

This page
is legacy
content.



Check out the current
u s e n i x
Web site.

OSDI '02 Paper [[OSDI '02 Tech Program In](#)

Pp. 299-314 of the *Proceedings* ☐

Secure routing for structured peer overlay networks

To appear in the [Fifth Symposium on Operating
Implementation \(OSDI 2002\)](#)

[Miguel Castro](#)¹, [Peter Druschel](#)², [Ayalvadi Ganesh](#)
[Dan S. Wallach](#)²

¹*Microsoft Research Ltd., 7 J J Thomson Avenue,
{[mcastro](#), [ajg](#), [antr](#)}@microsoft.com*

²*Rice University, 6100 Main Street, MS 132, Houston, TX 77005
{[druschel](#), [dwallach](#)}@cs.rice.edu*

Abstract:

Structured peer-to-peer overlay networks provide a suitable infrastructure for a wide range of large-scale, decentralized applications, including distributed storage, distributed communication, and content distribution. These overlays are resilient: they can route messages correctly even when a large fraction of the nodes in the network partitions. But current overlays are not secure: malicious nodes can prevent correct message delivery. This security problem is particularly serious in open peer-to-peer systems where autonomous parties without pre-existing trust relations interact over shared resources. This paper studies attacks aimed at preventing

in structured peer-to-peer overlays and presents defenses. We describe and evaluate techniques that allow nodes to join and leave the routing state, and to forward messages securely in the presence of malicious nodes.

1 Introduction

Structured peer-to-peer (p2p) overlays like CAN [16], Tapestry [21] provide a self-organizing substrate for large-scale applications. These systems provide a powerful platform for a variety of decentralized services, including network storage and application-level multicast. Structured overlays allow each node to object in a probabilistically bounded, small number of hops, requiring per-node routing tables with only a small number of entries. The systems are scalable, fault-tolerant and provide efficient routing.

However, to fully realize the potential of the p2p paradigm, the system must be able to support an open environment where nodes with conflicting interests are allowed to join. Even in a large scale, it may be unrealistic to assume that none of the nodes have been compromised by attackers. Thus, structured overlays are vulnerable to a variety of security attacks, including the case where a subset of nodes act maliciously. Such nodes may mis-route, corrupt or delete routing information. Additionally, they may attempt to corrupt other nodes and corrupt or delete objects they are supposed to store in the system.

In this paper, we consider security issues in structured overlays. We describe attacks that can be mounted against such overlays, the techniques they support, and present the design of secure techniques to defend against these attacks. In particular, we identify *secure routing* as a key security issue. Combined with existing, application-specific security techniques, we design secure, decentralized applications upon structured overlays. We focus on (1) a secure assignment of node identifiers, (2) secure routing and (3) secure message forwarding. We present techniques to defend against these problems, and show how using these techniques, secure routing can be achieved efficiently despite up to 25% of malicious participating nodes. We show that the overhead of secure routing is acceptable and practical in the presence of malicious nodes.

The rest of this paper is organized as follows. Section 2 describes structured p2p overlays, specifies models and assumptions. Section 3 describes secure routing. Sections 3, 4 and 5 present attacks on and solutions for node identifiers to nodes, routing table maintenance and message forwarding respectively. Section 6 explains how the overhead of secure routing is minimized by using self-certifying data. Finally, Section 7 discusses related work and Section 8 provides conclusions.

2 Background, models and

In this section, we present some background on structured p2p overlay networks like CAN, Chord, Tapestry and Pastry. Space limitations prevent a detailed overview of each protocol. Instead, we describe the general structured p2p overlay networks that we use to keep track of any particular protocol. For concreteness, we also give a brief overview and point out relevant differences with the other protocols. Finally, we define secure routing and outline our solution.

Throughout this paper, most of the analyses and technical details of our abstract model, and should apply to other structured p2p overlay networks, otherwise noted. However, the security and performance of our solution is fully evaluated only in the context of Pastry; a full evaluation of other protocols is future work.

2.1 Routing overlay model

We define an abstract model of a structured p2p routing overlay. We capture the key concepts common to overlays like CAN, Chord, Tapestry and Pastry.

In our model, participating nodes are assigned uniformly random identifiers from a large *id space* (e.g., the set of 128-bit unsigned integers). Specific objects are assigned unique identifiers, called *keys*, from the same *id space*. Each key is mapped by the overlay to a unique node, called the *root*. The protocol routes messages with a given key to the root.

To route messages efficiently, each node maintains a *neighbor set* of other nodes and their associated IP addresses. Moreover, each node maintains a *neighbor set*, consisting of some number of nodes with identifiers close to its own in the *id space*. Since node *id* assignment is random, the *neighbor set* represents a random sample of all participating nodes.

For fault tolerance, application objects are stored at multiple nodes in the overlay. A *replica function* maps an object's key to a set of nodes, the set of *replica roots* associated with the replica key. The *replica function* is a function of participating nodes in the overlay. Each replica root is a node in the overlay.

Next, we discuss existing structured p2p overlay protocols and how they fit into our abstract model.

`\psfig{file=table.eps,width=1.0\textwidth}`

Figure 1: Routing table of a Pastry node with node identifier v in base 16, \square represents an arbitrary node identifier

`\psfig{file=ring.eps,width=1.0\textwidth}`

Figure 2: Routing a message from node v to node w

2.2 Pastry

Pastry nodeIds are assigned randomly with uniform distribution in the 128-bit id space. Given a 128-bit key, Pastry routes a message to the live node whose nodeId is numerically closest to the key. Each node keeps track of its neighbor set and notifies application of changes.

Node state: For the purpose of routing, nodeIds and keys are represented as a sequence of digits in base b (b is a configuration parameter).

A node's routing table is organized into r rows and c columns. Entries in row i of the routing table contain the IP address of the node whose nodeId shares the first i digits with the present node's nodeId. The i th digit of the node in column j of row i equals d_{ij} . The value of d_{ij} corresponds to the value of the $(i \cdot c + j)$ th digit of the local node's nodeId, or is empty. A routing table entry is left empty if no node with the given prefix is known. Figure 1 depicts an example routing table.

Each node also maintains a neighbor set (called a "leaf set") of n nodes with nodeIds that are numerically closest to the current node's nodeId, with $n/2$ larger and $n/2$ smaller nodeIds than the current node's nodeId. n is a constant for all nodes in the overlay, with a typical value of 20, where n is the number of expected nodes in the overlay. The neighbor set ensures reliable message delivery and is used to store application objects.

Message routing: At each routing step, a node seeks for a node in the routing table whose nodeId shares with the key a prefix one digit (or $\log_2 b$ bits) longer than the prefix that the key shares with the node's id. If no such node can be found, the message is routed to the node whose nodeId shares a prefix with the key as long as the current node's id is closer to the key than the present node's id. If no appropriate node is found in the routing table or neighbor set, then the current node is the message's final destination.

Figure 2 shows the path of an example message. Analysis shows that the number of routing hops is slightly below $\log_b N$, with N being the number of nodes around the mean. Moreover, simulation shows that the system is robust to crash failures.

To achieve self-organization, Pastry nodes must dynamically maintain their state, i.e., the routing table and neighbor set, in the presence of node failures. A newly arriving node with the new nodeId

by asking any existing Pastry node u to route a special message. The message is routed to the existing node u with node u then obtains the neighbor set from u and constructs rows from the routing tables of the nodes it encounters u to u . Finally, u announces its presence to the initial nodes which in turn update their own neighbor sets and routing tables. This overlay can adapt to abrupt node failure by exchanging information (`neighborSet`) among a small number of nodes.

2.3 CAN, Chord, Tapestry

Next, we briefly describe CAN, Chord and Tapestry, and their differences relative to Pastry.

