

Consistency-preserving Neighbor Table Optimization for P2P Networks*

Huaiyu Liu and Simon S. Lam

Dept. of Computer Sciences, Univ. of Texas at Austin, Austin, TX 78712
{huaiyu, lam}@cs.utexas.edu

Abstract

Constructing and maintaining consistent neighbor tables and optimizing neighbor tables to improve routing locality are two important issues in p2p networks. In this paper, we address the problem of preserving consistency while optimizing neighbor tables for p2p networks with node dynamics. We present a general strategy: identify a consistent subnet as large as possible and only replace a neighbor with a closer one if both of them belong to the subnet. We realize the general strategy in the context of hypercube routing. First, we present a join protocol that enables the identification of a large consistent subnet with very low cost when new nodes join. Next, we define an optimization rule to constrain neighbor replacements to preserve consistency, and present a set of optimization heuristics to optimize neighbor tables with low cost. The join protocol is then integrated with a failure recovery protocol. By evaluating the protocols through simulation experiments, we found our protocols and optimization heuristics to be effective, efficient, and scalable to a large number of network nodes.

1 Introduction

Structured peer-to-peer networks are being investigated as a platform for building large-scale distributed systems [10, 11, 13, 14, 17]. The primary function of these networks is object location, that is, mapping an object ID to a node in the network. For efficient routing, each node maintains neighbor pointers in a table, called its *neighbor table*. The design of protocols to construct and maintain “consistent” neighbor tables for network nodes that may join, leave, and fail concurrently and frequently is an important issue. (Consistency ensures that a network is fully connected, i.e., there exists a path from any node to any other node.) Another important issue is to optimize neighbor tables so that the average distance traveled for each hop (locality) is optimized. Various ideas have been proposed to optimize neighbor tables for improving routing locality [1, 2, 3, 8, 12].

An important problem that has not been addressed adequately is how to preserve consistency (and thus preserve

established reachability) while optimizing neighbor tables, when there are nodes that join, leave, or fail concurrently and frequently. We address the problem in this paper and present a general strategy: Identify a consistent subnet as large as possible, and only allow a neighbor to be replaced by a closer one if both of them belong to the subnet. To implement this strategy in a distributed p2p network, where there is no global knowledge, the following problems need to be addressed: (1) how to identify nodes that belong to such a consistent subnet with minimum cost, (2) how to expand the subnet when new nodes join, and (3) how to maintain consistency of the subnet when nodes leave or fail.

In this paper, we realize the general strategy in the context of the hypercube routing scheme that is used in several proposed systems [10, 13, 17] to achieve scalable routing. In this scheme, given *consistent* [6] and *optimal* (that is, they store nearest neighbors) neighbor tables, it is guaranteed to locate a nearby copy of an object with asymptotically optimal cost if the object exists [10].

In [6], we have proposed a join protocol for the hypercube routing scheme. We proved that when an arbitrary number of nodes join an initially consistent network using the join protocol, the network is consistent again after all joins have terminated. The protocol is later extended to construct K -consistent neighbor tables to improve system robustness [4]. Correctness of the join protocol relies on preserved reachability: once a node can reach another node, it always can thereafter. In order not to break established reachability when replacing neighbors (to optimize neighbor tables), one approach is to apply optimization algorithms without interfering with joins, that is, applying optimization algorithms when joins have terminated and the network is already consistent. However, in a distributed p2p network, where nodes keep joining, it is difficult, if not impossible, to identify a quiescent time period in which there is no node joining and which is long enough for optimizations. Executing optimization algorithms while nodes are joining, on the other hand, may result in an inconsistent network, since replacing neighbors arbitrarily may break established reachability of some source-destination pairs, and thus affect the correctness of the join protocol.

We observe that within a subnet that is already consistent, replacing any neighbor with another, when both of them be-

*Research sponsored by NSF grant no. ANI-0319168 and Texas Advanced Research Program grant no. 003658-0439-2001

long to the subnet, does not break consistency conditions and thus does not break established reachability. (See Section 2.2 for the definition of consistency.) Following the observation, we first extend our join protocol in [4] so that at any time, the set of nodes whose join processes have terminated (including the nodes in the initial network) form a consistent subnet. The extended join protocol leads to solutions to the first two problems mentioned before: (1) identifying whether a neighbor is in the consistent subnet or not can be easily achieved by recording the state of the neighbor to indicate whether its join process has terminated or not; (2) the consistent subnet is expanded whenever a node's join process terminates, by including the node. Next, we integrate the extended join protocol with our failure recovery protocol presented in [5]. (Node leave is treated as a special case of failure.) The failure recovery protocol always tries to repair a hole left by a failed neighbor with a qualified node that is in the consistent subnet, thus it naturally follows the general strategy and provides a solution to problem (3). Through extensive simulation experiments [5], we have shown that the failure recovery protocol is able to maintain 1-consistency and re-establish K -consistency in every experiment with failures, for $K \geq 2$.

Contributions of this paper are the following:

- We present a general strategy to preserve consistency while optimizing neighbor tables for p2p networks with node dynamics.
- We extend the join protocol in [4] and prove that with the extended protocol, at any time t , the set of initial nodes plus the set of nodes whose joins have terminated form a consistent subnet. The extended protocol enables easy identification of nodes in the consistent subnet, and the costs of protocol extensions are shown to be very low.
- We present an optimization rule. Optimization algorithms should be applied within the constraint of this rule to preserve consistency. To optimize neighbor tables with low cost, we present a set of heuristics that search for nearby neighbors by primarily using information carried by join protocol messages.
- We integrate the extended join protocol with our failure recovery protocol and evaluate the protocols and the optimization heuristics by simulation experiments.
- We show that the extended join protocol and the optimization heuristics can also be used for initializing a K -consistent and optimized network.

Among related work, both Pastry [13] and Tapestry [17] make use of hypercube routing. In Pastry, in addition to a neighbor table for hypercube routing, each node maintains a set of nearest nodes on the ID ring, which is actively maintained and ensures success of routing as well as object location. Pointers for hypercube routing, on the other hand, are used as shortcuts and maintained lazily. Therefore, how to preserve established reachability while optimizing neighbor

tables is not addressed. Tapestry's join and failure recovery protocols are based upon use of a lower-layer Acknowledged Multicast protocol supported by all nodes [2], which also relies on established reachability. An algorithm to locate k nearest neighbors for each table entry, $k \geq 1$, is also presented [2]. However, it is not addressed how to preserve established reachability when nearest neighbors are located and old neighbors are replaced. Thus it is not clear how optimization operations will interfere with the correctness of their join protocol.

