

Practical Robust Communication in DHTs Tolerating a Byzantine Adversary

Maxwell Young Aniket Kate Ian Goldberg Martin Karsten
University of Waterloo, ON, Canada
{m22young, akate, iang, mkarsten}@cs.uwaterloo.ca

Abstract

There are several analytical results on distributed hash tables (DHTs) that can tolerate Byzantine faults. Unfortunately, in such systems, critical operations such as data retrieval and message sending incur significant communication costs. For example, a simple scheme used in many Byzantine fault-tolerant DHT constructions of n nodes requires $O(\log^3 n)$ messages; this is likely impractical for real-world applications. Currently, the best known message complexity is $O(\log^2 n)$ in expectation; however, the corresponding protocol suffers from prohibitive costs owing to hidden constants in the asymptotic notation and to setup costs.

In this paper, we focus on reducing the communication costs against a computationally bounded adversary. We employ threshold cryptography and distributed key generation to define two protocols both of which are more efficient than existing solutions. In comparison, our first protocol is deterministic with $O(\log^2 n)$ message complexity and our second protocol is randomized with expected $O(\log n)$ message complexity. Further, both the hidden constants and setup costs for our protocols are small and no trusted third party is required. Finally, we present results from microbenchmarks conducted over PlanetLab showing that our protocols are practical for deployment under significant levels of churn and adversarial behaviour.

1. Introduction

The peer-to-peer (P2P) paradigm is a popular approach to providing large-scale decentralized services. However, the lack of admission control in many such systems makes them vulnerable to malicious interference [24], [50], [55]. This is a practical concern since large-scale P2P systems are in existence today such as the Azureus DHT [16] and the KAD DHT [51], both of which see more than one million users per day. In addition to file sharing, there are proposals for using P2P systems to protect archived data [18], mitigate the

impact of computer worms [3] and re-implement the Domain Name System [54]; such applications would likely benefit from increased security.

There are a number of results on P2P systems that can provably tolerate Byzantine faults [4]–[6], [17], [22], [28], [37], [46]. This includes the Sybil attack [15] although for ease of exposition, we refer to a Byzantine adversary throughout this work; our results can also be used in conjunction with some proposals specific to the Sybil attack (see the survey of [53]). To date, the majority of results pertain to distributed hash tables (DHTs). A common technique in DHTs that tolerate adversarial faults is the use of *quorums* which are sets of peers such that a minority of the members suffer adversarial faults. A quorum replaces an individual peer as the atomic unit. Adversarial behavior can be overcome by majority action allowing for communication between correct peers; we call this *robust communication*. Since critical operations such as data queries are performed in concert by members of a quorum, robust communication must be efficient.

Several protocols using quorums have been proposed; however, there is a common theme in the way such quorums are utilized. A message m originating from a sending peer p traverses a sequence of quorums Q_1, Q_2, \dots, Q_ℓ until a destination peer is reached. A typical example is a query for content where the destination is a peer q holding a data item. Initially p notifies its own quorum Q_1 that it wishes to transmit m . Each peer in Q_1 forwards m to all peers in Q_2 . A peer in Q_2 determines the correct message by majority filtering on all incoming messages and, in turn, sends to all peers in the next quorum. This forwarding process continues until the quorum Q_ℓ holding p is reached. Assuming that a majority of peers in each quorum are correct, transmission of m is guaranteed. Unfortunately, this simple protocol is costly. If all quorums have size s and the path length is ℓ , then the message complexity is $\ell \cdot s^2$. Typically, $s = \Theta(\log n)$ and, as in Chord [52], $\ell = O(\log n)$ which gives a $O(\log^3 n)$ message complexity which is

likely prohibitively expensive for practical values of n .

Saia and Young [46] give a randomized protocol which provably achieves $O(\log^2 n)$ messages in expectation. While communication between two quorums incurs an expected constant number of messages, the analysis in [46] yields a prohibitively large constant. Furthermore, with probability $1 - o(1)$ some peers will incur $\omega(1)$ message complexity (see appendix for details). The protocol also employs a link architecture between peers requiring the use of a Byzantine agreement protocol. Finally, maintenance and asynchronicity issues remain unresolved.

Therefore, while results exist on the feasibility of robust communication, work on the practicalities has lagged behind. This dearth presents an impediment to the deployment of such systems and we seek to address this outstanding problem.

