function:  $p_n(x) = 1/(n \log x)$  if  $x \in [\frac{1}{n}, 1]$ . It contacts a node  $n'$  (using the routing protocol itself) which manages the value  $r + (M_a - m_a)x$  (wrapped around) in its hub. Finally, it attempts to make  $n'$  its neighbor. As a practical consideration, we set a fan-in limit of  $2k$  links per node. We will refer to a network constructed according to the above algorithm as a **ValueLink** network.

Under the assumption that node ranges are uniform, we can prove (see [18]) that the expected number of routing hops for routing to any value within a hub is  $\mathcal{O}(\frac{1}{k} \log^2 n)$ . Since inter-hub routing can take at most one hop, the number of hops taken for routing is at most  $\mathcal{O}(\frac{1}{k} \log^2 n)$  as well.<sup>4</sup> This guarantee is based upon Kleinberg’s analysis of small-world networks [14].

Unfortunately, the “uniform node ranges” assumption can be easily violated for many reasons. For example, explicit load-balancing would cause nodes to cluster closely in parts of the ring which are popular. In the Section 4, we present a novel distributed histogram maintenance scheme based on light-weight random sampling to provide efficient routing even with highly non-uniform ranges.

## Caching

For many applications, there can be significant locality in the generated data-items (incremental updates, for example) as well as queries (popular searches, for example.) Mercury provides hooks for the application so that it can insert its own specific caching behavior into the protocol. Essentially, Mercury allows an application to specify additional long-distance links that represent cached destinations as an addendum to the routing table. When looking for the neighbor closest to the destination, Mercury also considers nodes present in the cache.

## 3.4 Node Join and Leave

While the above describes the steady-state behavior of Mercury, it does not address how nodes join or leave the

<sup>4</sup>For a query, we count the number of routing hops to reach the first value in the range it covers.

system. This section describes the detailed protocol used by nodes during join and departure.

Recall that each node in Mercury needs to construct and maintain the following set of links: *a)* successor and predecessor links within the attribute hub, *b)*  $k$  long-distance links for efficient intra-hub routing and *c)* one cross-hub link per hub for connecting to other hubs. The cross-hub link implies that each node knows about at least one representative for every hub in the system. In order to recover during node departures, nodes keep a small number (instead of one) of successor/predecessor and cross-hub links.

**Node Join.** Like most other distributed overlays, an incoming Mercury node needs information about at least one (or at most a few) node(s) already part of the routing system. This information can be obtained via a match-making server or any other out-of-band means. The incoming node then queries an existing node and obtains state about the hubs along with a list of representatives for each hub in the system. Then, it randomly chooses a hub to join and contacts a member  $m$  of that hub. The incoming node installs itself as a predecessor of  $m$ , takes charge of half of  $m$ ’s range of values and becomes a part of the hub.

To start with, the new node copies the routing state of its successor  $m$ , including its long-distance links as well as links to nodes in other hubs. At this point, it initiates two maintenance processes: firstly, it sets up its own long-distance links by routing to newly sampled values generated from the harmonic distribution (as described above.) Secondly, it starts random-walks on each of the other hubs to obtain new cross-hub neighbors distinct from his successor’s. Note that these processes can be run lazily, as they are not essential for correctness, and only affect the efficiency of the routing protocol.

**Node Departure.** When nodes depart, the successor/predecessor links, the long-distance links and the inter-hub links within Mercury must be repaired. To repair successor/predecessor links within a hub, each node maintains a short list of contiguous nodes further clockwise on the ring than its immediate successor. When a node’s successor departs, that node is responsible for finding the next node along the ring and creating a new successor link.

A node’s departure will break the long-distance links of a set of nodes in the hub. These nodes establish new long-distance links to replace the failed ones. Nodes which are not directly affected by this departure do not take any action. The departure of several nodes, however, can distort the distribution of links of nodes which are not affected directly. To repair the distribution, nodes periodically re-construct their long-distance links using recent estimates of node counts. Such repair is initiated only when the number of nodes in the system changes dramatically (by a factor of 2 – either by addition or departure).<sup>5</sup>

Finally, to repair a broken cross-hub link, a node considers the following three choices: *a)* it uses a backup cross-hub link for that hub to generate a new cross-hub neighbor (using a random walk within the desired hub), or *b)* if such a backup is not available, it queries its successor and predecessor for their links to the desired hub, or *c)* in the worst case, the node contacts the match-making (or bootstrap server)

<sup>5</sup>Intuitive justification: routing performance is only sensitive to the *logarithm* of the number of nodes.

to query the address of a node participating in the desired hub.

## 4. EFFICIENCY IN THE FACE OF NON-UNIFORMITY

The Mercury protocol we have described thus far is largely a derivative of previous structured overlay protocols. We have shown that it can provide efficient (logarithmic) routing when the responsibility of handling various attribute values is uniformly distributed to all nodes within a hub. However, as alluded to in Section 3.2, the desire to balance routing load can create a highly non-uniform distribution of ranges.

We begin this section by analyzing why such non-uniform range distributions conflict with the goal of efficient routing. We find that Kleinberg’s basic small-world network result makes certain assumptions which are non-trivial to satisfy in a distributed setting when node ranges in a network are non-uniform. We then present a novel algorithm which ensures efficient routing even when the assumptions are violated.

We then tackle non-uniformity in two other dimensions: query selectivity, and data popularity. We show how the core of the algorithm for efficient routing under non-uniform range distributions can be re-used to optimize query performance given non-uniformity in query selectivity and data popularity.

### 4.1 Small-world Networks

Let  $G$  represent a circle on  $n$  nodes. Define *node-link distance*  $d_n(a, b)$  between two nodes  $a$  and  $b$  as the length of the path from  $a$  to  $b$  in the *clockwise* direction. The objective is to find “short” routes between any pair of nodes using distributed algorithms. Kleinberg [14] showed that if each node,  $A$ , in  $G$  constructs one additional “long-link” in a special manner, the number of expected hops for routing between any pair of nodes becomes  $\mathcal{O}(\log^2 n)$ . Each node  $A$  constructs its link using the following rule:  $A$  generates an integer  $x \in (0, n)$  using the harmonic distribution, viz.,  $h_n(x) = 1/(n \log x)$ , and establishes a link to the node  $B$  which is  $x$  links away in the clockwise direction from  $A$ . The routing algorithm for each node is to choose the link which takes the packet closest to the destination with respect to the node-link distance. Symphony [18] extends this result by showing that creating  $k$  such links reduces the routing hop distance to  $\mathcal{O}(\frac{1}{k} \log^2 n)$ .

Creating the long-links appears deceptively straightforward. However, it may be difficult and expensive ( $\mathcal{O}(x)$ ) for a node  $A$  to determine which node,  $B$ , is  $x$  hops away from it. Contacting node  $B$  would be simpler if we could easily determine what value range  $B$  was responsible for. This would allow the use of any existing long-links to contact this node more efficiently and reduce the number of routing hops to  $\mathcal{O}(\log^2 n)/k$ .

