

Fault-tolerant Routing in Peer-to-peer Systems

James Aspnes^{*}
Yale University,
New Haven, CT 06520, USA.
aspnes@cs.yale.edu

Zoë Diamadi[†]
Yale University,
New Haven, CT 06520, USA.
diamadi@cs.yale.edu

Gauri Shah[‡]
Yale University,
New Haven, CT 06520, USA.
gauri.shah@yale.edu

ABSTRACT

We consider the problem of designing an overlay network and routing mechanism that permits finding resources efficiently in a peer-to-peer system. We argue that many existing approaches to this problem can be modeled as the construction of a random graph embedded in a metric space whose points represent resource identifiers, where the probability of a connection between two nodes depends only on the distance between them in the metric space. We study the performance of a peer-to-peer system where nodes are embedded at grid points in a simple metric space: a one-dimensional real line. We prove upper and lower bounds on the message complexity of locating particular resources in such a system, under a variety of assumptions about failures of either nodes or the connections between them. Our lower bounds in particular show that the use of inverse power-law distributions in routing, as suggested by Kleinberg [5], is close to optimal. We also give heuristics to efficiently maintain a network supporting efficient routing as nodes enter and leave the system. Finally, we give some experimental results that suggest promising directions for future work.

1. INTRODUCTION

Peer-to-peer systems are distributed systems without any central authority and with varying computational power at each machine. We study the problem of locating resources in such a large network of heterogeneous machines that are subject to crash failures. We describe

^{*}Supported by NSF grants CCR-9820888 and CCR-0098078.

[†]Supported in part by ONR grant N00014-01-1-0795.

[‡]Supported by NSF grants CCR-9820888 and CCR-0098078.

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.

PODC '02 Monterey, California USA

Copyright 2001 ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

how to construct distributed data structures that have certain desirable properties and allow efficient resource location.

Decentralization is a critical feature of such a system as any central server is a vulnerable point of failure and also wastes the power of the clients. Equally important is scalability: the cost borne by each node must not depend too much on the network size and should ideally be proportional, within polylogarithmic factors, to the amount of data the node seeks or provides. Since we expect nodes to arrive and depart at a high rate, the system should be resilient to failures. Furthermore, disruptions to parts of the data structure should self-heal to provide self-stabilization.

Our approach provides a hash table-like functionality, based on keys that uniquely identify the system resources. To accomplish this, we map resources to points in a metric space from the keys' hash values. We construct and maintain a random graph linking these points and use greedy routing to traverse its edges to find data items. The principle we rely on is that failures leave behind yet another (smaller) random graph, ensuring that the system is robust even in the face of considerable damage. Another compelling advantage of random graphs is that they eliminate the need for global coordination. Thus, we get a fully-distributed, egalitarian, scalable system with no bottlenecks. We measure performance in terms of the number of messages sent by the system for a search operation. Given the growing storage capacity of machines, we are less concerned with minimizing the storage at each node, but the space requirements are still small: the information stored at a node consists only of a network address and location in the metric space for each neighbor.

The rest of the paper is organized as follows. Section 2 explains our abstract model in detail, and Section 3 describes some of the existing peer-to-peer systems. We prove our results for routing in Section 4, followed by a discussion on how to construct the random graph in Section 5. We present experimental results in Section 6 and conclusions and future work in Section 7.

2. OUR APPROACH

The idea underlying our approach consists of three basic parts: (1) embed resources as points in a metric space, (2) construct a random graph by appropriately

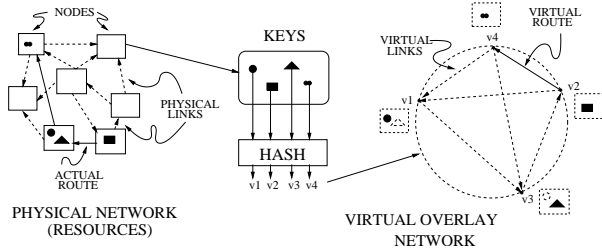


Figure 1: An example of the metric-space embedding.

linking these points, and (3) efficiently locate resources by routing greedily along the edges of the graph. Let R be a set of resources spread over a large, heterogeneous network N . For each resource $r \in R$, $owner(r)$ denotes the node in N that provides r and $key(r)$ denotes the resource's key. Let K be the set of all possible keys. We assume a hash function $h : K \rightarrow V$ such that resource r maps to the point $v = h(key(r))$ in a metric space (V, d) , where V is the point set and d is the distance metric as shown in Figure 1. The hash function is assumed to populate the metric space evenly. Note that via this resource embedding, a node n is mapped onto the set $V_n = \{v \in V : \exists r \in R, v = h(key(r)) \wedge (owner(r) = n)\}$, namely the set of metric-space points assigned to the resources the node provides.

Our next step is to carefully construct a directed random graph from the points embedded in V . We assume that each newly-arrived node n is initially connected to some other node in N . Each node n generates the outgoing edges for each vertex $v \in V_n$ independently. An edge $(v, u) \in V_n \times V_m$ simply denotes that n knows that m is the network node that provides the resource mapped to u ; hence, we can view the graph as a virtual overlay network of information. Node n constructs each edge by executing the search algorithm to locate the resource that is mapped to the sink of that edge. If the metric space is not populated densely enough, the choice of a sink may result in a vertex corresponding to an absent resource. In that case, n chooses the neighbor present closest to the original sink. Moving to nearby vertices will introduce some bias in the edge distribution but the magnitude of error does not appear to be large. A more detailed description of the graph construction is given in Section 5.

Having constructed the overlay network of information, we can now use it for resource location. At any time t , let $R^t \subseteq R$ be the set of available resources and I^t be the corresponding overlay network. A request by node n to locate resource r at time t is served in a simple, localized manner: n calculates the metric-space point v that corresponds to r , and a request message is then routed over I^t from the vertex in V_n that is closest to v to v itself. Each node needs only local information, namely its set of neighbors in I^t , to participate in the resource location. Routing is done greedily by forwarding the message to the node mapped to a metric-space point as close to v as possible. The problem of resource loca-

tion is thus translated into routing on random graphs embedded in a metric space.

The advantage of using random graphs is that they are robust against failures: a node-induced subgraph of a random graph is generally still a random graph; therefore, the disappearance of a vertex will still allow routing over the structure. Further, embedding the graph in a metric space has the very important property that the only information needed to locate a resource is the location of its corresponding metric-space point. That location is both permanent, in the sense of being unaffected by disruption of the data structure, and easily computable by any node that seeks the resource. So, while the pattern of links between nodes may be damaged or destroyed by failure of nodes or of the underlying communication network, the metric space forms an invulnerable foundation over which to build the ephemeral parts of the data structure.

3. CURRENT PEER-TO-PEER SYSTEMS

Most peer-to-peer systems in widespread use are not scalable. Napster's [8] approach of using a central server has the weaknesses of a vulnerable single point of failure and wasted client computational power. Gnutella [1] floods the network to locate a resource. Flooding creates a trade-off between overloading every node in the network for each request and cutting off searches before completion. While the use of super-peers [7] ameliorates the problem somewhat in practice, it does not improve performance in the limit.