1.1. Our Contributions

We improve over all previously known results involving communication between quorums [5], [17], [37], [46]. We summarize our main results below:

Theorem 1. *In the computational setting, for an adversary that controls up to an $\epsilon < 1/3$ -fraction of any quorum of size at most s , there are two protocols for achieving robust communication of a message m to a set of peers $D \subseteq Q_i$ for some quorum Q_i over a path of length ℓ . Our Robust Communication Protocol I (RCP-I) has the following properties:*

- The total message complexity (number of messages sent and received) and the message complexity of the sending peer is each at most $2 \cdot s + 4 \cdot s \cdot (\ell - 2) + |D|$.
- The message complexity of every non-sending peer along the lookup path is at most 4.
- The latency (number of roundtrip communication rounds) is at most $2 \cdot (\ell - 2) + 2$.

For our Robust Communication Protocol II (RCP-II):

- The expected total message complexity and the expected message complexity of the sending peer is each at most $2 \cdot s + \frac{(\ell-2)}{(1-\epsilon) \cdot c} + (\ell - 2) + |D|$.
- The expected message complexity of a non-sending peer on the lookup path is at most $\frac{2}{(1-\epsilon) \cdot c \cdot s}$.
- The expected latency is at most $\frac{(\ell-2)}{(1-\epsilon) \cdot c} + 2$.

Here, the constant $c > 0$ is the probability that the response time of a correct peer is at most Δ .

Using the Chord-based construction of [17], the message complexity of RCP-I is $O(\log^2 n)$ and for RCP-II it is $O(\log n)$ in expectation. We tolerate a

large fraction of adversarial peers; strictly less than a $1/3$ -fraction compared to the roughly $1/4$ -fraction in [46]. Our use of a distributed key generation (DKG) scheme allows for security *without* a trusted party or costly updating of public/private keys outside of each quorum. This obviates the need for a trusted third party. To the best of our knowledge, this is the first use of DKG in a Byzantine-tolerant P2P setting.

Finally, we provide microbenchmark results involving two quorums using PlanetLab. Our experimentation demonstrates that our protocols perform well under significant levels of churn and faulty behaviour. In particular, for a 10^5 -node system with a path length of $\ell = 20$, our results imply RCP-I and RCP-II complete in under 4 seconds and 5 seconds, respectively.

2. Related Work

We begin by summarizing a number of related applied results for achieving Byzantine fault tolerance. We then give an overview of a number of theoretical results in the area.

State machine replication (SMR) is a standard method for implementing highly fault-tolerant services [48]. Services are replicated over multiple servers providing a high-integrity distributed system.

A large body of literature exists on implementing Byzantine protocols. While P2P systems do not align perfectly with the state machine replication (SMR) paradigm [48], the large body of literature on Byzantine fault-tolerant replication is relevant to our work. Early work by Reiter [41] yielded protocols for Byzantine agreement and atomic broadcast. Our first protocol shares some common features with the multicast protocol of [41], yet we differ significantly since in the P2P domain we must contend with issues of scalability, churn, and spurious requests aimed at consuming resources. More recently, Castro and Liskov [11] demonstrated efficient Byzantine fault-tolerant SMR; however, this seems unsuitable for a P2P setting due to scaling issues. Several other Byzantine fault-tolerant systems have been implemented such as SINTRA [10], FARSITE [2], the Query/Update protocol [1] and the HQ system [13]; however, scalability issues make the use of these protocols in a P2P setting unlikely.

Two *implemented* large-scale Byzantine fault tolerant storage architectures are OceanStore [32] and Rosebud [44]. The latter scales up to tens of thousands of nodes and handles changing membership. However, with only a single Byzantine node per replication group, Rosebud incurs significant overhead. In contrast, our protocols perform efficiently with 10% of

the peers being Byzantine. Rosebud relies on a *configuration service* (CS) which tracks system membership, ejects faulty nodes, and handles new nodes. The CS, implemented over a set of nodes, introduces a potential bottleneck and possibly a point of attack; similarly, a “primary tier” of replicas is used in OceanStore. In comparison, our protocol is completely decentralized and no special set of nodes is required.