In systems like Symphony, this problem is solved by approximating the hop distance of any node. Since Symphony places nodes randomly along its routing hub, it makes the assumption that all nodes are responsible for ranges of approximately the same size,  $r$ . By simply multiplying  $r$  by  $x$  and adding to the start of the values range of  $A$ , Symphony is able to estimate the start of the range that  $B$  is responsible for. Unfortunately, this technique does not work when not all nodes are responsible for the same range size of values, i.e., when ranges are highly non-uniform in size.

Non-uniform range sizes, however, are exactly what we use in order to provide load balancing. Specifically, our load balancing algorithm, which we elaborate in Sec. 4.4, balances load by increasing the density of nodes in the portions of the value range that are heavily loaded. Accordingly, the value ranges for nodes in heavily loaded regions are smaller than the value ranges for nodes in lightly loaded regions.

To enable the efficient creation of long-links over an overlay with non-uniform range sizes, we introduce a novel random sampling algorithm (Sec. 4.2), and use this sampling algorithm to create an estimate of the density of nodes in different parts of the routing hub, i.e., a histogram of the distribution of nodes (Sec. 4.2.1). This allows us to easily map from  $x$  (the desired length of a long-link) to the start of the value range for  $B$  (the node at distance  $x$ ). This mapping, in turn, enables us to construct the long-distance links of Section 3.3 despite non-uniform node ranges.

### 4.2 Random Sampling

Maintaining state about a uniformly random subset of global participants in a distributed network, in a scalable, efficient and timely manner is non-trivial. In the context of our system, the naïve approach of routing a sample request message to a randomly generated *data-value* works well only if node ranges are uniformly distributed. Unfortunately, as already explained, this assumption is easily violated.

Another obvious approach is to assign each node a random identifier (by using a cryptographic hash, for example) and route to a randomly generated identifier to perform sampling. However, in order for the sampling process to be efficient, we need a routing table for delivering messages to node identifiers. Another approach is to use protocols like Ransub [15] which are specifically designed for delivering random subset information. Unfortunately, both these approaches require incurring the overhead of maintaining a separate overlay – one which may not be well suited for efficient data-value routing.

Mercury’s approach for sampling is novel – we show that the hub overlay constructed by Mercury in a randomized manner is an *expander* [19] with a high probability. An expander has the property that random walks over the links of such a network converge very quickly to the stationary distribution of the random walk. Since the hub overlay graph is regular, the stationary distribution is the uniform distribution. We state the lemma in a semi-rigorous manner.<sup>6</sup>

**LEMMA 1.** *Let  $G$  be a circle on  $n$  nodes with  $\mathcal{O}(\log n)$  additional links per node generated using the harmonic probability distribution (as described in Section 4.1). Let  $\Pi_\infty$  denote the stationary distribution of a random walk on  $G$  and let  $\Pi_t$  denote the distribution generated by the random walk after  $t$  steps. Then, with high probability,  $d_1(\Pi_t, \Pi_\infty) < \epsilon$  for  $t > \mathcal{O}(\log^c(n/\epsilon))$  for small constants  $c$ , where  $d_1$  denotes the statistical or  $L_1$  distance between two distributions. (See [20] for rigorous definitions.)*

This leads to a very simple algorithm for performing random sampling: send off a **sample-request** message with a small (e.g.,  $\log n$  hop) Time-To-Live (TTL). Every node along the path selects a random neighbor link and forwards it, decrementing the TTL. The node at which the TTL expires sends back a sample. Notice that this algorithm uses

<sup>6</sup>The proof is omitted for reasons of space, and will be available in a related tech-report.

only local information at every stage in the sampling process and adapts easily to a highly dynamic distributed overlay. In addition, these messages could be piggy-backed on any existing keep-alive traffic between neighbors to reduce overhead. Our simulations (Sec. 5.1) show that Mercury can indeed perform near-perfect uniform random sampling using a TTL of  $\log n$ .

We now describe three important ways in which we utilize random sampling in our system viz., to maintain node-count histograms, for estimating the selectivity of queries and for effective load balancing.

#### 4.2.1 Maintaining Approximate Histograms

This section presents the mechanism used by nodes for maintaining histograms of any system statistic (e.g., load distribution, node-count distribution<sup>7</sup>, etc.) The basic idea is to sample the distribution *locally* and exchange these estimates throughout the system in an epidemic-style protocol.

Let  $N_d$  denote the “local”  $d$ -neighborhood of a node - i.e., the set of all nodes within a distance  $d$  *ignoring the long distance links*. Each node periodically samples nodes  $\in N_d$  and produces a local estimate of the system statistic under consideration. For example, if the node-count distribution is being measured, a node’s local estimate is  $(M_a - m_a)|N_d|/(\sum_{k \in N_d} |r_k|)$  where  $r_k$  is the range of a node  $k$  and  $m_a, M_a$  are the minimum and maximum attribute values for the attribute  $a$ . In our experiments, we use  $d = 3$ .

In addition, a node periodically samples  $k_1$  nodes uniformly at random using the sampling algorithm described in Section 4.2. Each of these nodes reports back its local estimate and the most recent  $k_2$  estimates it has received. As time progresses, a node builds a list of tuples of the form: `{node_id, node_range, time, estimate}`. (The timestamp is used to age out old estimates.) Each of these tuples represent a point on the required distribution - stitching them together yields a piecewise linear approximation.

$k_1$  and  $k_2$  are parameters of the algorithm which trade-off between overhead and accuracy of the histogram maintenance process. In Section 5, we show through simulations that setting each of  $k_1$  and  $k_2$  to  $\log(n)$  is sufficient to give reasonably accurate histograms for sampling population distribution.

If the system needs to generate an average or histogram of node properties, the collected samples can be used exactly as they are collected. However, if the desire is to generate an average or histogram of properties around the routing hub, a minor modification is needed. Namely, in order to generate unbiased node-count histograms, the samples received are weighted differently; samples reporting lower densities are given higher weight to account for the fact that there would be less nodes to produce low density samples.

Once a histogram is constructed, long distance links are formed as follows: first, the number of nodes  $n$  in the system is estimated. For each long-distance link, a value  $n_i$  between  $[1, n]$  is generated using the harmonic distribution. This represents the number of nodes that must be skipped along the circle (in the clockwise direction, let’s say) to get to the desired neighbor. The histogram is then used to estimate a value  $v_i$  that this desired neighbor will be responsible for. Finally, a join message is sent to this *value*  $v_i$  which will get routed to the desired neighbor using the existing routing network.

<sup>7</sup>Number of nodes responsible for a given range of values.

### 4.3 Query Selectivity

Recall that a query  $Q$  is sent to only one of the attribute hubs in  $\mathcal{A}_Q$ . Also a query  $Q$  is a *conjunction* of its predicates each of which can have varying degree of selectivity. For example, some predicate might be a *wildcard* for its attribute while another might be an exact match. Clearly, a wildcard predicate will get flooded to every node within its attribute hub. Thus, the query  $Q$  should be sent to that hub for which  $Q$  is *most selective* to minimize the number of nodes that must be contacted.