Tapestry is very similar to Pastry but differs in its approach to routing and in how it manages replication. In Tapestry, nodes and namespaces are not aware of each other. When a node's routing table has an entry for a node that matches a key's i th digit, the node routes the message to the node with the next higher value in the i th digit, modulo the size of the routing table. This procedure, called *surrogate routing*, maps keys to nodes as long as the node routing tables are consistent. Tapestry does not maintain a separate neighbor set, although one can think of the lowest populated entries in the routing table as a neighbor set. For fault tolerance, Tapestry maintains a set of random keys, yielding a set of replicated objects. The expected number of routing hops in Tapestry is $O(\log n)$.

Chord uses a 160-bit circular id space. Unlike Pastry, Chord only routes in clockwise direction in the circular id space. In Chord, instead of a routing table in Pastry, Chord nodes maintain a routing table consisting of pointers to other live nodes (called a "finger table"). The finger table of node u refers to the live node with the smallest id greater than u . The first entry points to u 's successor, and subsequent entries point to nodes at repeatedly doubling distances from u . Each node also maintains pointers to its predecessor and to its k successors in the successor list. The successor list represents the neighbor set in our model. A routing function maps an object's key to the nodeIds in the neighbor set, i.e., replicas are stored in the neighbor set of the key's predecessor. The expected number of routing hops in Chord is $O(\log n)$.

CAN routes messages in a d -dimensional space, where d is the dimension of the routing table with 2^d entries and any node can be reached in d hops on average. The entries in a node's routing table are chosen from a set of 2^d nodes.

dimensional space. CAN's neighbor table duals as both neighbor set in our model. Like Tapestry, CAN's replication keys for storing replicas at diverse locations. Unlike Pastry, CAN's routing table does not grow with the network size. The number of hops grows faster than $\log n$ in this case.

Tapestry and Pastry construct their overlay in a Internet-like manner to reduce routing delays and network utilization. In these systems, routing entries can be chosen arbitrarily from an entire segment of the network, increasing the expected number of routing hops. The paper [17] discusses initializing the routing table to refer to nodes that are close to the topology and have the appropriate nodeId prefix. This approach simplifies routing [17]. However, it also makes these systems vulnerable to attacks shown in Section 4.

The choice of entries in CAN's and Chord's routing table is based on proximity. The CAN routing table entries refer to specific neighbors in the d -dimension, while the Chord finger table entries refer to nodes in the key space. This makes proximity routing harder but it protects against attacks that exploit attacking nodes' proximity to their victims.

2.4 System model

The system runs on a set of n nodes that form an overlay network as described in the previous section. We assume a bound f on the number of nodes that may be faulty. Faults are modeled using the Byzantine failure model, i.e., faulty nodes can behave arbitrarily. All faulty nodes necessarily be operating as a single conspiracy. The system is partitioned into independent coalitions, which are disjoint sets of nodes C_1, \dots, C_k (where $k \geq 1$). When $f \leq |C_i|$, all faulty nodes may corrupt C_i and cause the most damage to the system. We model the coalition structure into multiple independent coalitions by setting $f \leq |C_i|$. Nodes in C_i work together to corrupt the overlay but are unaware of nodes in other coalitions. We studied the behavior of the system with f ranging from 1 to $n/2$ for different failure scenarios.

We assume that every node in the p2p overlay has a static IP address that can be contacted. In this paper, we ignore nodes with dynamic IP addresses, and nodes behind network address translation (NAT). Modern p2p overlays can be extended to address these concerns by using techniques similar to traditional network hosts.

The nodes communicate over normal Internet connections. We consider two types of communication: *network-level*, where nodes communicate without routing through the overlay, and *overlay-level*, where nodes communicate through the overlay using one of the protocols discussed in Section 2.1. Nodes may use cryptographic techniques to prevent adversaries from

network-level communication between correct nodes. control over network-level communication to and from can compromise overlay-level communication that is. Adversaries may delay messages between correct nodes. message sent by a correct node to a correct destination no faulty nodes is delivered within time $\frac{1}{\epsilon}$ with proba

2.5 Secure routing

Next, we define a secure routing primitive that can be techniques to construct secure applications on structured sections show how to implement the secure routing primitive network models that we described in the previous section.

The routing primitives implemented by current structured best-effort service to deliver a message to a replica root. With malicious overlay nodes, the message may be dropped or be delivered to a malicious node instead of a legitimate node. these primitives cannot be used to construct secure applications. inserting an object, an application cannot ensure that the legitimate, diverse replica roots as opposed to faulty nodes. roots. Even if applications use cryptographic methods, malicious nodes may still corrupt, delete, deny access to all replicas of an object.

To address this problem, we define a secure routing primitive *primitive ensures that when a non-faulty node sends a message reaches all non-faulty members in the set of nodes with high probability.* \mathcal{R} is defined as the set of nodes that the set of replica keys associated with \mathcal{R} , a live root node replica key. In Pastry, for instance, \mathcal{R} is simply a set of numerically closest to the key. Secure routing ensures eventually delivered, despite nodes that may corrupt, and (2) the message is delivered to all legitimate replica nodes that may attempt to impersonate a replica root.

Secure routing can be combined with existing security maintain state in a structured p2p overlay. For instance, stored on the replica roots, or a Byzantine-fault-tolerant BFT [4] can be used to maintain the replicated state. Since the replicas are initially placed on legitimate replica roots, message reaches a replica if one exists. Similarly, secure build other secure services, such as maintaining file metadata distributed storage utility. The details of such services paper.

Implementing the secure routing primitive requires the securely assigning nodeIds to nodes, securely maintaining the routing table, and securely forwarding messages. Secure nodeId assignment cannot choose the value of nodeIds assigned to the nodes. Without it, the attacker could arrange to control all replication and mediate all traffic to and from a victim node.

Secure routing table maintenance ensures that the fraction of faulty nodes that appear in the routing tables of correct nodes does not exceed a small fraction of faulty nodes in the entire overlay. Without this, correct message delivery, given only a relatively small number of replicas. Finally, secure message forwarding ensures that at least one message to a key reaches each correct replica root for the key value. Sections [3](#), [4](#) and [5](#) describe solutions to each of these problems.

3 Secure nodeId assignment

The performance and security of structured p2p overlays rely on the fundamental assumption that there is a uniform random distribution of nodeIds. This cannot be controlled by an attacker. This section discusses how the attacker violates this assumption, and how this problem can be solved.

3.1 Attacks

Attackers who can choose nodeIds can compromise the security of the overlay, without needing to control a particularly large number of nodes. For example, an attacker may partition a Pastry or Chord overlay into two complete and disjoint neighbor sets. Such attackers mediate all traffic to victim nodes by carefully choosing nodeIds. For example, an attacker can choose every entry in a victim's routing table and neighbor set in a Chord overlay. At that point, the victim's access to the object is completely mediated by the attacker.

Attackers who can choose nodeIds can also control access to objects. An attacker can choose the closest nodeIds to all replica keys for a given object, thus controlling all replica roots. As a result, the attacker can corrupt, or deny access to the object. Even when attacked nodes are not corrupt, they may still be able to mount all the attacks above (as long as there is a large number of legitimate nodeIds easily. This is known as a "sybil" attack.

Previous approaches to nodeId assignment have either chosen nodeIds randomly by the new node [[5](#)] or compute nodeIds by hashing the node's IP address [[20](#)]. Neither approach is secure because an attacker can choose nodeIds that are not necessarily random, or can choose nodeIds that hash to a desired interval in the nodeId space. Particularly, even modest attackers will have more potential IP addresses than

there are likely to be nodes in a given p2p network.

3.2 Solution: certified nodeIds

One solution to securing the assignment of nodeIds is a central, trusted authority. We use a set of trusted certificates to assign nodeIds to principals and to sign *nodeId certificates*. A nodeId is bound to the public key that speaks for its principal and we ensure that nodeIds are chosen randomly from the id space to prevent forging nodeIds. Furthermore, these certificates give the network infrastructure, suitable for establishing encrypted and authenticated communication between nodes.

