

# SybilGuard: Defending Against Sybil Attacks via Social Networks

Haifeng Yu

Michael Kaminsky  
Intel Research Pittsburgh

Phillip B. Gibbons

Abraham Flaxman  
Carnegie Mellon University

{haifeng.yu,michael.e.kaminsky,phillip.b.gibbons}@intel.com

abie@cmu.edu

## ABSTRACT

Peer-to-peer and other decentralized, distributed systems are known to be particularly vulnerable to *sybil attacks*. In a sybil attack, a malicious user obtains multiple fake identities and pretends to be multiple, distinct nodes in the system. By controlling a large fraction of the nodes in the system, the malicious user is able to “out vote” the honest users in collaborative tasks such as Byzantine failure defenses. This paper presents *SybilGuard*, a novel protocol for limiting the corruptive influences of sybil attacks. Our protocol is based on the “social network” among user identities, where an edge between two identities indicates a human-established trust relationship. Malicious users can create many identities but few trust relationships. Thus, there is a disproportionately-small “cut” in the graph between the sybil nodes and the honest nodes. SybilGuard exploits this property to bound the number of identities a malicious user can create. We show the effectiveness of SybilGuard both analytically and experimentally.

## Categories and Subject Descriptors

C.2.4 [Computer-Communication Networks]: Distributed Systems—*Distributed applications*; C.2.0 [Computer-Communication Networks]: General—*Security and protection (e.g., firewalls)*

## General Terms

Security, Design, Algorithms, Experimentation

## Keywords

Sybil attack, sybil identity, SybilGuard, social networks

## 1. INTRODUCTION

As the scale of a decentralized distributed system increases, the presence of malicious behavior (e.g., Byzantine failures) becomes the norm rather than the exception. Most designs against such malicious behavior rely on the assumption that a certain fraction of the nodes in the system are honest. For example, virtually all protocols for tolerating Byzantine failures assume that at least  $2/3$  of the nodes are honest. This makes these protocols vulnerable to *sybil attacks* [9], in which a malicious user takes on multiple identities and pretends

to be multiple, distinct nodes (called *sybil nodes* or *sybil identities*) in the system. With sybil nodes comprising a large fraction (e.g., more than  $1/3$ ) of the nodes in the system, the malicious user is able to “out vote” the honest users, effectively breaking previous defenses against malicious behaviors. Thus, an effective defense against sybil attacks would remove a primary practical obstacle to collaborative tasks on peer-to-peer (p2p) and other decentralized systems. Such tasks include not only Byzantine failure defenses, but also voting schemes in file sharing, DHT routing, and identifying worm signatures or spam.

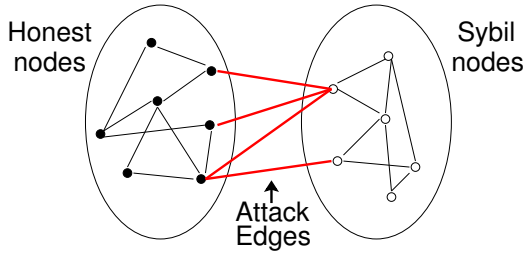
**Problems with using a central authority.** A trusted central authority that issues and verifies credentials unique to an actual human being can control sybil attacks easily. For example, if the system requires users to register with government-issued social security numbers or driver’s license numbers, then the barrier for launching a sybil attack becomes much higher. The central authority may also instead require a payment for each identity. Unfortunately, there are many scenarios where such designs are not desirable. For example, it may be difficult to select/establish a single entity that every user worldwide is willing to trust. Furthermore, the central authority can easily be a single point of failure, a single target for denial-of-service attacks, and also a bottleneck for performance, unless its functionality is itself widely distributed. Finally, requiring sensitive information or payment in order to use a system may scare away many potential users.

**Challenges in decentralized approaches.** Defending against sybil attacks without a trusted central authority is much harder. Many decentralized systems today try to combat sybil attacks by binding an identity to an IP address. However, malicious users can readily harvest (steal) IP addresses. Note that these IP addresses may have little similarity to each other, thereby thwarting attempts to filter based on simple characterizations such as common IP prefix. Spammers, for example, are known to harvest a wide variety of IP addresses to hide the source of their messages, by advertising BGP routes for unused blocks of IP addresses [19]. Beyond just IP harvesting, a malicious user can *co-opt* a large number of end-user machines, creating a *botnet* of thousands of compromised machines spread throughout the Internet. Botnets are particularly hard to defend against because nodes in botnets are indeed distributed end users’ computers.

The first investigation into sybil attacks [9] proved a series of negative results, showing that they cannot be prevented unless special assumptions are made. The difficulty stems from the fact that resource-challenge approaches, such as computation puzzles, require the challenges to be posed/validated simultaneously. Moreover, the adversary can potentially have significantly more resources than a typical user. Even puzzles that require human efforts, such as CAPTCHAs [23], can be reposted on the adversary’s web site to be solved by other users seeking access to the site. Furthermore, these challenges must be performed directly instead of trusting

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

SIGCOMM’06, September 11–15, 2006, Pisa, Italy.  
Copyright 2006 ACM 1-59593-308-5/06/0009 ...\$5.00.



**Figure 1: The social network with honest nodes and sybil nodes. Note that regardless of which nodes in the social network are sybil nodes, we can always “pull” these nodes to the right side to form the logical network in the figure.**

someone else’s challenge results, because sybil nodes can vouch for each other. A more recent proposal [4] suggests the use of network coordinates [17] to determine whether multiple identities belong to the same user (i.e., have similar network coordinates). Despite its elegance, a malicious user controlling just a moderate number of network positions (e.g., tens in practice) can fabricate network coordinates and thus break the defense. Finally, reputation systems based on historical behaviors of nodes are not sufficient either, because the sybil nodes can behave nicely initially, and later launch an attack. Typically, the damage from such an attack can be much larger than the initial contribution (e.g., the damage caused by throwing away another user’s backup data is much larger than the contribution of storing the data). In summary, there has been only limited progress on how to defend against sybil attacks without a trusted central authority, and the problem is widely considered to be quite challenging.

**SybilGuard: A new defense against sybil attacks.** This paper presents *SybilGuard*, a novel decentralized protocol that limits the corruptive influence of sybil attacks, including sybil attacks exploiting IP harvesting and even some sybil attacks launched from botnets outside the system. Our design is based on a unique insight regarding *social networks* (Figure 1), where identities are nodes in the graph and (undirected) edges are human-established trust relations (e.g., friend relations). The edges connecting the honest region (i.e., the region containing all the honest nodes) and the sybil region (i.e., the region containing all the sybil identities created by malicious users) are called *attack edges*. Our protocol ensures that the number of attack edges is independent of the number of sybil identities, and is limited by the number of trust relation pairs between malicious users and honest users.

The basic insight is that if malicious users create too many sybil identities, the graph becomes “strange” in the sense that it has a small *quotient cut*—i.e., a small set of edges (the attack edges) whose removal disconnects a large number of nodes (all the sybil identities) from the rest of the graph. On the other hand, we will show that social networks do not tend to have such cuts. Directly searching for such cuts is not practical, because we would need to obtain the global topology and verify each edge with its two endpoints. Even if we did know the global topology, the problem of finding cuts with the smallest quotient (the *Minimum Quotient Cut* problem) is known to be NP-hard.

Instead, SybilGuard relies on a special kind of verifiable random walk in the graph and intersections between such walks. These walks are designed so that the small quotient cut between the sybil region and the honest region can be used against the malicious users, to bound the number of sybil identities that they can create. We will show the effectiveness of SybilGuard both analytically and

experimentally.

The next section more precisely defines our system model and the sybil attack. Section 3 provides an overview of SybilGuard. Sections 4 and 5 elaborate on SybilGuard in depth. The effectiveness of SybilGuard is shown experimentally in Section 6. Finally, Section 7 discusses related work and Section 8 draws conclusions.

## 2. MODEL & PROBLEM FORMULATION

This section formalizes the desirable properties and functions of a defense system against sybil attacks. We begin by defining our system model. The system has  $n$  honest human beings as *honest users*, and one or more malicious human beings as *malicious users*. By definition, a user is distinct. Each honest user has a single (honest) *identity*, while each malicious user has one or more (malicious) *identities*. To unify terminology, we simply refer to all the identities created by the malicious users as *sybil identities*. Identities are also called *nodes*, and we will from now on use “identity” and “node” interchangeably. All malicious users may collude, and we say that they are all under the control of an *adversary*.

Nodes participate in the system to receive and provide service (e.g., file backup service) as peers. Because the nodes in the system may be honest or sybil, a defense system against sybil attacks aims to provide a mechanism for a node  $V$  to decide whether or not to *accept* or *reject* another node  $S$ . Accepting  $S$  means that  $V$  is willing to receive service from and provide service to  $S$ . Ideally, the defense system should guarantee that  $V$  accepts only honest nodes. Because such an idealized guarantee is challenging to achieve, we aim at providing the following guarantees that, while weaker, are still sufficiently strong to be useful.

**Bounding the number of sybil groups.** The first guarantee is based on defining an *equivalence relation* among accepted nodes. The equivalence relation partitions all accepted nodes into equivalence classes, called *equivalence groups*. Notice that nodes that are rejected do not belong to any equivalence groups. An equivalence group that includes one or more sybil nodes is called a *sybil group*. The defense system provides a guaranteed bound on the *number* of sybil groups, without necessarily knowing which groups are sybil.

Such notion of equivalence groups was also implicitly used by Bazzi and Konjevod [4], where they define (implicit) equivalence classes according to network coordinates. In their scheme, all nodes are accepted, and those nodes with similar network coordinates (e.g., nodes within the same university campus) are considered equivalent. Thus, the number of sybil groups is simply the number of distinct network locations that the adversary controls.