The problem of estimating the selectivity of a query has been very widely studied in the database community. The established canonical solution is to maintain approximate histograms of the number of database records per bucket. In our case, we want to know the number of nodes in a particular bucket. Each node within a hub can easily gather such an histogram for its own hub using the histogram maintenance mechanism described above. In addition, using its inter-hub links, it can also gather histograms for other hubs efficiently. These histograms are then used to determine the selectivity of a subscription for each hub. We see in Section 5.4 that even with a very conservative workload, this estimation can reduce a significant amount of query flooding.

### 4.4 Data Popularity and Load Balancing

When a node joins Mercury, it is assigned responsibility for some range of an attribute. Unfortunately, in many applications, a particular range of values may exhibit a much greater popularity in terms of database insertions or queries than other ranges. This would cause the node responsible for the popular range to become overloaded. One obvious solution is to determine some way to partition the ranges in proportion to their popularity. As load patterns change, the system should also move nodes around as needed.

We leverage our approximate histograms to help implement load-balancing in Mercury. First, each node can use histograms to determine the average load existing in the system, and, hence, can determine if it is relatively heavily or lightly loaded. Second, the histograms contain information about which parts of the overlay are lightly loaded. Using this information, heavily loaded nodes can send probes to lightly loaded parts of the network. Once the probe encounters a lightly loaded node, it requests this lightly loaded node to gracefully leave its location in the routing ring and re-join at the location of the heavily loaded node. This leave and re-join effectively increases the load on the neighboring (also likely to be lightly-loaded) nodes and partitions the previous heavy load across two nodes.

Let the average load in the system be denoted by  $\bar{L}$ . Define the *local load* of a node as the average of load of itself, its successor and its predecessor. A node is said to be lightly loaded if the ratio of its local load to  $\bar{L}$  is less than  $\frac{1}{\alpha}$  and heavily loaded if the ratio is greater than  $\alpha$ . This definition ensures that if a node is lightly loaded, its neighbors will be lightly loaded with a high probability. If this is not the case (when the ratio of neighbor loads is  $> \alpha$ ), the lighter neighbor performs a load balance with the heavier one to equalize their loads. It is easy to show<sup>8</sup> that the leave-rejoin protocol described above *decreases* the variance of the load distribution at each step and bounds the maximum load imbalance

<sup>8</sup>We omit the proof for reasons of space. The idea is simply that variance reduction ‘near’ the heavier node is larger than the variance increase ‘near’ the lighter node.

in the converged system by a factor of  $\alpha$ , provided  $\alpha \geq \sqrt{2}$ . By tolerating a small skew, we prevent load oscillations in the system.

Over time, the leaves and re-joins result in a shift in the distribution of nodes to reflect the distribution of load. However, this shift in node distribution can have significant implications. Many of the properties of Mercury’s routing and sampling rely on the harmonic distance distribution of the random long-links. When nodes move to adjust to load, this distribution may be changed. However, our technique for creating long-links actually takes the node distribution into account explicitly as stated previously.

We note that a similar leave-join based load balancing mechanism has been proposed concurrently in [13] and [8]. However, [13] and [8] do not handle skewed node range distributions. Because it exploits our random sampling mechanism (from Sec. 4.2), our load balancing mechanism works even in the presence of skewed node ranges.

## 5. EVALUATION

This section presents a detailed evaluation of the Mercury protocol using simulations. We implemented a simple discrete event-based simulator which assigns each application level hop a unit delay. To reduce overhead and enable the simulation of large networks, the simulator does not model any queuing delays or packet loss on links. The simplified simulation environment was chosen for two reasons: first, it allows the simulations to scale to a large (up to 50K) number of nodes, and secondly, this evaluation is not focused on proximity routing. Since our basic design is similar in spirit to Symphony and Chord, we believe that heuristics for performing proximity-based routing (as described in [10]) can be adapted easily to Mercury.

Our evaluation centers on two main features of the Mercury system: 1) scalable routing for queries and data-records, and 2) balancing of routing load throughout the system. We begin with an evaluation of our core routing mechanisms – random sampling and histogram maintenance. We then study the impact of these mechanisms on the overall routing performance under various workloads. Finally, we present results showing the utility of caching and query selectivity estimation in the context of Mercury.

Except for query selectivity estimation, most of our experiments focus on the routing performance of data *within* a single routing hub. Hence,  $n$  will denote the number of nodes within a hub. Unless stated otherwise, every node establishes  $k = \log n$  intra-hub long-distance links. For the rest of the section, we assume without loss of generality that the attribute under consideration is a `float` value with range  $[0, 1]$ . Each node in our experiments is thus responsible for a value interval  $\subset [0, 1]$ .

In what follows, **NodeLink** denotes the ideal small-world overlay, i.e., long distance links are constructed using the harmonic distribution on node-link distance. **ValueLink** denotes the overlay when the harmonic distribution on value-distance is used (Section 3.3). **HistoLink** denotes the scenario when links are created using node-count histograms (see Section 4.) Note that the performance of the **ValueLink** overlay is representative of the performance of a plain DHT (e.g., Chord, Symphony) under the absence of hashing and in the presence of load balancing algorithms which preserve value contiguity.

For evaluating the effect of non-uniform node ranges on

our protocol, we assign each node a range width which is inversely proportional to its popularity in the load distribution. Such a choice is reasonable since load balancing would produce precisely such a distribution – more nodes would participate in a region where load is high. The ranges are actually assigned using a Zipf distribution. In particular, data values near 0.0 are most popular and hence a large number of nodes share responsibility for this region, each taking care of a very small node range. For reference, in our simulator setup, these are also the nodes with lowest numeric IDs.

### 5.1 Random-Walk Based Sampling

The goal of our random-walk based sampling algorithm is to produce a uniform random sample of the nodes in the system. We measure the performance of our algorithm in terms of the statistical distance (alternatively called  $L_1$  distance) of the perfect uniform distribution from the distribution obtained via the random walks. For these experiments, nodes are assigned ranges using a highly-skewed Zipf distribution ( $\alpha = 0.95$ ). In each sampling experiment, we pick a node at random and record the distribution of the samples taken by  $kn$  random walks starting from this node. If our sampling algorithm is good, the random walks should hit each node roughly  $k$  times. Note that the parameter  $k$  is just for evaluating the distribution obtained – the protocol does not use it in any manner.

