

Improving Lookup Performance over a Widely-Deployed DHT

Daniel Stutzbach, Reza Rejaie
Department of Computer & Information Science
University of Oregon
{agthorr,reza}@cs.uoregon.edu

Abstract—During recent years, Distributed Hash Tables (DHTs) have been extensively studied by the networking community through simulation and analysis, but until recently were not adopted by popular P2P applications. Therefore, it was not feasible to examine DHT behavior in practice. Recently, the popular eMule file-sharing software incorporated a Kademlia DHT, called Kad, into their software. The success of Kad appears to have triggered other P2P applications (notably BitTorrent) to also incorporate Kademlia.

In this paper, we empirically study the performance of the key DHT function, lookup, over this DHT which has more than one million simultaneous users. We develop new analysis for predicting performance, characterize the accuracy of Kad’s routing tables, and empirically determine the optimum level of parallel lookup and the level of replication needed to cope with lookup inconsistencies. Our results show that compared to the current eMule client, overhead can be significantly reduced while simultaneously increasing performance.

I. INTRODUCTION

Distributed Hash Tables (DHTs) present an elegant distributed solution to deterministically map items to locations. They provide a *structured* approach to Peer-to-Peer (P2P) file-sharing applications since their item-to-location mapping can be used to (i) place data on specific peers and (ii) efficiently locate the data by identifying its corresponding peer. During the past few years, the potential of DHTs has motivated a wealth of research on different aspects of DHTs including the design of new DHTs [1]–[5], performance evaluation and improvement [6], [7], and the development of a wide range of DHT-based distributed applications [8], [9]. Despite a great deal of attention from the research community, DHTs have not been widely deployed in practice until recently. In the absence of any large scale deployment, all the previous studies on DHTs relied only on simulation, theoretical analysis, and limited-scale experiments. Therefore, the behavior of DHTs in practice has not been examined and thus is not well understood.

In practice, the dynamics of peer participation, or *churn*, can affect the accuracy of routing tables at each peer, and thus the performance of lookup operations in a DHT. More specifically, some entries in the routing table of individual peers might be missing or stale. Therefore, each peer does not have the expected connectivity to other peers. The inaccuracy of routing tables in turn affects the efficiency and consistency of lookup operations that are conducted by individual clients.

For example, a lookup may either take more than the ideal number of hops or could map to inconsistent endpoints.

There are two class of solutions to address the effect of churn on DHTs: (i) *DHT-based*: DHTs can incorporate various techniques to actively improve their resiliency to churn by increasing the degree of redundancy or frequency of updates for the routing table at each peer. (ii) *Client-based*: Alternatively, a client operating over an inaccurate DHT can improve its lookup efficiency by conducting lookup in parallel or cope with lookup inconsistency by active replication of its content.

Previous studies have examined both DHT-based [10], [11] and client-based [4], [5] solutions as well as the interactions and trade-offs between them [7]. All of the previous studies have used either simulation, analysis, or small-scale experiments to study these issues. However, churn, and thus the degree of inaccuracy in routing tables, is a real-world phenomenon generated by user behavior. Churn’s impact on the performance of lookup operations has not been empirically studied due to the limited deployment of DHTs. Given the limited understanding of churn characteristics, it is unclear how well simulation-based analysis of DHTs represent observed behavior in practice. Section VI discusses the related work in more detail.

This paper presents a measurement-based characterization of routing table inaccuracy and its impact on lookup performance in a widely deployed DHT, namely *Kad*. Kad is an open, Kademlia-based [4] DHT (with more than 1 million concurrent users) that has been recently deployed by the popular eMule¹ file-sharing application to improve efficiency of search in the face of growing population. Section II presents an overview of Kademlia and Kad.

To study the inaccuracy of routing tables in a DHT, Section III first establishes an analytical framework to quantify the effect of routing table richness on the performance of lookup. To our knowledge, this is the first analysis to show the direct benefits on lookup performance of Kademlia’s *k*-buckets. Then, we characterize both freshness and completeness of routing tables in Kad through detailed and representative measurement. To explain the observed behavior, we carefully examine the implementation of eMule² and identify

¹eMule began as an open-source alternative for the eDonkey unstructured network.

²There is no written specification that describes the Kad protocol so our explanations are based on our reading of the source code.

the underlying policies for routing table update and redundancy management that cause the observed behavior. Since we are dealing with a deployed DHT system, we are unable to explore DHT-based solutions to improve the accuracy of routing tables. Instead, we explore the design space of client-based solutions to improve both efficiency and consistency of lookup operations over DHT in practice.

Section IV examines client-based solutions to improve lookup performance using kLookup, a tool we developed to conduct a realistic lookup from *any* source ID to *any* destination ID without having local access to the designated peers for these IDs. Further leveraging the iterative lookup scheme in Kad, kLookup implements a spectrum of parallel crawling techniques, enabling us to study various parallel crawling techniques through real-world experiment over Kad and identify major design trade-offs.

Finally, Section V characterizes the degree of inconsistency in lookup operations in Kad and relates it to the underlying reasons of inaccurate routing tables. We then explore replication techniques to cope with lookup inconsistency and find that keeping three copies in the network dramatically improves lookup consistency.

Our main contributions can be summarized as follows:

- An analytical framework for computing the average performance of lookups
- The surprising result that redundancy in routing tables, such as Kademia’s k -buckets, directly improves mean lookup performance by reducing hop count
- kFetch, a tool for extracting the routing table from Kad peers
- kLookup, a parameterized tool for performing lookups over Kad using a variety of lookup algorithms
- Characterizations of the completeness and freshness of the routing tables in Kad
- Demonstrating experimentally that parallel lookup can reduce hop count, as predicted by our framework
- Experiments finding the optimum degree of parallelism ($\alpha = 3$) for use over Kad
- Experiments finding the optimum degree of replication ($c = 3$) to overcome routing table inconsistencies in Kad
- Methodology that can be applied to studying and tuning parameters for any widely deployed DHT.

While this study is inevitably centered around Kad, our analysis, methodologies, tools and findings are mostly applicable to other DHTs with proper adjustment. To address wider applicability of our work, we briefly discuss how some issues can be pursued in the context of other DHTs. Our extensive examination of eMule’s implementation also revealed several bugs in their implementation, some which will be fixed in the next revision.

II. BACKGROUND

Since the Kad network which we use for our empirical observations is based on Kademia, we first present some background on Kademia. Like most DHTs, peers in Kademia each have an identifier that is assigned either uniformly at

random or via a cryptographic hash. To determine the distance between two peers, Kademia uses a unique “XOR metric”, the bitwise XOR of their identifiers. For example, the distance between 0100 and 0111 is 0011 (or 3).

Kademia belongs to the general class of prefix-matching DHTs, such as Pastry [3] and Tapestry [12]. At the high-level, these DHTs work in the same way. A *lookup* consists of a sequence of *lookup steps*. The first step consults the client’s routing table for the target ID, which is guaranteed to have a route where the high-order b bits match. The next step is to consult that peer, which is guaranteed to have a route where the first $2b$ bits match. The process continues until no next route can be found, indicating that the closest peer to the ID has been reached. We can view the distance between two identifiers as the number of bits which must be matched, or $\log_2 n$ in the typical case, where n is the network population size. The expected number of steps is $\frac{\log_2 n}{b}$. We call b the *symbol size*, and in basic Kademia, $b = 1$. Section III examines the utility of different choices for b .