Some of these first-generation systems have inspired the development of more sophisticated ones like CAN [11], Chord [12] and Tapestry [2]. CAN partitions a d -dimensional metric space into *zones*. Each key is mapped to a point in some zone and stored at the node that owns the zone. Each node stores $O(d)$ information, and resource location, done by greedy routing, takes $O(dn^{1/d})$ time. Chord maps nodes to identities of m bits placed around a *modulo* 2^m *identifier circle*. Resources are stored at existing *successor* nodes of the nodes they are mapped to. Each node maintains a routing table of size m and uses greedy routing to give an $O(m)$ delivery time. Tapestry uses Plaxton's algorithm [10], a form of suffix-based, hypercube routing, as the routing mechanism: in this algorithm, the message is forwarded deterministically to a node whose identifier is one digit closer to the target identifier. To this end, each node maintains $O(\lg n)$ pieces of information and delivery time is also $O(\lg n)$.

Although these systems seem vastly different, there is a recurrent underlying theme in the use of some variant of an overlaid metric space in which the nodes are embedded. The location of a resource in this metric space is determined by its key. Each node maintains some information about its neighbors in the metric space, and routing is then simply done by forwarding packets to neighbors *closer* to the target node with respect to the metric. It is this inherent common structure that leads to similar results for the performance of these networks.

In this paper, we describe a general setting for such overlay metric spaces, although most of our results are obtained from one-dimensional spaces.

In general, the fault-tolerance properties of these systems are not well-defined. Each system provides a repair mechanism for failures but makes no performance guarantees till this mechanism kicks in. For large systems, where nodes appear and leave frequently, resilience to repeated and concurrent failures is a desirable and important property. Our experiments show that with our overlay space and linking strategies, the system performs reasonably well even with a large number of failures.

4. ROUTING

In this section, we present our lower and upper bounds on greedy routing.

4.1 Tools

Some of our upper bounds will be proved using a well-known upper bound of Karp *et al.*[3] on probabilistic recurrence relations. We will restate this bound as Lemma 1, and then show how a similar technique can be used to get *lower bounds* with some additional conditions in Theorem 2.

LEMMA 1 ([3]). *The time $T(X_0)$ needed for a non-increasing real-valued Markov chain $X_0, X_1, X_2, X_3 \dots$ to drop to 1 is bounded by*

$$T(X_0) \leq \int_1^{X_0} \frac{1}{\mu_z} dz, \quad (1)$$

when $\mu_z = E[X_t - X_{t+1} : X_t = z]$ is a nondecreasing function of z .

This bound has a nice physical interpretation. If it takes one second to jump down μ_x meters from x , then we are traveling at a rate of μ_x meters per second during that interval. When we zip past some position z , we are traveling at the average speed μ_x determined by our starting point $x \geq z$ for the interval. Since μ is nondecreasing, using μ_z as our estimated speed underestimates our actual speed when passing z . The integral computes the time to get all the way to zero if we use μ_z as our instantaneous speed when passing position z . Since our estimate of our speed is low (on average), our estimate of our time will be high, giving an upper bound on the actual expected time.

We would like to get lower bounds on such processes in addition to upper bounds, and we will not necessarily be able to guarantee that μ_z , as defined in Lemma 1, will be a nondecreasing function of z . But we will still use the same basic intuition: The average speed at which we pass z is at least the minimum average speed of any jump that takes us past z . We can find this minimum speed by taking the minimum over all $x > z$; unfortunately, this may give us too small an estimate. Instead, we choose an upper bound U on “short” jumps, compute the minimum speed of short jumps of at most U for all x between z and $z + U$, and handle the (hopefully rare) long jumps of more than U by conditioning

against them. Subject to this conditioning, we can define an upper bound m_z on the average speed passing z , and use essentially the same integral as in (1) to get a lower bound on the time. Some additional tinkering to account for the effect of the conditioning then gives us our real lower bound, which appears in Theorem 2 below, as Inequality (8).

THEOREM 2. *Let X_0, X_1, X_2, \dots be Markov process with state space S , where X_0 is a constant. Let f be a non-negative real-valued function on S such that, for all t ,*

$$\Pr[f(X_t) - f(X_{t+1}) \geq 0 : X_t] = 1. \quad (2)$$

Let U and ϵ be constants such that for any $x > 0$,

$$\Pr[f(X_t) - f(X_{t+1}) \geq U : X_t = x] \leq \epsilon. \quad (3)$$

Let

$$\tau = \min\{t : f(X_t) = 0\}. \quad (4)$$

For each x with $f(x) > 0$, let $\mu_x > 0$ satisfy

$$\mu_x \geq E[f(X_t) - f(X_{t+1}) : X_t = x, f(X_t) - f(X_{t+1}) < U]. \quad (5)$$

Now define

$$m_x = \sup\{\mu_x : x \in S, f(x) \in [z, z + U]\} \quad (6)$$

and define

$$T(x) = \int_0^{f(x)} \frac{1}{m_z} dz \quad (7)$$

Then

$$E[\tau] \geq \frac{T(X_0)}{\epsilon T(X_0) + (1 - \epsilon)} \quad (8)$$

4.2 Lower bound on greedy routing

We will now show a lower bound on the expected time taken by greedy routing on a random graph embedded in a line, where each node in the graph has expected outdegree at most ℓ , and the probability that a node at position x is connected to positions $x - \Delta_1, x - \Delta_2, \dots, x - \Delta_k$ depends only on the set $\Delta = \{\Delta_1, \dots, \Delta_k\}$ and not on x and is independent of the choice of outgoing edges for other nodes.¹

We consider to two variants of the greedy routing algorithm. Without loss of generality, we assume that the target of the search is labeled 0. In *one-sided greedy routing*, the algorithm never traverses a link that would take it past its target. So if the algorithm is currently at x and is trying to reach 0, it will move to the node $x - \Delta_i$ with the smallest non-negative label. In *two-sided greedy routing*, the algorithm chooses a link that minimizes the distance to the target, without regard to which side of the target the other end of the link is. In the two-sided case the algorithm will move to a node $x - \Delta_i$ whose label has the smallest absolute value, with ties broken arbitrarily. One-sided greedy routing

¹We assume that nodes are labeled by integers, and identify each node with its label to avoid excessive notation.

can be thought of as modeling algorithms on a graph with a boundary when the target lies on the boundary, or algorithms where all links point in only one direction (as in Chord).

Our results are stronger for the one-sided case than for the two-sided case. With one-sided greedy routing, we show a lower bound of $\Omega(\ln n / (\ell \ln \ln n))$ on the time to reach 0 from a point chosen uniformly from the range 1 to n that applies to any link distribution. For two-sided routing, we show a lower bound of $\Omega(\ln n / (\ell^2 \ln \ln n))$, with some constraints on the distribution. We conjecture that these constraints are unnecessary, and that $\Omega(\ln n / (\ell \ln \ln n))$ is the correct lower bound for both models.