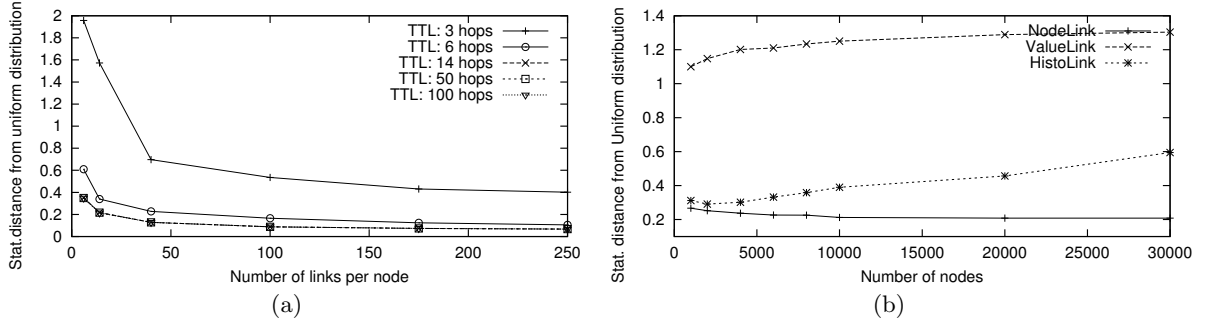
Figure 3(a) plots the accuracy of the sampling process as the degree of the graph and the TTL for the random-walks is varied. The underlying overlay we consider is a perfect small-world network (**NodeLink**). We find that, over a certain threshold ( $\log n$ ), the TTL of the random-walks does not influence sampling accuracy. Also, the sampled distribution is almost perfectly random for graph degrees  $c \log n$ , where  $c$  is a small constant. In practice, we found that, for routing purposes, sufficiently accurate histograms are obtained even for  $c = 1$ .

Figure 3(b) shows how the construction of the underlying network affects sampling accuracy. We see that the **NodeLink** and **HistoLink** overlays perform much better than the **ValueLink** (a vanilla DHT without hashing and in the presence of load balancing) overlay. These effects are explained using Figure 4 which plots the distribution of long-distance links. As described earlier, in our experiments, nodes with the lowest identifiers (responsible for values near 0.0) are the most popular while nodes at the other end of the value range are the least popular.

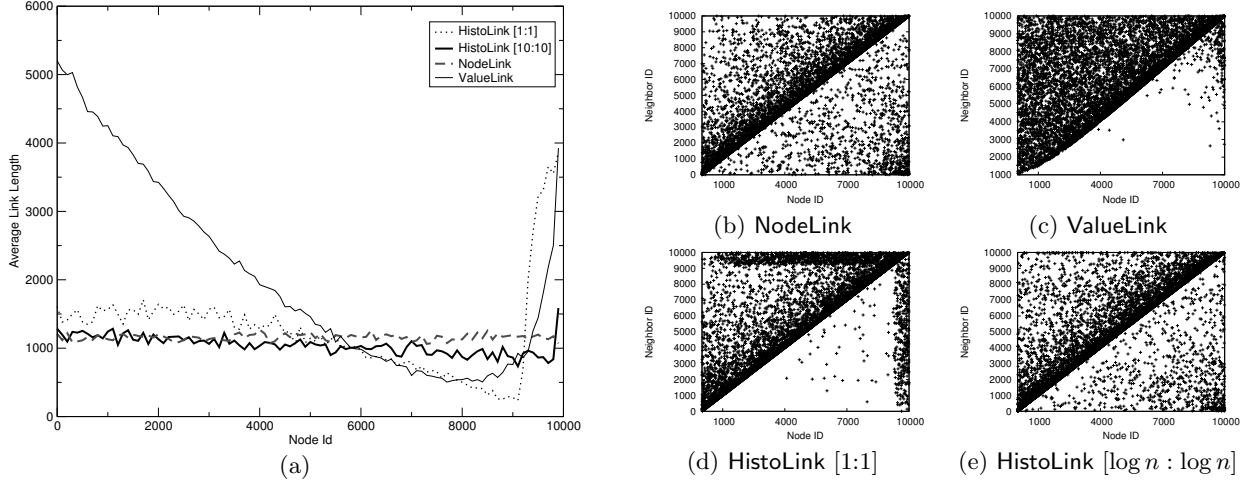
Recall that, in a **ValueLink** overlay, nodes construct their links by routing to *values* generated using a harmonic distribution. However, in this case node ranges are not uniformly distributed – in particular, nodes near the value 1.0 (i.e., nodes with higher IDs) are less popular, so they are in charge of larger range values. Hence, the long-distance links they create tend to *skip* over less nodes than appropriate. This causes all the links (and correspondingly, the random walks) to crowd towards the least popular end of the circle. The **HistoLink** overlay offsets this effect via the maintained histograms and achieves sampling accuracy close to that achieved by the optimal **NodeLink** overlay.

Each **sample-request** message travels for TTL hops and hence obtaining one random sample generates TTL additional messages in the overall system. However, all these messages are sent over existing long-distance links. Thus,





**Figure 3: Accuracy of random-walk based sampling.** Figure (a) shows the effect of the degree of the graph. ( $n = 10000$ ;  $\log n = 14$ ; NodeLink overlay.) Figure (b) shows the effect of link structure.



**Figure 4: Distribution of long-distance links.** Plot (a) summarizes the average link length as a function of the Node ID for the different link construction algorithms. Plots (b)-(e) provide more detail about the links constructed by each algorithm. In plots (b)-(e), the Y-axis plots the ID of the  $k = \log n$  neighbors for each node on the X-axis. Nodes are consecutively arranged on the hub circle. Number of nodes = 10000. For HistoLink,  $[k_1 : k_2]$  means  $k_1$  nodes were queried per round each giving  $k_2$  estimate reports; 5 exchange rounds were performed.

they can be easily multiplexed and piggy-backed (by simply appending the IP address of the requesting node) over the regular keep-alive pings that might be sent between neighbors. Also, if the samples are uniformly distributed over all nodes, each node receives  $\mathcal{O}(1)$  sample requests per sampling period. Hence, we conclude that the overhead of the sampling method is very small.

## 5.2 Node-Count Histograms

In this section, we evaluate the accuracy of the node-count histograms obtained by nodes using the mechanism described in Section 4.2. These histograms, introduced in Section 4.2.1, provide an estimate of the total number of nodes in the system and help in establishing the long-distance links correctly.

We measure the accuracy of the obtained histogram in terms of its *distance* from the “true” histogram under the  $L_1$  norm. Figure 5(a) plots the average accuracy of the histogram (the average is taken over all nodes) as the parameters for the histogram maintenance process are varied. In this experiment, 10 rounds of exchanges were performed. We see that the error is consistently small and decreases

rapidly as the number of the nodes queried per round increases.<sup>9</sup> Although not obvious from the graph, the same pattern is observed when the number of reports queried from each node is increased. These effects are expected, since the number of samples received by each node per round grows linearly with either of these parameters.