Like conventional CAs, ours can be offline to reduce the need for signing keys. They are not involved in the regular operation of the network. Nodes with valid nodeId certificates can join the overlay, route, and leave repeatedly without involvement of the CAs. As with a conventional CA, the CA's public keys must be well known, and can be installed in the software itself, as is done with current Web browsers.

The inclusion of an IP address in the certificate deserves further discussion. p2p protocols, such as Tapestry and Pastry, measure the distance to nodes and choose routing table entries that minimize this distance. If multiple legitimate nodeId certificates could freely swap their IP addresses, controls, it might be able to increase the fraction of an attacker's node in a node's routing table. By binding the nodeId to an IP address, we prevent an attacker to move nodeIds across nodes. We allow multiple nodeIds per IP address because the IP addresses of nodes may change over time. Otherwise, attackers could deny service by hijacking valid nodeIds.

A downside of binding nodeIds to IP addresses is that if a node's IP address changes, either as a result of dynamic address assignment or organizational network changes, then the node's old nodeId becomes invalid. In p2p systems where IP addresses are allowed to change, nodeId swapping attacks may be unavoidable.

Certified nodeIds work well when nodes have fixed network addresses, as in Chord, Pastry, and Tapestry. However, it might be harder to implement in assignment in CAN. CAN nodeIds represent a zone in the identifier space. It is split in half when a new node joins [16]. Both the network and the nodeId of the joining node change during this process.

3.2.1 Sybil attacks

While nodeId assignment by a CA ensures that nodeIds are unique, it is also important to prevent an attacker from easily obtaining multiple nodeId certificates. One solution is to require an attacker to

certificates, via credit card or any other suitable mechanism. The cost of an attack grows naturally with the size of the network. nodeIds certificates cost \$20, controlling 10% of an overlay network costs \$2,000 and the cost rises to \$2,000,000 with 1,000,000 nodes. An attack is even higher; it costs an expected \$20,000 to find a particular point in the id space in an overlay with 1,000,000 nodes. If attacks economically expensive, these fees can also fund the network.

Another solution is to bind nodeIds to real-world identities and money. In practice, different forms of CAs are suitable. An identity-based CA is the preferred solution in a virtual world organization that already maintains employment or membership identity checks. In an open Internet deployment, a monetary CA is suitable because it avoids the complexities of authentication.

None of the known solutions to nodeId assignment are suitable for a network that is very small. For small overlay networks, where all members of the network are trusted not to cheat. Only in a critical mass, where it becomes sufficiently hard for an attacker to acquire resources to control a significant fraction of the overlay network, is it allowed to join.

3.3 Rejected: distributed nodeId generation

The CAs represent points of failure, vulnerable to both central and distributed attacks. Also, for some p2p networks, it may be cumbersome to require nodes to pay money or prove their real-world identities. Therefore, we reject the idea to construct secure p2p overlays without requiring centralized identity checks. Unfortunately, fully decentralized nodeId generation has fundamental security limitations [10]. None of the known techniques can ultimately prevent a determined attacker from acquiring nodeIds.

However, several techniques may be able to, at a minimum, delay an attacker which an attacker can acquire nodeIds. One possible solution is to require prospective nodes to solve crypto puzzles [15] to gain access to the network. This approach has been taken to address a number of distributed denial of service attacks [13,8]. Unfortunately, the cost of solving a crypto puzzle is proportional to the slowest legitimate node, yet the puzzle must be slow enough to slow down an attacker with access to many fast machines. The effectiveness of any such technique.

For completeness, we briefly describe here one relatively simple technique to generate certified nodeIds in a completely distributed manner. The idea is to require new nodes to generate a key pair such that the SHA-1 hash of the public key has the first k bits zero. The number of operations required to generate such a key pair is $O(2^k)$. This technique allows cryptography allow the nodes to use a secure hash of the public key. This hash should be computed using SHA-1 with a digest length of k bits.

MD5 to avoid reducing the number of random bits in $ipadd$ that they performed the required amount of work to use a nodeId with information that would allow others to reuse their work. This is necessary to achieve the desired level of security.

It is also possible to bind IP addresses with nodeIds to exploit network locality. The idea is to require nodes to find a string $ipadd$ such that $SHA-1(SHA-1(ipadd) \parallel nodeId)$ is close to zero. Nodes would be required to present such an $ipadd$ to be accepted by others.

Finally, it is possible to periodically invalidate nodeIds. A node entity broadcast to the overlay a message supplying a new salt for the hash computations. This makes it harder for an attacker to reuse nodeIds over time and to reuse nodeIds computed for one overlay. However, it requires legitimate nodes to periodically update and communication to maintain their membership in the overlay.

4 Secure routing table maintenance

We now turn our attention to the problem of secure routing table maintenance mechanisms are used to create and maintain neighbor sets for joining nodes, and to maintain them. A routing table and neighbor set should have an average of 10^6 entries that point to nodes controlled by the attacker (controlled nodes). Attackers can increase the fraction of bad entries by supplying bad entries which reduces the probability of routing successfully.

Preventing attackers from choosing nodeIds is necessary, but it is not sufficient as illustrated by the two attacks discussed below. Solutions to this problem.

4.1 Attacks

The first attack is aimed at routing algorithms that use nodeIds to improve routing efficiency: attackers modify the fraction of bad routing table entries. For example, a nodeId assumed allows an attacker to control communication to the node it controls. When a correct node $nodeId$ sends a probe to establish a connection with a certain nodeId, an attacker can intercept the probe and reply closest to $nodeId$ reply to it. If the attacker controls enough nodes on the Internet, it can make nodes that it controls appear close to the target.

the probability that they are used for routing. The attacker (with a maximal fraction of colluding nodes) is small even if

This attack can be ruled out by a more restrictive constraint: nodeId certificates bind IP addresses to nodeIds (see Section 4.1). Faulty nodes can only observe messages that are sent to them, so this attack is prevented. But note that a rogue ISP or a set of offices around the world could easily perform this attack on all routers appropriately. The attack is also possible if the attacker, in the indirection that the attacker can control, e.g., mobile IP.

The second attack does not manipulate proximity information, but the fact that it is hard to determine whether routing updates from overlay protocols like Tapestry and Pastry. Nodes receive updates when they join the overlay and when other nodes join, and they update from other nodes in their routing table periodically to account for delays. In these systems, attackers can more easily suppress updates from always point to faulty nodes. This simple attack causes routing table entries to increase toward one as the bad routing updates. More precisely, routing updates from correct nodes point to a node with probability at least $\frac{1}{2}$ whereas this probability can be as low as $\frac{1}{2}$ for updates from faulty nodes. Correct nodes receive updates from a node with probability at most $\frac{1}{2}$ and from faulty nodes with probability at least $\frac{1}{2}$. Therefore, the probability that a routing table entry is updated by a correct node is at least $\frac{1}{2}$, which is greater than $\frac{1}{2}$. This causes the fraction of faulty entries to decrease in subsequent update, causing the fraction of faulty entries to decrease.

Systems without strong constraints on the set of nodeIds that can fill each table slot are more vulnerable to this attack. Pastry and Tapestry have constraints at the top levels of routing tables. This flexibility allows to determine if routing updates are unbiased but it allows an attacker to exploit network proximity to improve routing performance. To avoid this, impose strong constraints on nodeIds in routing tables. For example, the closest nodeIds to some point in the id space. This makes it hard to exploit proximity to improve performance but it is good for security. If an attacker chooses the nodeIds they control, the probability that a routing table entry is closest to a point in the id space is $\frac{1}{2}$.

4.2 Solution: constrained routing

To enable secure routing table maintenance, it is important to have strong constraints on the set of nodeIds that can fill each slot. For example, the entry in each slot can be constrained to be the closest node to a point in the id space as in Chord. This constraint can be enforced independent of network proximity information, which

attackers.

The solution that we propose uses two routing tables: proximity information for efficient routing (as in Pastry) constrains routing table entries (as in Chord). In normal operation, the proximity table is used to forward messages to achieve good performance, and the constrained routing table is used only when the efficient routing technique fails. We describe how to detect when routing fails.