As IP is also a prefix-matching protocol, we borrow some terminology from IP to describe Kademia routing tables. Each line in a Kademia routing table is labeled with a subnet address and mask. When performing a lookup for a key, the most-specific routing table entry with a matching subnet is used, just as in IP routing. In this paper, the familiar “slash-notation” specifies masks.

In Kademia, the routing table contains one line per address bit, with increasingly specific masks. The subnet addresses are the same as the peer hosting the routing table. For example, consider a Kademia network using 4-bit identifiers³ and a particular peer with the address 0000. The routing lines are: /0, 0/1, 00/2, 000/3, 0000/4. Because more-specific routes are preferred, the effective routing lines are: 1/1, 01/2, 001/3, 0001/4. Put another way, the /0 line will only contain 1/1 addresses since any 0/1 address would go in one of the more specific lines. The routing table structure can be viewed as a binary tree, as shown in Figure 1(a).

Globally, the routing tables in all the Kademia peers form one large binary tree, with each peer containing a fraction ($O(\log n)$) of it. During a lookup, each routing step pivots to a different peer which is one bit closer to the target, guaranteeing that the lookup requires at most $O(\log n)$ steps.

For redundancy purposes, each routing line (or node in the binary tree) points to a list, called a k -bucket, of k matching contacts. Each contact includes the Kademia ID, IP address, and port of the remote peer. Thus, each lookup step has a choice of k different contacts. Section III examines some of the consequences for choosing different values of k . We note that k -buckets could be adapted for use in most varieties of DHT.

Kademia makes use of parallel routing to speed up lookups, as do EpiChord [13] and Accordion [5]. By sending out α lookup requests at a time, a client can avoid stalling by trying

³In practice no DHT would use such a small identifier space, but it’s more tractable for illustrative purposes.

to reach a departed peer and also increase the probability of finding low-latency peers. Section IV examines using different values of α in Kad.

Kademlia uses iterative routing, where the client is responsible for the entire lookup process. At each step, the client sends a *lookup request* to the next-hop peer and waits for a *lookup reply*. The reply lets the client know what the next hop is. Iterative routing contrasts with recursive routing, where the lookup request is forwarded automatically from one peer to another. While it has been shown that recursive routing typically has lower latency [14], iterative routing has several useful practical properties:

Fate-Sharing: Lookups cannot be lost through the departure of an intermediate carrying the lookup request.

Debugging: Iterative routing is easier to debug since information at each step is reported back to the client performing the lookup.

Compartmentalization: Iterative routing allows route table structure and lookup efficiency to be improved independently in a deployed network. Our tool, kLookup, uses this division to evaluate a variety of lookup techniques directly over the existing Kad network, as shown in Sections IV and V.

Route Table Extraction: It is possible to download the routing table of any peer, which we make use of in our tool, kFetch, described in Section III-B.

In summary, the key properties of Kademlia (and, thus, Kad) are as follows:

- Routing by prefix-matching
- Redundancy in routing tables (k -buckets)
- Parallel routing
- Iterative routing

Of these, redundancy, parallel routing, and iterative routing could each be added directly to most varieties of DHT. For example, EpiChord is a variant of Chord with parallel routing. Prefix-matching is an intrinsic property of Kademlia's design, which it shares with a number of other DHTs such as Pastry and Tapestry.

Kad is a Kademlia-based network for file-sharing, composed of eMule clients. It has approximately 1 million simultaneous users, plus many more firewalled peers who utilize the Kad DHT for lookups but do not participate in the DHT structure. For each file an eMule client is sharing, it computes the hash of each word in the filename, and publishes information about itself and the file to the peers responsible for the hashes. When an eMule user enters a keyword search, eMule computes the hash of the first keyword and initiates a lookup for the hash. The lookup returns a set of endpoints to which the client submits the full keyword list. Those peers process the query and return a set of matching results.

III. INACCURACY OF ROUTING TABLES

In this section, first we establish an analytical framework to examine the effect on lookup performance of adding redundancy to routing tables. This framework provides insight on

key trade-offs between different ways to increase the richness of routing tables and provides a formula for computing the average performance. Second, we empirically characterize the degree of routing table accuracy in Kad and identify the underlying reasons for inaccuracies. These characteristics help us explain the observed lookup performance in Section IV.

A. Benefits of Redundancy: An Analytical Framework

Every DHT has some structure that determines which neighbors a peer can choose from based on their identifiers. For example, in basic Kademlia a peer must have a neighbor with a different high-order ID bit, a neighbor with a matching first bit and a different second bit, a neighbor with the first two bits matching and a different third bit, etc. We call each address-mask pair a *bucket* (following the Kademlia terminology) and the neighbor information stored in the bucket a *contact*. In the base case, the DHT contains just enough information to perform the lookup in $\log_2 n$ expected hops. In prefix-matching systems such as Kademlia, this implies a symbol size of $b = 1$ and one contact per bucket.

A DHT can enrich the routing table structure beyond this base by either *adding more buckets* or *adding more contacts per bucket*. By adding more buckets, a DHT can guarantee that a larger number of bits will be improved at each step, thereby decreasing the number of hops for a lookup. For example, Pastry [3] uses a default symbol size of $b = 4$ which guarantees 4 bits will be improved at each step. Tables in Chord can also be enriched in this way [7].

Adding more contacts per bucket is used to guard against churn, an approach employed by DHTs such as Kademlia [4] and Tapestry [12]. By having other contacts handy, a peer can more quickly repair its routing table when a failure is detected. Also, as observed in [4], with heavy-tailed session times, storing backups and only evicting unresponsive peers implicitly leads to a set of peers with good uptime characteristics. Finally, multiple contacts allows for the use of parallel routing. A bucket with k contacts is called a k -bucket.

To examine the benefits and costs of the above two approaches for enriching routing tables, we analyze their impact in the context of Kademlia. Our analysis also directly applies to other prefix-matching systems such as Pastry and Tapestry, where we can quantify the improvement at each step in terms of the number of matching bits. For other DHTs that use a different basic geometry, our analysis could be adapted by modifying the formulas to reflect the appropriate distance metric.

There are two different approaches for adding more buckets to a routing table, both of which improve the number of lookup hops from $\log_2 n$ to $\log_2 b$:

Discrete Symbols: With this approach, illustrated in Figure 1(b), each interior node points to $b - 1$ buckets and an additional interior node. When searching a routing table, a peer begins by checking the first b bits. If all of them match the peer's ID, then it proceeds to the next b bits (*i.e.*, the next interior node). Otherwise, it proceeds immediately to the appropriate bucket. Using Discrete