To understand why bounding the number of sybil groups is sufficient in some scenarios, imagine that we are maintaining replicas of a file that has been digitally signed for authenticity. Our goal is to ensure that not all replicas are placed on sybil nodes. If the defense system guarantees that the number of sybil groups is at most some value  $g$ , then placing the file on nodes from  $g + 1$  different equivalence groups will ensure at least one good copy of the file. Another example is replicating a file that is not signed. As long as we obtain the file from  $2g + 1$  nodes from  $2g + 1$  different equivalence groups, the majority is guaranteed to have the correct file.

**Bounding the size of sybil groups.** In some other scenarios, only bounding the number of sybil groups is not effective. Unavoidably, the bound on the number of sybil groups depends on how “powerful” the adversary is. For example, the adversary can always “bribe” or even threaten honest users to act maliciously and thus force the defense system to accept more sybil groups. As a result, one may want a pessimistic estimation of the number of sybil groups  $g$ . On the other hand, even when  $g$  is only moderately large (e.g., 100),

maintaining  $g + 1$  replicas is wasteful.

To be more effective, a defense system may further bound the number of nodes accepted into each of the  $g$  sybil groups. If the number of nodes in each sybil group (or the *size* of the sybil group) is at most  $w$ , then a node will accept at most  $g \cdot w$  sybil nodes. To see the benefits of bounding both the number and size of the sybil groups, consider our running example of replicating unsigned and signed files. Suppose we use a simple assignment that maps replicas to random nodes. If  $g \cdot w$  is smaller than the number of honest nodes  $n$ , then from Chernoff bounds [14], the probability of having a majority of the replicas on honest nodes (as required for unsigned files) approaches 1.0 exponentially fast with the number of replicas. Similarly, as long as  $g \cdot w$  is not much larger than  $n$ , the probability of having at least one replica on an honest node (as required for signed files) also approaches 1.0 exponentially fast.

Choosing roughly uniformly random nodes as replicas is not difficult in most decentralized distributed systems. For example, DHT-based systems (such as those based on Chord [22]) typically place replicas on a random set of nodes. It may appear that instead of choosing uniformly random nodes, we could avoid the need for bounding sybil group sizes by instead choosing uniformly random equivalence groups (and then picking a random node from each chosen group). However, such a design would cause severe load imbalance under heterogeneous group sizes, which is the case, for example, in the network coordinates approach. Moreover, for DHT-based systems, the design would completely disrupt DHT routing.

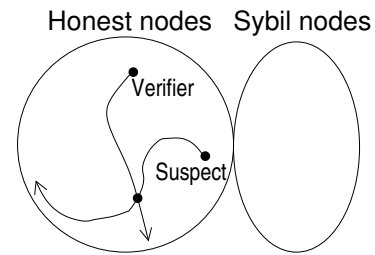
**Side-effects on honest nodes.** As side-effects of bounding the number and size of sybil groups, the defense system may both (mistakenly) reject some honest nodes and (mistakenly) consider two or more distinct honest nodes as equivalent. For example, as noted above, all honest nodes in the same university campus may be considered equivalent in the network coordinates approach.

**Summary of SybilGuard functionalities.** SybilGuard is completely decentralized and all functionalities are with respect to a given node. SybilGuard guarantees that an honest node accepts, and also is accepted by, most other honest nodes (except a few percent in our later simulation) with high probability. Thus, an honest node can successfully obtain service from, and provide service to, most other honest nodes. SybilGuard also guarantees that with high probability, an honest node only accepts a bounded number of sybil nodes. Notice that since SybilGuard is decentralized, the set of accepted nodes by node  $V_1$  can be different from those accepted by node  $V_2$ . However, the difference should be small since both  $V_1$  and  $V_2$  should accept most honest nodes with high probability.

SybilGuard further enables a node  $V$  to partition the accepted nodes (by  $V$ ) into equivalence groups such that only a certain number of those groups contain sybil nodes. Notice that if the application only wants to bound the number of sybil nodes accepted, the notion of equivalence groups does not need to be visible to the application. It is possible for two distinct honest users to be mistakenly considered by SybilGuard to belong to the same equivalence group. This does not affect their ability to receive service. As for providing service, the application may prevent them from, for example, both storing replicas of the same file. As argued in [4], as long as there are a sufficiently large number of equivalence groups, this will not likely result in wasted resource capacity.

### 3. SYBILGUARD OVERVIEW

**Social network and attack edges.** SybilGuard leverages the existing human-established trust relationships among users to bound both the number and size of sybil groups. All honest nodes and sybil nodes in the system form a *social network* (see Figure 1). An



**Figure 2: Verifier accepts the suspect because their random routes intersect. SybilGuard leverages the facts that (1) the average honest node’s random route is highly likely to stay within the honest region and (2) two random routes from honest nodes are highly likely to intersect within  $w$  steps.**

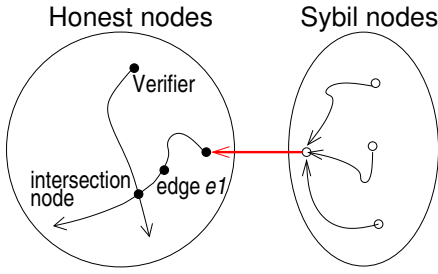
undirected edge exists between two nodes if the two corresponding users have strong social connections (e.g., colleagues or relatives) and trust each other not to launch a sybil attack. If two nodes are connected by an edge, we say the two users are *friends*. Notice that here the edge indicates strong trust, and the notion of friends is quite different from friends in other systems such as online chat rooms. An edge may exist between a sybil node and an honest node if a malicious user (Malory) successfully fools an honest user (Alice) into trusting her. Such an edge is called an *attack edge* and we use  $g$  to denote the total number of attack edges. The authentication mechanism in SybilGuard ensures that regardless of the number of sybil nodes Malory creates, Alice will share an edge with at most one of them (as in the real social network). Thus, the number of attack edges is limited by the number of trust relation pairs that the adversary can establish between honest users and malicious users. While the adversary has only limited influence over the social network, we do assume it may have full knowledge of the social network.

The degree of the nodes in the social network tends to be much smaller than  $n$ , so the system would be of little practical use if nodes only accepted their friends. Instead, SybilGuard bootstraps from the given social network a protocol that enables honest nodes to accept a large fraction of the other honest nodes. It is important to note that SybilGuard does not increase or decrease the number of edges in the social network as a result of its execution.

**Random routes and route intersection.** SybilGuard uses a special kind of random walks, called *random routes*, in the social network. In a standard random walk, at each hop, the current node flips a coin on the fly and selects a (uniformly) random edge to direct the walk. In random routes, each node uses a pre-computed random permutation as a one-to-one mapping from incoming edges to outgoing edges. As a result, two random routes entering an honest node along the same edge will always exit along the same edge (called the *convergence property*). Furthermore, the outgoing edge uniquely determines the incoming edge as well; thus the random routes can be back-traced (called the *back-traceable property*). Of course, these properties can be guaranteed only for the portions of a route that do not contain sybil nodes. Sybil nodes may deviate from any aspect of the protocol.

In the simplest form of SybilGuard, each node performs a random route (starting from itself)<sup>1</sup> of a certain length  $w$  (e.g.,  $w$  is roughly 2000 for the one-million node topology in our later experiments). These random routes form the basis of SybilGuard whereby an honest node (called the *verifier*) decides whether or not to accept another node (called the *suspect*). In particular, the verifier only

<sup>1</sup>In the full protocol, each node performs multiple random routes.



**Figure 3:** All random routes traversing the same edge merge.

accepts a suspect whose random route intersects with the verifier’s random route (see Figure 2). Because of the limited number of attack edges, with appropriate  $w$ , the verifier’s route will remain entirely within the honest region with high probability. (An exception is a verifier with a nearby attack edge; our redundancy techniques discussed in Section 4.4 will address such nodes.)

**Bounding the number and size of sybil groups.** To intersect with the verifier’s random route, a sybil node’s random route must traverse one of the attack edges (whether or not the sybil nodes follow the protocol). Suppose there were only a single attack edge (as in Figure 3). Based on the convergence property, the random routes from sybil nodes must merge completely once they traverse the attack edge. Thus, all of these routes will have the same intersection node with the verifier’s route; furthermore, they enter the intersection node along the same edge (edge  $e1$  in the figure). The verifier thus considers all of these nodes to be in the same equivalence group, and hence there is only a single sybil group. In the more general case of  $g$  attack edges, the number of sybil groups is bounded by  $g$ .

SybilGuard further bounds the size of equivalence groups (and hence of sybil groups) within the length of the random routes  $w$ . From the back-traceable property, we know there can be at most  $w$  distinct routes that (i) intersect with the verifier’s random route at a given node, and (ii) enter the intersection node along a given edge (e.g., along edge  $e1$  in Figure 3). Specifically, the  $i$ th such route,  $i = 1, \dots, w$ , traverses the given edge in its  $i$ th hop. Thus, the verifier accepts exactly one node for each of the  $w$  hop numbers at a given intersection point and a given edge adjacent to the intersection point. In summary, there are many equivalence groups, but only  $g$  are sybil and each has at most  $w$  nodes.

**Guarantees on honest nodes.** For honest nodes, we will show that with appropriate  $w$ , (i) an honest node’s random route intersects with the verifier’s route with high probability, and (ii) such an honest node will never compete for the same hop number with any other node (including sybil nodes). Thus, the average honest node will be accepted with high probability.