We modified Pastry to use this solution. We use the normal Pastry routing table and an additional *constrained* Pastry routing table. The constrained routing table of a node with identifier id , the slots of which contain any nodeId that shares the first k digits with id . The entry in the i th slot of the constrained routing table, the entry in the i th slot of the normal routing table points to the closest nodeId to a point $id + 2^i$ in the domain. If a nodeId n shares the first k digits with id , it has the value 2^i in the i th slot of the constrained routing table if n has the same remaining digits as id .

Pastry's message forwarding works with the constrained routing table. The same would be true with Tapestry. The initialization and maintenance of the routing table were modified to account for the constrained routing table.

All overlay routing algorithms rely on a *bootstrap node* to initialize the state of a newly joining node. The bootstrap node is responsible for sending a message using the nodeId of the joining node as the key. If the bootstrap node is faulty, it can completely corrupt the view of the overlay network of the joining node. Therefore, it is necessary to use a set of diverse bootstrap nodes to ensure that with very high probability, at least one of the bootstrap nodes' nodeId certificates makes the task of choosing such a bootstrap node. A node cannot forge nodeIds.

A newly joining node, n , picks a set of bootstrap nodes. It uses each bootstrap node to route using its nodeId as the key. Then, non-faulty bootstrap nodes use the forwarding techniques (described in Section 5.2) to obtain the closest live nodeId to the joining node. Node n collects the proposed neighbor sets from each bootstrap node, and picks the "closest" live nodeIds from each neighbor set (where the definition of closest is protocol-specific).

The locality-aware routing table is initialized as before. The difference is that the entry in the i th slot of the constrained routing table is initialized by sending a message along the route to the nodeId. The difference is that the entry in the i th slot of the constrained routing table picks the entry with minimal network delay from the entries in the i th slot of the normal routing table for each routing table slot.

Each entry in the constrained routing table can be initialized by sending a message along the route to the nodeId. This is similar to what is done in Chord. The problem is that the entry in the i th slot of the constrained routing table picks the entry with minimal network delay from the entries in the i th slot of the normal routing table for each routing table slot.

with `maxOutgoingEdges` (recall that `maxOutgoingEdges` controls the number of columns in `RoutingTable` and `Pastry`). To reduce the overhead, we can use `maxOutgoingEdges` to limit the number of entries in the constrained routing table. That is, by induction, the constrained routing tables of the nodes point to entries that are close to the desired point `nodeId`. The nodes use the constrained routing tables from the nodes in its neighbor set and update the constrained routing table. From the set of candidates to `nodeId`, it picks the `nodeId` that is closest to the desired point `nodeId`. At the end of this process, `nodeId` informs the nodes in its neighbor set.

We exploit the symmetry in the constrained routing table to update their routing tables to reflect `nodeId`'s arrival: `nodeId` chooses a set of candidates for each entry to determine which candidate is closest to `nodeId`. It updates the constrained routing table entries to point to `nodeId`. It informs those candidates.

To ensure neighbor set stabilization in the absence of `nodeId`, `nodeId` informs the members of its neighbor set whenever it changes its neighbor set. It retransmits this information until its receipt is acknowledged.

5 Secure message forwarding

The use of certified `nodeIds` and secure routing table maintenance in the constrained routing table (and neighbor set) has an advantage: it prevents random entries that point to nodes controlled by the attacker. The constrained routing table is not sufficient because the attacker can still prevent the probability of successful delivery by simply not forwarding the message. The algorithm. The attack is effective even when `nodeId` is secure. This section describes an efficient solution to this problem.

5.1 Attacks

All structured p2p overlays provide a primitive to send a message. In the absence of faults, the message is delivered to the root node. The average of `maxOutgoingEdges` routing hops. But routing may fail if any node on the route between the sender and the root are faulty; faulty nodes may drop the message, route the message to the wrong place, or prevent the message from being delivered. Therefore, the probability of routing successfully between a fraction `fraction` of the nodes is faulty is only: $(1 - \text{fraction})^{\text{maxOutgoingEdges}}$, where `fraction` is the fraction of the nodes that are faulty.

The root node for a key may itself be faulty. As discussed in the previous section, we tolerate root faults by replicating the information associated with the key to multiple nodes -- the *replica roots*. Therefore, the probability of

correct replica root is only: . The value of α is in CAN, in Chord, and in Pastry.

We ran simulations of Pastry to validate this model. The probability of success is slightly lower than the probability in the analytical simulations (because the number of Pastry hops is slightly higher than the average [3]) but the error was below 2%.

Figure 3 plots the probability of routing to a correct replica (using the model) for different values of α , β , and γ . The probability is quite fast when α or β increase. Even with only 10% of the nodes being correct, the probability of successful routing is only 65% when $\alpha = 0.1$ and $\beta = 0.1$ in Pastry overlay.

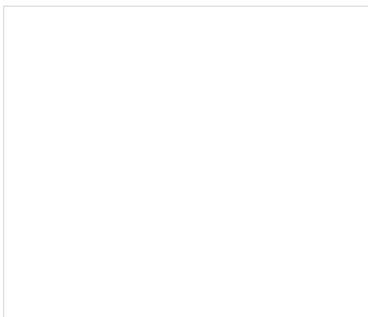


Figure 3: Probability of routing to a correct replica.

In CAN, Pastry, and Tapestry, applications can reduce the number of hops by increasing the value of α or β . Fewer hops increase the probability of routing correctly. For example, the probability of successful routing to a correct replica in 100,000 nodes is 65% in Pastry when $\alpha = 0.1$ and 75% when $\alpha = 0.2$. However, increasing α also increases the cost of routing table maintenance; a successful routing to a correct replica requires an impractically large value of α . Chord, with $\alpha = 0.1$, which results in a low probability of success, e.g., 42% under the same conditions.

5.2 Solution: detect faults, use diversity

The results in Figure 3 show that it is important to design a secure routing primitive. We want a *secure routing primitive* that takes a key and ensures that with very high probability at least one message reaches each correct replica root for the key. The question is how to do this efficiently.

Our approach is to route a message efficiently and to determine if routing worked. We only use more expensive routing if the failure test returns positive. In more detail, our second approach routes a message efficiently to the root of the destination key using the routing table. Then, it collects the prospective set of replica roots, chooses a prospective root node and applies the *failure test* to that node. If the prospective replica roots are accepted as the correct replica roots, message copies are sent over diverse routes toward the destination key so that with high probability each correct replica root is reached. We describe how to implement the failure test. Then we discuss why we rejected an alternate approach called iterative routing.

5.2.1 Routing failure test

The failure test takes a key and a set of prospective replica roots. It returns negative if the set of roots is likely to be correct and positive otherwise. Of course, routing can fail without the use of prospective replica roots. The sender detects this by sending a message. If it does not receive a response before a timeout, the failure test returns positive triggering the use of redundant routing.

Detecting routing failures is difficult because a coalition of nodes can pretend to be the legitimate replica roots for a given key. Since the correct replica roots are determined by the structure of the overlay, a node who is not a replica root must rely on overlay routing to determine the correct replica roots. If a message is routed by a faulty node, the adversary can choose a different replica root set that contain only nodes it controls. Furthermore, it is possible that the adversary just happens to legitimately control the correct replica roots. This problem is common to all structured p2p overlays.

The routing failure test is based on the observation that the density of nodeIds per unit of "volume" in the id space is greater for faulty nodes than for live nodes. The test works by comparing the density of nodeIds in the neighbor set of the sender with the density of nodeIds close to the destination key. We describe the test in detail only in the context of Pastry to simplify the presentation; the generalization to other overlays is straightforward. Overlays that distribute replica keys for a key uniformly can use this check by comparing the density at the sender with the density at the destination key between each replica key and its root's nodeId.