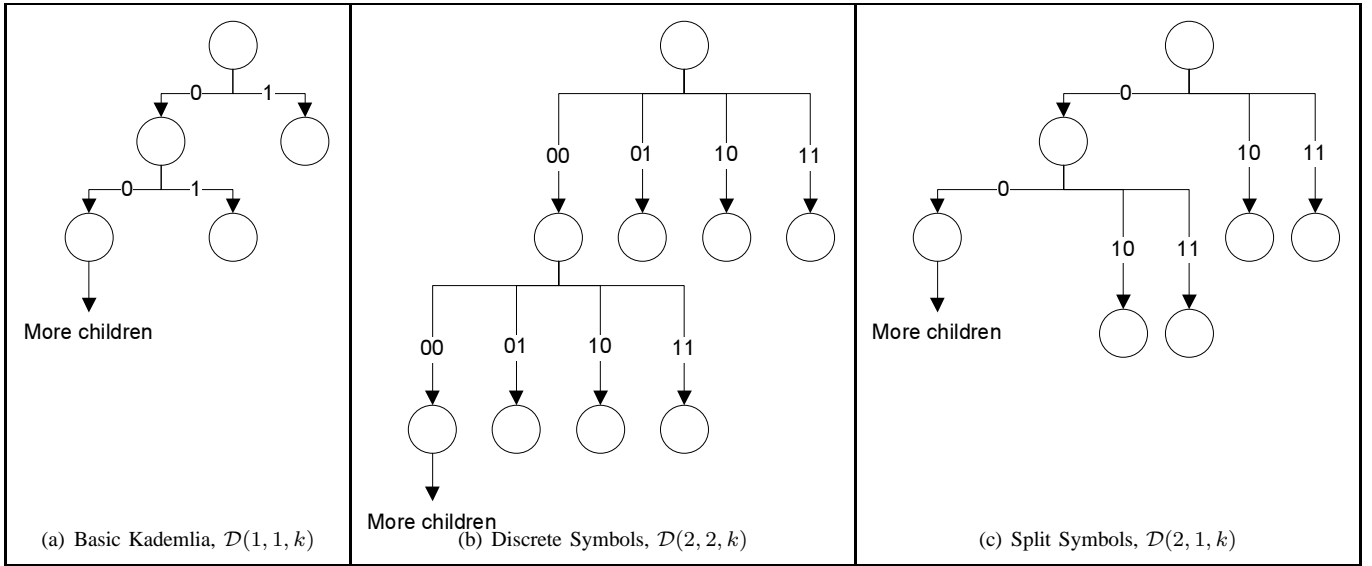


Fig. 1. Routing Table Structures

Symbols increases the routing table size from $\log_2 n$ rows of one k -bucket each to $\log_{2^b} n$ rows of $2^b - 1$ k -buckets each. This is the approach used in Kademlia and Pastry.

Split Symbols: With this approach, illustrated in Figure 1(c), each interior node points to 2^{b-1} buckets and an additional interior node. When searching a routing table, a peer begins by checking the first single bit. If it matches the peer's ID, then it proceeds to the next bit (*i.e.*, the next interior node). Otherwise, it examines the next b bits and proceeds to the appropriate bucket. Using Split Symbols increases the routing table size for $\log_2 n$ rows of one k -bucket each to $\log_2 n$ rows of 2^{b-1} k -buckets each. This is the approach used in Kad.

To compare and contrast these approaches for organizing routing table contacts, we create a general framework for analyzing their performance. We define $\mathcal{D}(b, r, k)$ as a system which uses b -bit symbols with r -bit resolution and k -buckets. $\mathcal{D}(1, 1, k)$ is the basic Kademlia approach, $\mathcal{D}(b, b, k)$ is the Discrete Symbol approach, and $\mathcal{D}(b, 1, k)$ is the Split Symbol approach used in Kad. Each routing table has $\log_{2^r} n$ rows of $2^b - 2^{b-r}$ k -buckets, for a total size of $k(2^b - 2^{b-r}) \log_{2^r} n$ contacts. Normalizing by a factor of $\log_2 n$ yields a normalized size of $k \frac{2^b - 2^{b-r}}{r}$.

Prior work has concerned itself exclusively with the worst-case scenario where the selected contact will not match any *additional* bits of the target identifier. For example, consider searching for the key 111 in the routing table of peer 000 with the base $b = 1$ system. The peer looks in the bucket with the prefix 1, and returns a contact which we know matches the first bit of the key. However, that contact could be any of the peers 100, 101, 110, or 111. In other words, there's a $\frac{1}{2}$ chance of improving at least 1 extra bit, a $\frac{1}{4}$ chance of improving at least 2 extra bits, and so on. More precisely, the probability

of improving at least δ bits is:

$$Pr[X \geq \delta] = \frac{1}{2^\delta} \quad (1)$$

Therefore, the average-case is better than the worst-case given in prior work. In particular, the key insight is that large buckets ($k > 1$) improve the probability of *randomly* finding a contact with more matching bits since there are more options to choose from. As we will show, the average number of bits improved increases logarithmically with k , making the performance boost of increasing k comparable to the performance boost of increasing b . Generally, for a k -bucket the probability of improving by at least δ extra symbols is:

$$F(\delta, r, k) = Pr[X \geq \delta] = 1 - \left(1 - \frac{1}{2^r} \right)^k \quad (2)$$

and the probability of improving *exactly* δ symbols is:

$$f(\delta, r, k) = Pr[X = \delta] = F(\delta, r, k) - F(\delta + 1, r, k) \quad (3)$$

The key question is: how many bits improve on average? This is equal to the product of r and the mean value of $f(\delta, k)$ as a function of k as follows:

$$m(r, k) = r \frac{\sum_{\delta=0}^{\infty} \delta \cdot f(\delta, r, k)}{\sum_{\delta=0}^{\infty} f(\delta, r, k)} \quad (4)$$

While we were unable to find a simple closed form for $m(r, k)$, it can be computed numerically without difficulty. For $r = 1$, $m(1, k)$ asymptotically approaches $\log_2 k + 0.3327$, somewhat exceeding this value for low k . Significantly, $m(1, 1) = 1$, indicating that on average $r = 1$ systems improve one extra bit per step even with no redundancy! In other words, a basic $\mathcal{D}(1, 1, 1)$ system on average performs a lookup in *half as many hops* as reported by previous work.

For a Discrete Symbol configuration, $\mathcal{D}(b, b, k)$, the number of bits improved on average is $b + m(b, k)$. For a Split

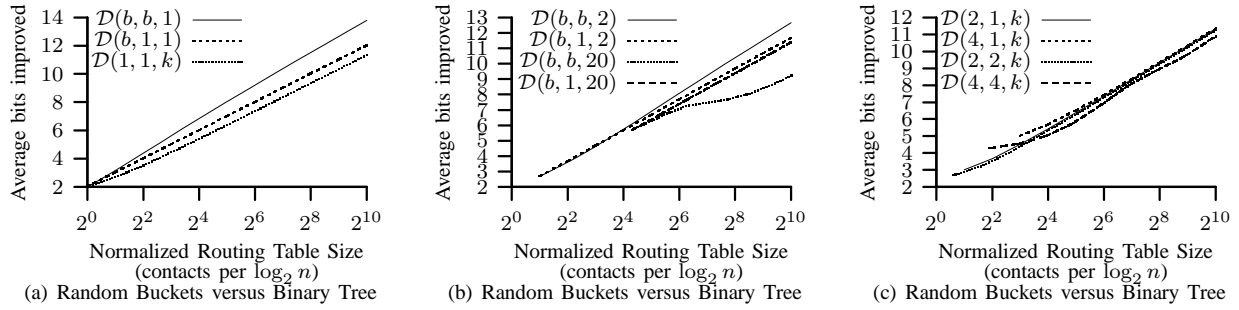


Fig. 2. Relative performance of different routing table structures

Symbol configuration, $\mathcal{D}(b, 1, k)$, the number of bits improved on average is $b + m(1, k)$. In other words, while the Split Symbol approach does use more routing table space, it has the advantage that it can leverage a random improvement of a single bit; the Discrete Symbol approach must randomly improve by b bits to make use of the improvement.

To compare the different approaches, first consider the three extreme cases: $\mathcal{D}(1, 1, k)$ (pure Redundancy), $\mathcal{D}(b, b, 1)$ (pure Discrete Symbols), and $\mathcal{D}(b, 1, 1)$ (pure Split Symbols). Figure 2(a) presents the performance of each approach as a function of the normalized routing table size. Split Symbols and Redundancy have nearly identical performance, while Discrete Symbols performs slightly better.