Both Rodrigues, Kouznetsov and Bhattacharjee [43] and Rodrigues, Liskov and Shriram [45] give *proposals* for applying the SMR approach on a large scale; the latter describes a P2P system. However, both works rely on a CS and neither provides empirical results or discusses the details of secure data retrieval and message passing. Wang *et al.* [56] design and implement a routing scheme that tolerates Byzantine faults and demonstrates good performance. However, they require both a certificate authority (CA) and a special set of nodes, called a neighborhood authority, similar to a CS.

In summary, work on practical Byzantine fault-tolerance has focused on ensuring consistency and availability of replicas in the SMR paradigm; as noted in [43], this does not necessarily characterize the P2P domain. In many cases, the reliance on a special set of nodes introduces a potential bottleneck to the system; in contrast, our use of DKG and threshold cryptography removes this bottleneck. Finally, both empirically validated systems and proposals for larger-scale applications are either not aimed at or do not directly address the issue of robust communication in a peer-to-peer-like environment.

There are several theoretical results on Byzantine fault-tolerant DHTs [5], [17], [22], [37]. These results make use of *quorums*, which are sets of $\Theta(\log n)$ peers such that a majority of the peers in a quorum are correct. Awerbuch and Scheideler show how to maintain quorums [4]–[6]. Saia and Young [46] demonstrate more efficient robust communication but, as discussed earlier, several issues remain unresolved.

Castro *et al.* [24], Halo [29], and Salsa [36] handle Byzantine faults by routing along multiple diverse routes. The proposal in [24] requires a CA whereas we do not rely on any trusted third party. In both [29] and [36], the guarantees are unclear against an adversary who owns a large IP-address space or targets identifiers over time as described in [5]. In contrast, defenses for quorum-based protocols are known [4]–[6]. Such an adaptive adversary could potentially compromise the “knuckle” nodes in [29] or the global contacts used in [36]; in contrast, defenses for quorum-based protocols are known [4]–[6].

There are several other works relating to issues of

security in P2P networks. Finally, the ShadowWalker system [34] addresses the issue of anonymity and routes securely using the notion of multiple “shadows” which are similar to a quorum; however, our protocols differ significantly. The Brahms system [9] allows for uniform sampling of peers despite a Byzantine adversary. The Fireflies architecture [28] allows each peer to remain informed of live members in the system despite malicious attacks; however, its applicability likely extends only to single hop overlays such as in work by Gupta *et al.* [21] and secure routing in multi-hop networks is not treated.

3. System Model

Each peer p is assumed to have a unique identifier, p_{ID} , and a network address, p_{addr} . Byzantine peers are also referred to as *faulty* or *adversarial*; all other peers are called *correct*. A fraction of the correct peers may crash due to a system failure or leave the DHT gracefully. We model such peers as having *crashed*.

We adopt an asynchronous communication model with unbounded message delivery time. However, for liveness in DKG and in our second protocol, we use a *weak synchrony* assumption by Castro and Liskov [12].

Peers p and q are said to communicate directly if each has the other in its routing table. The target of m is a set of peers D within a single quorum; m may be a data item request and D may consist of a single peer or multiple peers depending on how data is stored.

3.1. The Quorum Topology

There are several different approaches to how quorums are created and maintained [5], [37], [46]; we refer the reader to [17] for a detailed explanation. Despite these different approaches, we may view the setup of quorums as a graph where nodes correspond to quorums and edges correspond to communication capability between quorums; we refer to this as the *quorum topology*. Figure 1 illustrates how quorums can be linked in a DHT such as Chord. Peers will likely have different views of the network and hence membership lists for Q_i may differ for two peers; however, such issues can be overcome (see [17]). We assume the following four simple invariants are true:

- 1) *Goodness*: each quorum has size $\Theta(s)$ for $s = \Omega(\log n)$ and possesses at most an ϵ -fraction of Byzantine peers for $\epsilon < 1/3$.
- 2) *Membership*: every peer belongs to at least one quorum.