In general, we assume that each node is connected to its immediate neighbors; the simplest way to model this is to require that ± 1 appear in Δ .

We will begin by developing machinery that will be useful in the proofs of both the one-sided and two-sided lower bounds.

4.2.1 Link sets: notation and distributions

First we describe some notation for Δ sets. Write each Δ as

$$\{\Delta_{-s}, \dots, \Delta_{-2}, \Delta_{-1} = -1, \Delta_1 = 1, \Delta_2, \dots, \Delta_t\},$$

where $\Delta_i < \Delta_j$ whenever $i < j$. Each Δ is a random variable drawn from some distribution on finite sets; the individual Δ_i are thus in general *not* independent. Let Δ^- consist of the s negative elements of Δ and Δ^+ consist of the t positive elements. Formally define $\Delta_{-i} = -\infty$ when $i > s$ and $\Delta_i = +\infty$ when $i > t$.

For one-sided routing, we make no assumptions about the distribution of Δ except that $|\Delta|$ must have finite expectation and Δ always contains 1. For two-sided routing, we assume that Δ is generated by including each possible δ in Δ with probability p_δ , where p is symmetric about the origin (i.e., $p_\delta = p_{-\delta}$ for all δ), $p_1 = p_{-1} = 1$, and p is unimodal, i.e. nonincreasing for positive δ and nondecreasing for negative δ .² We also require that the events $[\delta \in \Delta]$ and $[\delta' \in \Delta]$ are pairwise independent for distinct δ, δ' .

4.2.2 The aggregate chain S^t

For a fixed distribution on Δ , the trajectory of a single initial point X^0 is a Markov chain X^0, X^1, X^2, \dots , with $X^{t+1} = s(X^t, \Delta^t)$, where Δ^t determines the outgoing links from the node reached at time t and s is a *successor function* that selects the next node $X^{t+1} = X^t - \Delta_i^t$ according to the routing algorithm. Note that the chain is Markov, because the presence of ± 1 links guarantees that no node ever appears twice in the sequence, and so each new node corresponds to a new choice of links.

From the X^t chain we can derive an *aggregate chain* that describes the collective behavior of all nodes in some range. Each state of the aggregate chain is a contiguous sets of nodes whose labels all have the same sign; we define the sign of the state to be the common

²These constraints imply that $p_0 = 1$; formally, we imagine that 0 is present in each Δ but is ignored by the routing algorithm.

sign of all of its elements. For one-sided routing each state is either $\{0\}$ or an interval of the form $\{1 \dots k\}$ for some k . For two-sided routing the states are more general. The aggregate states are characterized formally in Lemma 4.

Given a contiguous set of nodes S and a set Δ , define

$$S_{\Delta_i} = \{x \in S : s(x, \Delta) = x - \Delta_i\}.$$

The intuition is that S_{Δ_i} consists of all those nodes for which the algorithm will choose Δ_i as the outgoing link. Note that S_{Δ_i} will always be a contiguous range because of the greediness of the algorithm. Now define, for each $\sigma \in \{-, 0, +\}$:

$$S_{\Delta_i \sigma} = \{x \in S_{\Delta_i} : \text{sgn } s(x, \Delta) = \sigma\}.$$

Here we have simply split S_{Δ_i} into those nodes with negative, zero, or positive successors.

For any set A and integer δ write $A - \delta$ for $\{x - \delta : x \in A\}$.

We will now build our aggregate chain by letting the successors of a range S be the ranges $S_{\Delta_i \sigma} - \Delta_i$ for all possible Δ, i , and σ . As a special case, we define $S^{t+1} = \{0\}$ when $S^t = \{0\}$; once we arrive at the target, we do not leave it. For all other S^t , we let

$$\Pr [S^{t+1} = S_{\Delta_i \sigma}^t - \Delta_i : \Delta] = \frac{|S_{\Delta_i \sigma}^t|}{|S^t|}, \quad (9)$$

and define the unconditional transition probabilities by averaging over all Δ .

Lemma 3 shows that moving to the aggregate chain does not misrepresent the underlying single-point chain:

LEMMA 3. *Let X^0 be drawn uniformly from the range S^0 . Let Y^t be a uniformly chosen element of S^t . Then for all x and t , $\Pr[X^t = x] = \Pr[Y^t = x]$.*

Lemma 4 justifies our earlier characterization of the aggregate state spaces:

LEMMA 4. *Let $S^0 = \{1 \dots n\}$ for some n . Then with one-sided routing, every S^t is either $\{0\}$ or of the form $\{1 \dots k\}$ for some k ; and with two-sided routing, every S^t is an interval of integers in which every element has the same sign.*

The advantage of the aggregate chain over the single-point chain is that, while we cannot do much to bound the progress of a single point with an arbitrary distribution on Δ , we can show that the size of S^t does not drop too quickly given a bound ℓ on $\mathbb{E}[|\Delta|]$. The intuition is that each successor set of size $a^{-1}|S^t|$ or less occurs with probability at most a^{-1} , and there are at most 3ℓ such sets on average.

LEMMA 5. *Let $\mathbb{E}[|\Delta|] \leq \ell$. Then for any $a \geq 1$, in either the one-sided or two-sided model,*

$$\Pr [|S^{t+1}| \leq a^{-1}|S^t|] \leq 3\ell a^{-1}. \quad (10)$$

Another way to write (10) is to say that $\Pr [\ln |S^t| - \ln |S^{t+1}| \geq \ln a] \leq 3\ell a^{-1}$, which will give the bound (3) on the probability of large jumps when it comes time to apply Theorem 2.

4.2.3 Boundary points

Lemma 5 says that $|S^t|$ seldom drops by too large a ratio at once, but it doesn't tell us much about how quickly $|S^t|$ drops in short hops. To bound this latter quantity, we need to get a bound on how many subranges S^t splinters into through the action of $s(\cdot, \Delta)$. We will do so by showing that only certain points can appear as the boundaries of these subranges in the direction of 0.

For fixed Δ , define for each $i > 0$

$$\beta_i = \left\lceil \frac{\Delta_i + \Delta_{i+1}}{2} \right\rceil$$

and

$$\beta_{-i} = \left\lfloor \frac{\Delta_{-i} + \Delta_{-i-1}}{2} \right\rfloor.$$

Let β be the set of all finite β_i and β_{-i} .

LEMMA 6. Fix S and Δ and let β be defined as above. Suppose that S is positive. Let $M = \{\min(S_{\Delta i \sigma}) : S_{\Delta i \sigma} \neq \emptyset\}$ be the set of minimum elements of subranges $S_{\Delta i \sigma}$ of S . Then M is a subset of β and contains no elements other than

1. $\min(S)$,
2. Δ_i for each $i > 0$,
3. $\Delta_i + 1$ for each $i > 0$, and
4. At most one of β_i or $\beta_i + 1$ for each $i > 0$,

where the last case holds only with two-sided routing.

If S is negative, the symmetric condition holds for $M = \{\max(S_{\Delta i \sigma}) : S_{\Delta i \sigma} \neq \emptyset\}$.