SybilGuard partitions the honest nodes in the system into at most  $z$  different equivalence groups, where  $z$  is the sum of the degrees of the  $w$  nodes on the verifier’s route. While  $z$  can still be far from  $n$ , note that  $z$  can easily be much larger than the number of different equivalence groups needed in practice (e.g., when choosing  $g + 1$  different equivalence groups for placing replicas).

Our SybilGuard design leverages the following three important facts to bound the number of sybil nodes: (i) social networks tend to be *fast mixing* (defined in the next section), which necessarily means that subsets of honest nodes have good connectivity to the rest of the social network, (ii) too many sybil nodes (compared to the number of attack edges) disrupts the fast mixing property, and (iii) the verifier is itself an honest node, which breaks symmetry. We will elaborate on these aspects later.

## 4. SYBILGUARD DESIGN

With the preceding high-level sketch in mind, this section provides the detailed design of SybilGuard, explains the insights, and also formally argues about its properties.

### 4.1 Social Network

Consider the social network defined in the previous section. Each pair of friends shares a unique symmetric secret key (e.g., a shared password) called the *edge key*. The edge key is used to authenticate messages between the two friends (e.g., with a Message Authentication Code). Because only the two friends need to know the edge key, key distribution is easily done out-of-band (e.g., via phone calls). A node can also revoke an edge key unilaterally simply by discontinuing use of the key and discarding it.

Because of the nature of the social network and the strong trust associated with the notion of friends in SybilGuard, we expect node degrees to be relatively small and will tend not to increase significantly as  $n$  grows. As a result, a user only needs to invoke out-of-band communication a small number of times. In order to prevent the adversary from increasing the number of attack edges ( $g$ ) dramatically by compromising high-degree honest nodes, each honest node (before compromised) voluntarily constrains its degree within some constant (e.g., 30). Doing so will not affect the guarantees of SybilGuard as long as the social network remains fast mixing. On the other hand, researchers have shown that even with rather small constant node degrees, social networks (or more precisely, small-world topologies) are fast mixing [6, 11].

A node informs its friends of its IP address whenever its IP address changes, to allow continued communication via the network. This IP address is used only as a hint. It does not result in a vulnerability even if the IP address is wrong, because authentication based on the edge key will always be performed. If DNS and DNS names are available, nodes may also provide DNS names and only update the DNS record when the IP address changes.

### 4.2 Limiting the Number of Attack Edges

The effectiveness of SybilGuard relies on there being a limited number of attack edges ( $g$ ). There are several ways the adversary might attempt to increase  $g$ :

- The malicious users establish social trust and convince more honest users in the system to “be their friends” in real life. But this is quite difficult to do on a large scale.
- A malicious user (Malory) who managed to convince an honest user (Alice) to be her friend creates many sybil nodes, and then tries to convince Alice to also be friends with these sybil nodes. But Alice only has a single edge key corresponding to the edge between Alice and Malory. As a result, all messages authenticated using that edge key will be considered by Alice to come from the same edge. Thus the number of attack edges remains unchanged.
- The adversary compromises a single honest node with degree  $d$ . Because  $d$  was already constrained (before the node is compromised) within some constant by the user,  $g$  can be increased by at most some constant. On the other hand, the adversary will not be able to create further attack edges from the node because adding an edge to another honest user requires out-of-band verification by that user. When a user drops and then makes new friends, it is possible for the adversary with access to the old edge keys to “resurrect” dropped edges and hence further increase  $g$ . However, we expect such effect to be negligible in practice and if necessary, can be prevented by requiring out-of-band confirmation when deleting edges.



**Figure 4: Two routes of length 3. Sharing an edge necessarily means that one route starts after the other.**

- The adversary compromises a small fraction of the nodes in the system. This will not likely increase  $g$  excessively due to the reasons above.
- The adversary compromises a large fraction of the nodes in the system. Here the system has already been subverted, and the adversary does not even need to launch a sybil attack. SybilGuard will not help here.
- The adversary compromises a large number of computers (i.e., creates a botnet), only some of which belong to the system. The increase in  $g$  is upper bounded by some constant times the number of compromised computers which already belong to the system. The increase is not affected by the total size of the botnet. Although acquiring a botnet with many nodes may be relatively easy (e.g., in the black market), acquiring a botnet containing many nodes that are already in the system is more challenging.

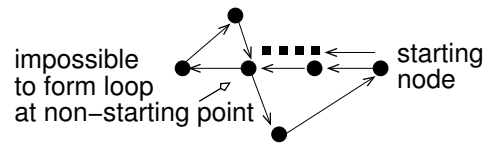
In summary, SybilGuard is quite effective in limiting the number of attack edges, as long as not too many honest users are compromised. Relatively speaking, SybilGuard is more effective defending against malicious users than defending against compromised honest users that belong to the system. This is because a malicious user must make real friends in order to increase the number of attack edges, while compromised honest users already have friends.

### 4.3 Random Routes

Starting from here, the rest of Section 4 assumes a static social network where all nodes are online—we will discuss user and node dynamics in Section 5. SybilGuard relies on the convergence and back-traceable properties in random routes to bound the number and size of sybil groups. Here, we elaborate on how to achieve these properties and their implications.

For random routes, each node uses a randomized routing table to choose the next hop. A node  $A$  with  $d$  neighbors uniformly randomly chooses a permutation “ $x_1, x_2, \dots, x_d$ ” among all permutations of  $1, 2, \dots, d$ . If a random route comes from the  $i$ th edge,  $A$  uses edge  $x_i$  as the next hop. It is possible that  $i = x_i$  for some  $i$ . The routing table of  $A$ , once chosen, will never change (unless  $A$ ’s degree changes—see Section 5). Using such a randomized routing table introduces some correlation in the random choices if a random route visits the same node multiple times. It is possible that random routes become repeated loops due to this; however, later we will explain intuitively and also demonstrate experimentally why this is unlikely.

For random routes in the honest region, these routing tables give us the following properties. First, once two routes traverse the same edge along the same direction, they will merge and stay merged (i.e., the convergence property). Using a permutation as the routing table further guarantees that the random routes are back-traceable. In other words, it is impossible for two routes to enter the same node along different edges but exit along the same direction. With the above properties, if we know that a random route of a certain length  $w$  traverses a certain edge  $e$  along a certain direction in its  $i$ th hop, the entire route is uniquely determined. In other words, there can be only one route with length  $w$  that traverses  $e$  along the given direction at its  $i$ th hop. In addition, if two random routes ever share an edge in the same direction, then one of them must start in the middle of the other (Figure 4).



**Figure 5: A loop can form only at the starting point of a route.**

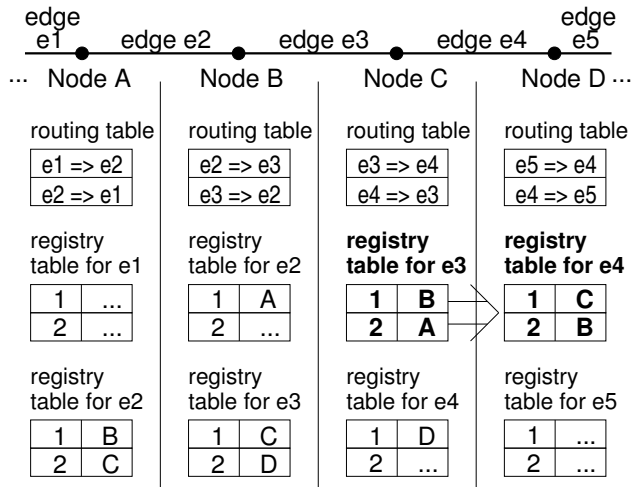
### 4.4 Problematic Routes and Redundancy

A random route is *problematic* if either (i) it traverses some **edge** in the same direction more than once (i.e., a *loop*), or (ii) it enters the sybil region. Note that a route traversing the same **node** more than once may or may not be a loop. Because of the use of routing tables, loops will repeatedly visit many nodes, reducing the “effective” length of the route and the probability of route intersection. On the other hand, random routes that go into the sybil region fall under the control of the adversary. If a verifier uses such a route, it may accept an unbounded number of sybil nodes.

Because the routing table is a permutation, if a random route ever traverses the same edge twice in the same direction, the first edge in the route must be the first edge that is traversed twice. In other words, loops can only form at the starting node (Figure 5). If a loop is formed, the random route must have come back to the starting point, and the starting point must have decided to forward the route along the first edge. Also notice that the smallest loop has three hops, otherwise it is impossible for the route to traverse the same edge (via the same direction) twice. More concretely, consider a simplified scenario where all nodes have the same degree  $d$ . At the second hop, the route will return to the starting point with probability  $1/d$ . At the third hop, if a loop is formed, the starting point must have decided to forward the route along the same edge as the first hop. Thus, a loop is formed at the third hop with probability  $1/d^2$ . As the route proceeds, the chance of repeating the first hop edge will usually become smaller and smaller. In fact, in a fast mixing graph, after a small number of hops a random walk is equally likely to be traversing any edge in a given hop. This provides an intuition as to why loops are unlikely. As for the probability of a random route extending to the sybil region, we will later formally argue (Theorem 1) why this probability is also likely to be small. Finally, Section 6 will provide concrete experimental results demonstrating that problematic random routes are relatively rare.

An effective way to further avoid problematic random routes is to use redundancy. In SybilGuard, a node with degree  $d$  performs  $d$  random routes, one along each of its edges. Now imagine that a verifier  $V$  tries to decide whether to accept a suspect  $S$ . Those routes that are loops can still be used, because they do not compromise security—they are simply less “effective.” We can also safely use all routes from  $S$  regardless of whether they extend to the sybil region: If  $S$  is an honest node, then using all routes simply increases the probability of some route intersecting with  $V$ ’s routes. On the other hand, if  $S$  is a sybil node, then all of  $S$ ’s routes still need to cross the attack edges before intersecting with  $V$ ’s routes that are in the honest region. Because of the convergence property, we can easily see that this will not compromise SybilGuard’s guarantees (Section 3).