In Pastry, the set of replica roots for a key is a subset of the neighbor set of the root node, called the key's *root neighbor set*. Each correct replica root is at an average numerical distance, $\frac{1}{2}$, between consecutive nodes in the neighbor set. The neighbor set of $\frac{1}{2}$ contains $\frac{1}{2}$ live nodes: $\frac{1}{2}$, the nodes with nodeIds less than $\frac{1}{2}$'s, and the $\frac{1}{2}$ nodes with the closest nodeIds to $\frac{1}{2}$. To test a prospective root neighbor set, $\frac{1}{2}$, we compare the density of nodeIds in the neighbor set of $\frac{1}{2}$ with the density of nodeIds in the neighbor set of $\frac{1}{2}$.

1. all nodeIds in \mathcal{N}_i have a valid nodeId certificate
is the middle one, and the nodeIds satisfy the de
2. the average numerical distance, \bar{d} , between co
satisfies: $\bar{d} \leq \frac{1}{2} \log \frac{1}{\epsilon}$

If \mathcal{N}_i satisfies both conditions, the test returns negative
positive. The test can be inaccurate in one of two way
when the prospective root neighbor set is correct, or it
when the prospective set is incorrect. We call the prob
and the probability of false negatives ϵ . The paramete
between ϵ and δ . Intuitively, increasing δ decreases ϵ

Assuming that there are N live nodes with nodeIds un
id space (which has length 2^L), the distances bet
approximately independent exponential random variab
 $\frac{1}{N}$. The same holds for the distances between consecut
that can collude together but the mean is $\frac{1}{N}$. It
and δ are independent of ϵ . They only depend on the
fraction of colluding nodes because faulty nodes only
nodes that they collude with.

Under these assumptions, we have derived the followi
and δ (see detailed derivation in the Appendix):

$$\epsilon = \frac{1}{2} \log \frac{1}{\delta} \quad \text{and} \quad \delta = \frac{1}{2} \log \frac{1}{\epsilon}$$

$$\epsilon = \frac{1}{2} \log \frac{1}{\delta} \quad \text{and} \quad \delta = \frac{1}{2} \log \frac{1}{\epsilon}$$

These expressions can be used to compute ϵ and δ nu
the following closed-form upper bounds for ϵ and δ :

$$\epsilon \leq \frac{1}{2} \log \frac{1}{\delta} \quad \text{and} \quad \delta \leq \frac{1}{2} \log \frac{1}{\epsilon}$$

where n is the number of distance samples used to compute \bar{d} , and n' . The

The analysis shows that \bar{d} and \bar{d}' are independent of n . The test's accuracy can be improved by increasing the number of distance samples used to compute the means. It is easy to increase the number of samples used to compute \bar{d} by augmenting the mechanism that is already used to compute neighbor sets. This mechanism propagates nodeIds through the network from a neighbor set to the other members of the set; it propagates nodeIds further but we omit the details due to lack of space. The number of samples used to compute \bar{d} because of this mechanism is described below. Therefore, we keep $n = 100$.

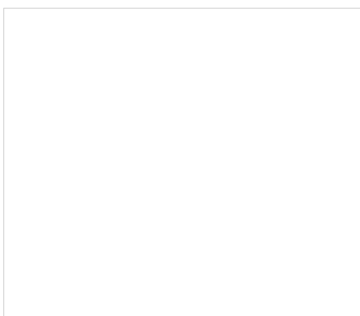


Figure 4: Routing failure test: probability of routing failure (P_r) and negatives (P_n). The predicted curves are indistinguishable from the simulation means. The upper bounds are not tight.

We ran several simulations to evaluate the effectiveness of the test. The simulations ran in a system with 100,000 random nodes. The values of \bar{d} and \bar{d}' for different values of n with $n' = 100$ at the sender is shown in Figure 4. The number of root neighbors is shown in Figure 5. The predicted values computed numerically, the upper bounds, and the simulation means are shown in Figure 6. The predicted curves match the simulation means almost exactly but the upper bounds are not very tight. The upper bounds obtained when $n = 100$, which is equal to $n' = 100$ with $n' = 100$.

5.2.1.1 Attacks:

There are several attacks that could invalidate the analysis failure test. First, the attacker can collect nodeIds certified by the overlay, and use them to increase the density of a prospective root neighbor set. Second, the attacker can include both nodeIds of nodes and correct nodes in a prospective root neighbor set. Both attacks decrease the probability that messages reach all correct replica roots. To counter in overlays that distribute replica keys over multiple roots have no detailed knowledge about the nodeIds of nodes.

These attacks can be avoided by having the sender communicate with root neighbors to determine if they are live and if they have not been omitted from the prospective root neighbor set. To improve the test, the prospective root returns to the sender a message with a list with the secure hashes of the neighbor sets reported by the prospective root neighbors, and the set of nodeIds (not in the prospective root neighbor set) are used to compute the hashes in this list. The sender can then verify that the hashes are consistent with the identifiers of the prospective root neighbors. For each prospective root neighbor the corresponding neighbor set is

In the absence of faults, the root neighbors will confirm the test. The sender can perform the density comparison immediately. For this happens with probability $\frac{1}{2^k}$, where k is the number of root neighbors in the prospective root neighbor set, the routing failure test must be run after a round of communication before the density check can be run. We consider this as a strategy to deal with this case. Currently, we consider the case where the prospective root neighbors don't agree and use redundant communication to make it worthwhile investing some additional communication for the routing.

In addition to these attacks, there is a *nodeId suppression* attack that is unavoidable and significantly decreases the accuracy of the test. The attacker suppress nodeIds close to the sender by leaving the overlay.

Similarly, the attacker can suppress nodeIds in the root neighbor set. This increases the probability of error. Furthermore, the attacker can alternate between honest and dishonest nodes have no way of detecting in which mode the sender is operating.

We ran simulations to compute the minimum error probability for the test without nodeId suppression attacks for different values of k . The error probability increases fast with k and it is higher than $\frac{1}{2^k}$ for $k > 1$. The number of samples at the sender. The nodeId suppression attack increases the probability of error for large percentages of compromised nodes. The probability of error is higher than 0.001 for $k > 1$ even if the sender is honest. Figures 5 and 6 show the results without and with nodeId suppression attacks, respectively.



Figure 5: Routing failure test: minimum error rate without nodeId suppression attacks and varying number of samples.

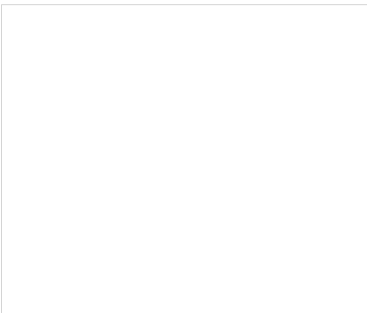


Figure 6: Routing failure test: minimum error rate with nodeId suppression attacks and varying number of samples.

These results indicate that our routing failure test is not perfect. However, fortunately we can trade off an increase in the number of samples to achieve a lower error rate by using redundant routing to disambiguate false positives. We can also achieve the minimum error rate that can be achieved for a target number of samples, and different numbers of samples at the sender. This is shown in Figure 7 with nodeId suppression attacks.

Figure 7: Routing failure test: probability for a false negative rate of 0.001 with no attacks and varying number of s

The results show that the test is not meaningful for this nodeId suppression attacks. However, setting `nodeIdSuppression` sender enables the routing failure test to achieve the target value of 0.001 and with `nodeIdSuppression` without nodeId suppression attacks the value of 0.001 is achieved routing is required 12% of the time.

5.2.2 Redundant routing

The redundant routing technique is invoked when the routing failure test is positive. The idea is simply to route copies of the message over multiple diverse routes to each replica key. If enough copies of the message are sent over diverse routes to each replica key, all correct replica roots will receive a copy of the message with high probability.