The rest of this paper is organized as follows. In Section 2, we briefly review the hypercube routing scheme, K -consistency, our original join protocol [4], and our theoretical foundation of protocol design and proofs. In Section 3, we present our general strategy for consistency-preserving optimization, extend the join protocol following the strategy, and present an optimization rule and a set of optimization heuristics. Correctness of the extended join protocol is proved and scalability of the protocol is analyzed. In Section 4, we evaluate the effectiveness of optimization heuristics by conducting simulation experiments in which nodes may join and fail concurrently and frequently. In Section 5, we explain how to initialize a K -consistent and optimized network. We conclude in Section 6.

2 Foundation

2.1 Hypercube routing scheme

In this section, we briefly review the hypercube routing scheme used in PRR [10], Pastry [13], and Tapestry [17]. Consider a set of nodes. Each node has a unique ID, which is a fixed-length random binary string. A node's ID is represented by d digits of base b , e.g., a 160-bit ID can be represented by 40 Hex digits ($d = 40$, $b = 16$). Hereafter, we will use $x.ID$ to denote the ID of node x , $x[i]$ the i th digit in $x.ID$, and $x[i-1]...x[0]$ a suffix of $x.ID$. We count digits in an ID from right to left, with the 0th digit being the *rightmost* digit. See Table 1 for notation used throughout this paper. Also, we will use "network" instead of "hypercube routing network" for brevity.

Notation	Definition
$\langle V, \mathcal{N}(V) \rangle$	a hypercube network: V is the set of nodes in the network, $\mathcal{N}(V)$ is the set of neighbor tables
$[\ell]$	the set $\{0, \dots, \ell-1\}$, ℓ is a positive integer
d	the number of digits in a node's ID
b	the base of each digit
$x[i]$	the i th digit in $x.ID$
$x[i-1]...x[0]$	suffix of $x.ID$; denotes empty string if $i = 0$
$x.table$	the neighbor table of node x
$j \cdot \omega$	digit j concatenated with suffix ω
$N_x(i, j)$	the set of nodes in (i, j) -entry of $x.table$, also referred as the (i, j) -neighbors of node x
$N_x(i, j).prim$	the primary (i, j) -neighbor of node x

Table 1. Notation

Given a message with destination node ID, $z.ID$, the objective of each step in hypercube routing is to forward the

message from its current node, say x , to a next node, say y , such that the suffix match between $y.ID$ and $z.ID$ is at least one digit longer than the match between $x.ID$ and $z.ID$.¹ If such a path exists, the destination is reached in $O(\log_b n)$ steps on the average and d steps in the worst case, where n is the number of network nodes. Figure 1 shows an example path for routing from source node 21233 to destination node 03231 ($b = 4, d = 5$). Note that the ID of each intermediate node in the path matches 03231 by at least one more suffix digit than its predecessor.

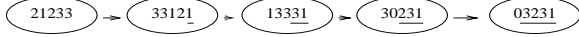


Figure 1. An example hypercube routing path

To implement hypercube routing, each node maintains a *neighbor table* that has d levels with b entries at each level. Each table entry stores link information (IDs and IP addresses) to nodes whose IDs have the entry's required suffix, defined as follows. (Hereafter, we will use “neighbor” or “node” instead of “node’s ID and IP address” whenever the meaning is clear from context.) Consider the table in node x . The *required suffix* for entry j at level i , $j \in [b]$, $i \in [d]$, referred to as the (i, j) -entry of $x.table$, is $j \cdot x[i-1] \dots x[0]$. Any node whose ID has this required suffix is said to be a **qualified node** for the (i, j) -entry of $x.table$. Nodes stored in the (i, j) -entry of $x.table$ are called the (i, j) -neighbors of x , denoted by $N_x(i, j)$. Ideally, these neighbors are chosen from qualified nodes for the entry according to some proximity criterion [10], with the nearest one designated as the *primary* (i, j) -neighbor. Furthermore, node x is said to be a *reverse* (i, j) -neighbor of node y if y is an (i, j) -neighbor of x . Each node also keeps track of its reverse-neighbors.

Note that node x has the required suffix for each $(i, x[i])$ -entry, $i \in [d]$, of its own table. For routing efficiency, we fill each node’s table such that $N_x(i, x[i]).prim = x$ for all $x \in V$, $i \in [d]$. Figure 2 shows an example neighbor table. The string to the right of each entry is the required suffix for that entry. An empty entry indicates that there does not exist a node in the network whose ID has the entry’s required suffix. For clarity, IP addresses are not shown in Figure 2.

Neighbor table of node 21233 ($b=4, d=5$)

^	01233	10233	0233	31033	033	22303	03	01100	0
11233	11233	21233	1233	03133	133	13113	13	33121	1
21233	21233	^	2233	21233	233	00123	23	12232	2
^	31233	03233	3233	^	333	21233	33	21233	3
level 4	level 3	level 2	level 1	level 0					

Figure 2. An example neighbor table

2.2 K -consistent networks

Constructing and maintaining consistent neighbor tables is an important design objective for structured peer-to-peer

¹In this paper, we follow PRR [10] and use suffix matching, whereas other systems use prefix matching. The choice is arbitrary and conceptually insignificant.

networks. We defined consistency for a hypercube routing network as follows [6]: A network, $\langle V, \mathcal{N}(V) \rangle$, is **consistent** if and only if the following conditions hold: (i) For every table entry in $\mathcal{N}(V)$, if there exists at least one qualified node in V , then the entry stores at least one qualified node; (ii) otherwise, the entry is empty. In a consistent network, any node x can reach any other node y using hypercube routing in k steps, $k \leq d$; more precisely, there exists a neighbor sequence (**path**), (u_0, \dots, u_k) , $k \leq d$, such that u_0 is x , u_k is y , and $u_{i+1} \in N_{u_i}(i, y[i])$, $i \in [k]$.

If nodes may fail frequently in a network, a natural approach to improve robustness is to store in each table entry multiple qualified nodes. For this approach, we generalized the definition of consistency to K -consistency as follows [4]. A network, $\langle V, \mathcal{N}(V) \rangle$, is **K -consistent** if and only if the following conditions hold: (i) For every table entry in $\mathcal{N}(V)$, if there exist H qualified nodes in V , $H \geq 0$, then the entry stores at least $\min(K, H)$ qualified nodes; (ii) otherwise, the entry is empty. For $K \geq 1$, K -consistency implies consistency (in particular, 1-consistency is the same as consistency).