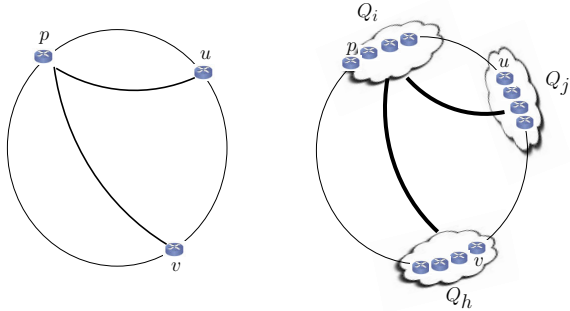


Figure 1. (Left) Three peers on a DHT ring where p links to u and v . (Right) An example of a quorum topology in a DHT ring where $p \in Q_i$, $u \in Q_j$ and $v \in Q_h$. Thick lines signify inter-quorum links.

- 3) *Intra-Quorum Communication*: every peer can communicate directly to all other members of its quorums.
- 4) *Inter-Quorum Communication*: if Q_i and Q_j share an edge in the quorum topology, then $p \in Q_i$ may communicate directly with any member of Q_j and the converse is true.

These four invariants are standard in the sense that previous works on quorums in robust DHTs systems ensure they hold with probability nearly equal to 1. For example, results for maintaining the goodness invariant in DHTs are known [4]–[6]. In terms of the membership invariant, there exist quorum topologies where a peer may belong to several different quorums simultaneously [17], [37]. Finally, to the best of our knowledge, no implementation of a quorum topology exists; this represents another gap between theory and practice. A number of challenges remain in bridging this gap and such an endeavor is outside the scope of this current work. However, the literature suggests that, with the proper deployment, maintaining these four invariants in real-world DHTs is plausible.

3.2. Assumptions

The adversary is assumed to have full knowledge of the network topology and control all faulty peers, which forms a constant fraction of all nodes in the system. In concert with the goodness invariant, strictly less than $1/3$ of the peers in any quorum can be faulty. These peers may collude and coordinate their attacks. Our adversary is computationally bounded with a security parameter κ and it has to do 2^κ computation to break the security of the Gap Diffie-Hellman (GDH) problem [23] in an appropriate group.

Our protocols guarantee successful transmission of a message; however, feasibility is not enough. Our protocols must be efficient, both in terms of (1) the costs incurred by correct peers for legitimate network

operations and (2) the costs incurred by adversarial behavior. The latter concern is crucial since it does no good to provide solutions that allow the adversary to easily launch costly attacks. We first discuss the cryptographic techniques for gaining efficiency and then elaborate on points (1) and (2).

3.3. Threshold Cryptography

We use threshold cryptography to authenticate messages. The idea behind an (η, t) -threshold scheme is to distribute a secret key among η parties in order to remove any single point of failure. Any subset of more than t parties can jointly reconstruct the secret key or perform the required computation securely in the presence of a Byzantine adversary which controls up to t parties. We use threshold signatures to authenticate the communication between quorums.

Threshold Signatures: In an (η, t) -threshold signature scheme, a signing (private) key k is distributed among η parties by a trusted dealer using a verifiable secret sharing protocol [25] or by a completely distributed approach using a DKG protocol [39]. Along with private key shares k_i for each party, the distribution algorithm also generates a verification (public) key K and the associated public key shares \hat{K} . To sign a message m , any subset of $t + 1$ or more parties use their shares to generate the signature shares σ_i . Any party can combine these signature shares to form a message-signature pair $S = (m, \sigma) = [m]_k$ that can be verified using the public key K ; however, this does not reveal k . We refer to a message-signature pair S as a signature. It is also possible to verify σ_i using the public key shares \hat{K} . We assume that no computationally bounded adversary that corrupts up to t parties can forge a signature $S' = (m', \sigma')$ for a message m' . Further, malicious behavior by up to t parties cannot prevent generation of a signature.

Many threshold signature schemes have been constructed in the literature. However, the threshold signature schemes other than [7], [19], [49] are impractical for a variety of reasons. Out of these three practical schemes, the threshold DSS scheme [19] requires a significant amount of interaction among all the parties for every signature generated, while removing the requirement of a trusted dealer is a difficult multiparty computation problem in the threshold RSA signature scheme [49]. The threshold version [7] of the Boneh-Lynn-Shacham (BLS) signature scheme [23] avoids all of the above problems. Its key generation does not mandate a trusted dealer. The signature generation protocol does not require any interaction among the signing parties or any zero-knowledge proofs. Further,

the BLS signature size and generation algorithm are more efficient than RSA and DSS signatures. Therefore, to authenticate the communication between the quorums, we use the threshold BLS signature scheme. This scheme uses the concept of bilinear pairing in the elliptic curve cryptography (ECC) setting and its security is based on difficulty of solving GDH problem (Refer to [7] for a detailed description).