PROOF. Consider some subrange $S_{\Delta i \sigma}$ of S . If $S_{\Delta i \sigma}$ contains $\min(S)$, the first case holds. Otherwise: (a) if $S_{\Delta i \sigma} = S_{\Delta i 0}$, the second case holds; (b) if $S_{\Delta i \sigma} = S_{\Delta i +}$, the third case holds; (c) if $S_{\Delta i \sigma} = S_{\Delta i -}$, the fourth case holds, with $\min(S_{\Delta i -}) = \beta_{i-1}$ if $\Delta_{i-1} + \Delta_i$ is odd, and either β_{i-1} or $\beta_{i-1} + 1$ if $\Delta_{i-1} + \Delta_i$ is even, depending on whether the tie-breaking rule assigns β_{i-1} to $S_{\Delta(i-1)+}$ or $S_{\Delta i -}$. \square

We will call the elements of M *boundary points* of S .

4.2.4 Bounding changes in $\ln |S^t|$

Now we would like to use Lemmas 5 and Lemma 6 to get an upper bound on the rate at which $\ln |S^t|$ drops as a function of the Δ distribution.

The following lemma is used to bound a sum that arises in Lemma 8.

LEMMA 7. Let $c \geq 0$. Let $\sum_{i=1}^n x_i = M$ where each $x_i \geq 0$ and at least one x_i is greater than c . Let B be the set of all i for which x_i is greater than c . Then

$$\frac{\sum_{i \in B} x_i \ln x_i}{\sum_{i \in B} x_i} \geq \ln \left(\max \left(c, \frac{M}{n} \right) \right). \quad (11)$$

PROOF. If $\frac{M}{n} < c$, we still have $x_i > c$ for all $i \in B$, so the left-hand side cannot be less than $\ln c$. So the interesting case is when $\frac{M}{n} > c$.

Let B have b elements. Then $\sum_{i \notin B} x_i < (n-b)c$ and $\sum_{i \in B} x_i \geq M - (n-b)c = M - nc + bc$. Because $x_i \ln x_i$ is convex, its sum over B is minimized for fixed $\sum_{i \in B} x_i$ by setting all such x_i equal, in which case the left-hand side of (11) becomes simply $\ln(x_i)$ for any $i \in B$.

Now observe that setting all x_i in B equal gives $x_i = \frac{M - nc + bc}{b} = \frac{M - nc}{b} + c \geq \frac{M - nc}{n} + c = \frac{M}{n}$. \square

LEMMA 8. Fix $a > 1$, and let $S = S^t$ be a positive range with $|S| \geq a$. Define β as in Lemma 6. Let $S' = [\min(S) + \lceil a^{-1}|S| \rceil - 1, \max(S) - 1]$. Then

$$\begin{aligned} \mathbb{E} [\ln |S^t| - \ln |S^{t+1}| : S^t, \ln |S^t| - \ln |S^{t+1}| < \ln a] \\ \leq \ln \frac{1}{1-a^{-1}} + \ln \mathbb{E}[1 + Z : S^t] \end{aligned} \quad (12)$$

where $Z = 2|\Delta \cap S'|$ with one-sided routing and $Z = 2|\Delta \cap S'| + |\beta \cap S'|$ with two-sided routing.

PROOF. Call a subrange $S_{\Delta i \sigma}$ *large* if $|S_{\Delta i \sigma}| > a^{-1}|S|$ and *small* otherwise. Observe that $\lceil a^{-1}|S| \rceil \geq 2$, implying any large set has at least two elements.

For any large $S_{\Delta i \sigma}$, $\max(S_{\Delta i \sigma}) \geq \min(S) + \lceil a^{-1}|S| \rceil - 1$. Similarly $\min(S_{\Delta i \sigma}) \leq \max(S) - 1$. So any large $S_{\Delta i \sigma}$ intersects S' in at least one point.

Let $T = \{T_1, T_2, \dots, T_k\}$ be the set of subranges $S_{\Delta i \sigma}$, large or small, that intersect S' . It is immediate from this definition that $\bigcup T \supseteq S'$ and thus $\sum |T_j| \geq |S'|$.

Using Lemma 6, we can characterize the elements of T as follows:

1. There is at most one set T_j that contains $\min(T_j)$.
2. There is at most one set T_j that has $\min(T_j) = \Delta_i$ for each Δ_i in S' .
3. There is at most one set T_j that has $\min(T_j) = \Delta_i + 1$ for each Δ_i in S' .
4. With two-sided routing, there is at most one set T_j that has $\min(T_j) = \beta_i$ or $\min(T_j) = \beta_i + 1$ for each β_i in S' . Note that there may be a set whose minimum element is $\beta_i + 1$ where $\beta_i = \min(S') - 1$, but this set is already accounted for by the first case.

Thus T has at most $1 + Z = 1 + 2|\Delta \cap S'|$ elements with one-sided routing and at most $1 + Z = 1 + 2|\Delta \cap S'| + |\beta \cap S'|$ elements with two-sided routing.

Conditioning on $|S^{t+1}| > a^{-1}|S|$, $|S^{t+1}|$ is equal to $|S_{\Delta i \sigma}|$ for some large $S_{\Delta i \sigma}$ and thus for some large $T_j \in T$. Which large T_j is chosen is proportional to its size, so for fixed T , we have

$$\begin{aligned} \mathbb{E}[\ln |S^{t+1}| : T] &= \frac{\sum_{j=1}^{|T|} |T_j| \ln |T_j|}{\sum_{j=1}^{|T|} |T_j|} \\ &\geq \ln \left(\max \left(a^{-1}|S|, \frac{|\bigcup T|}{|T|} \right) \right) \\ &\geq \ln \left(\frac{|S'|}{|T|} \right), \end{aligned}$$

where the first inequality follows from Lemma 7.

Now let us compute

$$\begin{aligned} & \mathbb{E}[\ln |S^t| - \ln |S^{t+1}| : S^t] \\ & \leq \ln |S^t| - \mathbb{E}[\ln |S^t| - \ln |T| : S^t] \\ & = \ln \frac{|S^t|}{|S^t|} + \mathbb{E}[\ln |T| : S^t] \\ & \leq \ln \frac{1}{1-a^{-1}} + \ln \mathbb{E}[|T| : S^t]. \end{aligned}$$

In the last step we use $\mathbb{E}[\ln |T| : S^t] \leq \ln \mathbb{E}[|T| : S^t]$, which follows from the concavity of \ln and Jensen's inequality. \square

4.2.5 Putting the pieces together

We now have all the tools we need to prove our lower bound.