For the case of Split Symbols ($\mathcal{D}(b, 1, 1)$), the b -bit symbols guarantee an improvement of b bits in the worst case, plus an additional $m(1, 1) = 1$ bits on average, for a total of exactly $b + 1$ bits, dividing by the normalized size yields $\frac{b+1}{2^b-1}$. This is the slope of the Split Symbols ($\mathcal{D}(b, 1, 1)$) line in Figure 2(a).

For the case of Discrete Symbols ($\mathcal{D}(b, b, 1)$), the b -bit symbols again guarantee an improvement of b bits in the worst case, plus an additional $m(r, 1)$ bits on average. However, $m(r, 1)$ asymptotically approaches 0 for large r . As a point of reference, for Pastry's typical value of $b = r = 4$, the average improvement is 4.27 bits per step, roughly a 4% reduction in the mean number of lookup hops⁴ compared to that reported by the Pastry authors [3]. The average improvement divided by the normalized size is $b \frac{b+m(r,1)}{2^b-1}$.

For the case of large buckets and 1-bit symbols ($\mathcal{D}(1, 1, k)$), the 1-bit symbols guarantee an improvement of 1 bit in the worst case, plus an additional $m(1, k)$ bits on average, for a total of $1 + m(1, k)$. Dividing by the normalized size results in $\frac{1+m(1,k)}{k}$. As a point of reference, for the value of $k = 20$ suggested in the Kademia paper [4], the average improvement is 5.7 bits per step rather than 1 bit per step, resulting in a 60% reduction in the mean number of hops!

Can performance be improved by using a mixture of large buckets and large symbols? Not really. Figures 2(b) and 2(c) plot several other permutations of $\mathcal{D}(b, r, k)$. Figure 2(b) holds k constant and varies b , while Figure 2(c) holds b constant and varies k . For small values of k (e.g., 2) with varying b ,

both Discrete Symbols and Split Symbols have performance in between their regular performance and $\mathcal{D}(1, 1, k)$. For moderate values (e.g., 20) of k , the performance of Split Symbols is virtually identical to $\mathcal{D}(1, 1, k)$, while the performance of Discrete Symbols plummets (as seen in Figure 2(b)). Because Discrete Symbols can't make good use of randomness, the k -redundancy imposes a cost with little benefit on lookup performance.

In summary, increasing the symbol size offers a constant-factor improvement to worst-case performance, while using k -buckets unexpectedly offers comparable average-case improvement. Moreover, k -buckets offer other advantages:

- Reduced implementation complexity
- Lower maintenance bandwidth; fewer restrictions on acceptable contacts allows for more contacts to be acquired passively
- Better resistance to churn by accumulating high-quality contacts

While our framework is motivated by our study of Kad, it applies to any prefix-matching DHT and could be extended to other DHTs that can accommodate different symbol or bucket sizes, such as Chord. In the following section, we use the formulas we have developed to compute a bound on the performance of Kad and empirically examine how close the actual performance is to our model's prediction.

B. Accuracy of Routing Tables in Kad

Our goal is to explore the structure and redundancy (i.e., b and k) of routing tables in Kad by examining eMule source code, and then empirically study the impact of churn on the accuracy of routing tables in Kad.

Analysis of Accuracy/Redundancy: Close examination of the eMule 0.46a source code along with empirical validations revealed that Kad is based on Kademia with a bucket size of 10 contacts ($k = 10$), and 3.25-bit Split Symbols⁵ which means that Kad is a $\mathcal{D}(3.25, 1, 10)$ system. Therefore, according to Equation 4 the mean number of improved bits per step is 6.98 in Kad. Additionally, as a special case Kad always

⁵The $\frac{1}{4}$ bit is due to the fact that Kad uses unbalanced subtrees. Each interior node has branches with labels 0, 1000, 1001, 101, 110, and 111. The 0 branch leads to the next interior node; the other branches lead to k -buckets. The average improvement per step is $3.25 + m(1, k)$ bits.

⁴The expected number of hops is equal to $\log_B n$ where B is the average number of bits of improvement.

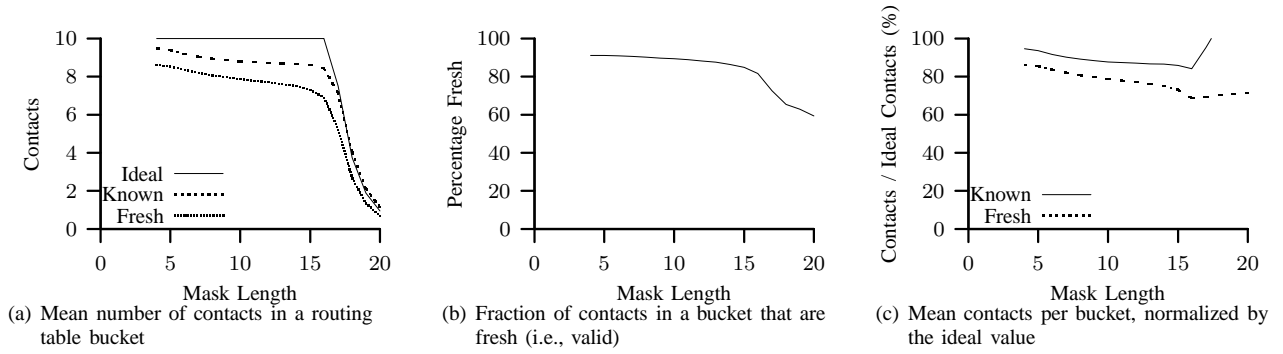


Fig. 3. Mean observed routing table bucket properties, as a function of how specific the bucket’s mask is

improves at least the 4 most significant bits on the first step, or 7.73 on average for a maximally distant node⁶. Thus, the expected number of hops for a maximally distant node in Kad is $\frac{\log_2(n)-7.73}{6.98} + 1$. To estimate the number of hops from this equation, we determine n later in this subsection.

Characteristics of Accuracy: We examine two dimensions of inaccuracy in the routing tables in Kad as follows:

Completeness: represents whether the routing table contains the appropriate number of entries, given the size of the Kad network.

Freshness: captures the number of contacts in the routing table that are still active (*i.e.*, do not point to a departed peer).

Toward this end, first we need to estimate the size of the Kad network (n). In our previous work [15], we developed a parallel peer-to-peer overlay crawler, called *Cruiser*. Crawling the *entire* Kad network takes too long due to the large size of routing tables at each peer and the large number of peers in the network. Because churn occurs while the crawler runs, a long crawl results in an inflated population count as it records a large number of short-lived peers which were not simultaneously present. However, since Kad identifiers are selected uniformly at random, any subnet of the ID space is a representative sample of the total population, and a subnet can be crawled much faster than crawling the whole network. Multiplying the size of the subnet by the number of such subnets yields an estimate of the population size. By taking the mean over many such samples, we can get a good estimate for n . Moreover, our formulas are not sensitive to minor fluctuations in population size as they are based on $\log_2 n$.