Distributed Key Generation (DKG): In absence of a trusted party in the P2P paradigm, we use a DKG scheme to generate the (distributed) private key. An (η, t) -DKG protocol allows a set of η nodes to construct a shared secret key k such that its shares k_i are distributed over the nodes and no coalition of fewer than t nodes may reconstruct the secret; no trusted dealer is required. There is also an associated public key K and a set of public key shares \hat{K} for verification.

The protocol in [31] is the first DKG for an asynchronous setting; therefore, it is uniquely suitable for deployment in a P2P network. Along with a Byzantine adversary, this protocol also tolerates crash failures. For a quorum of size $s = \eta$, with t Byzantine nodes and f correct nodes that can crash, the DKG protocol requires that $s \geq 3t + 2f + 1$. In our case, this *security threshold* holds due to the goodness invariant in Section 3.1. The DKG protocol allows for system dynamics without changing the system public key K . Notably, the message complexity of a batch of peers (say set P) all joining and/or all leaving the quorum is the same as for a single peer joining/leaving the quorum, while the bit complexity increases only linearly with $|P|$ (see [31, Sec. 6]); for efficiency, we batch such operations during our analysis in Section 5.2. The DKG protocol also considers *mobile adversary* [38] and provides proactive security using share renewal and share recovery protocols.

3.4. Spamming Attacks

A critical concern is that the adversary may launch spurious communications aimed at consuming resources; we refer to such behavior as *spamming*. For example, a malicious peer may initiate a number of data retrieval requests [50], [55]. Here the situation is more dire since the impact of such attacks is multiplied by the group action in a quorum-based system.

Ultimately, there is no perfect defence against an adversary with the resources to initiate massive spamming attacks and this is not our focus. In such circumstances, spamming amounts to a denial-of-service (DoS) attack which forces correct peers to quit the system; unfortunately, there appears to be no adequate remedy. On the other hand, an extended cuckoo rule

exists for maintaining the goodness invariants against limited DoS attacks [6] and this result is compatible with our proposal. Regardless, handling massive spamming or DoS attacks is a challenging problem that falls beyond the scope of this current work. Rather we show that our protocols do not afford the adversary an advantage in launching such attacks. Our goal is to prevent the adversary from forcing a peer to perform expensive operations with impunity. For any operation initiated by a spammer p , this can be accomplished by either (A) placing the bulk of the cost of executing said operation on p or (B) making the detection of spamming inexpensive. As we will show in Section 4, our protocol RCP-1 in Section 4.1 employs principle (A) while our protocol RCP-II in Section 4.2 employs principle (B).

In addition to cryptographic techniques, we assume a *rule set* to reduce the impact of spamming attacks as introduced by Fiat *et al.* [17]. A rule set defines acceptable behavior in a quorum; for example, the number of data lookup operations a peer may execute per duration of time, or tit-for-tat behavior for uploads/downloads. Such rules are known to everyone within a quorum and can be implemented at the software level or agreed upon by quorum members. Requests from a peer that deviates from the rule set are ignored by the other members of its quorum.

3.5. Efficiency (Not Feasibility) Through Cryptography: The Prove-and-Verify Scenario

We now discuss the merits of employing cryptographic techniques. In the presence of Byzantine peers, no single peer can be trusted. Quorums are employed to overcome this trust deficit through majority action. Using the simple protocol outlined in Section 1, messages are guaranteed to reach the intended recipient. *Therefore, quorums allow for robust communication without the need for cryptographic techniques.*

However, spamming attacks still pose a critical problem. For example, a group of Byzantine peers may pretend to be a quorum and initiate requests. Therefore, simply obeying a request because it *appears* to come from a quorum does not prevent spamming. A standard fix is that a quorum responds only to requests that are “proven” to be legitimate. Yet, there is a cost to proving legitimacy; we explore this to motivate our protocols. First, we expand on the utility of a quorum topology in proving legitimacy. We then show how cryptographic techniques improve the efficiency of this task.