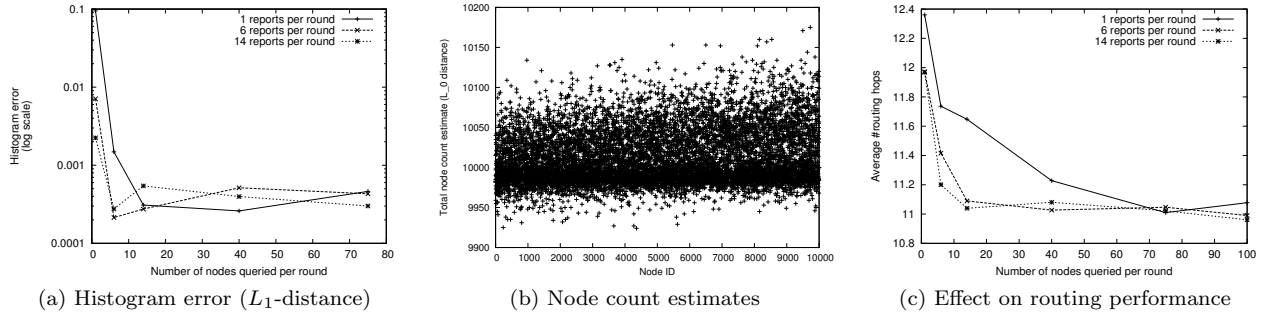
Figure 5(b) plots the overall node-count estimates produced by each node in a system of  $n = 10000$  nodes. The experiment was run for 10 exchange rounds, with  $k_1 = k_2 = \log n$ . We see that the estimates are very tightly clustered around the actual node-count value of 10000.

During each round of the histogram maintenance process, each node queries  $k_1$  randomly sampled nodes and receives  $k_2$  estimate samples from each node. The messaging overhead per round per node is thus proportional to  $k_1 k_2$ .

## 5.3 Routing Performance

We now present an evaluation of the overall routing performance of Mercury. This factors in the effects of the ran-

<sup>9</sup>The graph does show some fluctuations, but their magnitudes are tiny (result of experimental variations).



**Figure 5: Accuracy of sampled histograms.** ( $n = 10000$ ) Figure (a) shows the effect of changing parameters on average histogram error. Figure (b) shows the node count estimates gathered by each node. Figure (c) shows effect of changing parameters on overall routing performance.

dom sampling and histogram maintenance strategies. We present the performance of the basic protocol with and without route caching optimizations, discuss the effect of skewed node-range distributions and validate our claim that the protocol using histograms achieves near-optimal routing delays. As before, we concentrate on routing within a single hub. In each of the following experiments, nodes establish  $k = \log n$  long-distance links within a hub.

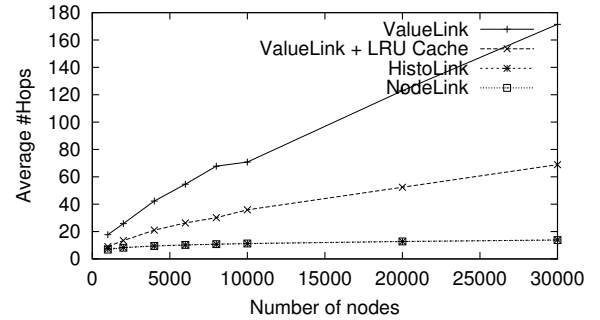
We experiment with two different data workloads – *uniform* and *Zipf*. The Zipf workload is high-skewed and is generated using the distribution  $x^{-\alpha}$  where  $\alpha = 0.95$ . Notice that this means that the attribute values near 0.0 are the most popular and those around 1.0 are the least popular. We also show the performance of two types of route caching policies, viz., LRU replacement and a direct-mapped cache.<sup>10</sup> Our objective here is not to find the best possible policy for our workload. Rather, our aim is to show the ease with which application-specific caching can co-exist fruitfully with Mercury routing. In our caching experiments, each node keeps a cache of  $\log n$  route entries.

Figure 6 shows the performance of Mercury when node ranges are uniformly distributed. The Y-axis shows the average number of hops taken by a data-item to reach its destination (node where it is stored) in the hub. Although these graphs show results for HistoLink overlay, the performance of NodeLink and ValueLink is very similar, as expected.

We see that, for uniform node ranges, the number of routing hops scales logarithmically (with very low constant factors) as the number of nodes increases, irrespective of the workload used. Thus, Mercury can provide low end-to-end routing delays to applications even for a large number of nodes. With caching enabled, there is a significant reduction in hop count. While this is easy to see for a skewed workload, the reduction for a uniform workload results from the fact that a cache effectively increases Mercury’s routing table size. We believe that caching is an important optimization which Mercury can easily incorporate into its basic protocol.

### Effect of Non-Uniform Ranges

Figure 7 compares the performance of the protocol with and without approximate histograms to guide the selection of the long-distance links. In this experiment, the node-range



**Figure 7: Effect of non-uniform node ranges on the average number of routing hops.** As workload, we use the Zipf distribution with  $\alpha = 0.95$ .

distribution and the data distribution are Zipf-skewed. For histogram maintenance in this experiment, we used 5 exchange rounds, where each node queried  $\log n$  nodes per round asking each for  $\log n$  estimate reports.

As explained in Section 5.1, the naïve ValueLink overlay (vanilla DHT in the presence of load balancing) creates links which *skip* the crowded and popular region (see Figure 4.) Hence, packets destined to these nodes take circuitous routes along the circle rather than taking short cuts provided by the long-distance links. Although caching ameliorates the effect, the performance is still much worse as compared to the optimal NodeLink overlay.

On the other hand, we see that the performance of the HistoLink overlay is nearly the same as that of the optimal NodeLink overlay. Again, looking at Figure 4, we find that node-count histograms enable nodes to establish a correct link distribution (corresponding to the NodeLink overlay) quickly using very low overheads.

Figure 5(c) shows the effect of histogram accuracy on the overall routing performance. We see that as the parameters  $k_1$  and  $k_2$  in the histogram maintenance process increase, the routing performance improves as expected. We note that this influence is limited (note the scale of the graph) since it is directly dependent on the accuracy of the obtained histograms (see Figure 5(a).)

### 5.4 Estimating Query Selectivity

To evaluate the usefulness of forwarding queries to the most selective attribute hubs, we set up an experiment with

<sup>10</sup>For an  $n$ -entry cache, there is one entry for each of the  $(1/n)^{\text{th}}$  region of the attribute space.

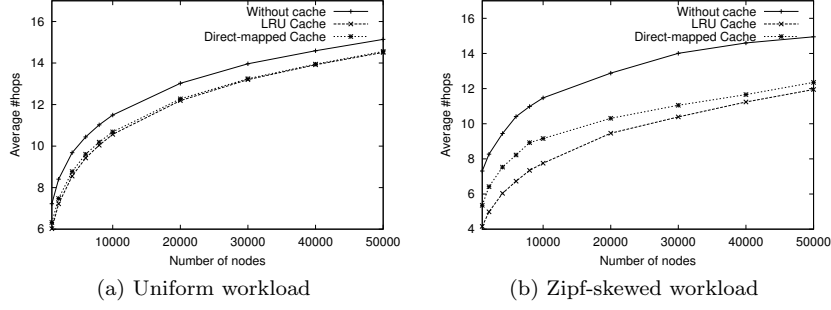


Figure 6: Performance of the basic Mercury protocol for various workloads.