On the other hand, if a route from  $V$  extends to the sybil region,  $V$  will not be able to bound the number of sybil nodes using that route.  $V$  uses the following technique to mask the misleading effects of routes extending to the sybil region. For each of  $V$ ’s routes, as long as at least one route from  $S$  intersects that route from  $V$ , that route from  $V$  *accepts*  $S$ . If at least a threshold  $t$  of  $V$ ’s routes accept  $S$ ,  $V$  *accepts*  $S$ . The parameter  $t$  involves the following tradeoff: if  $t$  is too small, then  $V$  may have a large probability of having more than  $t$  routes enter the sybil region; if  $t$  is too large, then  $V$  may have



**Figure 6: Maintaining the registry tables. In order to simplify this example,  $w = 2$ , each node has exactly two edges, and the routing tables are carefully chosen. The node names in the registry tables stand for the nodes’ public keys.**

trouble accepting other honest nodes if more than  $(d-t)$  routes from  $V$  enter the sybil region and if the sybil nodes prevent intersection from happening. A simplified analysis [27] shows that  $t = d/2$  tends to provide a good tradeoff, and this effectively becomes majority voting.

## 4.5 Secure and Decentralized Design for Random Routes and Their Verification

The previous sections explained the basics of random routes. In the actual SybilGuard protocol, these routes are performed in a completely decentralized way. The two local data structures (registry tables and witness tables) described in this section are the only data structures that each node needs to maintain. Also, propagating these tables to direct neighbors is the only action each node needs to take in order to perform random routes.

**Registration.** In SybilGuard, each node  $S$  with degree  $d$  must perform  $d$  random routes of  $w$  hops each and remember these routes. To prevent  $S$  from “lying” about its routes, SybilGuard requires  $S$  to register with all  $w$  nodes along each of its routes. A node  $Q$  along the route permits  $S$  to register only if  $S$  is one of the nodes that are within  $w$  hops “upstream”. When the verifier  $V$  wants to verify  $S$ ,  $V$  will ask the intersection point (between  $S$ ’s route and  $V$ ’s route) whether  $S$  is indeed registered.

In this registration process, each node needs to use a “token” that cannot be easily forged by other nodes. Note that the availability of such tokens does not solve the sybil attack problem by itself, because a malicious user may have many such tokens. A node will be accepted based on its token. The token must be unforgeable to prevent the adversary from stealing the token of an honest node (unless the node is compromised). Our initial design of SybilGuard used a node’s IP address as its token and the node simply registered its IP address. This design assumed no IP spoofing, and was mainly suited for users with static or slowly changing IP addresses.

Our current design of SybilGuard uses public key cryptography for the tokens. This improved design does not rely on the stability of IP addresses, and is secure even under IP spoofing. Each honest node has a locally generated public/private key pair. Notice that these public and private keys have no connection with the edge keys (which are secret symmetric keys). Malicious nodes may create

as many public/private key pairs as they wish. We use the private key of each node as the unforgeable token, while the public key is registered along the random routes as a proof of owning the token. Note that we do not intend or need to solve the public key distribution problem, because we are not concerned with associating public keys to, for example, human beings or computers. The only property SybilGuard relies on is that the private key is unforgeable and its possession can be verified. To perform the registration in a secure and completely decentralized manner, SybilGuard uses registry tables and witness tables, as described next.

**Registry tables.** Each node  $I$  maintains a registry table for each of its edges (Figure 6). The  $i$ th entry in the registry table for edge  $e$  lists the public key of the node whose random route enters  $I$  along  $e$  at its  $i$ th hop. For example, consider the registry table on  $C$  for edge  $e_3$  in Figure 6. Here, one of  $B$ ’s random routes is  $B \rightarrow$  (via edge  $e_3$ )  $C \rightarrow$  (via edge  $e_4$ )  $D$ . In other words, in the first hop of this random route,  $B$  enters  $C$  via edge  $e_3$ . Thus the first entry in the registry table is  $B$ ’s public key. Similarly, the second entry is  $A$ ’s public key. As a result, the registry table has  $w$  entries that are the public keys of the  $w$  “upstream” nodes along the direction of edge  $e_3$  from  $C$ .

Suppose that according to  $C$ ’s routing table,  $e_4$  is the outgoing direction corresponding to  $e_3$  (as in Figure 6).  $C$  will forward its registry table for  $e_3$  to its neighbor  $D$  along  $e_4$ , via a secure channel established using the edge key for  $e_4$ .  $D$  then populates its registry table for  $e_4$  by shifting the registry table from  $C$  downward by one entry and adding  $C$ ’s public key as the new first entry.

As shown in Figure 6, this simple design will ultimately register each node’s public key with all nodes on its  $d$  random routes. The protocol does not have to proceed in synchronous rounds, and nodes in the system may start with empty registry tables. The overhead of the protocol is small as well. Even with one million nodes, if we were to use  $w = 2000$  (already pessimistic given our simulation results), then a registry table is roughly 256KB when using 1024-bit public keys. For a node with 10 neighbors, the total data sent is 2.56MB. A further optimization is to store cryptographically secure hashes of the public keys in the registry table instead of the actual public keys. With each hashed key being 160-bit, the total data sent by each node would be roughly 400KB. Finally, it is important to notice that registry table updates are needed only when social trust relationships change (Section 5). Thus, we expect the bandwidth consumption to be quite acceptable.

**Witness tables.** Registry tables ensure that each node registers with the nodes on its random routes. Each node, on the other hand, also needs to know the set of nodes that are on its random routes. This is achieved by each node maintaining a witness table for each of its edges. The  $i$ th entry in the table contains the public key (or its hash, if we use the above optimization) and the IP address of the node encountered at the  $i$ th hop of the random route along the edge. The public key will later be used for intersection and authentication purposes, while the IP address will be used as a hint to find the node. If the IP address is stale or wrong, it will have the same effect as the intersection node being offline. (Offline nodes are addressed in Section 5.1.)

The witness table is propagated and updated in a similar fashion as the registry table, except that it propagates “backward” (using the reverse of the routing table). In this way, a node will know the  $w$  “downstream” nodes along the direction of each of its edges, which is exactly the set of nodes that are on its random routes. Different from registry tables, witness tables should be updated when a node’s IP address changes (even with a static social network). But this updating can be done lazily, given the optimizations described below in the verification process.

**Verification process.** For a node  $V$  to verify a node  $S$ ,  $V$  needs to perform an intersection between each of its random routes and all of  $S$ 's random routes. To do this,  $S$  sends all of its witness tables to  $V$ , together with  $S$ 's public key. The communication overhead in this step can be reduced using standard optimizations such as Bloom Filters [14] to summarize the nodes in witness tables.

For each of  $V$ 's witness tables,  $V$  performs an intersection with all of  $S$ 's tables, and determines the (hashed) public key of the first intersection point  $X$  (if any) on  $V$ 's route.  $V$  then contacts  $X$  using the recorded IP address in the witness table as a hint.  $V$  authenticates  $X$  by requiring  $X$  to sign each message sent, using its private key. If hashed keys are used,  $X$  also sends its public key, which  $V$  hashes and compares with the stored hash, before authenticating  $X$ . If  $X$  cannot be found using the recorded IP address,  $V$  will try to obtain  $X$ 's IP address from nearby nodes in the witness table. They will likely have  $X$ 's more up-to-date IP address because they are near  $X$ . Because  $V$  will always authenticate  $X$  based on  $X$ 's public key, this does not introduce a vulnerability.

$V$  then checks with  $X$  whether  $S$ 's public key is indeed present in one of  $X$ 's registry tables. The entry number is not relevant. If it is present, then that route from  $V$  accepts  $S$ . If at least half of  $V$ 's routes accept  $S$ ,  $V$  accepts  $S$  (i.e.,  $S$ 's public key). Finally, when interacting with  $S$ ,  $V$  always authenticates  $S$  by requiring  $S$  to sign every message sent, using its private key.

**Key revocation.** A node can easily revoke its old public/private key pair by unilaterally switching to a new public/private key pair, and then using the new public key in its registry table and witness table propagation. The old public key in registry and witness tables will be overwritten by the new public key.

**Sybil nodes.** We described the protocol for the case where all nodes behave honestly. A sybil node may not follow the protocol and may arbitrarily manipulate the registry tables and witness tables. SybilGuard is still secure against such attacks. To understand why and obtain intuition, it helps to consider the set of all registry table entries on all honest nodes in the system. For simplicity, assume that all honest nodes have the same degree  $d$ . Thus there are altogether,  $n \cdot d \cdot w$  registry table entries in the system.

Consider a malicious node  $M$  and a single attack edge connecting an honest node  $A$  with  $M$ . Clearly,  $M$  can propagate to  $A$  an arbitrary registry table, thus polluting the  $w$  entries in  $A$ 's registry table. Suppose  $A$  next forwards the registry table to  $B$ , who shifts the table downward and adds  $A$  as the first entry. Thus  $w - 1$  entries in  $B$ 's registry table are polluted. Continuing this argument, we see that a single attack edge enables  $M$  to control  $w + (w - 1) + \dots + 1 \approx w^2/2$  entries system-wide. With  $g$  attack edges and even when  $gw$  approaches  $n$ , the total number of polluted entries ( $gw^2/2$ ) is still less than half of the total number of entries ( $n \cdot d \cdot w$ ). This provides some intuition why the number of accepted sybil nodes is properly bounded even though the adversary may not follow the SybilGuard protocol.