The issue is how to ensure that routes are diverse. One approach is to select multiple members of the sender's neighbor set to forward the message to each replica keys. This technique is sufficient in overlays that route messages uniformly over the id space (e.g., CAN and Tapestry). However, in overlays that choose replica roots in the neighbor set of the sender (e.g., Pastry) because the routes all converge on the key's root. For these overlays, we developed a technique called *redundant routing*. It sends multiple copies of the message toward the destination key until it reaches the key's root in its neighbor set. Then it uses the detailed routing information it has about the portion of the id space around the destination key to ensure that correct replica roots receive a copy of the message.

To simplify the presentation, we only describe in detail how this technique works in Pastry. If a correct node `node` sends a message to a destination key and the routing failure test is positive, it does the following:

(1) α sends β messages to the destination key γ . Each message is sent to a different member of α 's neighbor set; this causes the messages to take different routes. All messages are forwarded using the constrained shortest path and include a nonce.

(2) Any correct node that receives one of the messages from α and is in α 's neighbor set returns its nodeId certificate and the nonce to α .

(3) α collects in a set S the β nodeId certificates that are closest to α on the left, and the γ closest to α on the right. Only the nonces are added to S and they are first marked *pending*.

(4) After a timeout or after all β replies are received, α sends a message in S to each node marked *pending* in S and marks them as received.

(5) Any correct node that receives this list forwards it to all nodes in its neighbor set that are not in the list, or it sends a confirmation if there are no such nodes. This may cause steps 2 to 4 to be repeated.

(6) Once α has received a confirmation from each of its neighbors, it has executed three times, it computes the set of replica roots.

If the timeout is sufficiently large and correct nodes have at least one neighbor in each half of their neighbor set¹, the probability of reaching the root of γ is approximately equal to the probability that at least one message is forwarded over a route with no faults to a correct node or root in its neighbor set. Assuming independent routes,

$$P_{\text{success}} = 1 - \prod_{i=1}^n (1 - p_i)$$

where binom is the binomial distribution [6] with 0 success and n trials, and p is the probability of routing successfully in each trial is p . The probability of success for this technique depends on n and is independent of the number of nodes in the network.

We also ran simulations to determine the probability of reaching the root of γ with our redundant routing technique. Figure 8 plots the probability measured in the simulator for 100,000 trials against the analytic model. The analytic model matches the results well for $n \geq 10$. The results show that the probability of reaching the root of γ for $n = 10$ is approximately 0.9. Therefore, this technique combined with

achieve a reliability of approximately 0.999 for

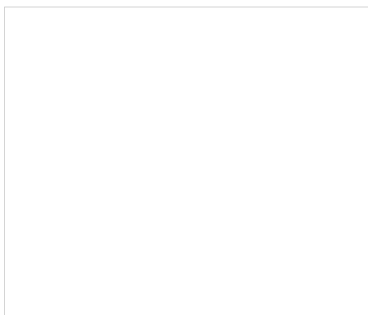


Figure 8: Model and simulation results for reaching all correct replica roots using redundant routing with a neighbor set anycast.

We studied several versions of redundant routing that all achieved a high level of success but perform differently. For example, the simplest version ensures that the nodeId certificates in the routing table belong to live nodes by signing nonces by periodically signing their clock readings. With synchronized clocks, and no signatures are necessary for verifying source addresses. We are still exploring the design space for ways it is possible to improve performance significantly by sending multiple probes at a time and using a version of the routing failure test that the current approach would also work well when reading self-certificates.

5.2.3 Putting it all together: performance

The performance of Pastry's secure routing primitive depends on the routing failure test, the cost of redundant routing, and the probability that redundant routing can be avoided. This section presents the cost of redundant routing and probability. For simplicity, we assume that all failures are detected, the number of probes used by redundant routing is equal to the number of samples at the source for routing (see Section 5.2.1), the number of nodes in the overlay is N .

The cost of the routing failure test when it returns negative is τ delay and α messages. The total number of bytes sent is

$$N \cdot \text{IdSize} \cdot \alpha \cdot \text{HashSize} + N \cdot \text{IdCertSize}$$

Using PSS-R [1] for signing nodeId certificates with modulus for the node public keys, the nodeId certificate the extra bandwidth consumed by the routing failure test with and 2.9 KB with (plus the space used). When the test returns positive, it adds the same number the extra delay is the timeout period.

The cost of redundant routing depends on the value of α when all of the root's neighbor set is added to in the redundant routing adds extra message delay messages. The total number of bytes in these messages

$$\alpha \cdot \text{IdSize} \cdot \text{IdCertSize} \cdot \text{SigSize} \cdot \text{IdSize}$$

Using PSS-R for signing nonces, the signed nonce size bandwidth consumed in this case is 22 KB with the space used up by message headers).

Under attack redundant routing adds a delay of at most the expected number of extra messages is less than where is the expected number of neighbor set that is added to in the first iteration. The number of messages is less than 451 with and and . The total number of bytes sent under attack value except that the sender sends an additional lists and the number of messages increases. This is an and and 1 KB with and (plus headers).

The probability of avoiding redundant routing is given by $P_{\text{avoid}} = (1 - \alpha)^{\alpha \cdot \text{IdSize} \cdot \text{IdCertSize} \cdot \text{SigSize} \cdot \text{IdSize}}$ is the probability that the overlay routes the message to the probability that there are no faulty nodes in the neighbor set is the false positive rate of the routing failure test. We assume that routing tables have an average of random assumption holds for the locality-aware routing table discussed in Section 4 and it holds for the constrained attacks. We do not have a good model of the effect of

aware routing table but we believe that they are very hard to find ().

The parameters α and β , should be set based on the deployment scenario. α can be expressed as the probability that all correct nodes in the overlay have the message. The overlay size and the assignment of values to α and β implicitly define a bound on α . If this bound is exceeded, redundant routing is invoked. For example, we saw that with $\alpha = 0.9$ and $\beta = 0.9$, redundant routing is invoked 12% of the time for this value of α .

One can trade off security for improved performance by decreasing α to reduce the cost of the routing function and to increase β . For example, consider the following two scenarios: (1) with $\alpha = 0.9$ and $\beta = 0.9$, and (2) with $\alpha = 0.8$ and $\beta = 0.9$. Figure 9 plots the probability of avoiding redundant routing for these two scenarios for different values of α . Without redundant routing, the probability of avoiding redundant routing is only 0.5% of the time in scenario (1) and 0.4% in scenario (2). When the fraction of faulty nodes is small, the routing function significantly improves performance by avoiding the cost of redundant routing.

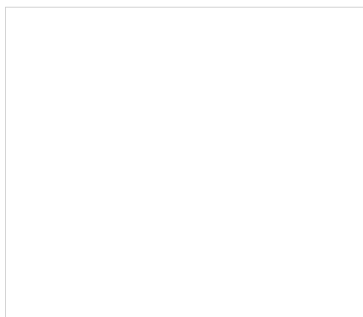


Figure 9: Probability of avoiding redundant routing for two scenarios: (1) with $\alpha = 0.9$ and $\beta = 0.9$ and (2) with $\alpha = 0.8$ and $\beta = 0.9$.

5.2.4 Rejected: checked iterative routing

An alternative to redundant routing is iterative routing [19]: the sender starts by looking up the next hop in the routing table, setting a variable current_node to point to this node; then, the sender sends the message to current_node and updates current_node to point to the returned value. The process repeats until the sender reaches the root of the destination key.

Iterative routing doubles the cost relative to the more efficient recursive routing, but it may increase the probability of routing successfully. We let the sender pick an alternative next hop when it fails to find a next hop. This is not a strong defense against an attacker who can corrupt the next hop. However, iterative routing can be augmented with a hop test, whether the next hop in a route is correct.

Hop tests are effective in systems like Chord or Pastry where the routing table because each routing table entry should contain information about a specific point in the id space. One can use a mechanism for density checking that we used for the routing failure test. We compare the average distance between consecutive nodeIds close to the nodeId in the routing table to the distance between the nodeId in the routing table and the next hop.

We ran simulations to compute the false positive and false negative rates for our approach with different values of k (these rates are independent of k). The minimum error for this hop test ($k=1$) is equal to 0.3 for $k=1$ and 256 samples to compute the mean at the sender.