Given a Kad overlay subnet as an input (*e.g.*, 0x5cd/12), *Cruiser* walks the DHT structure to capture a snapshot of all the active peers with IDs in the specified subnet. For example, it can capture a /10 subnet with roughly 1000 peers in around 3–4 minutes and a /12 subnet with roughly 250 peers in around one minute. During June of 2005, we captured the population size for several hundred randomly selected subnets with *Cruiser*. Our measurements reveal that the Kad network has a mean population size of approximately 980,000 concurrent peers. Given this estimated group size for the Kad

network, a lookup over Kad requires $\log_2 980,000 \approx 19.9$ bits improvement per step, and a lookup in Kad should take $\frac{19.9-7.73}{6.98} + 1 = 2.7$ hops, assuming perfect routing tables. This is significantly better than predicted by prior techniques of $\frac{\log_2(n)-4}{3} + 1 = 6.30$ hops. Correctly taking into account the effect of randomness alters the predicted performance by *more than a factor of two*.

kFetch: To study the accuracy of routing tables, we developed a new tool called *kFetch*. *kFetch* chooses a Kad peer at random, downloads its complete routing table, and identifies stale entries in the routing table by actively probing (*i.e.*, sending a lookup request) to each contact in the routing table. The routing table of the target peer should be downloaded quickly in order to minimize any error in our characterization due to on going churn (*i.e.*, a contact that is considered stale might be actually present in the network). Each Kad node only replies to a lookup request from participating peers to lookup an ID. There are two challenges to download a routing table efficiently: (*i*) the rate of requests (which are UDP messages) should be properly paced to rapidly download the table without causing any packet loss, and (*ii*) lookup messages should properly select requested IDs from a target peer to extract its table with the minimum number of messages. *kFetch* implements NewReno-style congestion control to determine the proper rate for issuing requests to the target node. Furthermore, *kFetch* sends requests for each possible routing table line, a strategy which can be used on any DHT that uses iterative routing. Additionally, for each contact learned, *kFetch* pings the contact to verify whether it is still present in the network, concurrently with continuing to download the routing table. To locate a peer at random, *kFetch* generates a random Kad identifier, then performs a Kad lookup to locate the peer closest to that Kad identifier.

Characteristics of Kad Tables: Using *kFetch*, we retrieved the routing tables of approximately 80,000 distinct Kad peers in June 2005 and examined their freshness and completeness. Figure 3(a) shows the mean number of contacts in each routing table bucket as a function of the bucket’s subnet mask. It also shows what fraction of these contacts are fresh (*i.e.*, the contact responded to our ping). The “Ideal” line indicates the average number of contacts in the bucket if we had global information, based on the population size, *i.e.*, $\min(10, \frac{n}{2^x})$

⁶The root node has a full 16 branches.

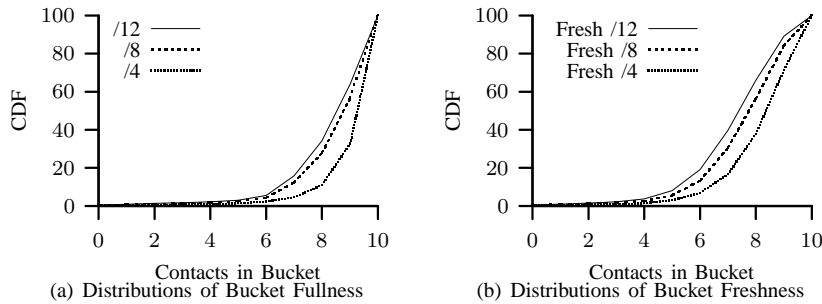


Fig. 4. Distributions of Bucket Contents

where x is the number of bits in the address mask and n is the population size. All three curves (Ideal, Known, and Fresh) curve off steeply when the mask length exceeds 16 bits, which is where the number of contacts is limited by the number of available contacts in the system. For shorter masks, on average each bucket has one or two empty slots and contains one stale contact. The mean number of empty slots is slightly higher as the mask length increases. Figure 3(b) shows the mean number of fresh contacts as a fraction of the number of contacts actually present. This shows that around 90% of entries are fresh for masks up to length around 16, then the fraction of fresh entries decreases (meaning that the number of stale entries increases). This is because the current implementation of eMule doesn't ping peers in buckets which are not at least 70% full. In fact, in Figure 3(c), where we examine the number of contacts relative to the ideal number, above /17 there are actually more stale contacts than active peers anywhere in that subnet! As each peer is up, it gradually accumulates stale contacts in these buckets which are expunged too slowly. As a consequence virtually every lookup in Kad necessarily ends with timeouts to stale peers even though the closest peer has already been contacted! As this routing table maintenance problem can trivially be corrected by pinging all buckets instead of just those mostly full, in the remainder of this paper we consider only the time required to hear from the closest responsive peer.

From Figure 3(a), we see that on average there are 1.5 empty slots plus 1 stale contact per bucket. We could plug $k = 10 - 1.5 - 1 = 7.5$ into our formula, but first we must validate that most buckets are close to the average state. If the variance is very high, for example if 85% of buckets had 10 entries and the other 15% were completely empty, then using the average would introduce considerable error. Towards this end, Figures 4(a) and 4(b) present the CDF of the number of contacts and fresh contacts across all observed buckets for masks /4, /8, and /12. They show that for both completeness and freshness, nearly all buckets are close to the average value. Therefore, we may use the average value for the purposes of our computations without introducing considerable error.

Using an average of 1.5 empty slots plus 1 stale contact per bucket, we have an effective bucket size of $k = 10 - 1.5 - 1 = 7.5$. This increases the expected hop count to $\frac{\log_2(n) - 7.33}{6.58} +$

$1 = 2.91$ hops. This is still significantly better than the value predicted by prior analytical techniques of 6.30 hops.

Note that we are unable to change the routing tables in the entire Kad network. Therefore, we explore client-based alternatives to improve efficiency and consistency of lookup in Kad and evaluate their performance in the following two sections, respectively.

IV. LOOKUP EFFICIENCY

We turn our attention to client-based approaches to improve the performance of iterative lookup over a DHT that has inaccurate routing tables. While incomplete buckets will degrade performance as described in the previous section, stale contacts can dramatically increase latency by causing timeouts to occur. Since the timeout interval is typically set to at least a few round-trip times, it can easily exceed the desired time for the entire lookup.

To improve performance despite inaccurate routing tables, clients (*i.e.*, end-points) can perform parallel lookup. While parallel lookup has traditionally been used exclusively with iterative DHTs, Jinyang Li *et al.* [5] present a technique for performing parallel lookup on a recursive DHT.

In a parallel lookup, a client simultaneously manages multiple lookup requests to different peers and performs the lookup process based on the information obtained from all requests. The inherent redundancy in the information collected by parallel lookup reduces the problem of hitting stale contacts, improving lookup performance at the cost of higher network overhead (*i.e.*, a larger number of requests per lookup). Parallel lookup has two other significant advantages. First, lookup requests facilitate populating or passively updating the routing tables, which in turn reduces the bandwidth requirement for explicit updates, as shown in [7]. Second, during each step of the lookup process, parallelism increases the number of contacts searched, increasing the probability of finding a contact closer to the target (*i.e.*, with more matching bits) and thus decreasing the number of hops needed to reach the target.