## 4.6 Designing the Length of Random Routes

A critical design choice in SybilGuard is  $w$ , the length of the random routes. The value of  $w$  must be sufficiently small to ensure that (i) a verifier's random route remains entirely within the honest region with high probability; and (ii) the size of sybil groups is not excessively large. On the other hand,  $w$  must be sufficiently large to ensure that routes will intersect with high probability.

In the following, we provide some analytical assurance that having  $w = \Theta(\sqrt{n} \log n)$  will likely satisfy the above requirements simultaneously. Our results are for random walks instead of the random routes used in SybilGuard—considering random walks allows us to leverage the well-established theory on such walks. Our

full paper [27] explains how these results likely apply to random routes, which will be further confirmed in our later experiments.

We first study the probability that a random walk starting from a random honest node enters the sybil region of the topology.

**THEOREM 1.** *For any connected and non-bipartite social network, the probability that a length- $w$  random walk starting from a uniformly random honest node will ever traverse any of the  $g$  attack edges is upper bounded by  $gw/n$ . In particular, when  $w = \Theta(\sqrt{n} \log n)$  and  $g = o(\sqrt{n}/\log n)$ , this probability is  $o(1)$ .*

We leave the proof to our full paper [27]. The condition of “connected and non-bipartite” on the social network serves to exclude theoretical corner cases. As long as the network has any cycle with an odd number of edges, the network is non-bipartite. The actual likelihood, as shown in our later experiments, is much better than the above pessimistic theoretical bound of  $gw/n$ .

We should point out that the above theorem provides only an “average” guarantee for all honest nodes. Those honest nodes that are closer to attack edges are likely to have a larger probability of walking into the sybil region. Our later simulation results, however, will show that using the redundancy techniques from Section 4.4 will give most nodes a high probability of success.

The next property we would like to show is that  $w = \Theta(\sqrt{n} \log n)$  is likely to be sufficiently large for routes from an honest verifier and an honest suspect to intersect with high probability. Such a property for random walks has been rigorously proved [3, 15] in several other contexts, and thus we only give a high-level review. First, we need to provide some informal background. With a length- $w$  random walk, clearly the distribution of the ending point of the walk depends on the starting point. However, for connected and non-bipartite graphs, the ending point distribution becomes independent of the starting point when  $w \rightarrow \infty$ . This distribution is called the *stationary distribution* of the graph. The *mixing time*  $T$  of a graph quantifies how fast the ending point of a random walk approach the stationary distribution. In other words, after  $\Theta(T)$  steps, the node on the random walk becomes roughly independent of the starting point. If  $T = \Theta(\log n)$ , the graph is called *fast mixing*.

Many randomly-grown topologies are fast mixing, including social networks (or more specifically, small-world topologies) [6, 11]. Thus, a walk of  $\Theta(\sqrt{n} \log n)$  steps contains  $\Theta(\sqrt{n})$  independent samples drawn roughly from the stationary distribution. When the verifier's and the suspect's walks remain in the honest region, both walks draw  $\Theta(\sqrt{n})$  independent samples from roughly the same distribution. It follows from the generalized Birthday Paradox [3, 15] that they intersect with probability  $1 - o(1)$ .

## 4.7 Locally Determining the Appropriate Length of Random Routes

Because SybilGuard is decentralized, each node needs to locally determine  $w$ . Directly setting  $w = \Theta(\sqrt{n} \log n)$  requires the knowledge of  $n$ . This is challenging because we must exclude sybil nodes when estimating  $n$ , which requires running SybilGuard with an appropriate  $w$ .

Instead, to locally determine  $w$ , a node  $A$  first performs a short random walk (e.g., 10 hops), ending at some node  $B$ . Because the random walk is short, with high probability, it stays in the honest region and  $B$  is an honest node. Next  $A$  and  $B$  conceptually both perform random routes to determine how long the two routes need to be to intersect. In practice,  $A$  and  $B$  should have already performed random routes along all directions, thus  $B$  simply needs to hand over one of its witness tables to  $A$ . It is important here to use a standard random walk (instead of a random route) to choose  $B$ , otherwise  $A$ 's random route will always intersect with  $B$  within a small number of

hops. Also, our later simulation will show that even a walk as short as 3 hops suffices to obtain good estimations on  $w$  in a million-node social network.

The intuition behind the above design is that in fast mixing graphs, a random walk of short length is sufficient to approach the stationary distribution. Thus,  $B$  is just a random node drawn from the stationary distribution, and the procedure yields a random sampling of  $w$ . The sampling, however, is biased because the stationary distribution is not necessarily a uniform distribution and  $B$  is more likely to be a higher-degree node than a lower-degree node. On the other hand, notice that if we start a random walk from a uniformly random node  $C$ , then after  $\Theta(T)$  steps ( $T$  being the mixing time), the walk will be at a node roughly drawn from the stationary distribution. Thus the needed route length for two routes (starting from  $A$  and  $C$ , respectively) to intersect is at most  $\Theta(T) + w$ . Since  $w = \Theta(\sqrt{n \log n})$  and  $T = \Theta(\log n)$ , we can safely ignore the term of  $\Theta(T)$ , which will be further confirmed in our later experiments.

Finally, node  $A$  obtains multiple such samples using the above procedure, and calculates the median  $m$  of the samples (see Section 6 for the number of samples needed). It then sets  $w = 2.1m$ , where the constant 2.1 is derived from our analysis of Birthday Paradox distributions [27]. The analysis proves that multiplying the median by 2.1 is sufficient to ensure a collision probability of 95%, regardless of  $n$ . Note that when  $B$  is itself a sybil node or the random route from either  $A$  or  $B$  enters the sybil region, the adversary controls that particular sample. Thus, using the median sample to estimate  $w$  is much more robust than directly using the 95th percentile.

## 5. SYBILGUARD UNDER DYNAMICS

Our protocol so far assumes that the social network is static. In decentralized distributed systems, a typical user first downloads and installs the software (i.e., the user is *created*). The node corresponding to the user may then freely *join* or *leave* the system (i.e., become *online* and *offline*) many times. Finally, the user may decide to uninstall the software and never use it again (i.e., the user is *deleted*). Node join/leave tends to be much more frequent than user creation/deletion. For example, dealing with frequent node join/leave (or “churn”) is often a critical problem faced by DHTs.

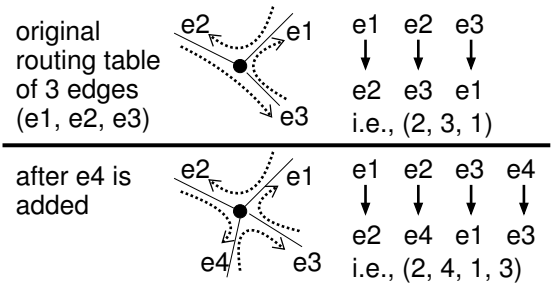
SybilGuard is designed such that it needs to respond only to user creation/deletion, and *not* to node churn. The social network in this paper always includes all users/nodes that have been created and not yet deleted. In other words, many of the nodes in the graph can be offline at any given time.

### 5.1 Dealing with Offline Nodes

In SybilGuard, a node communicates with other nodes only when (i) it tries to verify another node, and hence needs to contact the intersection nodes of the random routes, and (ii) it propagates its registry and witness tables to its neighbors.

For the first scenario, because both the verifier  $V$  and the suspect  $S$  perform multiple random routes (Section 4.4), there will likely be multiple intersections. In fact, even a single route from  $V$  and a single route from  $S$  may still have multiple intersections. The verification can be done as long as a majority of  $V$ ’s routes have at least one intersection point online.

For propagating registry and witness tables, note that this occurs when a random route changes, due to user creation/deletion or edge creation/deletion in the social network. Witness table propagation may also be needed when IP addresses change, but such updating can be performed lazily (Section 4.5). Previous studies [5] on p2p systems show that despite high node churn rate, user creation/deletion occurs only infrequently and the average user lifetime is roughly a



**Figure 7: Incremental maintenance of routing tables. The example assumes that  $d = 3$  and  $k = 2$ . Note that after edge  $e4$  is added, only routes entering via edge  $e2$  need to be redirected.**

year. Similarly, people make and lose social trust relations in real life over months-long time horizons. Thus, the system can afford to take days to completely propagate a new registry or witness table, waiting for nodes to come online. In the case of a new user, prior to becoming a full participant, she can always use the system via a friend as a proxy. As an optimization, a simple *lookahead* routing table design [27] may further help to bypass some offline nodes. For a given node and a given edge adjacent to the node, the lookahead routing table (established in a secure way) records which nodes the route should traverse on the next  $k$  hops.

In the process of propagating/updating registry and witness tables, the social network may change again. Thus, it is helpful to consider it as a decentralized, background stabilization process. This means that if the topology were to stop changing, then the registry and witness tables would eventually stabilize to a consistent state for this (now static) topology.

### 5.2 Incremental Routing Table Maintenance

When users and edges are added or deleted in the social network, the routing tables must be updated as well. Adding a new node can be considered as first adding a node with no edges and then successively adding its edges one by one. Deleting a node can be considered similarly. Thus we only need to discuss edge creation and deletion.

We first explain how  $A$  updates its routing table when a new edge is added between  $A$  and  $B$ . Suppose  $A$ ’s original degree is  $d$  and its original routing table is the permutation “ $x_1, x_2, \dots, x_d$ ”. A trivial way to update  $A$ ’s routing table would be to pick a new random permutation of “ $1, 2, \dots, d, d+1$ ” that is unrelated to “ $x_1, x_2, \dots, x_d$ ”. Doing so, however, would affect/redirect many routes, and incur unnecessary overhead in updating registry and witness tables.