The error is high because there is a single sample at the sender. Our simulations indicate that iterative lookups using Pastry with this hop check improve the probability of routing successfully. For example, the probability of routing successfully with $k=1$ and 256 samples to compute the mean at the sender is 0.3 to 0.4 . But it adds an extra 2.7 hops to each route on average for false positives.

We tried to increase the number of samples by having the sender store the routing table row during each iterative routing step with the required slot. Unfortunately, this performs worse than recursive routing because the attacker can combine good and bad routing information to increase the high average density.

We also tried to combine checked iterative routing with the recursive routing technique that we described before. We used checked iterative routing with neighbor set anycast messages in the hope that the improved routing would result in an improvement over recursive routes. But there was no visible improvement. Iterative routes are less independent than the recursive routes. Routing failure test combined with redundant routing is not a good idea for implementing secure routing.

6 Self-certifying data

The secure routing primitive adds significant overhead to the routing process.

In this section, we describe how the use of secure routing using *self-certifying data*.

The reliance on secure routing can be reduced by storing an overlay, i.e., data whose integrity can be verified by the client, to use efficient routing to request a copy of an object. After receiving the object, it can check its integrity and resort to secure routing if the integrity check fails or there was no response within a timeout.

Self-certifying data does not help when inserting new objects, when verifying that an object is not stored in the overlay, or the secure routing primitive to ensure that all correct replicas have the object. Similarly, node joining requires secure routing. Nevertheless, this can eliminate the overhead of secure routing in common cases.

Self-certifying data has been used in several systems. PAST uses a cryptographic hash of a file's contents as the key during insertion of the file, and PAST [18] inserts signed files into the overlay.

The technique can be extended to support mutable objects and group guarantees. One can use a system like PAST to store *group descriptors* that identify the set of hosts responsible for a group of objects. Group descriptors can be used as follows. At object creation, the object uses secure routing to insert a group descriptor that identifies the object. The descriptor contains the public keys of the object's replica holders and it is signed by the object creator.

The replica group can run a Byzantine-fault-tolerant replication primitive (BFT) [4] and the initial group membership is the set of hosts in the key. In this setting, read and write operations can be performed. A client uses efficient routing to retrieve a group descriptor for the object, checks the descriptor's signature; if correct, it uses the descriptor to authenticate the replica holders and to invoke a replication primitive. If the client fails to retrieve a valid descriptor or if it fails to authenticate the replica holders, it uses the secure routing primitive to obtain a correct group descriptor or that the object does not exist. This procedure provides linearizability guarantees (linearizability [11]) for reads and writes with a failure test in the common case.

Changing the membership of the group that is responsible for a group of objects is not trivial; it requires securely inserting a new group descriptor and ensuring that clients can reliably detect stale group descriptors. This technique allows groups to change membership while maintaining linearizability guarantees. Each group of hosts that stores replicas of an object has a private/public key pair associated with the group. When the membership changes, each host in the new membership generates a new key pair. The hosts in the old membership use their old keys to sign the new group descriptor containing the new keys, and then delete the old keys.

If this operation is performed by a quorum of replica holders, the number of faulty group members is exceeded [4],

will not be able to collude to pretend they are the current owner and thus form the quorum necessary to authenticate themselves.

Group descriptors can be authenticated by following a chain of descriptors with an owner signature and has signatures of a quorum of nodes from subsequent membership change. The chain can be shown to the owner from the owner or, alternatively, replicas can use proofs of ownership to avoid the need for chaining signatures.

7 Related work

Sit and Morris [19] present a framework for performing secure lookups in networks. Their adversarial model allows for nodes to store arbitrary contents, but assumes that nodes cannot interfere with each other. They then present a taxonomy of possible attacks. At the root of the taxonomy are lookup, routing table maintenance, and network partitioning. They then discuss security risks. They also discuss issues in higher-level storage, where nodes may not necessarily maintain the same data. Finally, they discuss various classes of attacks, including rapidly joining and leaving the network, and nodes to send bulk volumes of data to overload a victim (e.g., distributed denial of service attacks).

Dingledine *et al.* [9] and Douceur [10] discuss addressing the problem of a large number of potentially malicious nodes in the system. In a system with a central authority to certify node identities, it becomes difficult to verify whether you can trust the claimed identity of someone you are communicating with before communicating. Dingledine proposes to address this by using a system including the use of micro-cash, that allow nodes to be anonymous.

Bellovin [2] identifies a number of issues with Napster and Gnutella, and how difficult it might be to limit Napster and Gnutella. He argues that they can leak information that users might consider private. He also discusses queries they issue to the network. Bellovin also expresses concern about a "push" feature, intended to work around firewalls, which could be used for distributed denial of service attacks. He considers Napster to be more secure against such attacks, although it requires a central server.

It is worthwhile mentioning a very elegant alternative to routing table maintenance and forwarding that we rejected. The idea is to have a node by a group of diverse replicas as suggested by Lybman. The replicas are coordinated using a state machine replication algorithm that can tolerate Byzantine faults. BFT can replicate arbitrary state. It can replicate Pastry's routing table maintenance and forwarding. Additionally, the algorithm in [14] provides strong coordination for overlay routing and maintenance.

However, there are two disadvantages: the solution is not fault-tolerant to Byzantine faults, and it is less resilient than the solution that we presented in [1]. It is more expensive because it requires an agreement protocol between all replicas. If replicas should be geographically dispersed to reduce the impact of Byzantine faults, agreement latency will be high. Additionally, each node can have less than $\frac{1}{2}$ of its nodes faulty. This bound on the number of faulty nodes per group results in a relatively low probability of successful routing. The probability that a replica group with k replicas is correct is $\sum_{i=0}^{k-1} \binom{k}{i} p^i (1-p)^{k-i}$, where p denotes the binomial distribution with k successes, p the probability of a success $\frac{1}{2}$. For example, the probability that a replica group with $k=32$ nodes compromised and 32 replicas is less than 93×10^{-6} . The probability of routing correctly with 100,000 nodes in the Pastry overlay is compromised is $\sum_{i=0}^{k-1} \binom{k}{i} p^i (1-p)^{k-i}$.

8 Conclusions

Structured peer-to-peer overlay networks have previously been used as a model for nodes; any node accessible in the network was assumed to follow the protocol. However, if nodes are malicious and do not follow the protocol, it is possible for a small number of nodes to compromise the network and applications built upon it. This paper has presented techniques for secure node joining, routing table maintenance, and message forwarding in structured p2p overlays. These techniques can be combined with existing techniques to create a more robust network in the presence of malicious participants. A routing protocol of efficient proximity-aware routing in the common case and a costly redundant routing technique only when the test fails to detect interference by an attacker. Moreover, we show how the overhead can be reduced by using self-certifying application data. The network can tolerate up to 25% malicious nodes while providing good performance if the fraction of compromised nodes is small.

Acknowledgments

We wish to thank Robert Morris, Rodrigo Rodrigues, shepherd David Wetherall and the anonymous referee for their comments. We also wish to thank Adam Stubblefield for many discussions about the assignment problem. This work was supported in part by NSF (003604-0079-2001) and NSF (CCR-9985332).

Bibliography

1 M. Bellare and P. Rogaway.

The exact security of digital signatures- How to
In *Advances in Cryptology - EUROCRYPT 96*,
Science, Vol. 1070. Springer-Verlag, 1996.

2 Steve Bellovin.

Security aspects of Napster and Gnutella.

In *2001 Usenix Annual Technical Conference*,
2001.

Invited talk.

3 Miguel Castro, Peter Druschel, Y. Charlie Hu, and
Exploiting network proximity in peer-to-peer overlay
Technical Report MSR-TR-2002-82, Microsoft

4 Miguel Castro and Barbara Liskov.

Practical byzantine fault tolerance.