To examine different lookup strategies, we developed a new tool, called *kLookup*, which performs a lookup from *any* source ID to *any* destination ID without requiring local access to those peers. To emulate a lookup from a source ID, *kLookup* takes the following steps. First, it uses a local Kad routing table to locate the peer closest to the source ID (*i.e.*, the source

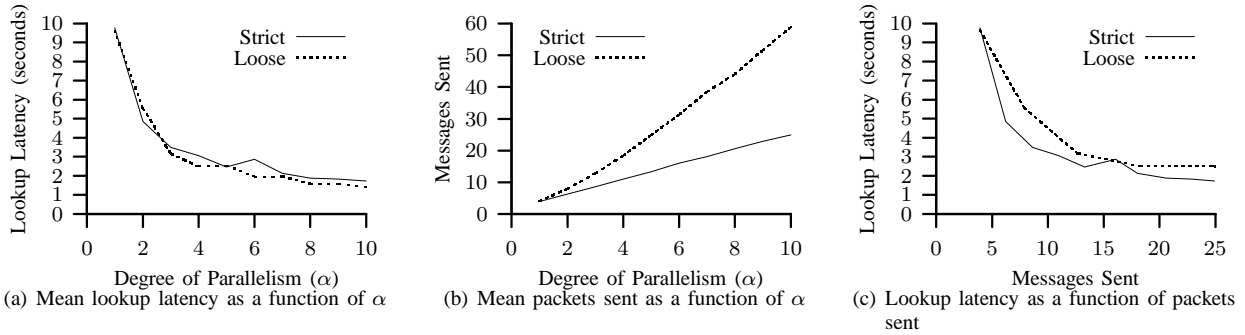


Fig. 5. Strict versus Loose Lookup Parallelism

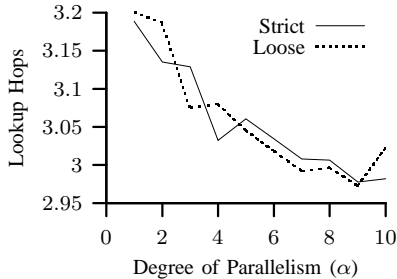


Fig. 6. Parallel lookup reduces hop count

peer), then it extracts the routing table of the source peer using kFetch. Finally, it performs a lookup to the destination ID using the routing table of the source peer. kLookup implements an adjustable degree of parallelism (α) with the following two classes of parallel lookup techniques: (i) Strict Parallel lookup and (ii) Loose Parallel lookup.

Strict Parallel Lookup: In this approach, a client begins a lookup by sending lookup requests to the α best known contacts. Similar to the window-based congestion control in TCP, a client restricts the number of requests in-flight to α . A new request is issued only when a pending request times out or a response is received. The resulting overhead is limited to a factor of α . The downside of the strict approach is that when a client sends a packet to a departed contact, it must wait for a timeout to occur before giving up. In the meantime, the degree of parallelism is effectively reduced by one. However, a timeout is typically set to at least a few round-trip times which is on the order of the desired time for the entire lookup. Thus, in the strict approach, α roughly determines the number of timeout events a client can experience without experiencing a significant latency penalty. Kademlia uses this approach.

Loose Parallel Lookup: Parallel lookup can be performed in a looser fashion by allowing more than α requests in flight. In this approach, a client can issue a lookup request to a contact that is among the top α contacts as soon as such a contact is identified, even if this lookup request increases the number of pending requests beyond α . For example, if $\alpha = 3$, the lookup begins by sending 3 lookup requests. If the first response contains 3 better contacts (which is likely), 3 more

requests are sent immediately. While this approach appears to be significantly more expensive than strict parallel lookup, it incurs only modest additional overhead since later responses from the same step are less likely to contain better contacts (*i.e.*, each time a packet is sent, the bar has been raised). The advantage of this looser approach is the ability to quickly abandon lookups that are likely to time out. This approach is used by eMule.

We evaluated the performance of both types of parallel lookup techniques under varying degrees of parallelism. Using kLookup, we captured several hundred lookups for different values of α for both strict and loose parallelism. Each lookup used a unique, randomly selected source and a unique, randomly selected destination. In our evaluation, we examine three metrics:

Hops: The number of hops from the source to the destination

Latency: The duration from the start of the lookup to when a response is received by the final destination, which is a function of the number of hops and the time spent waiting for responses and timeouts

Messages Sent: The overhead used to perform the lookup

Earlier, we mentioned that increasing α can reduce the number of lookup hops by providing more opportunities to randomly improve extra bits. Figure 6 provides empirical support by showing that the mean number of hops decreases slightly as α increases⁷. Furthermore, the hop count for $\alpha = 1$ is around 3.2, which is close to our predicted value of 2.9.

Since the number of hops is as expected, the next question is: how much latency is introduced by timing out while trying to contact stale contacts? Figure 5(a) compares the latency of the two approaches for several values of α . The first observation is that the latency for $\alpha = 1$ is very high—close to 10 seconds. Using a value of $\alpha = 3$ dramatically reduces the latency, with diminishing returns for larger α . Second, Figure 5(a) reveals that the loose approach is just barely quicker for constant α . The greatest advantage of loose parallelism is that it is significantly less likely to get stuck waiting for timeouts to occur. However, as we show in Sec-

⁷This figure is particularly noisy due to the small y -axis scale. The general downward trend is nevertheless visible.

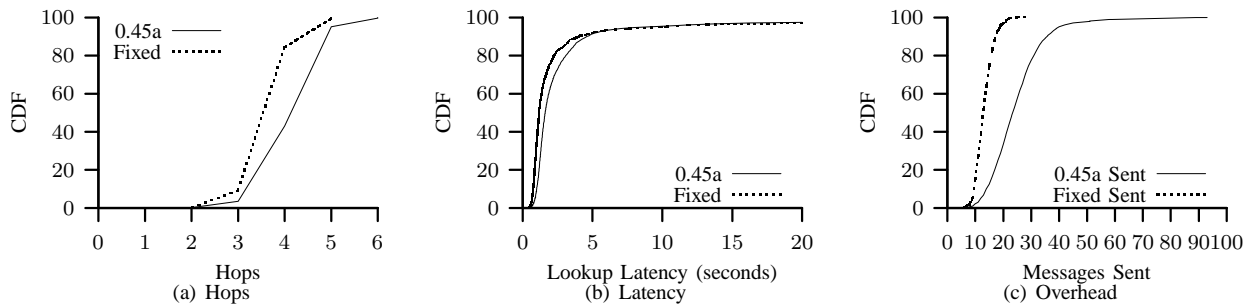


Fig. 7. Efficiency metrics across many lookups

tion III, few contacts in Kad are stale. This explains why loose parallelism does not show much performance improvement for this network.

To examine the overhead of parallelism, Figure 5(b) shows the number of packets sent as a function of α for the two approaches. In both cases, the overhead increases roughly linearly with α , with loose approach generating roughly twice as many messages as the strict approach. Given that for fixed α the performance of strict and loose parallelism are quite similar, this suggests that strict parallelism is the better choice for the current Kad network. To more directly compare the two, Figure 5(c) factors out α by plotting the performance as a function of the overhead. From this, we see that asymptotically the performance of strict and loose parallelism are surprisingly similar. A large number of messages represents the upper-bound on performance: no amount of increased parallelism will significantly improve it. At the low-end, the two perform the same since the two approaches result in identical behavior for the special case $\alpha = 1$. However, the sweet-spot for strict parallelism ($\alpha = 3$) is significantly better than the sweet-spot for loose parallelism.

In summary, these observations show that strict parallelism with $\alpha = 3$ is a good choice for the current Kad network. Higher values of α and loose parallelism substantially increase overhead without much change in performance. Also, our findings in Figure 6 provide evidence of the correctness of our analysis in Section III.

Comparing with eMule: As part of creating kLookup, we also attempted to exactly reimplement eMule 0.46a’s lookup algorithm. We validated this mode of kLookup by extending tcpdump to decode Kad packets and performing lookups for the same key using kLookup and eMule itself to verify their similarity. In the process of implementing eMule’s lookup algorithm, we discovered a few bugs which significantly degrade its⁸ efficiency.