2.3 Join protocol

In [4], we presented a join protocol for the hypercube routing scheme and proved that it constructs and maintains K -consistent neighbor tables for an arbitrary number of concurrent joins. Here we briefly review the protocol design.

In designing and proving the correctness of the protocol for nodes to join a network $\langle V, \mathcal{N}(V) \rangle$, we made the following assumptions: (i) $V \neq \emptyset$ and $\langle V, \mathcal{N}(V) \rangle$ is a K -consistent network, (ii) each joining node, by some means, knows a node in V initially, (iii) messages between nodes are delivered reliably, and (iv) there is no node leave or node failure during the joins. Then, tasks of the join protocol are to update neighbor tables of nodes in V and to construct tables for the joining nodes so that after the joins, the network is K -consistent again.

Each node in the network maintains a state variable named *status*, which begins in *copying*, then changes to *waiting*, *notifying*, and *in_system* in that order. A node in status *in_system* is called an *S-node*; otherwise, it is a *T-node*. Each node also stores, for each neighbor in its table, the neighbor’s state, which can be *S* indicating that the neighbor is an S-node or *T* indicating that it is not yet.

In status *copying*, a joining node, say x , copies neighbor information from S-nodes to fill in most entries of its table level by level. It copies level-0 neighbor information from the node it knows in V (an S-node), say g_0 , and finds an S-node g_1 among the level-0 neighbors of g_0 such that $g_1.ID$ shares the rightmost digit with $x.ID$. x then copies level-1 neighbors from g_1 , and finds an S-node g_2 that shares the rightmost two digits with it, and so on. When after coping level- $(i-1)$ neighbors, x cannot find an S-node that shares the rightmost i digits with it, $i \geq 1$, x changes status to

waiting. In this status, x tries to “attach” itself to the network, i.e., to find an S-node, say y , that shares at least the rightmost $i - 1$ with x and stores x as a neighbor. When x is attached, its status becomes *notifying*. Then, x seeks and notifies nodes that share the rightmost j digits with it, where j is the lowest level that x is stored in y ’s table (the attach-level of x , as defined in [4]). Lastly, when it finds no more node to notify, x changes status to *in_system* and becomes an S-node.

Figure 3 presents the protocol messages. In particular, *JoinWaitMsg* is the message that a joining node sends out to request for attachment. It is worth pointing out that when a node, y , receives a *JoinWaitMsg* from some joining node, y processes the message and replies immediately if y is already an S-node; otherwise, y saves the message to be processed later when it becomes an S-node. That is, a joining node is always stored as a neighbor by an S-node first.

CpRstMsg, sent by x to request a copy of receiver’s neighbor table.
CpRlyMsg(x .table), sent by x in response to a *CpRstMsg*.
JoinWaitMsg, sent by x to notify receiver of the existence of x and request the receiver to store x , when x .status is *waiting*.
JoinWaitRlyMsg(r , i , x .table), sent by x in response to a *JoinWaitMsg*, when x .status is *in_system*. $r \in \{\text{negative}, \text{positive}\}$, i : an integer.
JoinNotiMsg(i , x .table), sent by x to notify receiver of the existence of x , when x .status is *notifying*. i : an integer.
JoinNotiRlyMsg(r , Q , x .table, f), sent by x in response to a *JoinNotiMsg*.
 $r \in \{\text{negative}, \text{positive}\}$, Q : a set of integers, $f \in \{\text{true}, \text{false}\}$.
SpeNotiMsg(x , y), sent or forwarded by a node to inform receiver of the existence of y , where x is the initial sender.
SpeNotiRlyMsg(x , y), response to a *SpeNotiMsg*.
InSysNotiMsg, sent by x when x .status changes to *in_system*.
RvNghNotiMsg(y , s), sent by x to notify y that x is a reverse neighbor of y , $s \in \{T, S\}$.
RvNghNotiRlyMsg(s), sent by x in response to a *RvNghNotiMsg*, $s = S$ if x .status is *in_system*; otherwise $s = T$.

Figure 3. Join protocol messages

2.4 C-set tree

C-set tree is a conceptual foundation for guiding our protocol design and reasoning about K -consistency [4, 6]. To introduce C-set trees, we first present the notion of *notification set of x regarding V* , denoted by V_x^{Notify} [4]. Suppose a set of nodes W join a K -consistent network $\langle V, \mathcal{N}(V) \rangle$ and $x \in W$. Intuitively, V_x^{Notify} is the set of nodes in V that need to update their tables if x were the only node that joins $\langle V, \mathcal{N}(V) \rangle$.

Intuitively, a C-set tree organizes nodes in V that need to update their tables as well as nodes in W into a tree structure, if the notification sets regarding V (*noti-sets*, in short) of all nodes in W are the same. Generally, the noti-sets of all nodes in W may not be the same. Then, nodes in W with the same noti-set belong to the same C-set tree and the C-set trees for all nodes in W form a forest. Each C-set tree in the forest can be treated separately in proving protocol correctness. In the rest of this subsection, we focus on a single C-set tree, i.e., we assume that the noti-sets of the joining

nodes are the same. (Formal definitions for C-set trees are presented in [4, 6].)

Given V , W and K , the structure of the C-set tree is determined, which we call a *C-set tree template*. For example, suppose $W = \{30633, 41633, 10533\}$ ($b = 8, d = 5$) and $V = \{02700, 14263, 62332, 72413\}$. The corresponding C-set tree template is shown in Figure 4(a). Here we assume $K = 1$ to simplify illustration. In this example, noti-set of the joining nodes is the set of nodes in V with suffix 3, denoted by V_3 . Observe that the joining nodes introduce new suffixes to the network. For each new suffix, there is a corresponding C-set, and all C-sets plus set V_3 form a tree according to their suffixes.

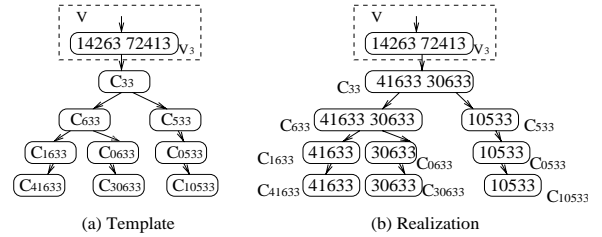


Figure 4. C-set tree example

The task of the join protocol is to construct and update neighbor tables such that paths are established between nodes; *conceptually* nodes are filled into each C-set. For instance, in the above example, when 14263 stores a node with suffix 33, say node 30633, in its (1, 3)-entry, then conceptually 30633 is filled into C_{33} . We call the C-set tree realized at the end of all joins a *C-set tree realization*. Figure 4(b) shows one possible realization of the template in Figure 4(a). At the end of joins, we check whether some correctness conditions [4] are satisfied by the C-set tree realization. If they are, then neighbor tables of nodes in $V \cup W$ are guaranteed to be K -consistent.