In *Proceedings of the Third Symposium on Operating
Implementation (OSDI'99)*, New Orleans, Louisiana

5 Ian Clarke, Oskar Sandberg, Brandon Wiley, and
Freenet: A distributed anonymous information system
In *Workshop on Design Issues in Anonymity and
320*, July 2000.

ICSI, Berkeley, California.

6 Thomas H. Cormen, Charles E. Leiserson, and
Introduction to Algorithms.

MIT Electrical Engineering and Computer Science

7 Frank Dabek, M. Frans Kaashoek, David Karger, and
Stoica.

Wide-area cooperative storage with CFS.

In *Proc. ACM SOSP'01*, Banff, Canada, October

8 Drew Dean and Adam Stubblefield.

Using client puzzles to protect TLS.

In *10th Usenix Security Symposium*, pages 1-8,
2001.

9 Roger Dingledine, Michael J. Freedman, and David
Accountability measures for peer-to-peer systems

In *Peer-to-Peer: Harnessing the Power of Disruptive Technologies*, Morgan Kaufmann and Associates, November 2000.

10

John R. Douceur.

The Sybil attack.

In *Proceedings for the 1st International Workshop on Peer-to-Peer Systems (IPTPS '02)*, Cambridge, Massachusetts, March 2002.

11

M. P. Herlihy and J. M. Wing.

Axioms for Concurrent Objects.

In *Proceedings of 14th ACM Symposium on Principles of Distributed Computing Languages*, pages 13-26, January 1987.

12

A. Herzberg, M. Jakobsson, S. Jarecki, H. Krawczyk.

Proactive public key and signature systems.

In *Proc. of the 1997 ACM Conference on Computer and Communications Security*, 1997.

13

Ari Juels and John Brainard.

Client puzzles: A cryptographic defense against denial of service.

In *Internet Society Symposium on Network and Security (NDSS '99)*, pages 151-165, San Diego, California, April 1999.

14

Nancy Lynch, Dahlia Malkhi, and David Ratajczak.

Atomic data access in content addressable networks.

In *Proceedings for the 1st International Workshop on Peer-to-Peer Systems (IPTPS '02)*, Cambridge, Massachusetts, March 2002.

15

Ralph C. Merkle.

Secure communications over insecure channels.

Communications of the ACM, 21(4):294-299, April 1978.

16

Sylvia Ratnasamy, Paul Francis, Mark Handley, Richard Kretzschmar, and Scott Shenker.

A scalable content-addressable network.

In *Proc. ACM SIGCOMM'01*, San Diego, California, August 2001.

17

Antony Rowstron and Peter Druschel.

Pastry: Scalable, distributed object location and high performance peer-to-peer systems.

In *Proc. IFIP/ACM Middleware 2001*, Heidelberg, Germany, November 2001.

18

Antony Rowstron and Peter Druschel.

Storage management and caching in PAST, a large-scale peer storage utility.

In *Proc. ACM SOSP'01*, Banff, Canada, October 2001.

19

Emil Sit and Robert Morris.

Security considerations for peer-to-peer distributed storage.

In *Proceedings for the 1st International Workshop on Peer-to-Peer Security (IPTPS '02)*, Cambridge, Massachusetts, March 2002.

20

Ion Stoica, Robert Morris, David Karger, M. Frans Kaashoek, and Balakrishnan.

Chord: A scalable peer-to-peer lookup service for Internet applications.

In *Proc. ACM SIGCOMM'01*, San Diego, California, August 2001.

21

Ben Y. Zhao, John D. Kubiatowicz, and Anthony D. Joseph.

Tapestry: An infrastructure for fault-resilient wide-area network services.

Technical Report UCB//CSD-01-1141, U. C. Berkeley, 2001.

Appendix

This appendix describes an analytic model for the probability of routing failures in the presence of node failures. We consider the routing failure test with n nodes and m negatives in the routing failure test.

We assume that there exist n nodeIds distributed uniformly over an interval of length L . If L is large and we look at a small, arbitrarily chosen location on this interval (for some x), the number of nodeIds is well approximated in distribution by a Poisson distribution with mean $\lambda = nL^{-1}x$.

In particular, the inter-point distances are approximately independent random variables with mean L/n .

Let X denote the exponential distribution with mean L/n . Let Y denote the exponential distribution with mean L/m , where m is the number of faulty nodes. Suppose X_1, \dots, X_n are independent identically distributed random variables from one of these two distributions and we are required to determine the distribution they are drawn from, e.g., X_1 can be drawn from X or Y . In the context of Pastry and we are trying to determine if the set of nodes is only faulty nodes. An optimal hypothesis test is based on comparing the likelihood ratio to a threshold; by writing down the likelihood ratio test is equivalent to comparing the sample mean, denoted \bar{X} , to a threshold.

We are in a situation where θ is unknown but we have n samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$ (i.e., the samples that we collect from the nodeIds close to \mathbf{x} in the space). We propose the following hypothesis test: choose τ , for some constant τ , and accept/reject H_0 by comparing $\hat{\theta}$ to this threshold. We now compute the true positive probability, $\text{TPR}(\tau)$, and the false negative probability, $\text{FNR}(\tau)$, for this test.

Denote θ_i by θ_i and assume without loss of generality $\theta_i \geq 0$, define

$$\theta_i = \frac{\sum_{j=1}^n \mathbf{x}^{(j)} \cdot \mathbf{x}^{(i)}}{\|\mathbf{x}^{(i)}\|^2}$$

and note that the θ_i are iid random variables. Let $\theta_i \sim \text{Exp}(\lambda)$ denote iid exponential random variables with mean $1/\lambda$. Then, the event that $\hat{\theta} \geq \tau$ is equivalent to the event that $\theta_i \geq \tau$ for some i . Thus,

$$\text{TPR}(\tau) = \mathbb{P}(\theta_i \geq \tau \text{ for some } i) = 1 - \mathbb{P}(\theta_i < \tau \text{ for all } i)$$

where we write \mathbb{P}_θ to denote probabilities when the θ_i are iid $\text{Exp}(\lambda)$. Recalling that θ_i has the gamma distribution with shape parameter n and rate parameter λ , we can rewrite the above as

$$\text{TPR}(\tau) = 1 - \int_0^\tau \frac{\lambda^n}{\Gamma(n)} x^{n-1} e^{-\lambda x} dx$$

$$= 1 - \frac{\lambda^n}{\Gamma(n)} \int_0^\tau x^{n-1} e^{-\lambda x} dx$$

where we used the change of variables $y = \lambda x$ and the equality $\Gamma(n) = (n-1)!$. This expression can be used to compute $\text{TPR}(\tau)$ numerically.

We now derive a simple closed-form expression for $\text{TPR}(\tau)$. The Chernoff bound shows that $\text{TPR}(\tau)$ decays exponentially in the sample size n . We now derive the exact exponential rate of decay. For arbitrary τ , we have the bound that

Now, if λ has an exponential distribution with mean $\frac{1}{\lambda}$

and $\lambda > 0$ for $\lambda > 0$. Thus, for all $\lambda > 0$

The tightest upper bound is obtained by minimising the

λ . The minimum is attained at $\lambda = 1$. So

the bound,

We can derive an expression for the false negative prob-

lines. Now, the X_i are iid with distribution λ , i.e., they

with mean $\frac{1}{\lambda}$, and we are interested in the event

happens, then we fail to reject the hypothesis that the

where we write λ_i to denote probabilities when the

λ_i . In this case, λ_i has the same distribution as λ

distribution as λ , and we obtain using (1) that

This allows us to compute λ numerically and by comb

the following closed-form upper bound

Footnotes

... set¹

The neighbor set size should be chosen to ens

Generated by Sitaram Iyer (ssiyer@cs.rice.edu)

*This paper was originally
published in the
Proceedings of the 5th
Symposium on Operating
Systems Design and
Implementation,
December 9–11, Boston,
MA, US
Last changed: 6 Nov. 2002
aw*

[Technical Program](#)
[OSDI '02 Home](#)
[USENIX home](#)