As part of our study, we wanted to compare the performance of eMule’s current lookup algorithm with and without the bugs, in the hope that it will be of use to the eMule developers. Again, we example performance in terms of hops and latency

⁸Our results are based on eMule version 0.46a, the most recent version available at the time of our study. We have been corresponding with the eMule developer team regarding these discoveries, and at least some of the bugs we report will be corrected in 0.46b.

and overhead in terms of the number of messages. These experiments are based on more than one-thousand experiments of using kLookup from unique and randomly selected sources and destinations. With the bugs fixed, eMule’s lookup algorithm is $\alpha = 3$ with loose parallelism.

Figure 7(a) presents a CDF of the number of hops to perform a lookup. The mean value is 3.59, somewhat worse than our ideal value of 2.91. Without bugs, the number of hops drops to 3.08, which is closer. Figure 7(b) shows the latency of the two versions. In both cases, there is a significant tail (not shown) out to around 70 seconds. We see that the fixed version improves by around 1 second in most cases. The most striking difference however is in the overhead, as shown in 7(c). The fixed version uses roughly half as many messages on average.

V. LOOKUP CONSISTENCY

Ideally, each peer in a DHT is responsible for a certain part of the DHT geometry and all lookups for any point in that space lead to that peer. In practice, peer churn causes two types of routing tables inaccuracies:

- 1) Peers may not yet have points to a recently arrived peer
- 2) Peers may have extraneous pointers to a recently departed peer

When routes are inconsistent, it may not be possible to find information. The extent of these problems is determined by how frequently the DHT checks its pointers, known as route stabilization, compared to the rate of churn in the system. One approach to minimize these problems is to increase the frequency of route stabilization. However, this significantly increases the bandwidth required for route maintenance.

An alternative approach is to map each identifier to the set of the c closest peers, rather than to just the single closest peer. The publishing operation performs a regular lookup, then searches the surrounding area to find the closest c . The search operation does the same, and as long as the two find any peer in common the search will succeed. Kademia [4] takes this approach as a basic principle; however, it can be used in any DHT. For example, DHash [8] implements this technique over Chord. The parameter c must be chosen, based on knowledge of the degree of routing table inaccuracy, to guarantee with high likelihood that multiple lookups will be able to find peers in common.

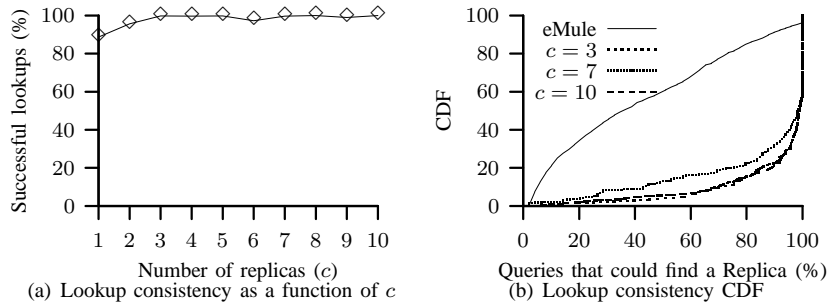


Fig. 8. Lookup Consistency

The key question is: what is the right value of c to guarantee a certain level of reliability p ? In the following subsection, we use empirical techniques to answer this question for Kad.

A. Achieving lookup consistency

To explore lookup consistency, we extended kLookup to locate the c closest points after its regular search has completed. To get an empirical measure for p , we use kLookup to perform 50 lookups to the same key, each from a different starting point in the Kad network. The first lookup is a publish operation which returns a set of peers to publish on. The following lookups to the same key are query operations, which return a set of peers to receive the actual query. Computing the fraction of queries that successfully find one of the target peers yields an empirical measure of the consistency, p , for that experiment. We perform the lookups as concurrently as possible to limit the effects of peer departure and arrival, to narrow the focus to issues of inconsistent routing tables. For these experiments, we used strict parallelism with $\alpha = 3$. We conducted this experiment 20 times for each value of c in the interval $[1, 10]$ (*i.e.*, 1000 lookups per value of c).

We observed that for $c = 1$, the consistency is only 89%, meaning that 11% of the time queries fail to find the same “closest” peer as a publisher. To explore how many replicas are needed, Figure 8(a) plots p as a function of c . For the value $c = 3$, the consistency is over 99.9% across the twenty 50-lookup trials. For $c = 2$, the consistency is in between, at around 96%.

The above values are for finding *any* of the replicas. However, another issue regarding consistency is how effectively all of the replicas can be found. If one replica can always be located, but the others cannot be, then lookups will fail if the one easy-to-locate replica fails. Therefore, for each replica we compute the number of lookups that found it, and plot it as a CDF in Figure 8(b). An ideal curve would be a vertical line at $x = 100\%$, indicating that every query found every replica. The Figure shows that the performance for the nearby-replication method is indeed good, with roughly 50% of queries able to find every replica, and 80–90% of queries able to find 80% of the replicas.

In summary, our results show that locating the three closest nodes after finding the closest peer is an effective way to cope with routing table inconsistencies. More importantly, we show

that even routing table inconsistencies can be a considerable problem in practice with more than 11% of lookups failing when no replication is used.

Comparing with eMule: Currently, eMule uses a fuzzy algorithm which selects several peers as part of the endpoint set, but which are not necessarily the closest. In addition to our experiments for different values of c , we also conducted more than 60 experiments using eMule’s algorithm⁹. We found that eMule’s approach produces 19 replicas on average and queries succeeds 99.9% of the time. While robust, this is 6.3 times more replicas than simply using $c = 3$. Furthermore, Figure 8(b) shows the CDF of the number of replicas each lookup found. The performance is substantially worse than the nearest- c approach, with many replicas being found by only a few queries. For example, 50% of replicas could be found less than one-third of the time, compare to just 3% for $c = 3$. Additionally, some replicas were not found by *any* queries.

VI. RELATED WORK

Early work on DHTs focused on introducing new DHTs [1]–[4] that each achieved $O(\log n)$ lookup hops using $O(\log n)$ state. Initially, it was difficult to directly compare the performance of these DHTs, as each DHT has several tunable parameters which might cause it to perform better or worse under different loads. For example, under low churn a DHT with a large routing table will perform better since it can achieve faster lookups and route maintenance is inexpensive. The same DHT will perform poorly under heavy churn.

Several studies [7], [10], [16]–[19] have attempted to address the issue of DHT performance under churn, in most cases using a simple Poisson model for session length. However, several measurement studies of peer-to-peer systems [20]–[24] show that session times follow a heavy-tailed distribution. In this study, we conduct experiments using the real Kad network, *i.e.*, under real churn.

Gummadi *et al.* [6] showed that DHTs can be broken into two components: geometry (or structure) and lookup strategy. Some DHT geometries provide greater routing flexibility than others in terms of neighbor selection or route selection. For example, in CAN a peer’s neighbors are precisely defined by the geometry, while in Chord there are 2^{i-1} options for

⁹For these experiments we also used eMule’s lookup techniques.

the i^{th} neighbor, providing Chord substantially more flexibility in selecting neighbors. Their results show that more flexible systems, such as Chord and Kademlia, can achieve better performance. We utilize their division between geometry and lookup to study the lookup behavior in light of the geometry of the deployed Kad network.