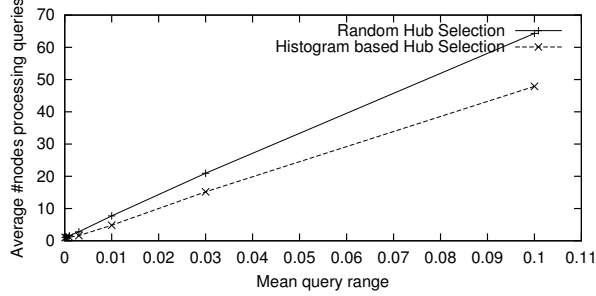


Figure 8: Nodes required to answer a query for varying query ranges, with and without selectivity estimation.

3 attribute hubs. Our workload is motivated by the distributed multi-player game application we describe in Section 6. The attribute hubs correspond to three dimensions of virtual space. Each query contained 3 range predicates – one for each attribute. Such a query specifies a cuboid region of the virtual space. The range-size of each predicate was Gaussian-distributed, while the range position within the attribute space was Zipf-distributed. The node-range distribution within each hub is skewed.

Figure 8 plots the average number of nodes processing a query for different sizes of queries. The query size is measured as the average fraction of the value range that a query covers in a *single* dimension. We call this the query range. The plot shows that, even with our conservative workload, choosing a hub based on the selectivity estimates results in up to 25-30% reduction in the degree of flooding of a query. With workloads exhibiting wildcards, much higher reductions would be expected. Because we lack a good model for wildcard workloads, however, we refrain from speculating about the exact degree of reduction that might be achieved.

## 5.5 Load Balancing

For evaluating the efficiency of load balancing achieved by Mercury’s load balancing algorithm, we conduct the following experiment: In a system of 1000 nodes, we assign load to each node using a Zipf distribution with varying values of the initial skew (Zipf parameter). The system is said to be load-balanced when  $\frac{1}{\Delta} \leq \text{load}/\text{avg\_load} \leq \Delta$  holds for all nodes in the system.

In the experiment, we run multiple *rounds* of the load balancing algorithm, until the system is load-balanced. Each round consists of the following steps:

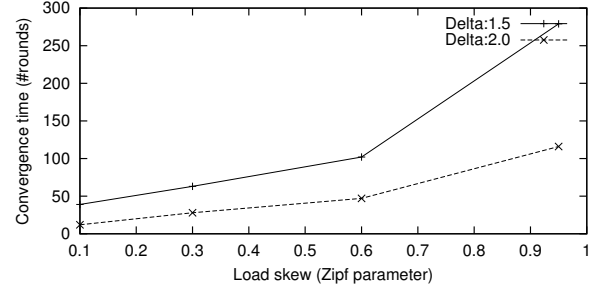


Figure 9: Rounds taken to achieve load balance as a function of the initial skew of the load. Delta is the degree of load balance sought.

1. Each node samples its neighbors, to determine the local node-count. This requires one round-trip.
2. Each node runs one round of the histogram maintenance algorithm. (Recall that a round of the histogram maintenance algorithm involves sending  $\log n$  probes in parallel, each of which must traverse  $1 + \log n$  hops.)
3. Nodes check their histograms to determine if they are heavily loaded. If a node is heavily loaded, it sends a probe to a lightly loaded node. This probe must traverse  $\log n$  hops.
4. Lightly loaded nodes leave and re-join the network. To re-join, the lightly loaded nodes must establish new long links. The link establishment messages traverse  $1 + \log n$  hops, in expectation.

Figure 9 plots the number of rounds of load-balancing required to achieve load balance. We see that Mercury can load-balance to within a factor of  $\Delta = 2$  within 100 rounds despite heavy skews in the workload (Zipf with  $\alpha = 0.95$ ). In practical terms, consider an overlay with 10000 nodes, and a 50 ms delay between nodes. The time to complete one round of load-balancing is the product of the number of hops traversed by messages in the load balancing algorithm<sup>11</sup>, and the inter-node delay. Thus the time to complete one round is  $50 * (4 + 3 \log n)$  ms. The time to load-balance the entire overlay is then  $100 * 50 * (4 + 3 \log n)$  ms, or about 220 seconds, which is reasonable for such a large network.

<sup>11</sup>Since the messages in step 2 are sent in parallel, we count the number of hops once, rather than multiplying by the number of messages. Similarly for step 4.

## 6. DISTRIBUTED APPLICATION DESIGN

Previous sections have demonstrated that Mercury provides scalable range-based lookups. In this section, we describe how the range query support provided by Mercury can also be used as a *building block* for distributed applications. Specifically, we consider a multiplayer game, and demonstrate how multi-attribute range queries can be fruitfully applied to solve the game’s distributed state maintenance problem.

### 6.1 Distributed State Maintenance

One of the difficulties in designing distributed multiplayer games is managing the game state. Game state includes such information as where a player is located, the resources (such as ammunition) he has, and how healthy he is. For a distributed game, this state must be available at multiple, perhaps physically distributed, machines. This is necessary so that the players in the game (who may be running the game on different nodes) have accurate views of the game. Thus, the central challenge is providing a way for the nodes in the game to update each others’ view of the game state.

While distributed state maintenance, in its full generality, is a large and difficult problem, a few observations about games in particular enable us to simplify the problem. First, observe that a node in a game is interested in only a small subset of the entire game state. Second, note that the objects belonging to this subset are not arbitrary, but are related to each other. For example, in most multiplayer games, a player is primarily interested in entities geographically nearby (such as within the player’s room). Third, note that games do not require strong consistency guarantees. In particular, even centralized games trade-off strict consistency to provide interactive response times.

These observations suggest a publish-subscribe [3] architecture for managing state updates. Specifically, a game node registers a “subscription” describing the objects which it wishes to keep updated. Also, whenever a node changes the state of a game object, the node creates a “publication” which is delivered to the other nodes interested in that object. Note that the subscription essentially implements “area-of-interest” filtering (which we describe in further detail in Section 6.4).

A key requirement of this design is a flexible subscription language which allows the game to express its object-subscriptions precisely. If the language is coarse, the game may end up receiving a large number of irrelevant updates which will have to be filtered locally. But a language may also be so flexible as to preclude scalable routing of object updates. Fortunately, we can use the multi-attribute range query primitive, which Mercury implements scalably, as the subscription language for distributed games.

In the remainder of this section, we describe the publish-subscribe system we have built on top of the Mercury routing layer; describe *Caduceus*, our simple multiplayer distributed game; compare publish-subscribe to other possible solutions, and present some basic performance measurements.

### 6.2 Implementing Publish-Subscribe with Mercury

To explain the implementation of publish-subscribe with Mercury, we begin with the application programmer’s inter-

face. We then explain how Mercury is used to implement this API.