THEOREM 9. *Let G be a random graph whose nodes are labeled by the integers. Let Δ_x for each x be a set of integer offsets chosen independently from some common distribution, subject to the constraint that -1 and $+1$ are present in every Δ_x , and let node x have an outgoing edge to $x-\delta$ for each $\delta \in \Delta_x$. Let $\ell = \mathbb{E}[|\Delta|]$. Consider a greedy routing trajectory in G starting at a point chosen uniformly from $1 \dots n$ and ending at 0.*

With one-sided routing, the expected time to reach 0 is

$$\Omega\left(\frac{\ln^2 n}{\ell \ln \ln n}\right) \quad (13)$$

With two-sided routing, the expected time to reach 0 is

$$\Omega\left(\frac{\ln^2 n}{\ell^2 \ln \ln n}\right), \quad (14)$$

provided Δ is generated by including each δ in Δ with probability p_δ , where (a) p is unimodal, (b) p is symmetric about 0, and (c) the choices to include particular δ, δ' are pairwise independent.

PROOF. Let $S^0 = \{1 \dots n\}$.

We are going to apply Theorem 2 to the sequence S^0, S^1, S^2, \dots with $f(S) = \ln |S|$. We have chosen f so that when we reach the target, $f(S) = 0$; so that a lower bound on τ gives a lower bound on the expected time of the routing algorithm. To apply the theorem, we need to show that (a) the probability that $\ln |S|$ drops by a large amount is small, and (b) that the integral in (7) is large.

Let $a = 3\ell \ln^3 n$. By Lemma 5, for all t , $\Pr[|S^{t+1}| \leq a^{-1}|S^t|] \leq 3\ell a^{-1} = \ln^{-3} n$, and thus $\Pr[\ln |S^t| - \ln |S^{t+1}| \geq \ln a] \leq \ln^{-3} n$. This satisfies (3) with $U = \ln a$ and $\epsilon = \ln^{-3} n$.

For the second step, Theorem 2 requires that we bound the speed of the change in $f(S)$ solely as a function of $f(S)$. For one-sided routing this is not a problem, as Lemma 4 shows that $f(S)$, which reveals $|S|$, characterizes S exactly except when $|S| = 1$ and the lower bound argument is done. For two-sided routing, the situation is more complicated; there may be some

S^t which is not of the form $\{1 \dots |S^t|\}$ or $\{0\}$, and we need a bound on the speed at which $\ln |S^t|$ drops that applies equally to all sets of the same size.

It is here (and only here) that we use our conditions on Δ for two-sided routing. Suppose that each δ appears in Δ with probability p_δ , that these probabilities are pairwise-independent, and that the sequence p is symmetric and unimodal. Let $\hat{\beta} = \{\text{absceil}(\frac{x+y}{2}) : x, y \in \Delta, x \neq y\}$, where $\text{absceil}(z)$, the *absolute ceiling* of z , is $\lceil z \rceil$ when $z \geq 0$ and $\lfloor z \rfloor$ when $z \leq 0$. Observe that $\hat{\beta} \supseteq \beta$; in effect, we are counting in $\hat{\beta}$ all midpoints of pairs of distinct elements of δ without regard to whether the elements are adjacent. For each k , the expected number of distinct pairs x, y with $x+y=z$ and $x, y \in \Delta$ is at most $b_k = \sum_{i=-\infty}^{\infty} p_{k-i} p_i$; this is a convolution of the non-negative, symmetric, and unimodal p sequence with itself and so it is also symmetric and unimodal. It follows that for all $0 \leq k < k'$, $b_k \geq b_{k'}$, and similarly $b_{-k} \geq b_{-k'}$.

Now for the punch line: for each $\delta \neq 0$, $q_\delta = b_{2\delta - \text{sgn } \delta} + b_{2\delta}$ is an upper bound on the expected number of distinct pairs x, y that put δ in β , which is in turn an upper bound on $\Pr[\delta \in \beta]$, and from the unimodularity of b we have that $q_\delta \geq q_{\delta'}$ and $q_{-\delta} \geq q_{-\delta'}$ whenever $0 < \delta < \delta'$. Though q grossly overcounts the elements of β (in particular, it gives a bound on $\mathbb{E}[|\beta|]$ of ℓ^2), its ordering property means that we can bound the expected number of elements of β that appear in some subrange of any positive S^t by using q to bound the expected number of elements that appear in the corresponding subrange of $\{1 \dots |S^t|\}$, and similarly for negative S^t and $\{-1 \dots -|S^t|\}$. Because p_i already satisfies a similar ordering property, we can thus bound the number of elements of both Δ and β that hit a fixed subrange of S^t given only $|S^t|$.

Which we now proceed to do. For convenience, formally define $q_i = 0$ for one-sided routing. For each integer $i > 0$ let $A_i = \{k \in \mathbf{Z} : a^i - 1 \leq k < a^{i+1} - 1\} = \{k \in \mathbf{Z} : \lfloor \ln_a k + 1 \rfloor = i\}$. Let $\gamma_i = \sum_{k \in A_i} 2p_i + q_i$. Note that $\gamma_i \geq 2\mathbb{E}[|A_i \cap \Delta|]$ for one-sided routing and $\gamma_i \geq 2\mathbb{E}[|A_i \cap \Delta|] + \mathbb{E}[|A_i \cap \beta|]$ for two-sided routing. Observe also that $\sum_{i=0}^{\infty} \gamma_i$ is at most 2ℓ for one-sided routing and by $2\ell + \ell^2$ for two-sided routing.

Consider some $S = S^t$. If $|S| \geq a$, then by Lemma 8 we have

$$\begin{aligned} & \mathbb{E}[\ln |S^t| - \ln |S^{t+1}| : S^t, \ln |S^t| - \ln |S^{t+1}| < \ln a] \\ & \leq \ln \frac{1}{1-a^{-1}} + \ln \mathbb{E}[1 + Z : S^t], \end{aligned} \quad (15)$$

where $Z = 2|\Delta \cap S^t|$ with one-sided routing and $Z = 2|\Delta \cap S^t| + |\beta \cap S^t|$ with two-sided routing, and $S^t = [\min(S) + \lceil a^{-1}|S| \rceil - 1, \max(S) - 1]$. By the monotonicity of p_i and q_i for positive i , $\ln \mathbb{E}[1 + Z]$ is at most

$$\mu_{\ln |S|} = \ln \frac{1}{1-a^{-1}} + \ln \left(1 + \sum_{i=\lceil a^{-1}|S| \rceil - 1}^{|S|-1} 2p_i + q_i \right), \quad (16)$$

provided $|S| \geq a$. For $|S| < a$, set $\mu_{\ln |S|} = \ln a$.

Let us now compute m_z , as defined in (6). For $z < \ln a$, $m_z = \ln a$. For larger z , observe that $m_z =$

$\sup \{m_{\ln |S|} : e^z \leq |S| < ae^z\}$. Now if $e^z \leq |S| < ae^z$, then the bounds on the sum in (16) both lie between $\lceil a^{-1}e^z \rceil - 1$ and $ae^z - 1$, so that

$$\begin{aligned} m_z &\leq \ln \frac{1}{1-a^{-1}} + \ln \left(1 + \sum_{i=\lceil a^{-1}e^z \rceil - 1}^{\lfloor ae^z - 1 \rfloor} 2p_i + q_i \right) \\ &\leq \ln \frac{1}{1-a^{-1}} + \ln(1 + \gamma_{z'} + \gamma_{z'+1} + \gamma_{z'+2}), \end{aligned}$$

where $z' = \lfloor z/\ln a \rfloor - 1$.