Jinyang Li *et al.* developed a performance-versus-cost framework (PVC) for comparing different DHTs [7]. Their key observation is that for a given bandwidth usage, there is a minimum lookup latency that can be achieved over the entire space of DHT parameters, and vice versa. In PVC, they simulate each DHT using a wide variety of parameters and plot the best lookup latency each DHT can achieve within a given bandwidth constraint. This allows them to compare how different DHTs make the performance-versus-cost trade-off under a given load. They show that using large routing tables with infrequent stabilizations and parallel lookup achieves a better balance than other approaches, culminating in their later development of the Accordion DHT [5]. However, PVC can only draw conclusions about how well the DHTs respond to the simulated workload. While their work is useful for drawing inferences about design trade-offs, our work is aimed at optimizing tunable parameters in a DHT that is already deployed.

In summary, prior work on DHTs has been driven by analysis, simulation, and limited experiments. In each case, a model is used to approximate or estimate real-world behavior. This paper presents experiments on a deployed DHT that has approximately one million real users, and develops tools and techniques for improving its performance.

VII. CONCLUSIONS AND FUTURE WORK

This paper examines lookup performance over the Kad DHT network. We analytically derive new formulas for the expected hop count, taking into account random improvements, and demonstrate that Kademlia's use of k -buckets leads to significantly better performance than previously reported. We present new tools, kFetch and kLookup, to characterize the accuracy of routing tables in Kad, examine the impact of accuracy on efficiency and consistency of the lookup operation, and experimentally verify our analysis. Furthermore, we explore two types of parallel lookup techniques and their impact on lookup efficiency and also examine the degree of replication needed to cope with routing consistency. While some of our empirical results are specific to Kad, our analysis applies to other prefix-matching DHTs such as Pastry and Tapestry and could be modified to handle other DHT geometries.

In our future work, we plan to measure the bandwidth eMule uses for route maintenance, and conduct a similar study in which we explore ways to maintain higher quality routing information at lower cost. We also plan to use our recent measurement-based characterization of churn in peer-to-peer systems [25] to determine the number of replicas needed to guarantee the availability of a piece of data within the network. This will include a mathematical analysis of the trade-off between republishing the data more frequently to a few peers

versus publishing infrequently to many peers, followed by empirical experiments to validate our findings.

REFERENCES

- [1] I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, and H. Balakrishnan, "Chord: A Scalable Peer-to-peer Lookup Service for Internet Applications," in *SIGCOMM*, San Diego, CA, Aug. 01.
- [2] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker, "A Scalable Content-Addressable Network," in *SIGCOMM*, 2001.
- [3] A. Rowstron and P. Druschel, "Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems," in *IFIP/ACM International Conference on Distributed Systems Platforms (Middleware)*, Heidelberg, Germany, Nov. 2001, pp. 329–350.
- [4] P. Maymounkov and D. Mazieres, "Kademlia: A Peer-to-peer Information System Based on the XOR Metric," in *International Workshop on Peer-to-Peer Systems*, 2002.
- [5] J. Li, J. Stribling, R. Morris, and M. F. Kaashoek, "Bandwidth-efficient Management of DHT Routing Tables," in *Networked Systems Design and Implementation*, Boston, MA, May 2005.
- [6] K. Gummadi, R. Gummadi, S. Gribble, S. Ratnasamy, S. Shenker, and I. Stoica, "The Impact of DHT Routing Geometry on Resilience and Proximity," in *SIGCOMM*, Karlsruhe, Germany, Aug. 2003.
- [7] J. Li, J. Stribling, F. Kaashoek, R. Morris, and T. Gil, "A Performance vs. Cost Framework for Evaluating DHT Design Tradeoffs under Churn," in *INFOCOM*, Miami, FL, Mar. 2005.
- [8] F. Dabek, M. F. Kaashoek, D. Karger, R. Morris, and I. Stoica, "Wide-area cooperative storage with CFS," in *SOSP*, Banff, AB, Canada, Oct. 2001.
- [9] V. Ramasubramanian and E. G. Sirer, "The Design and Implementation of a Next Generation Name Service for the Internet," in *SIGCOMM*, Portland, OR, 2004.
- [10] S. Rhea, D. Geels, and J. Kubiatowicz, "Handling Churn in a DHT," in *USENIX*, 2004, pp. 127–140.
- [11] S. S. Lam and H. Liu, "Failure Recovery for Structured P2P Networks: Protocol Design and Performance Evaluation," in *SIGMETRICS*, New York, NY, June 2004.
- [12] B. Y. Zhao, L. Huang, J. Stribling, S. C. Rhea, A. D. Joseph, and J. D. Kubiatowicz, "Tapestry: A Resilient Global-Scale Overlay for Service Deployment," *IEEE Journal on Selected Areas in Communications*, vol. 22, no. 1, pp. 41–53, Jan. 2004.
- [13] B. Leong, B. Liskov, and E. D. Demaine, "EpiChord: Parallelizing the Chord Lookup Algorithm with Reactive Routing State Management," in *International Conference on Networks*, Nov. 2004.
- [14] F. Dabek, J. Li, E. Sit, J. Robertson, M. F. Kaashoek, and R. Morris, "Designing a DHT for Low Latency and High Throughput," in *NSDI*, Berkeley, CA, 2004.
- [15] D. Stutzbach and R. Rejaie, "Capturing Accurate Snapshots of the Gnutella Network," in *Global Internet Symposium*, Miami, FL, Mar. 2005, pp. 127–132.
- [16] S. Krishnamurthy, S. El-Ansary, E. Aurell, and S. Haridi, "A Statistical Theory of Chord under Churn," in *International Workshop on Peer-to-Peer Systems (IPTPS)*, Ithaca, NY, Feb. 2005.
- [17] D. Liben-Nowell, H. Balakrishnan, and D. Karger, "Analysis of the Evolution of Peer-to-Peer Systems," in *Principles of Distributed Computing*, Monterey, CA, July 2002.
- [18] J. Li, J. Stribling, T. M. Gil, R. Morris, and F. Kaashoek, "Comparing the performance of distributed hash tables under churn," in *International Workshop on Peer-to-Peer Systems*, 2004.
- [19] R. Mahajan, M. Castro, and A. Rowstron, "Controlling the Cost of Reliability in Peer-to-Peer Overlays," in *International Workshop on Peer-to-Peer Systems*, 2003.
- [20] S. Sen and J. Wang, "Analyzing Peer-To-Peer Traffic Across Large Networks," *IEEE/ACM Transactions on Networking*, vol. 12, no. 2, pp. 219–232, Apr. 2004.
- [21] K. P. Gummadi, R. J. Dunn, S. Saroiu, S. D. Gribble, H. M. Levy, and J. Zahorjan, "Measurement, Modeling, and Analysis of a Peer-to-Peer File-Sharing Workload," in *SOSP*, 2003.
- [22] J. Pouwelse, P. Garbacki, D. Epema, and H. Sips, "The BitTorrent P2P File-sharing System: Measurements and Analysis," in *International Workshop on Peer-to-Peer Systems (IPTPS)*, Ithaca, NY, Feb. 2005.
- [23] M. Izal, G. Urvoy-Keller, E. W. Biersack, P. A. Felber, A. A. Hamra, and L. Garcés-Erice, "Dissecting BitTorrent: Five Months in a Torrent's Lifetime," in *PAM*, Apr. 2004.

- [24] J. Chu, K. Labonte, and B. N. Levine, "Availability and Locality Measurements of Peer-to-Peer File Systems," in *ITCom: Scalability and Traffic Control in IP Networks II Conferences*, July 2002.
- [25] D. Stutzbach and R. Rejaie, "Characterizing Churn in Peer-to-Peer Networks," University of Oregon, Eugene, OR, Tech. Rep. 2005-03, May 2005.