Instead, SybilGuard uses an incremental maintenance algorithm where only routes entering  $A$  along a specific edge may be affected (Figure 7). This reduces the expected overhead on the network by a factor of almost  $d$ . In this algorithm, when a new edge is added to  $A$ ,  $A$  chooses a uniformly random integer  $k$  between 1 and  $d+1$ , inclusive. If  $k = d+1$ , then  $A$ ’s new routing table will be “ $x_1, x_2, \dots, x_d, d+1$ ”. If  $1 \leq k \leq d$ ,  $A$ ’s new routing table will be “ $x_1, x_2, \dots, x_{k-1}, d+1, x_{k+1}, \dots, x_d, x_k$ ”. In other words, we replace  $x_k$  (if exists) with  $d+1$ , and then append  $x_k$  to the end of the permutation. Similarly, for edge deletion, suppose  $A$ ’s original routing table is “ $x_1, x_2, \dots, x_d, x_{d+1}$ ”. Without loss of generality, assume that we are deleting edge  $d+1$ , and let  $k$  be such that  $x_k = d+1$ . If  $k = d+1$ , then  $A$ ’s new routing table is trivially “ $x_1, x_2, \dots, x_d$ ”. Otherwise the new routing table will be “ $x_1, x_2, \dots, x_{k-1}, x_{d+1}, x_{k+1}, \dots, x_d$ ”. In other words, we simply substitute  $x_k$  with  $x_{d+1}$ . For both insertion and deletion, only routes entering  $A$  via edge  $k$  are affected, and one



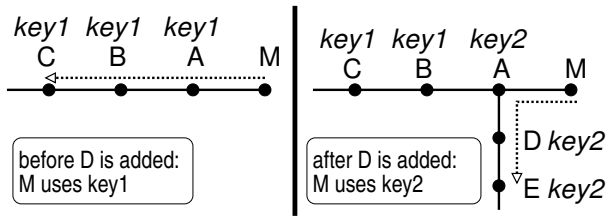


Figure 8: A potential attack by  $M$  during node dynamics.

can prove [27] that the resulting routing table is indeed a uniformly random permutation.

### 5.3 Attacks Exploiting Node Dynamics

This section shows that performing random routes along *all* directions (Section 4.4) actually is necessary for security and provides a defense against potential attacks under node dynamics. We first explain the potential attack scenario. Suppose each node were to perform only a single random route, and consider the example in Figure 8, where  $w = 3$ . Here  $M$  is malicious and the other nodes are honest.  $M$ 's random route is  $M \rightarrow A \rightarrow B \rightarrow C$ . Thus  $A$ ,  $B$ , and  $C$  record  $M$ 's public key  $key1$  in their registry tables. Now another honest node  $D$  joins, and establishes edges with  $A$  and  $E$ .  $A$  updates its routing table, and suppose that routes from  $M$  now go to  $D$  instead of  $B$ . Being malicious,  $M$  launches the attack by changing its public key to  $key2$ . Now  $A$ ,  $D$ , and  $E$  will record  $key2$  in their registry tables. At this point,  $key1$  is registered on  $w - 1$  nodes, while  $key2$  is registered on  $w$  nodes. Both of them are likely to be successfully verified with good probability.

The source of the above vulnerability is that when the routing table on  $A$  changes, the system needs to “revoke” the stale entry of  $key1$  from the registry tables on  $B$  and  $C$ , because  $M$ 's random route no longer passes through these nodes. Explicitly revoking stale entries would introduce considerable complexity because  $B$  and  $C$  may be offline. An alternative design would be to associate TTLs with table entries, which unavoidably introduces a trade-off between security and overheads to refresh expired entries.

SybilGuard prevents the above attack by having all nodes perform random routes along all directions. In particular, if  $D$  (with  $key3$ ) has a random route of  $D \rightarrow A \rightarrow B \rightarrow C$ , then  $key3$  will overwrite  $M$ 's  $key1$ . It is also possible that  $D$ 's route may not be  $D \rightarrow A \rightarrow B \rightarrow C$ . However, it is easy to show that the stale entries will always be overwritten by some node. To understand why, suppose that an entry in  $B$ 's registry table indicates that  $B$  is the  $i$ th hop in the random route of  $M$ . If this entry is stale, it means that  $B$  is no longer the  $i$ th hop in  $M$ 's route. From the back-traceable property of random routes, there must exist another node  $F$  somewhere, such that one of  $F$ 's routes visits  $B$  at the  $i$ th hop. Thus  $F$ 's public key will overwrite the stale entry on  $B$ . In other words, the back-traceable property ensures that for any registry table entry, there is one and exactly one “owner”. Under node dynamics, ownership may change and there may be temporary periods where a malicious user “owns” more entries than it should. However, after the system stabilizes, all entries will be “owned” by the right owner. Based on such observations, we can easily see that other similar attacks under node dynamics will be prevented by SybilGuard as well.

## 6. EVALUATION

This section uses simulation to evaluate the guarantees of SybilGuard. We choose to use simulation because it enables us to study large-scale systems. Because social networks tend to contain private information, there are only a limited number of publicly available

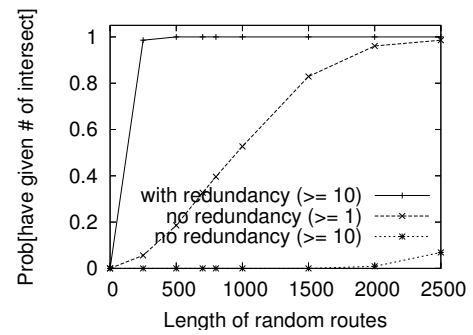


Figure 9: Probability of intersection. The legend “with redundancy” means that each node performs random routes along all directions, while “no redundancy” means performing a single random route. The legend “ $(\geq x)$ ” means that we are considering the probability of having at least  $x$  distinct intersections. SybilGuard corresponds to “with redundancy  $(\geq 10)$ ”.

social network datasets. Those that are publicly available [2, 1] are quite small, which prevents a thorough evaluation of probabilistic guarantees. Thus we use the widely accepted Kleinberg’s synthetic social network model [12] in our evaluation, which generalizes from the Watts-Strogatz model [25]. We use the model to instantiate three different graphs: a million-node graph with average node degree of 24, a 10000-node graph with average degree of 24, and a 100-node graph with average degree of 12. For space limitations, we leave to [27] a review of the model and the detailed parameters. We also focus on the million-node graph, and present only summary results for the other two graphs. All results below are for the million-node graph unless otherwise mentioned.

### 6.1 Results with No Malicious Users

We start by studying the basic behavior of SybilGuard when there are no malicious users. Without malicious users, the only property we are concerned with is whether an honest verifier accepts an honest suspect. This is affected by: (i) whether the random routes from the two nodes intersect; (ii) whether the random routes from the two nodes are loops (which will decrease the chance of intersection); (iii) whether there is at least one intersection node online; and (iv) whether the needed length of random routes is properly estimated.

**Probability of random routes being loops.** As discussed in Section 4.4, if a random route becomes a loop, then its effective length is reduced. Our simulation shows that 99.3% of the routes do not form loops in their first 2500 hops (while later we will show that the needed length of the routes is below 2000). Furthermore, with the redundancy technique in Section 4.4, all the nodes in our simulation have at least one route that is not a loop within their first 2500 hops. For the 10000-node topology, 99.7% of the routes do not form loops in their first 200 hops, which is above the needed route length. For the 100-node topology, 90% of the routes do not form loops in the first 50 hops, which is again above the needed route length.

As the results show that loops are quite rare, and also because they only impact effectiveness rather than security, we will not investigate them further. In all our results below, we do not distinguish loops from non-loops, and thus all the results will already capture the impact of random routes being loops.

**Probability of an honest node being successfully accepted.** We move on to study the probability of the verifier  $V$  accepting the suspect  $S$ . For  $V$  to accept  $S$ , their routes must intersect and at least one intersection must be online. We do not directly model nodes

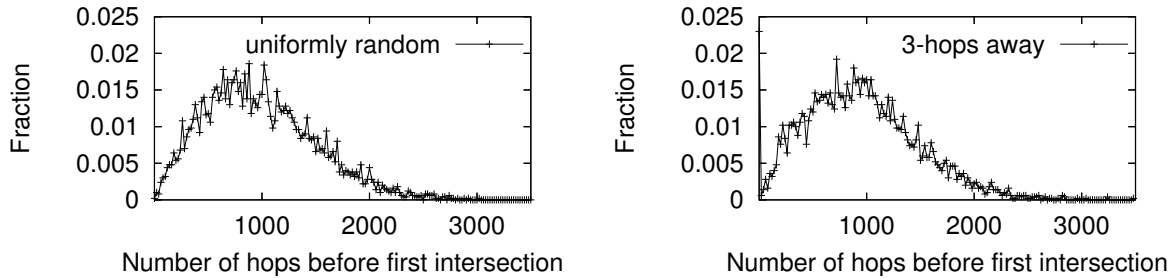


Figure 10: Probability distribution histogram for the number of hops needed before the first intersection.

being online or offline. Rather, we assume that as long as there are at least 10 intersections, the verification succeeds. Note that even when nodes are online only 20% of the time, the probability that at least one out of 10 intersections is online is already roughly 90%.