Finally, compute

$$\begin{aligned} T(\ln n) &= \int_0^{\ln n} \frac{1}{m_z} dz \\ &\geq \int_{\ln a}^{\ln n} \frac{1}{\ln \frac{1}{1-a^{-1}} + \ln(1 + \gamma_{z'} + \gamma_{z'+1} + \gamma_{z'+2})} dz \\ &\geq \sum_{i=0}^{\lfloor \ln n / \ln a \rfloor - 1} \frac{\ln a}{\ln \frac{1}{1-a^{-1}} + \ln(1 + \gamma_i + \gamma_{i+1} + \gamma_{i+2})}. \end{aligned}$$

To get a lower bound on the sum, note that

$$\sum_{i=0}^{\lfloor \ln n / \ln a \rfloor - 1} (\gamma_i + \gamma_{i+1} + \gamma_{i+2}) \leq 3 \sum_{i=0}^{\lfloor \ln n / \ln a \rfloor - 1} \gamma_i \leq 3 \sum_{i=0}^{\infty} \gamma_i$$

which is at most $L = 6\ell$ for one-sided routing and at most $L = 6\ell + 3\ell^2$ for two-sided routing. In either case, because $\frac{1}{c+\ln(1+x)}$ is convex and decreasing, we have

$$\begin{aligned} T(\ln n) &\geq \sum_{i=0}^{\lfloor \ln n / \ln a \rfloor - 1} \frac{\ln a}{\ln \frac{1}{1-a^{-1}} + \ln(1 + \gamma_i + \gamma_{i+1} + \gamma_{i+2})} \\ &\geq \sum_{i=0}^{\lfloor \ln n / \ln a \rfloor - 1} \frac{\ln a}{\ln \frac{1}{1-a^{-1}} + \ln \left(1 + \frac{L}{\lfloor \ln n / \ln a \rfloor} \right)} \\ &= \frac{\ln a \lfloor \ln n / \ln a \rfloor}{\ln \frac{1}{1-a^{-1}} + \ln \left(1 + \frac{L}{\lfloor \ln n / \ln a \rfloor} \right)}. \end{aligned} \quad (17)$$

We will now rewrite our bound on $T(\ln n)$ in a more convenient asymptotic form. We will ignore the 1 and concentrate on the large fraction. Recall that $a = 3\ell \ln^3 n$, so $\ln a = \Theta(\ln \ell + \ln \ln n)$. Unless ℓ is polynomial in n , we have $\ln n / \ln a = \omega(1)$ and the numerator simplifies to $\Theta(\ln n)$.

Now let us look at the denominator. Consider first the term $\ln \frac{1}{1-a^{-1}}$. We can rewrite this term as $-\ln(1 - a^{-1})$; since a^{-1} goes to zero as ℓ and n grow we have $-\ln(1 - a^{-1}) = \Theta(a^{-1}) = \Theta(\ell^{-1} \ln^{-3} n)$. It is unlikely that this term will contribute much.

Turning to the second term, let us use the fact that $\ln(1+x) \leq x$ for $x \geq 0$. Thus

$$\begin{aligned} \ln \left(1 + \frac{L}{\lfloor \ln n / \ln a \rfloor} \right) &\leq \frac{L}{\lfloor \ln n / \ln a \rfloor} \\ &= O \left(\frac{L(\ln \ell + \ln \ln n)}{\ln n} \right), \end{aligned}$$

and the bound in (17) simplifies to $\Omega(\ln^2 n / (L(\ln \ell + \ln \ln n)))$. We can further assume that $\ell = O(\ln^2 n)$, since otherwise the bound degenerates to $\Omega(1)$, and rewrite it simply as $\Omega(\ln^2 n / (L \ln \ln n))$.

For large L the approximation $\ln(1+x) \leq 1 + \ln x$ for $x \geq 0.59$ is more useful. In this case (17) simplifies to $T(\ln n) = \Omega(\ln n / \ln \ell)$, which has a natural interpretation in terms of the tree of successor nodes of some single starting node.

We are not quite done with Theorem 2 yet, as we still need to plug our T and ϵ into (8) to get a lower bound on $E[\tau]$. But here we can simply observe that $\epsilon T = O(\ln^{-1} n)$, so the denominator in (8) goes rapidly to 1. Our stated bounds are thus finally obtained by substituting $O(\ell)$ or $O(\ell^2)$ for L . \square

4.2.6 Possible strengthenings of the lower bound

Examining the proof of Theorem 9, both the ℓ^2 that appears in the bound (14) for two-sided routing and the extra conditions imposed on the Δ distribution arise only as artifacts of our need to project each range S onto $\{1 \dots |S|\}$ and thus reduce the problem to tracking a single parameter. We believe that a more sophisticated argument that does not collapse ranges together would show a stronger result:

CONJECTURE 10. *Let G , Δ , and ℓ be as in Theorem 9. Consider a greedy routing trajectory starting at a point chosen uniformly from $1 \dots n$ and ending at 0.*

Then the expected time to reach 0 is

$$\Omega \left(\frac{\ln^2 n}{\ell \ln \ln n} \right),$$

with either one-sided or two-sided routing, and no constraints on the Δ distribution.

We also believe that the bound continues to hold in higher dimensions than 1. Unfortunately, the fact that we can embed the line in, say, a two-dimensional grid is not enough to justify this belief; divergence to one side or the other of the line may change the distribution of boundaries between segments and break the proof of Theorem 9.

4.3 Upper Bounds

In this section, we present our upper bounds on the delivery time for settings that include failure of nodes and links. Note that we use the term node here to mean the point in the metric space and not the physical machine as implied in earlier sections. To simplify the theoretical analysis, we assume an ideal placement of one node at every point of the one-dimensional line. Each node is connected to its immediate neighbors and has ℓ multiple long-distance links chosen independently as follows: $\Pr[v \text{ is the } i\text{th neighbor of } u] = \frac{1}{d(u,v)} / \sum_{w \neq u} \frac{1}{d(u,w)}$, where $d(u,v)$ is the distance between u and v .

We first consider an idealized model with no failures. Kleinberg [5] proved that with n^d nodes embedded at grid points in a d -dimensional grid, with each node u

connected to its immediate neighbors and one long-distance neighbor v chosen with probability inversely proportional to $d(u, v)^d$, any message can be delivered in time polynomial in $\log n$ using greedy routing. While this result can be directly applied to our model with $d = 1$ and $\ell = 1$ to give a $O(\log^2 n)$ delivery time, we have a much simpler proof using Lemma 1, which we omit for lack of space.

The next interesting question to ask is if the performance improves by adding more long-distance links. We consider two different strategies for $\ell \in [1, \lg n]$ and $(\lg n, n^c]$, $c < 1$. In [6], Kleinberg uses a group structure to cover the first case with polylogarithmic links to get $O(\log n)$ delivery time. However, he uses a more complicated algorithm for routing while we obtain the same bound (for the case of a line) using only greedy routing.