3 Consistency-preserving Optimization

To date, correctness of proposed join protocols for the hypercube routing scheme [2, 4, 6] depends on preserved reachability, i.e., once a node can reach another node, it always can thereafter. Therefore, if optimization operations are to be performed, they should preserve reachability. There is a common operation in all optimization algorithms: replacing an old neighbor with a new one that is measured to be closer. However, if there is no constraint on such a replacement, it may break reachability of some source-destination pairs, affect correctness of the join protocol, and result in an *inconsistent* network after nodes join.

For example, suppose nodes 41633 (x) and 30633 (y) join a network concurrently with some other nodes. Let t_2 be the time that neighbor pointers along the path from x to y are completely established. Then x cannot reach y before time t_2 . If at some time t_1 , $t_1 < t_2$, some node that has stored y , say node 14263 (u), finds x to be closer and replaces y with x , then after the replacement, u cannot reach

y until time t_2 , as illustrated by Figure 5. In this case, reachability of pair (u, y) is not preserved by the optimization operation even if both join processes of x and y have terminated by time t_1 , since some nodes along the path from x to y may be still joining and neighbor pointers are still being established. Then, during the period $[t_1, t_2]$, joining nodes that are supposed to find out y through u will fail to do so and thus cannot construct their neighbor tables correctly. Even worse, the period may be arbitrarily long, if messages are delayed arbitrarily long in the network, or if reachability of some source-destination pair along the path from u to y is also broken.

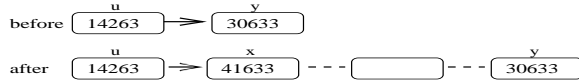


Figure 5. Paths before and after neighbor replacement

To construct and optimize neighbor tables without breaking established reachability when new nodes join a network, one possible approach is to first construct and update neighbor tables so that they are K -consistent, and then optimize neighbor tables after the joins. However, this approach is not practical in a distributed p2p network, since nodes keep joining and none of them is aware of any quiescent time period in which there is no node joining and which is long enough for optimization operations, if such a period exists.

3.1 Our strategy

We observe that for the hypercube routing scheme, within a subnet that is already consistent, replacing any neighbor with any other neighbor does not break consistency conditions if both neighbors belong to the consistent subnet. (Basically, consistency conditions require that for each table entry, if there exists qualified nodes in the subnet, then the entry is filled with at least such a node.) If a neighbor replacement does not break the consistency conditions, then after the replacement, nodes that are previously reachable via the old neighbor can now be reached via the new neighbor. This observation is also applicable to other structured p2p networks, such as the system proposed in [9].

When new nodes are joining a network, if we can identify a “core” of the network such that if we only consider the nodes in this core, their neighbor tables are consistent and they can reach each other, then we know that replacing a neighbor with a closer neighbor, both of which are in the core, is a safe operation and will not break established reachability. Note that before the joins start, the initial network is consistent and thus is the “core” of the network. However, if we optimize neighbor tables by only considering nodes in the initial network, the extent of optimization would be greatly limited. It is desired that after a node has joined the network, it becomes part of the core so that it can also be considered for optimization. It is also desired that when nodes fail, consistency of the core is maintained. To summarize, we present a general strategy for consistency-

preserving neighbor table optimization in presence of node dynamics.

A general strategy for consistency-preserving optimization: Identify a consistent subnet as large as possible; only allow a neighbor to be replaced by a closer one if both of them belong to the subnet; expand the consistent subnet after new nodes join; and maintain consistency of the subnet when nodes fail.

The join protocol in [4] guarantees that when a set of nodes join an initially K -consistent network, the network is K -consistent (and thus consistent) again after all join processes terminate. To implement the above strategy, we need another property from the join protocol: *at any time, the subnet consisting of all nodes whose join processes have terminated plus nodes in the initial network is consistent.* With this property, identifying nodes or neighbors that belong to the consistent subnet becomes easy: if the join process of a node has terminated, then it belongs to the subnet; otherwise, it is not. The property also ensures that the consistent subnet keeps growing when more join processes terminate. To maintain consistency of the subnet when nodes fail, a failure recovery protocol is needed to recover K -consistency.² The failure recovery protocol should always try to recover a hole left by a failed neighbor with a qualified node that is in the consistent subnet.

Recall that in our protocol design, when a node’s join process terminates, it becomes an S-node. (Nodes in the initial network are also S-nodes.) Hence, more specifically, our goals are to (1) design a join protocol so that at any time, the set of S-nodes form a consistent subnet, and (2) design a failure recovery protocol that recovers K -consistency of the subnet by repairing holes left by failed neighbors with qualified S-nodes. The failure recovery protocol presented in [5] naturally fits into the general strategy with minor extensions. Basically, it works in the following way. When a neighbor failure is detected by a node, a recovery process is initiated. The process always tries to repair a hole left by the failed neighbor with a qualified S-node, by searching in the node’s own neighbor table and querying the node’s neighbors. Only when it fails to find a qualified S-node will it repair the hole with a T-node. The failure recovery protocol has been shown to maintain consistency and re-establish K -consistency for networks with $K \geq 2$. Therefore, in this section, we focus on how to extend the join protocol in [4] to achieve goal (1).

3.2 Extended join protocol

To extend the join protocol, we first consider the basis of the proofs of protocol correctness. Proofs in [4] rely on the following properties of a network.

1. Once two S-nodes can reach each other, they always can thereafter.

² K -consistency provides redundancy in neighbor tables to ensure that a dynamically changing network remains fully connected.

2. Once a T-node can reach an S-node, it always can thereafter.
3. After a T-node, say x , is stored by another node, say y , x remains in the table of y when x is still a T-node.

If there is no table optimization involved during the joins, i.e., no neighbor in any entry would be replaced, the above properties hold trivially: once a path is established, the neighbor pointers from one hop to another along the path are always there and remain the same. When there are optimization operations that happen concurrently with joins, the above three properties must be preserved to ensure the correctness of the join protocol. To preserved property 3 is not difficult: we require that if a neighbor is still a T-node, it cannot be replaced even if another node is found to be closer than it. To preserve properties 1 and 2, goal (1) stated above needs to be achieved and neighbor replacement should be constrained to neighbors that are S-nodes.