The publish-subscribe system exposes a simple interface for applications. This API has three basic calls: `send_publication`, `register_interest`, and `unregister_interest`. The `send_publication` call is used to disseminate an object update via Mercury. The call takes the updated object as its argument. When an application invokes this call, the publish-subscribe system assumes responsibility for delivering the object update to other interested nodes. The `register_interest` is used to subscribe to updates. This call takes a multi-attribute range query as its argument. After an application invokes this call, the publish-subscribe system is responsible for delivering an update to the application whenever an object matching the query is added or modified. The `unregister_interest` call simply cancels a subscription.

Our implementation of publish-subscribe using Mercury is relatively straightforward. Publications are routed as regular data items in Mercury. Subscriptions are similar to Mercury queries, but with one important difference. Mercury queries can be discarded immediately after they are answered. Subscriptions, however, must match not only objects currently existing in the system, but must return future matching objects as well. Thus, subscriptions must be retained at the nodes which receive them. In our current implementation, subscriptions are hard state: they are retained until they are cancelled by the querying node. To provide robustness in the face of node failures, however, we plan to implement subscriptions as soft state.

### 6.3 Caduceus

Caduceus is a two-dimensional, multi-player, shooter game. Each player in the game has a ship and a supply of missiles. Players pursue each other around the two-dimensional space, and fire missiles when their opponents are in range. The goal is simply to kill as many opponents as possible. Figure 10(a) presents a screenshot of the game. At any given time, a player sees the other ships in her immediate vicinity, as defined by the game window. The virtual game world is larger than the window. Thus, there might, for example, be opponents located beyond any of the edges of the game window.

The state of the game is represented using two kinds of objects: ships and missiles. A ship consists of a location, velocity, and ancillary information such as fuel level. A missile is created whenever a ship fires shots. A missile consists of a location, velocity, and owner information.<sup>12</sup> The main loop of Caduceus, shown in Figure 10(b), is relatively simple.

### 6.4 Alternative Update Routing Architectures

As stated previously, a central problem in *distributed* multiplayer gaming is ensuring that all nodes have (weakly) consistent views of the game state. To provide this (weak) consistency, nodes send updates to other nodes whenever the game state has changed. To date, most games have used either broadcast messaging, or a centralized server, to deliver these updates. Unfortunately, both these designs scale poorly to a large number of nodes.

To improve scalability, researchers have proposed area-of-interest filtering [17, 24] schemes, which divide the world into

<sup>12</sup>The owner is used to credit the appropriate player when an opponent is killed.

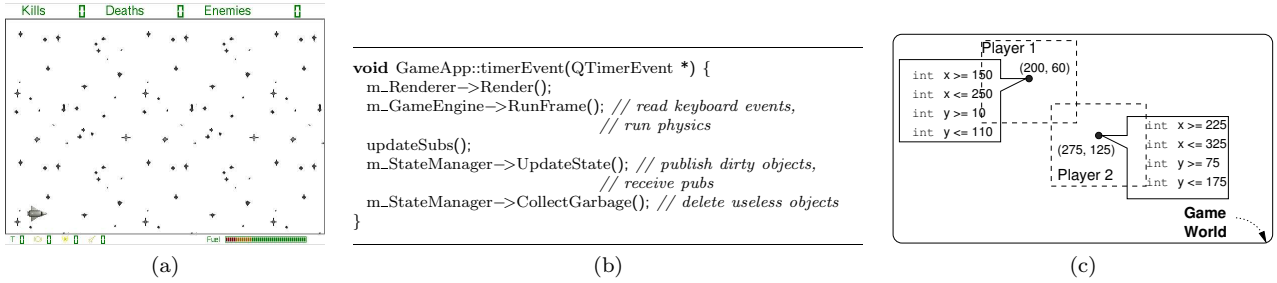


Figure 10: a) Screenshot of Caduceus b) Caduceus main loop c) Example subscriptions

a fixed set of regions and map them to IP multicast groups. However, the fixed regions result either in the delivery of a large number of irrelevant updates to clients, or in the maintenance of a large number of IP multicast groups at routers.

In contrast, Mercury’s subscription language is ideal for implementing area-of-interest filtering. In particular, the subscription language makes it easy to describe arbitrary physical regions. As an example, Figure 10(c) shows two nodes expressing their interests in the rectangular regions near them. Of interest is the fact that the regions do not, for example, need to fit a pre-defined tiling of the space. Note that while tiling the space, and assigning these tiles to different channels, would be possible for a simple two-dimensional game, it becomes far more difficult in games with irregular spaces, such as corridors, or which have to deal with visibility constraints such as horizons. It is also difficult for multicast group-based schemes to support arbitrary interests such as “the location of all my teammates”.

As shown in Figure 6(b), Mercury is able to handle 10000 nodes while keeping the number of routing hops below 8. Assuming that the average-case, one-way delay between nodes is about 20ms (e.g., they are all well connected and within the U.S. west coast), this results in an end-to-end delay of less than 160ms. We believe that game-specific caching algorithms could further improve the routing performance, making Mercury-based games scalable to thousands of nodes.

## 6.5 Performance Evaluation

We evaluate the performance of our system with two metrics: hop count, and message count. We run a varying number of players. The players move through the world according to a random waypoint model, with a motion time chosen uniformly at random from (1, 10) seconds, a destination chosen uniformly at random, and a speed chosen uniformly at random from (0, 360) pixels per second. The size of the game world is scaled according to the number of players. The dimensions are  $640n \times 480n$ , where  $n$  is the number of players. All results are based on the average of 3 experiments, with each experiment lasting 60 seconds. The experiments include the benefit of a  $\log n$  sized LRU cache at each node, but do *not* include the benefits of any long pointers.<sup>13</sup>

Table 1 summarizes the results. With respect to hop count, we find that the hop count increases only slightly

as we double the number of nodes. To evaluate Mercury’s messaging efficiency, we compare it to two alternatives. In the “broadcast messages” column of the table, we report the number of messages that would have been transmitted if every update were sent to every node (as was done in first-generation distributed games). In the “optimal messages” column, we report the number of messages required to exactly satisfy each node’s interests, without any control message overhead. We find that Mercury performs substantially better than a broadcast scheme (43% as many messages transmitted for 20 nodes), and that this performance difference increases when we increase the number of nodes, with Mercury using only 29% as many messages as broadcast for 40 nodes.

# of Players	Average Hops	Broadcast Messages	Mercury Messages	Optimal Messages
20	4.44	170000	74295	28154
40	4.61	695240	199076	58644

Table 1: Routing overheads for Caduceus, *without long pointers*.

## 7. CONCLUSION

In this paper, we have described the design and implementation of Mercury, a scalable protocol for routing multi-attribute range-based queries. Our contributions as compared to previous systems include support for multiple attributes and explicit load balancing. Mercury incorporates novel techniques to support random sampling of nodes within the system. Random sampling enables a number of lightweight approaches to performing load-balancing, node count estimation and query selectivity estimation. Our evaluation clearly shows that Mercury scales well, has low lookup latency and provides good load balancing properties.