Utility of the Quorum Topology: We compare two general scenarios in order to demonstrate the utility

of a quorum topology in proving requests legitimate. The first assumes that proofs and verifications are required to initiate operations; call this the *prove-and-verify scenario*. The second assumes no proof is required before acting (although, each peer may keep a record of misbehaving peers); call this the *passive scenario*. P2P systems often lack admission control and, if forced to leave the system, a Byzantine peer may simply rejoin the network with a new identity. In the worst case, perpetual and rapid rejoin operations result in a DoS attack. Therefore, we make the standard assumption that there is a cost for joining the network (for example, monetary costs as in [24] or CAPTCHAs as suggested in [36]). The best method of enforcing this cost is beyond the scope of this work.

Let τ denote the rate at which p can issue spurious requests before being forced to rejoin the system. In the passive scenario, a Byzantine peer p can contact *any* quorum Q_i by colluding with other faulty peers to obtain necessary routing information. Members of Q_i act on any request coming from p . Therefore, a correct peer may be required to maintain $O(n)$ records so that spam requests are ignored; clearly, this is far from ideal. Moreover, here τ is large due to the abundance of potential targets. In contrast, in a prove-and-verify system the members of Q_i must verify p 's proof before acting. Proof and verification may take different forms. For instance, constructions exist where two peers communicate only if their respective quorums are linked [17], [37]; that is, the quorum topology itself acts as proof. Verification occurs by having a quorum Q_i act on p 's request only if each peer in Q_i receives messages from a majority in Q_p . Here τ is smaller; however, there are still shortcomings to this method of proof and verification.

Efficiency in the Prove-and-Verify Scenario: We argue two things: (1) the form of proof discussed above is restrictive and (2) verification is expensive. First, the proof is restrictive since for Q_i and Q_j to communicate without sending through intermediary quorums, they must maintain links to one another; such maintenance is costly. Second, the verification process is expensive because when communication occurs from Q_i to Q_j , a correct peer $q \in Q_j$ must know to which peers in Q_i it must listen; this incurs more maintenance costs. These are two significant problems with existing schemes.

Cryptography allows us to improve asymptotically on the message complexity of verification. Under our protocols, each quorum has a public and private key established using DKG. Communication can occur between any two quorums that know and can verify each other's public key. Therefore, the form of proof is not

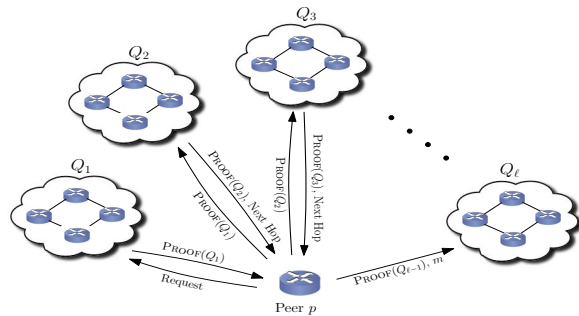


Figure 2. Our general robust communication scheme. At step $i = 1, \dots, \ell - 1$, peer p presents proof, $\text{PROOF}(Q_i)$, that quorum Q_i sanctions p 's action, and receives new proof from Q_{i+1} in addition to routing information for the next hop. At the final step ℓ , peer p sends $\text{PROOF}(Q_{\ell-1})$ and m .

as restricted by the quorum topology and we exploit this in RCP-II. Furthermore, verification is cheaper; using $O(s)$ messages in RCP-I or $O(1)$ expected messages in RCP-II. Of course, overhead is incurred by using cryptography. Message sizes increase by an additional $O(\kappa)$ bits and keys shares, but *not* the key itself, must be updated when membership changes. However, our experimental results in Section 5 show that this overhead is tolerable since the computation costs are significantly smaller than the network latency. Hence, cryptography provides a more efficient and flexible implementation of the prove-and-verify scenario.