We extend the join protocol to achieve goal (1) as follows. In short, a new status, *cset_waiting*, is inserted between *notifying* and *in_system*. When a joining node has finished its tasks and exited status *notifying*, it will not change to status *in_system* and become an S-node immediately. Instead, the node waits in status *cset_waiting* for some nodes that are joining concurrently and are likely to be in the same C-set with it (conceptually). When it is confirmed that all these nodes have exited status *notifying*, it changes status to *in_system*. (Pseudo-code of the extended join protocol is presented in [7].) The extensions ensure that when two nodes have both become S-nodes, paths between them (in both directions) have already been established.

- A new joining status, *cset_waiting*, is added after status *notifying*. Moreover, one more join protocol message, *SameCsetMsg(s)*, is introduced, where s is S if the sender is already an S-node and T otherwise.
- When a node, say x , receives a *JoinNotiMsg* or a *JoinNotiRlyMsg*, the message includes a copy of the sender's table. If x is in status *notifying* when it receives the message, and if from the copy of the sender's table, x finds a T-node, say y , that shares with x a suffix of length k , $k \geq x.att_level$, x saves y in set Q_{cset_wait} . ($x.att_level$ is the attach-level of x in the network [4], which is the lowest level x is stored in the table of the first S-node that stored x .)
- When a node in status *notifying* finds that it is not expecting any more *JoinNotiRlyMsg* or *SpeNotiRlyMsg*, it changes status to *cset_waiting*. It then sends a *SameCsetMsg(T)* to each node in set Q_{cset_wait} and waits for their replies. It also replies to each node in set Q_{cset_recv} (see discussion below) with a *SameCsetMsg(T)*. Each node that is in both Q_{cset_recv} and Q_{cset_wait} is then removed from Q_{cset_wait} .
- When a node, say x , receives a *SameCsetMsg(s)*, if it is already in status *in_system*, it sends a *SameCsetMsg(S)* back immediately if s is T (if s is S , x simply ignores

the message). If x is in status *cset_waiting*, it sends a *SameCsetMsg(T)* back immediately if it has not done so, and removes the sender from Q_{cset_wait} . If x is in any other status, x saves the sender into Q_{cset_recv} to reply later when x changes status from *notifying* to *cset_waiting*.

- When a node is in status *cset_waiting* and finds that Q_{cset_wait} is empty, it changes status to *in_system*.

The above extensions add extra delay into each join process. With the extra delay, a joining node will not become an S-node until it believes that nodes currently in the same C-set with it (conceptually) have all entered status *cset_waiting* or *in_system*. Since only after a node becomes an S-node can it store another joining node that has requested it for attachment (by sending a *JoinWaitMsg*), the above extensions ensure that only after a set of nodes in a parent C-set have all finished their joining tasks, will new joining nodes be attached to these nodes and filled into children C-sets. In the correctness proof [7], we show that when a new node is filled into a child C-set, neighbor pointers among the nodes that have been filled in ancestor C-sets have been established and those nodes already can reach each other.

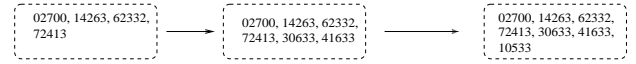


Figure 6. Evolution of consistent subnet

For instance, consider the example mentioned in Section 2.4 (the C-set tree template of which, assuming $K = 1$, is shown in Figure 4(a)). With the extended join protocol, the C-set tree is realized in the following way: only after C-set C_{33} is filled and nodes in it have all entered status *cset_waiting* or *in_system*, will new nodes (nodes other than those in C_{33}) be filled into the children C-sets, C_{633} and C_{533} , and so on.³ For example, for the realization as shown in Figure 4(b), it is realized as follows: only after nodes 41633 and 30633 (nodes in C_{33}) have entered status *cset_waiting* or *in_system*, will node 10533 be filled into C_{533} . Figure 6 shows the corresponding evolution of the consistent subnet.

3.3 Correctness and scalability of join protocol

We first present Theorem 1, which shows that when a set of new nodes join a network using the extended join protocol, at any time, all S-nodes at that time belong to a consistent subnet. This property guarantees that replacing a neighbor with another one is safe if both of them are S-nodes. Proof of Theorem 1 is based on the assumptions stated in Section 2.3. Proof details are presented in [7] and are omitted here due to space limitation.

³A node is a neighbor of itself and is stored in each entry whose required suffix is a suffix of its node ID. Therefore, after a node is filled into a C-set, it is automatically filled into descendant C-sets. For instance, when 41633 is filled into C_{33} , it is automatically filled into C_{633} , C_{1633} , and C_{41633} .

Theorem 1 Suppose a set of nodes, $W = \{x_1, \dots, x_m\}$, $m \geq 1$, join a K -consistent network $\langle V, \mathcal{N}(V) \rangle$ using the extended join protocol. Then at any time t , any node in set $S(t)$ can reach any other node in $S(t)$, where $S(t)$ is the set of S -nodes at time t .

Next, we demonstrate the scalability of the extended join protocol by analyzing communication costs of protocol extensions through simulation experiments. We implemented the extended join protocol in an event-driven simulator, and used the GT-ITM package [15] to generate network topologies. For a generated topology with a set of routers, overlay nodes (end hosts) were attached randomly to the routers. For the simulations reported in this paper, two topologies were used: a topology with 1056 routers to which 1000 overlay nodes were attached, and a topology with 2112 routers to which 4000 overlay nodes were attached. We simulated the sending of a message and the reception of a message as events, but abstracted away queueing delays. The end-to-end delay of a message from its source to destination was modeled as a random variable with mean value proportional to the shortest path length in the underlying network. For the 1056-router topology, end-to-end delays are in the range of 0 to 329 ms, with the average being 113 ms; for the 2112-router topology, end-to-end delays are in the range of 0 to 596 ms, with the average being 163 ms. In each experiment, we let m nodes join an initial network of n nodes, $m \gg n$. We set parameters b to be 16 and d to be 8.⁴

We first study the extra delay caused by the new status, *cset_waiting*. We define the **join duration** of a node to be the duration from the time the node starts joining to the time it changes status to *in_system*. Figure 7(a) plots the average join durations for 990 nodes joining an initial network of 10 nodes, as a function of K , for simulations using the original join protocol (presented in Section 2.3) and the extended join protocol, respectively. The underlying topology was the 1056-router topology. In each experiment, all joins started at exactly the same time. As shown in the figure, the average join durations for the extended protocol are only slightly longer than those for the original protocol, which indicates that the extra delay caused by waiting in status *cset_waiting* is small. The same conclusion can be drawn from Figure 7(b), where 1990 nodes joined an initial network of 10 nodes and the underlying topology is the 2112-router topology. Error-bars in Figure 7 show the minimum and maximum join durations observed from simulations using the extended join protocol.