In addition to providing high query-routing performance, Mercury provides a flexible range-based query primitive. Using this primitive, we are able to build an easy-to-use publish-subscribe facility for the maintenance of weakly consistent distributed state. We find that this facility is well suited for distributed state maintenance in distributed games. While we have only directly shown the suitability of Mercury for distributed games, we believe that the classes of applications that will benefit from our system include collaborative applications, such as shared whiteboards, distributed inventories and possibly sensing applications as well.

<sup>13</sup>We did not implement long-distance links because we were primarily interested in assessing the suitability of the Mercury service for distributed games. However, we were curious about performance.

## 8. ACKNOWLEDGEMENTS

We thank Justin Weisz for his initial work on adapting Caduceus to work with Mercury.

## 9. REFERENCES

- [1] BYERS, J., CONSIDINE, J., AND MITZENMACHER, M. Simple load balancing for distributed hash tables. *Second International Workshop on Peer-to-Peer Systems* (2003).
- [2] CABRERA, L. F., JONES, M. B., AND THEIMER, M. Herald: Achieving a Global Event Notification Service. In *Proceedings of the 8th IEEE Workshop on Hot Topics in Operating Systems* (Elmau, Germany, May 2001).
- [3] CARZANIGA, A., ROSENBLUM, D. S., AND WOLF, A. L. Design and Evaluation of a Wide-Area Event Notification Service. *ACM Transactions on Computer Systems* 19, 3 (Aug. 2001), 332–383.
- [4] CASTRO, M., DRUSCHEL, P., HU, Y. C., AND ROWSTRON, A. Exploiting network proximity in distributed hash tables. In *International Workshop on Future Directions in Distributed Computing (FuDiCo)* (June 2002), O. Babaoglu, K. Birman, and K. Marzullo, Eds., pp. 52–55.
- [5] CASTRO, M., DRUSCHEL, P., KERMARREC, A. M., NANDI, A., ROWSTRON, A., AND A., S. Splitstream: High-bandwidth multicast in a cooperative environment. In *Proceedings of the 19th Symposium on Operating System Principles* (Oct. 2003).
- [6] CASTRO M., ET. AL. SCRIBE: A large-scale and decentralized application-level multicast infrastructure. *IEEE Journal on Selected Areas in Communications (JSAC)* 20, 8 (Oct. 2002).
- [7] DABEK, F., KAASHOEK, M. F., KARGER, D., MORRIS, R., AND STOICA, I. Wide-area cooperative storage with CFS. In *Proceedings of the 18th Symposium on Operating System Principles* (Chateau Lake Louise, Banff, Canada, Oct. 2001).
- [8] GANESHAN, P., BAWA, M., AND GARCIA-MOLINA, H. Online Balancing of Range-Partitioned Data with Applications to Peer-to-Peer Systems. In *Conference on Very Large Databases (VLDB)* (2004). To appear.
- [9] GHOSH, B., LEIGHTON, F. T., MAGGS, B. M., MUTHUKRISHNAN, S., PLAXTON, C. G., RAJARAMAN, R., RICHA, A. W., TARJAN, R. E., AND ZUCKERMAN, D. Tight analyses of two local load balancing algorithms. In *Proceedings of the 27th ACM STOC* (1995), pp. 548–558.
- [10] GUMMADI, K., GUMMADI, R., RATNASAMY, S., GRIBBLE, S., SHENKER, S., AND STOICA, I. The Impact of DHT Routing Geometry on Resilience and Proximity. In *Proceedings of the ACM SIGCOMM '03* (Aug. 2003).
- [11] HARVEY, N. J. A., JONES, M. B., SAROIU, S., THEIMER, M., AND WOLMAN, A. Skipnet: A scalable overlay network with practical locality properties. In *Proceedings of the 4th USENIX Symposium on Internet Technologies and Systems* (Seattle, WA, Mar. 2003).
- [12] HEUBSCH, R., HELLERSTEIN, J., LANHAN, N., LOO, B. T., SHENKER, S., AND STOICA, I. Querying the Internet with PIER. In *Proceedings of the 29th International Conference on Very Large DataBases* (Sept. 2003).
- [13] KARGER, D., AND RUHL, M. Simple efficient load-balancing algorithms for peer-to-peer systems. *Third International Workshop on Peer-to-Peer Systems* (2004).
- [14] KLEINBERG, J. The Small-World Phenomenon: An Algorithmic Perspective. In *Proceedings of the 32th ACM STOC* (2000).
- [15] KOSTIC, D., RODRIGUEZ, A., ALBRECHT, J., BHIRUD, A., AND VAHDAT, A. Using random subsets to build scalable network services. In *Proceedings of the 4th USENIX Symposium on Internet Technologies and Systems* (Seattle, WA, Mar. 2003).
- [16] LI, X., KIM, Y.-J., GOVINDAN, R., AND HONG, W. Multi-dimensional range queries in sensor networks. In *Proceedings of the ACM Sensys 2003* (Nov. 2003).
- [17] MACEDONIA, M. R., ZYDA, M. J., PRATT, D. R., BRUTZMAN, D. P., AND BRAHAM, P. T. Exploiting reality with multicast groups: A network architecture for large-scale virtual environments. In *Proc. of the 1995 IEEE Virtual Reality Symposium (VRAIS95)* (Mar. 1995).
- [18] MANKU, G., BAWA, M., AND RAGHAVAN, P. Symphony: Distributed hashing in a small world. In *Proceedings of the 4th USENIX Symposium on Internet Technologies and Systems* (Seattle, WA, Mar. 2003).
- [19] MOTWANI, R., AND RAGHAVAN, P. *Randomized Algorithms*. Cambridge University Press, 1995.
- [20] RANDALL, D. Math 8213A - Rapidly Mixing Markov Chains. <http://www.math.gatech.edu/~randall/Course/lewis1.ps>, 2003.
- [21] RAO, A., LAKSHMINARAYANAN, K., SURANA, S., KARP, R., AND STOICA, I. Load Balancing in Structured P2P Systems. *Second International Workshop on Peer-to-Peer Systems* (2003).
- [22] RATNASAMY, S., FRANCIS, P., HANDLEY, M., KARP, R., AND SHENKER, S. A Scalable Content-Addressable Network. In *Proceedings of the SIGCOMM '01 Symposium on Communications Architectures and Protocols* (San Diego, California, Aug. 2001).
- [23] ROWSTRON, A., AND DRUSCHEL, P. Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems. In *IFIP/ACM International Conference on Distributed Systems Platforms (Middleware)* (Nov. 2001), pp. 329–350.
- [24] SINGHAL, S., AND CHERITON, D. Using projection aggregations to support scalability in distributed simulation. In *Proceedings of the 16th International Conference on Distributed Computing Systems* (1996).
- [25] STOICA, I., MORRIS, R., KARGER, D., KAASHOEK, F., AND BALAKRISHNAN, H. Chord: A scalable peer-to-peer lookup service for internet applications. In *Proceedings of the SIGCOMM '01 Symposium on Communications Architectures and Protocols* (2001).