Figure 9 plots the probability of  $V$  successfully accepting  $S$ , as a function of  $w$  (length of the random routes). For better understanding, we also include in Figure 9 two other curves for the cases where each node performs a single random route, and seeks either at least 1 or 10 intersections. The results show that in a million-node social network, even having a  $w$  as small as 300 yields a 99.96% probability of having at least 10 intersections. On the other hand, if we do not exploit redundancy, the needed length will be much larger. For our 10000-node topology,  $w = 30$  yields a 99.29% probability of having at least 10 intersections. For the 100-node topology,  $w = 15$  gives us a probability of 99.97%.

**Estimating the needed length of the routes  $w$ .** In SybilGuard, each node infers the needed length of the routes using the sampling technique described in Section 4.7. Using this technique, a node  $A$  first performs a short random walk ending at some node  $B$ . Then  $A$  and  $B$  both perform random routes to determine how long the routes need to be in order to intersect. Such estimation would be entirely accurate if (i)  $B$  were chosen uniformly randomly from all nodes in the system; and (ii) the number of samples were infinite. In practice, however, neither condition holds.

To gain insight into the impact of  $B$  not actually being a uniformly random node, Figure 10 depicts the distribution of the number of hops before intersection, comparing the case when  $B$  is chosen uniformly at random to the case when  $B$  is chosen using a 3-hop random walk from  $A$ . As the figure shows, the two distributions are quite similar. This will help to explain later the small impact of  $B$  not being uniformly random. Based on the distribution when  $B$  is chosen uniformly at random, we obtain an accurate  $w$  of 1906 needed for 95% of the pairs to intersect. This value of 1906 will be used as a comparison with SybilGuard’s estimated  $w$ .

To understand the error introduced by having only a finite number of samples, we study how the estimated  $w$  fluctuates and approaches 1906 as a node takes more and more samples. This experiment is repeated from multiple different nodes. In all cases, we observe that the estimated  $w$  always falls within  $1906 \pm 300$  after 30 samples. While after 100 samples, the estimated  $w$  always falls within  $1906 \pm 150$ . These results show that the estimated  $w$  is accurate enough even after a small number of samples. Even with only 30 samples and a worst case estimated  $w$  of 1606, Figure 9 still shows a close-to-100% intersection probability when using redundancy. On the other hand, because taking each sample only involves a 3-hop random walk and the transfer of a witness table, the overhead is quite small. Finally, since the number of users  $n$  changes slowly and  $w$  changes roughly proportionally to  $\sqrt{n} \log n$ , we do not expect  $w$  to change rapidly. Thus a node needs only to re-estimate  $w$ , for example, on

a daily basis. For our 10000-node topology, the accurate  $w$  is 197, and the estimated  $w$  falls within  $197 \pm 30$  after 35 samples. For the 100-node topology, the accurate  $w$  is 24, and the estimated  $w$  falls within  $24 \pm 7$  after 40 samples.

## 6.2 Results with Sybil Attackers

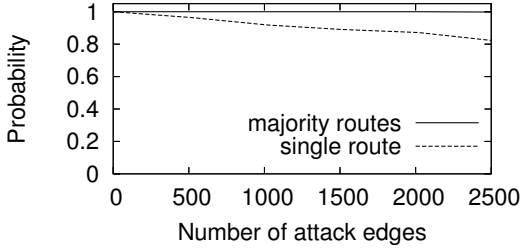
Next we study the behavior of SybilGuard when there are malicious users. In most security research, the term “malicious user” typically refer to a single malicious user who does not assume additional identities. In this paper, however, malicious users refer to powerful attackers who have the sophistication and computation power to launch sybil attacks. For clarity, we use “sybil attackers” to refer to these users in our evaluation. Each of these sybil attackers can potentially create an *unlimited* number of “malicious users”.

Sybil attackers influence the system by creating attack edges. There are clearly many possibilities regarding where the attack edges are in the graph, and we consider two extremes in our experiments. In `random`, we repeatedly pick uniformly random nodes in the graph as sybil attackers, until the total number of attack edges reaches a certain value. In `cluster`, we start from a “seed” node and perform a breadth-first search from the seed. Nodes encountered are marked as sybil attackers, until the total number of attack edges reaches a certain value. All our results below are based on `random` placement, unless explicitly mentioned. We have obtained all corresponding results for `cluster` as well, which are always slightly better but the difference is usually negligible. The reason for better results under `cluster` is that the random routes are more likely to cross attack edges under `random`.

For our experiments based on the million-node graph, we vary the number of attack edges  $g$  from 0 to 2500. When  $g = 2500$ , there are roughly 100 nodes marked as sybil attackers. It is crucial to understand that just having 100 sybil attackers in the system will not necessarily result in 2500 attack edges—on average, each attacker must be able to convince 25 real human beings to be his friend. The hardness of creating these social links is what SybilGuard relies on.

In the presence of sybil attackers, we are concerned with several measures of “goodness”: (i) the probability that an honest node accepts more than  $g \cdot w$  sybil nodes; (ii) the probability that an honest node accepts another honest node; and (iii) the impact of sybil nodes on estimating  $w$ .

**Probability of an honest node accepting more than  $g \cdot w$  sybil nodes.** Routes from an honest verifier  $V$  may enter the sybil region, and the adversary can then direct the routes to intersect with the routes of all sybil nodes. As explained in Section 4.4, SybilGuard uses redundant routes and majority voting to limit the influence of such problematic routes. The curve labeled “majority routes” in Figure 11 shows the probability that the majority of an honest node’s routes remain entirely in the honest region. Here we use  $w = 1906$  as obtained before (the same is true for all the following



**Figure 11: Probability of routes remaining entirely within the honest region.**

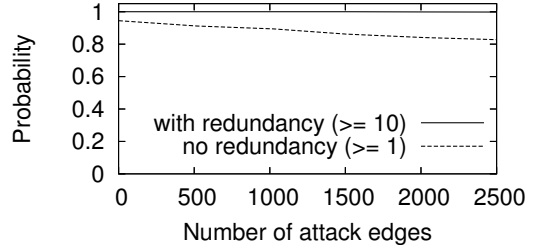
experiments). If a majority of the routes are in the honest region, then the remaining routes will not constitute a majority, and the adversary will not be able to fool the node into accepting more than  $g \cdot w$  sybil nodes. As we can see from the figure, the probability is always almost 100% before  $g = 2000$ , and only drops to 99.8% when  $g = 2500$ . This means that even with 2500 attack edges, only 0.2% of the nodes are not protected by SybilGuard. These are mostly nodes adjacent to multiple attack edges. In some sense, these nodes are “paying the price” for being friends of sybil attackers. For the 10000-node topology and the 100-node topology,  $g = 204$  and  $g = 11$  will result in 0.4% and 5.1% nodes unprotected, respectively. For better understanding, Figure 11 also includes a second curve showing the probability of a single route remaining entirely in the honest region.

**Probability of an honest node being successfully accepted.** In the presence of sybil nodes, the probability that an honest verifier  $V$  accepts another honest suspect  $S$  decreases. First, the routes from  $S$  may enter the sybil region, and the adversary can prevent these routes from intersecting with  $V$ ’s routes. The same is true for  $V$ ’s routes. Second, the presence of sybil nodes necessitates the technique of majority voting as in Section 4.4. This means that among the  $d$  routes from  $V$ , at least  $d/2$  routes need to successfully accept  $S$  before  $V$  can accept  $S$ .

To capture the worst case scenario, here we will assume that after a route (from  $V$  or  $S$ ) enters the sybil region, the rest of the route can no longer be used for verification/intersection. In some sense, the presence of sybil nodes “prunes” the routes. As in Section 6.1, we assume that a “pruned” route from  $V$  accepts  $S$  if it has at least 10 distinct intersections with  $S$ ’s “pruned” routes. Finally,  $V$  successfully accepts  $S$  if a majority of  $V$ ’s routes accept  $S$ .

Figure 12 presents the probability of  $V$  accepting  $S$ , as a function of the number of attack edges  $g$ . This probability is still 99.8% with 2500 attack edges, which is quite satisfactory. The case without using redundancy is much worse (even if we seek only a single intersection), demonstrating that exploiting redundancy is necessary. For our 10000-node topology and 100-node topology,  $g = 204$  and  $g = 11$  give probabilities of 99.6% and 87.7%, respectively. Notice that a 87.7% probability does not mean that 12.3% of the nodes will not be accepted by the system. It only means that given a verifier, 12.3% of the nodes will not be accepted by that verifier. Each honest node, on average, should still be accepted by 87.7% of the honest nodes (verifiers).

**Estimating the needed length of the routes  $w$ .** The final set of experiments seeks to quantify the impact of sybil nodes on the estimated  $w$ . Recall that to estimate  $w$ , a node  $A$  performs a short (3-hop in our experiments) random walk ending at some node  $B$ .  $A$  and  $B$  then both perform random routes to determine when the two routes intersect, which is used as a sample. The sample taken is *bad* (i.e., potentially influenced by the adversary) if any of the



**Figure 12: Probability of an honest node accepting another honest node (i.e., having at least a target number of intersections). The legends are the same as in Figure 9, and SybilGuard corresponds to “with redundancy ( $\geq 10$ )”.**

two routes or the short random walk enters the sybil region. Our simulation shows that the probability of obtaining bad samples roughly increases linearly with the number of attack edges  $g$ . Even when  $g$  reaches 2500, the fraction of bad samples is still below 20%. Since our estimation uses the median of the samples, these 20% bad samples will have only limited influence on the estimate for  $w$ . For our 10000-node topology and 100-node topology, the fraction of bad samples is always below 20% when  $g \leq 204$  and  $g \leq 11$ , respectively.

## 7. RELATED WORK

The sybil attack [9] is a powerful threat faced by any decentralized distributed system (such as a p2p system) that has no central, trusted authority to vouch for a one-to-one correspondence between users and identities. As mentioned in Section 1, the first investigation [9] into sybil attacks already proved a series of negative results.