Next, we study communication costs of the extended join protocol in terms of numbers of messages sent by a joining node. In [4], we have analyzed numbers of protocol messages sent by a joining node, for all message types except the one introduced in this paper, and showed that the communication costs are scalable to large networks. Hence, in

⁴In Tapestry, $b = 16$ and $d = 40$. In Pastry, $b = 16$ and $d = 32$. We found that the value of d is insignificant when $b^d \gg n$, where n is the number of nodes in a network.

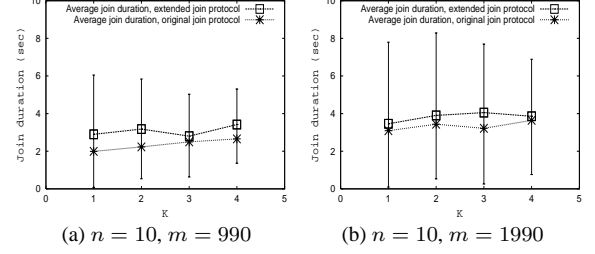


Figure 7. Join durations with/without protocol extensions

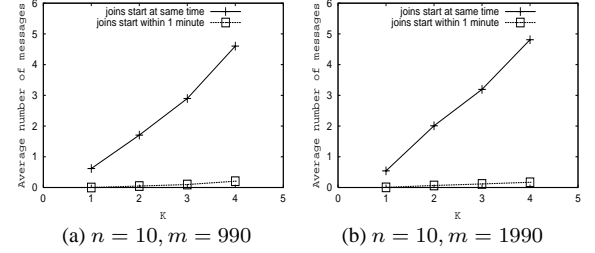


Figure 8. Average number of *SameCsetMsg*

this paper we only need to study the number of the new message (*SameCsetMsg*) sent by a joining node.

Figure 8 presents average numbers of *SameCsetMsg* sent by joining nodes as a function of K . The numbers are small in general, and increase when K increases. This is because when K increases, more neighbors are stored in each entry and thus each C-set tends to contain more nodes. By comparing the two curves in each diagram, we observe that in the simulations where joins did not start at exactly the same time, average numbers of *SameCsetMsg* were greatly reduced. Moreover, comparing Figure 8(a) and Figure 8(b), we see that with other parameters being the same, the average number of *SameCsetMsg* remained almost the same when the number of concurrent joins (m) was increased from 990 to 1990.

We conclude that the communication costs of the protocol extensions are very low and the extended join protocol is scalable to a large number of network nodes.

3.4 Optimization rule and heuristics

We now have an extended join protocol that expands the consistent subnet while nodes join a network, and a failure recovery protocol [5] that maintains consistency of the consistent subnet when nodes fail. To implement the general strategy (Section 3.1), we also need the following rule.

Optimization Rule When a node, x , intends to replace a neighbor, y , with a closer one, z , the replacement is only allowed when both y and z are S -nodes.

Recall that for each neighbor, a node stores the state of the neighbor. State S indicates that the neighbor is in status *in_system*, while state T indicates it is not yet. To implement the above rule, when x intends to replace y with z , it only does so when the states associated with both y and z are S . With the extended join protocol and the optimization rule, the three properties stated in Section 3.2 will be

preserved even when optimization operations happen concurrently with joins [7].

To optimize neighbor tables, an algorithm is needed to search for qualified nodes that are closer than current neighbors. We next present a set of heuristics to optimize neighbor tables when new nodes are joining a network and new tables are constructed. To search for closer neighbors with low cost, the heuristics are designed by primarily utilizing information carried in join protocol messages. Notice that whenever a closer neighbor is found for a table entry, it can be used to replace an old neighbor *only if* the replacement is allowed by the optimization rule.

Heuristic 1: Copy neighbor information from nearby nodes. Recall that in the *copying* status, a joining node, x , constructs most part of its neighbor table by copying neighbor information from other nodes (S-nodes). Suppose y is the node that x starts joining with. Instead of directly copying level-0 neighbors from y , x chooses the closest node from y 's neighbors, say g_0 , and copies level-0 neighbors from g_0 . If the level-0 neighbors of g_0 are close to g_0 , and g_0 and x are close to each other, then it is highly likely that these level-0 neighbors are also close to x [1]. To copy level-1 neighbors, x chooses a level-0 neighbor of g_0 that shares suffix $x[0]$ with it, say z , if such a node exists. Then from the level-1 neighbors of z (whose IDs all have suffix $x[0]$), x chooses the closest one to copy level-1 neighbors from, and so on.

Heuristic 2: Utilize protocol messages that include copies of neighbor tables. During status *waiting* and *notifying*, a joining node, x , sends out messages (*JoinWaitMsg* and *JoinNotiMsg*) to some nodes to notify them about itself. Replies to these messages all include copies of the neighbor tables of the senders. From each reply message, x searches for qualified nodes that are closer than some current neighbors for every table entry. Moreover, when x is in status *notifying*, a notification message sent by x includes a copy of $x.table$. The receiver of such a message also searches for closer nodes in $x.table$ to replace old neighbors.

Heuristic 3: Optimize neighbor tables when a node's join process terminates. When a joining node, x , changes status to *in_system*, it informs its reverse-neighbors (nodes that have stored x as a neighbor) as well as its neighbors that it becomes an S-node. These nodes then update the state of x to be *S* in their tables and try to optimize their table entries for which x is a qualified node. In addition to informing neighbors, x exchanges neighbor tables with its neighbors (not including reverse-neighbors) so that both x and its neighbors can optimize their tables at this time.