We prove that with $\ell \in [1, \lg n]$, we get an expected delivery time of $O(\lg^2 n/\ell)$. The main idea behind the proof is that the delivery of a message is divided into *phases*. This is an extension of the idea used in [5]. A message is said to be in phase j if the distance from the current node to the destination node is between 2^j and 2^{j+1} . There are at most $(\lg n + 1)$ such phases. We analyze the expected number of links present between any two phases and this number increases as ℓ increases. We prove that the expected time spent in each phase is at most $O(\lg n/\ell)$, thus giving a total upper bound on the delivery time as $O(\lg^2 n/\ell)$.

For $\ell \in (\lg n, n^c]$, we use a deterministic strategy. The location of each node is identified as a number to a base b . With $O(b \log_b n)$ links per node, routing is done by forwarding the message to a node with an identifier closer to the target identifier by one digit. With at most $O(\log_b n)$ digits, we get the same delivery time. This strategy is similar in spirit to Plaxton's algorithm [10].

THEOREM 11. *Choose an integer $b > 1$. With $\ell = (b - 1)\lceil \log_b n \rceil$, let each node link to nodes at distances $1x, 2x, 3x, \dots, (b - 1)x$, for each $x \in \{b^0, b^1, \dots, b^{\lceil \log_b n \rceil - 1}\}$. Then the expected delivery time $T(n) = O(\log_b n)$.*

4.3.1 Failure of Links

We get reasonable performance even with link failures. We assume that each long-distance link is present independently with probability p , but that each node is always connected to its immediate neighbors. This ensures that a message will always be delivered even if it takes a long time.

With $\ell \in [1, \lg n]$, we use the same idea of delivering messages in phases. Intuitively, the expected time in each phase increases inversely proportional to p and we get the following result.

THEOREM 12. *Let each node be connected to its immediate neighbors (at distance 1) and $\ell \in [1, \lg n]$ long-distance neighbors chosen independently with replacement, with probability inversely proportional to the distance between the nodes. Assume that the links to the immediate neighbors are always present. If the proba-*

bility of a long-distance link being present is p , then the delivery time is $O(\lg^2 n/p\ell)$.

A similar intuition works for $\ell \in (\lg n, n^c]$. If a link fails, then the node has to take a shorter long-distance link, which will not take the message as close to the target as the initial failed link. Clearly as p decreases, the message has to take shorter and shorter links which increases the delivery time.

THEOREM 13. *Let each node have $\ell = O(\log_b n)$, long-distance links to distances $b^0, b^1, b^2, \dots, b^{\lceil \log_b n \rceil}$. Assume that the links to the nearest neighbors are always present. If the probability of a link being present is p , then the expected delivery time $T(n) = O(b \lg n/p)$.*

4.3.2 Failure of Nodes

The analysis for node failures is not as simple as that for link failures because we lose the important property of independence between links of different nodes. It is no longer the case that if one node cannot communicate with some other node, it has a good chance of doing so by passing the message to its neighbor. We analyze the situation when a node forwards a message to its next best neighbor after it reaches a dead neighbor.

To prove our result, we again use the formulation of a message moving between phases to reach the target. The idea is that the jumps between phases are independent so once we move from phase j to phase $j-1$, further routing no longer depends on any nodes in phase j . We can condition on the number of nodes being alive in the lower phase and estimate the time spent in each phase. Intuitively, if a node is present with probability p , we would expect to wait for a time inversely proportional to p in anticipation of finding a node in the lower phase to jump to.

THEOREM 14. *Let the model be as in Theorem 12. and let each node be present with probability p . Then the expected delivery time $T(n) = O(\lg^2 n/p\ell)$.*

In contrast, it appears that our *deterministic* routing strategy can lead to very poor performance; we have not yet analyzed this situation formally.

5. CONSTRUCTION OF GRAPHS

As the group of nodes present in the network changes, so does the graph of the virtual overlay network. In order for our routing techniques to be effective, the graph must always exhibit the property that the likelihood of any two vertices v, u being connected is $\Omega(d(v, u)^{-1})$. We describe briefly a heuristic approach to construct and maintain a graph with such an invariant.

Since the choice of edges leaving each vertex is independent of the choices of other vertices, we can assume that points in the metric space are added one at a time. Let v be the k -th point to be added. Point v chooses the sinks of its outgoing edges according to the inverse-distance distribution and connects to them by running the search algorithm. If a desired sink u is not present, v connects to u 's closest, present neighbor. In effect,

each of the $k - 1$ points already present before v is surrounded by a basin of attraction, collecting probability mass in proportion to its length. Since we assume the hash function populates the metric space evenly, and because of absolute symmetry, the basin length L has the same distribution for all points. It is easy to see that with high probability, L will not be much smaller than its expectation: $\text{Prob}(L \leq c \cdot k^{-1}) = 1 - (1 - c \cdot k^{-1})^{k-1}$. A lower bound on the probability that the edge (v, u) is present is $c' \cdot k^{-1} \cdot d(v, f)^{-1}$, where f is the point in u 's basin that is the farthest from v .³ However, the bound holds only if u is amongst the $k - 1$ points added before v . Otherwise, the aforementioned probability is 0, which means that we need to amend our linking strategy so as to transfer probability mass from the former case to the latter one. We describe next how to accomplish this task.

Let u be a new point. We give earlier points the opportunity to obtain outgoing edges to u by having u (1) calculate the number of incoming edges it “should” have from points added before it arrived, and (2) choose such points according to an appropriate distribution α .⁴ If ℓ is the number of outgoing edges for each point, then ℓ will also be the expected number of incoming edges that u has to estimate in step (1). For graphs with a large number of points n , each point has roughly $1/n$ chance of ending at u because of symmetry. The number of links ending at u is thus distributed according to a Poisson distribution with rate ℓ , that is, the probability that u has k incoming edges is $\frac{e^{-\ell} \ell^k}{k!}$, and the expectation of the distribution is ℓ . After step (2) is completed by u , each chosen point v responds to u 's request by using a distribution β to choose one of its existing outgoing edges to replace with an edge to u . It is easy to show that if α and β are proportional to the inverse of the distance between points, the resulting graph has the property we want. This procedure can be repeated to allow for regeneration of links when a network node crashes.

We believe that our method will give a distribution close enough to the desired distribution for the routing algorithm to work. Unfortunately, the construction process interacts in a complicated way with the routing algorithm, which makes analysis difficult. We are currently conducting simulations to assess the effectiveness of constructing a graph using our heuristic.

There has recently been related work [9] on how to construct, with the support of a central server, random graphs with many desirable properties, such as small diameter and guaranteed connectivity with high probability. Although it is not clear what kind of fault-tolerance properties this approach offers if the central server crashes, or how the constructed graph can be used for efficient routing, it is likely that similar techniques could be useful in our setting.

³The constant c' has absorbed c and the normalizing constant for the distribution.

⁴All this can be easily calculated by u since the link probabilities are symmetric.