Bazzi and Konjevod [4] proposed using network coordinates [17] to foil sybil attacks, and a similar idea has also been explored for sensor networks [21]. The scheme relies on the assumption that a malicious user can have only one network position, defined in terms of its minimum latency to a set of beacons. However, with network coordinates in a  $d$ -dimensional space, an adversary controlling more than  $d$  malicious nodes at  $d$  different network positions can fabricate an arbitrary number of network coordinates, and thus break the defense in [4]. This is problematic because  $d$  is usually a small number (e.g.,  $< 10$ ) in practice. Moreover, a solution based on network coordinates fundamentally can only bound the number of sybil groups and not the size of the sybil groups.

Danezis *et al.* [8] proposed a scheme for making DHT lookups more resilient to sybil attacks. The scheme leverages the bootstrap tree of the DHT, where two nodes share an edge if one node introduced the other into the DHT. The insight is that sybil nodes will attach to the rest of the tree only at a limited number of nodes (or attack edges in our terminology). One can imagine defining a similar notion of equivalence groups here, which correspond to subtrees. The scheme can then properly bound the number of sybil groups. In comparison, SybilGuard exploits the graph property in social networks instead of the bootstrap tree, which helps to achieve much stronger properties. First, SybilGuard is able to further bound the size of sybil groups, which is not possible based on bootstrap trees. As a result, even with a single attack edge, the results in [8] deteriorate as the adversary creates more and more sybil nodes. Second, SybilGuard guarantees roughly  $\sqrt{n}$  equivalence groups to ensure sufficient diversity. A bootstrap tree can be in any shape and thus the number of equivalence groups can be rather small. Third, the sizes of different equivalence groups in SybilGuard are roughly the same. In the bootstrap tree approach the sizes can be quite different, which

can lead to significant load imbalance. Finally, compromising even a single node in the bootstrap tree will disconnect the tree, breaking the assumption of the scheme.

**Sybil attacks in sensor networks.** Sybil attacks have also been studied for sensor networks [16]. The solutions there, such as radio resource testing and random key predistribution, unfortunately do not apply to distributed systems in the wide-area. A sybil-related attack in sensor networks is the *node replication attack* [18], where a single compromised sensor is replicated indefinitely, by loading the node's cryptographic information into multiple generic sensor nodes. All these replicated nodes have the same ID (e.g., they all have to use the same secret key issued to the compromised sensor). The solution [18], which is based on simple random walk intersection, does not extend to sybil attacks because the sybil nodes do not necessarily share a single, verifiable ID.

**Sybil attacks in reputation systems.** In a reputation system, each user has a rating describing how well the user behaves. For example, eBay ratings are based on users' previous transactions with other users. Sybil attacks can create a large number of sybil nodes that collude to artificially increase a user's rating. Known defenses [7, 10, 20] against such attacks aim at preventing the sybil nodes from boosting a malicious user's rating (and attracting buyers, in the case of eBay). They cannot and do not aim to control the number or size of sybil groups. All the sybil nodes are able to obtain the same rating/reputation as the malicious user. Thus the sybil attack problem in reputation systems is fundamentally different from the one solved by SybilGuard.

In some other reputation systems such as Credence [24], users cast votes regarding the validity of shared files. The votes are then combined using a weighted average based on the ratings of the user. Sybil nodes are able to dramatically influence the average (even when applying the techniques from [7]), and thus Credence relies on a central authority to limit sybil nodes [24].

**Trust networks and random walks.** The social network in SybilGuard is one kind of trust network. Many previous works [7, 10, 24] use trust networks that are based on past successful transactions or demonstrated shared interest between users. The trust associated with our social network is much stronger, which is essential to the effectiveness of SybilGuard. Such a strong-trust social network is also leveraged by LOCKSS [13], where the verifier accepts all its direct social friends, as well as a proportional number of other nodes. The total number of nodes accepted (proportional to the degree of the verifier) can be orders of magnitude smaller than the system size. Because a node can only accept and thus use a limited number of other nodes in the system, LOCKSS is more suited for specific application scenarios such as digital library maintenance.

Trust propagation or transitive trust is a technique that researchers often use on trust networks [7, 10, 20, 24]. SybilGuard is more related to exploiting graph properties rather than trust propagation. Random walks have also been used to infer worm origin [26] by identifying nodes with a small number of incoming messages but with a large number of outgoing flows. Such techniques are not, however, applicable or related to sybil attacks.

## 8. CONCLUSION

This paper presented SybilGuard, a novel decentralized protocol for limiting the corruptive influences of sybil attacks, by bounding both the number and size of sybil groups. SybilGuard relies on properties of the users' underlying social network, namely that (i) the honest region of the network is fast mixing, and (ii) malicious users may create many nodes but relatively few attack edges. In all our simulation experiments with one million nodes, SybilGuard

ensured that (i) the number and size of sybil groups are properly bounded for 99.8% of the honest users, and (ii) an honest node can accept, and be accepted by, 99.8% of all other honest nodes. Currently we are working on obtaining real social network data to further validate SybilGuard.

## 9. ACKNOWLEDGMENTS

We thank David Andersen, Michael Freedman, Petros Maniatis, Adrian Perrig, Srinivasan Seshan, and the anonymous reviewers for many helpful comments on the paper.

## 10. REFERENCES

- [1] Center for Computational Analysis of Social and Organizational Systems (CASOS), 2006. [http://www.casos.cs.cmu.edu/computational\\_tools/data.php](http://www.casos.cs.cmu.edu/computational_tools/data.php).
- [2] International Network for Social Network Analysis, 2006. [http://www.insna.org/INSNA/data\\_inf.htm](http://www.insna.org/INSNA/data_inf.htm).
- [3] I. Abraham and D. Malkhi. Probabilistic quorums for dynamic systems. In *DISC*, 2003.
- [4] R. Bazzi and G. Konjevod. On the establishment of distinct identities in overlay networks. In *ACM PODC*, 2005.
- [5] W. J. Bolosky, J. R. Douceur, D. Ely, and M. Theimer. Feasibility of a serverless distributed file system deployed on an existing set of desktop PCs. In *ACM SIGMETRICS*, 2000.
- [6] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah. Gossip algorithms: Design, analysis and applications. In *IEEE INFOCOM*, 2005.
- [7] A. Cheng and E. Friedman. Sybilproof reputation mechanisms. In *ACM SIGCOMM Workshop on Economics of Peer-to-Peer Systems*, 2005.
- [8] G. Danezis, C. Lesniewski-Laas, M. F. Kaashoek, and R. Anderson. Sybil-resistant DHT routing. In *European Symposium On Research In Computer Security*, 2005.
- [9] J. Douceur. The Sybil attack. In *IPTPS*, 2002.
- [10] M. Feldman, K. Lai, I. Stoica, and J. Chuang. Robust incentive techniques for peer-to-peer networks. In *ACM Electronic Commerce*, 2004.
- [11] A. D. Flaxman. Expansion and lack thereof in randomly perturbed graphs. Manuscript under submission, 2006.
- [12] J. Kleinberg. The small-world phenomenon: An algorithm perspective. In *STOC*, 2000.
- [13] P. Maniatis, M. Roussopoulos, T. Giuli, D. S. H. Rosenthal, and M. Baker. The LOCKSS peer-to-peer digital preservation system. *ACM TOCS*, 23(1), 2005.
- [14] M. Mitzenmacher and E. Upfal. *Probability and Computing*. Cambridge University Press, 2005.
- [15] R. Morselli, B. Bhattacharjee, A. Srinivasan, and M. Marsh. Efficient lookup on unstructured topologies. In *ACM PODC*, 2005.
- [16] J. Newsome, E. Shi, D. Song, and A. Perrig. The Sybil attack in sensor networks: Analysis & defenses. In *ACM/IEEE IPSN*, 2004.
- [17] T. S. E. Ng and H. Zhang. Predicting internet network distance with coordinates-based approaches. In *IEEE INFOCOM*, 2002.
- [18] B. Parno, A. Perrig, and V. Gligor. Distributed detection of node replication attacks in sensor networks. In *IEEE Symposium on Security and Privacy*, 2005.
- [19] A. Ramachandran and N. Feamster. Understanding the network-level behavior of spammers. In *ACM SIGCOMM*, 2006.
- [20] M. Richardson, R. Agrawal, and P. Domingos. Trust management for the semantic web. In *International Semantic Web Conference*, 2003.
- [21] N. Sastry, U. Shankar, and D. Wagner. Secure verification of location claims. In *ACM Workshop on Wireless Security*, 2003.
- [22] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup service for internet applications. In *ACM SIGCOMM*, 2001.
- [23] L. von Ahn, M. Blum, N. J. Hopper, and J. Langford. CAPTCHA: Telling humans and computers apart. In *Eurocrypt*, 2003.
- [24] K. Walsh and E. G. Sirer. Experience with an object reputation system for peer-to-peer filesharing. In *USENIX NSDI*, 2006.
- [25] D. J. Watts and S. H. Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 393(6684), 1998.
- [26] Y. Xie, V. Sekar, D. Maltz, M. Reiter, and H. Zhang. Worm origin identification using random moonwalks. In *IEEE Symposium on Security and Privacy*, 2005.
- [27] H. Yu, M. Kaminsky, P. B. Gibbons, and A. Flaxman. SybilGuard: Defending against sybil attacks via social networks. Technical Report IRP-TR-06-01, Intel Research Pittsburgh, June 2006. Also available at <http://www.cs.cmu.edu/~yhf/sybilguard-tr.pdf>.