4 Experimental Results

We have integrated the extended join protocol with our failure recovery protocol and the optimization heuristics, under the constraint of the optimization rule. In this section, we validate our strategy for consistency-preserving op-

timization and evaluate the effectiveness of the heuristics through simulation experiments. To evaluate the optimization heuristics, we use a metric called p-ratio, defined below. Recall that the closest neighbor in an entry is called the primary-neighbor of that entry. For a table entry of a node, say x , suppose the primary-neighbor of the entry is y , and the closest node among all qualified nodes of the entry is z . We define **p-ratio** of the entry to be the ratio of the communication delay from x to y to the delay from x to z . A p-ratio of 1 indicates that y and z are of the same distance. If for every table entry in a network, p-ratio is 1, then the neighbor tables are optimal.

4.1 Optimization during joins

In each experiment where optimization happened concurrently with joins, we let m nodes join an initial K -consistent network of n nodes, $m \gg n$. Neighbor tables were then constructed, updated, and optimized according to the extended join protocol and the optimization heuristics. In the protocol implementation, an old neighbor is only replaced by a new neighbor if the distance of the new one is measured to be 10% shorter than the old one (plus that the replacement is allowed by the optimization rule). This is to prevent oscillation, since each end-to-end delay is modeled as a random number with a mean value proportional to the shortest path length in the underlying network. When all join processes had terminated, we checked whether K -consistency was maintained and calculated p-ratio for every table entry.

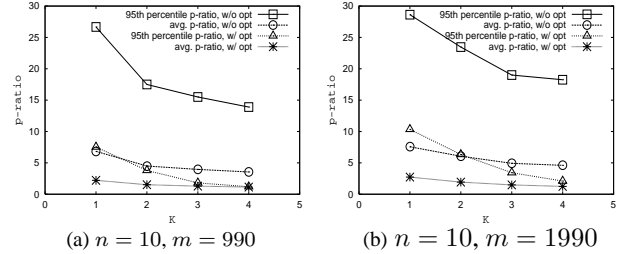


Figure 9. Effectiveness of optimization heuristics

Figures 9 presents results from experiments with $n = 10$ and $m = 990$, and from experiments with $n = 10$ and $m = 1990$. In each experiment, starting times of the joins were drawn randomly from range [0s, 60s] (i.e., all nodes joined within 1 minute). The results show that by primarily using information carried in join protocol messages, table entries can be greatly optimized. For instance, in Figure 9(a), without any optimization, the average p-ratio for $K = 1$ is more than 6.82, and the 95th percentile of p-ratio for $K = 1$ is 26.67 (i.e., 95% of p-ratios are no greater than 26.67); with the optimization heuristics, the values drop to 2.21 and 7.51, respectively. We also found that in every experiment, K -consistency was maintained after all joins had terminated, which demonstrates that our strategy preserves consistency and ensures correctness of the join protocol.

Results in Figure 9 also show that when K is increased, the average p-ratio decreases. The reason is that when K

becomes larger, more neighbors are stored in each table entry, thus more neighbor information is carried in protocol messages. Clearly, there is a tradeoff between the benefits and maintenance costs of K -consistency.⁵

4.2 Optimization with concurrent joins and failures

The extensions to the join protocol presented in this paper do not affect failure recovery actions, thus integrating the extended join protocol with the failure recovery protocol should not affect success of failure recoveries. On the other hand, since a substitute for a failed neighbor is searched locally (see Section 3.1), if neighbor tables have been optimized, the substitute node would not be too far away. Hence the average p-ratio would not be affected too much after a recovery action. Therefore, integration of the extended join protocol, the failure recovery protocol, and the optimization heuristics should be effective and stable in both consistency maintenance and neighbor table optimization.⁶ To demonstrate this, we conducted experiments with concurrent joins and failures as well as churn experiments.

Massive joins and failures We first conducted simulations in which massive number of joins and failures happened concurrently. Each experiment began with a K -consistent network, $\langle V, \mathcal{N}(V) \rangle$, which was constructed and optimized by the extended join protocol and optimization heuristics. Then, a set W of nodes joined and a set F of randomly chosen nodes failed. Join and failure events were generated according to a Poisson process at the rate of 10 events every second.

From the experiments, we found that K -consistency was maintained when all join and failure recovery processes had terminated, in every experiment with $K \geq 2$. This result indicates that our protocols are effective in consistency maintenance. Figure 10 presents results of average p-ratios at the end of the simulations. The lower curve presents results from simulations where 494 joins and 506 failures happened in a network that initially had 1000 nodes. The upper curve presents results from simulations where 968 joins and 1032 failures happened in a network that initially had 2000 nodes. As shown in the figure, even with massive joins and failures, the table entries were still optimized greatly: For $K \geq 2$, average p-ratios were less than 3.

Churn experiments We also investigated the impact of continuous node dynamics on protocol performance. To simulate node dynamics, Poisson processes with rates λ_{join} and λ_{fail} were used to generate join and failure events, respectively. We set $\lambda_{join} = \lambda_{fail} = \lambda$, which is said to be the *churn rate*. For each join event, a new node (T-node) was given a randomly chosen S-node to begin its join process. For each failure event, an S-node or a T-node was randomly chosen to fail and stay silent. Periodically in each

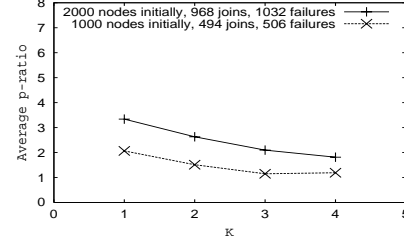


Figure 10. Optimization with massive joins and failures

experiment, we took snapshots of the neighbor tables of all S-nodes (the “core” of the network). For each snapshot, we calculated the average p-ratio as an indicator of how well table entries were optimized at the moment. We also checked whether consistency was maintained at each snapshot.

Figure 11 presents results from an experiment with $\lambda = 1$, i.e., join events were generated at a rate of 1 per second and so were the failure events. The initial K -consistent network of 2000 nodes, $K = 3$, was constructed and optimized by letting 1990 nodes join a network of 10 nodes. In the experiment, join and failure events were generated from the 1,000th second to the 4,000th second (simulated time). After that, no more join or failure events was generated and the experiment continued until all join, failure recovery, and optimization processes terminated. Snapshots were taken every 50 seconds. The lower curve in Figure 11(a) plots the average p-ratio for each snapshot. Although there were continuous joins and failures, neighbor tables remained optimized to a certain degree: The average p-ratio increased slightly at first, when joins and failures started to happen; it then remained below 2.3. (For comparison, the upper curve shows the average p-ratios from an experiment with the same simulation setup, in which no optimization heuristics were applied.) We also found that consistency was maintained at every snapshot, and K -consistency ($K = 3$) was recovered at the end of the simulation. Figure 11(b) plots the number of nodes in the network (T-nodes and S-nodes) versus the number of S-nodes for each snapshot. Note that the two curves are very close to each other, which demonstrates that at the given churn rate, the size of the subnet formed by S-nodes is consistently close to that of the entire network. It also demonstrates that with the given churn rate and the network size, our protocols can sustain a large stable “core” over the long term even when joins, failures, and neighbor table optimization happen concurrently.⁷

5 Network Initialization

To initialize a K -consistent and optimized network of n nodes, we can put any one of the nodes, say x , in V , and construct $x.table$ as follows. (Let $x.state(y)$ denote the state of neighbor y stored in the table of x .)