6. EXPERIMENTAL RESULTS

We simulated a network at the *application level* with $n = 2^{17}$ nodes. Each node has $\lg n = 17$ links chosen using the randomized rule explained in Section 4.3. Routing is done greedily by forwarding a message to the neighbor closest to its target node. In each simulation, the network is set up afresh and a fraction p of the nodes fail. We choose random source and destination nodes which have not failed and route a message between them. For each value of p , we ran 1000 simulations delivering 100 messages in each simulation, and averaged the number of delivery hops for successful searches and number of failed searches.

With node failures, a node may not be able to find a live neighbor that is closer to the target node than itself. We studied three possible strategies to overcome this problem as follows: (i) Terminate the search. (ii) Randomly choose another node, deliver the message to this new node and then try to deliver the message from this node to the original destination node (similar to the hypercube routing strategy as explained in [13]). (iii) Keep track of a fixed number (in our simulations, 5) of nodes through which the message is last routed and backtrack. When the search reaches a node from where it cannot proceed, it backtracks to the most recently visited node from this list and chooses the next best neighbor to route to. For all these strategies we note that once a node chooses its best neighbor, it does not send the message to any other link if it finds out that the best neighbor has failed.

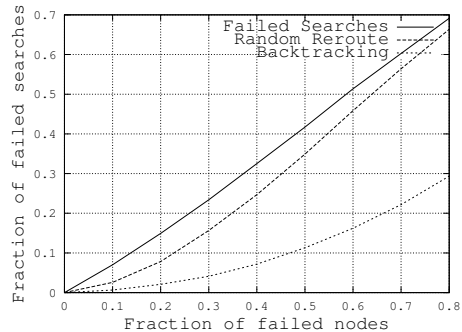


Figure 2: Fraction of failed searches

Figure 2 shows the fraction of messages that fail to be delivered versus the fraction of failed nodes. Figure 3 shows the number of hops for successful searches versus the fraction of failed nodes. It is very interesting to see how well the system behaves even with such a large number of failed nodes. In addition, backtracking gives a significant improvement in reducing the number of failures as compared to the other two methods, although it may take a longer time for delivery.

Our results may not be directly comparable to those of CAN[11] and Chord[12] since they use different simulators for their results. However, we see that we get results as good as theirs. Even if we just terminate the search, we get less than p fraction of failed searches with p fraction of failed nodes. Chord[12] has roughly

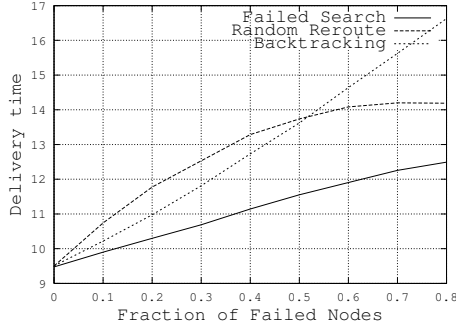


Figure 3: Delivery time

the same performance *after* their network stabilizes from some repair mechanism. Further, with backtracking we see that with 80% failed nodes, we still get less than 30% failed searches. These results are very promising and it would be very interesting to study backtracking analytically.

7. CONCLUSIONS AND FUTURE WORK

The following table summarizes our upper and lower bounds⁵:

Model	# Links ℓ	Upper Bound	Lower Bound
No failures	1	$O(\lg^2 n)$	$\Omega(\frac{\ln^2 n}{\ln \ln n})$
	$[1, \lg n]$	$O(\frac{\lg^2 n}{\ell})$	$\Omega(\frac{\ln^2 n}{\ell \ln \ln n})$
	$[\lg n, n^c]$	$O(\frac{\lg n}{\lg b})$	$\Omega(\frac{\lg n}{\lg \ell})$
$\Pr[\text{Link present}] = p$	$[1, \lg n]$	$O(\frac{\lg^2 n}{p\ell})$	-
	$[\lg n, n^c]$	$O(\frac{b \lg n}{p})$	-
$\Pr[\text{Node present}] = p$	$[1, \lg n]$	$O(\frac{\lg^2 n}{p\ell})$	-

We have shown that greedy routing in an overlay network organized as a random graph in a metric space can be a nearly optimal mechanism for searching a peer-to-peer system with low message complexity, even in the presence of many faults. We see this as an important first step in the design of efficient algorithms for such networks, but many issues still need to be addressed. Our results mostly apply to one-dimensional metric spaces like the line or a circle. One interesting possibility is whether similar strategies would work for higher-dimensional spaces, particularly ones in which some of the dimensions represent the actual physical distribution of the nodes in real space; good network-building and search mechanisms for this model might allow efficient location of nearby instances of a resource

⁵In the upper bound with $[\lg n, n^c]$ links, the number of links $\ell = O(b \log_b n)$. Also, the deterministic strategy used for links $\ell \in [\lg n, n^c]$, with link failures is slightly different than the one with no failures, and $\ell = O(\log_b n)$. In the lower bound column, the bound for $[1, \lg n]$ links is for one-sided routing.

without having to resort to local flooding (as in [4]). Another promising direction would be to study the security properties of greedy routing schemes, to see how they can be adapted to provide desirable properties like anonymity or robustness against Byzantine failures.

8. REFERENCES

- [1] GNUTELLA. <http://gnutella.wego.com>.
- [2] A. D. Joseph, J. Kubiawicz, and B. Y. Zhao. Tapestry: An infrastructure for fault-tolerant wide-area location and routing. Technical Report UCB/CSD-01-1141, University of California, Berkeley, Apr 2001.
- [3] R. M. Karp, E. Upfal, and A. Wigderson. The complexity of parallel search. *Journal of Computer and System Sciences*, 36(2):225–253, 1988.
- [4] D. Kempe, J. M. Kleinberg, and A. J. Demers. Spatial gossip and resource location protocols. In *Proceedings of 33rd Annual ACM Symposium on Theory of Computing*, pages 163–172, 2001.
- [5] J. Kleinberg. The small-world phenomenon: an algorithmic perspective. Technical Report 99-1776, Cornell University, Oct. 1999.
- [6] J. Kleinberg. Small-world phenomena and the dynamics of information. 2001.
- [7] MORPHEUS. <http://www.musiccity.com>.
- [8] NAPSTER. <http://www.napster.com>.
- [9] G. Panduragan, P. Raghavan, and E. Upfal. Building low-diameter P2P networks. In *Proceedings of 42nd Annual IEEE Symposium on the Foundations of Computer Science (FOCS)*, 2001.
- [10] C. Plaxton, R. Rajaram, and A. W. Richa. Accessing nearby copies of replicated objects in a distributed environment. In *Proceedings of the Ninth Annual ACM Symposium on Parallel Algorithms and Architectures (SPAA)*, June 1997.
- [11] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A scalable content-addressable network. In *Proceedings of the ACM SIGCOMM*, pages 161–170, 2001.
- [12] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishna. Chord: A scalable peer-to-peer lookup service for internet applications. In *Proceedings of SIGCOMM 2001*, pages 149–160, 2001.
- [13] L. Valiant. A scheme for fast parallel communication. *SIAM Journal on Computing*, 11:350–361, 1982.