- $N_x(i, x[i]).prim = x, x.state(x) = S, i \in [d]$.
- $N_x(i, j) = \emptyset, i \in [d], j \in [b]$ and $j \neq x[i]$.

⁷In [5], we have studied “sustainable churn rates” in detail.

⁵In [5], we had investigated the tradeoff in detail.

⁶In [5], we have shown that the integration of the original join protocol and the failure recovery protocol is effective and stable in consistency maintenance.

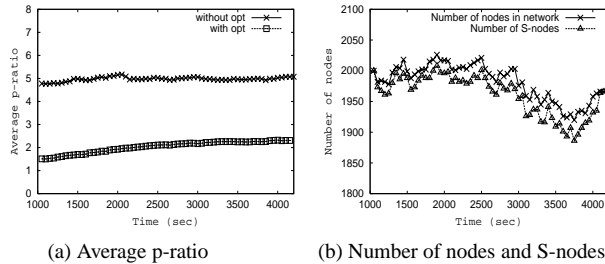


Figure 11. Churn experiment, $\lambda = 1$, $K = 3$

Next, let the other $n - 1$ nodes join the network concurrently. Each node is given x to start with and executes the extended join protocol with the optimization heuristics implemented. At the end of joins, a K -consistent network is constructed and table entries are optimized.

6 Conclusions

Constructing and maintaining consistent neighbor tables and optimizing neighbor tables to improve routing locality are two important issues in p2p networks. To construct and maintain consistent neighbor tables in presence of node dynamics, especially when new nodes are joining, it is desired that neighbor pointers remain unmodified once they are established so that new nodes are ensured to construct neighbor tables correctly following the pointers. On the other hand, to improve routing locality, it is desired that once closer neighbors are found, old neighbors that are farther away are replaced.

In this paper, we showed that the “divergence” between the two issues can be resolved by a general strategy: to replace a neighbor with a closer one only when they both belong to a consistent subnet. We realized the strategy in the context of hypercube routing. We first extended our join protocol in [4] so that the following property holds in a network: at any time, the set of S-nodes form a consistent subnet. This property enables both easy identification of a consistent subnet and expansion of the consistent subset whenever a join process terminates. Nevertheless, utilization of this property is not limited to consistency-preserving optimization.

The extended join protocol was then integrated with our failure recovery protocol and a set of optimization heuristics. The integrated protocols were evaluated through simulation experiments. We showed that our protocols are effective and efficient in maintaining K -consistency and scalable to a large number of network nodes. We also showed that by primarily using information carried in join protocol messages, neighbor tables can be greatly optimized. For p2p networks that have higher demand for optimality of neighbor tables, algorithms presented in [1, 2, 16] can be further applied with extra costs. No matter which algorithm is applied, it should be applied within the constraint of the optimization rule to preserve consistency.

References

- [1] M. Castro, P. Druschel, Y. C. Hu, and A. Rowstron. Exploiting network proximity in peer-to-peer overlay networks. In *Proc. of International Workshop on Future Directions in Distributed Computing*, June 2002.
- [2] K. Hildrum, J. D. Kubiatowicz, S. Rao, and B. Y. Zhao. Distributed object location in a dynamic network. In *Proc. of ACM Symposium on Parallel Algorithms and Architectures*, August 2002.
- [3] D. R. Karger and M. Ruhl. Finding nearest neighbors in growth-restricted metrics. In *Proc. of ACM Symposium on Theory of Computing*, May 2002.
- [4] S. S. Lam and H. Liu. Silk: a resilient routing fabric for peer-to-peer networks. Technical Report TR-03-13, Dept. of CS, Univ. of Texas at Austin, May 2003.
- [5] S. S. Lam and H. Liu. Failure recovery for structured p2p networks: Protocol design and performance evaluation. In *Proc. of ACM SIGMETRICS*, June 2004.
- [6] H. Liu and S. S. Lam. Neighbor table construction and update in a dynamic peer-to-peer network. In *Proc. of IEEE International Conference on Distributed Computing Systems (ICDCS)*, May 2003.
- [7] H. Liu and S. S. Lam. Consistency-preserving neighbor table optimization for p2p networks. Technical Report TR-04-01, Dept. of CS, Univ. of Texas at Austin, January 2004.
- [8] Y. Liu, Z. Zhuang, L. Xiao, and L. M. Ni. A distributed approach to solving overlay mismatching problem. In *Proc. of International Conference on Distributed Computing Systems*, March 2004.
- [9] P. Maymounkov and D. Mazieres. Kademlia: A peer-to-peer information system based on the xor metric. In *Proc. of International Workshop on Peer-to-Peer Systems*, March 2002.
- [10] C. G. Plaxton, R. Rajaraman, and A. W. Richa. Accessing nearby copies of replicated objects in a distributed environment. In *Proc. of ACM Symposium on Parallel Algorithms and Architectures*, June 1997.
- [11] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and Scott Shenker. A scalable content-addressable network. In *Proc. of ACM SIGCOMM*, August 2001.
- [12] S. Ratnasamy, M. Handley, R. Karp, and S. Shenker. Topologically-aware overlay construction and server selection. In *Proc. of IEEE INFOCOM*, June 2002.
- [13] A. Rowstron and P. Druschel. Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems. In *Proc. of IFIP/ACM International Conference on Distributed Systems Platforms*, November 2001.
- [14] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup service for internet applications. In *Proc. of ACM SIGCOMM*, August 2001.
- [15] E. W. Zegura, K. Calvert, and S. Bhattacharjee. How to model an internetwork. In *Proc. of IEEE Infocom*, March 1996.
- [16] H. Zhang, A. Goel, and R. Govindan. Incrementally improving lookup latency in distributed hash table systems. In *Proc. of SIGMETRICS*, June 2003.
- [17] 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), January 2004.