4. Our Robust Communication Protocols

We propose two robust communication protocols: RCP-I and RCP-II. Here we outline a general scheme in Figure 2 that is later refined to give our two protocols. Consider a sending peer p who wishes to send a message m to peer q . We assume m is associated with a key value which yields information necessary for distributed routing; that is, the next peer to which m should be forwarded is always known. Peer p notifies its quorum Q_1 that it is performing robust communication and receives $\text{PROOF}(Q_1)$. Peer p sends this to Q_2 as proof that p 's actions are legitimate; the form of this proof is discussed later. Depending on the scheme, one or more members of Q_2 examines the proof and, upon verifying it, sends to p : (1) routing information for Q_3 and (2) $\text{PROOF}(Q_2)$, that will convince Q_3 that p 's actions are legitimate. This continues iteratively until p contacts the quorum holding q and m is delivered. We employ the following concepts:

Quorum Public/Private Keys: Each quorum Q_i is associated with a (distributed) public/private key pair

(K_{Q_i}, k_{Q_i}) ; however, there are two crucial differences between how such a key pair is utilized here in comparison to traditional implementations. First, only those quorums linked to Q_i in the quorum topology, and not everyone in the network, need to know K_{Q_i} . Second, (K_{Q_i}, k_{Q_i}) is created using the DKG protocol and \widehat{K}_{Q_i} is the associated set of public key shares.

Individual Public/Private Key Shares: Each peer $p \in Q_i$ possesses a private key share $(k_{Q_i})_p$ of k_{Q_i} produced using DKG. Unlike the quorum public/private key pair of Q_i which must be known to all quorums to which Q_i is linked in the quorum topology, only the members of Q_i need to know the corresponding public key shares \widehat{K}_{Q_i} , which plays an important role in allowing members of Q_i to verify that the signature share sent to peer p is valid.

4.1. Robust Communication Protocol I

We now illustrate RCP-I for a peer p who wishes to send a message m . The path m takes through quorums is denoted by Q_1, \dots, Q_ℓ . We assume that $p \in Q_1$ and the target of the message is a set of peers $D \subseteq Q_\ell$.

Overview: We outline RCP-I; the pseudocode is given in Figure 3. Initially, the correct peers of Q_1 must acquiesce to p 's request. Peer p begins by sending $[p_{\text{ID}}|p_{\text{addr}}|\text{key}|ts_1]$ to all peers in its quorum Q_1 . The value key corresponds to the intended destination of m and ts_1 is a time stamp. The message m can also be sent, and its hash can be added inside the signature below; however, for simplicity, we assume m is sent only in the last step. Each correct peer $q \in Q_1$ then consults the rule set and sends its signature share to p if p is not in violation. Peer p interpolates these signature shares to generate the signature: $S_1 \leftarrow [p_{\text{ID}}|p_{\text{addr}}|\text{key}|ts_1]_{k_{Q_1}}$.

In each intermediate step $i = 2, \dots, \ell - 1$, p sends its most recent signature S_{i-1} and a new time stamp ts_i to each peer $q \in Q_i$ along the lookup path. Since Q_i is linked to Q_{i-1} in the quorum topology, each q knows the public key $K_{Q_{i-1}}$ to verify S_{i-1} . If S_{i-1} is verified and ts_i is valid, q sends back its signature share, $K_{Q_{i+1}}$ and the routing information. Peer p collects the shares to form S_i and majority filters on the routing information for Q_{i+1} . In terms of majority filtering, both group membership and the corresponding routing information are agreed upon using DKG. Finally, for Q_ℓ , p sends m along with $S_{\ell-1}$ to peers in the set D .

Share Corruption Attack: Note the following attack: a set of Byzantine peers $B \subsetneq Q_i$ send invalid shares to p and, therefore, p will fail to construct S_i . We refer to this attack as the *share corruption attack*. Here, the

RCP-I: SENDING PEER p

Initial Step:

- 1: $p \in Q_1$ sends the following request to all peers in Q_1 : $[p_{\text{ID}}|p_{\text{addr}}|\text{key}|ts_1]$
- 2: p interpolates all received signature shares to form: $S_{Q_1} \leftarrow [p_{\text{ID}}|p_{\text{addr}}|\text{key}|ts_1]_{k_{Q_1}}$

Intermediate Steps:

- 3: **for** $i = 2$ to $\ell - 1$ **do**
- 4: p sends $S_{Q_{i-1}}$ and ts_i to every peer in Q_i and requests a signature S_{Q_i} , public key $K_{Q_{i+1}}$ and routing information for Q_{i+1} .
- 5: p interpolates received signature shares to form $S_{Q_i} \leftarrow [p_{\text{ID}}|p_{\text{addr}}|\text{key}|ts_i]_{k_{Q_i}}$.
- 6: p verifies if S_{Q_i} is valid using K_{Q_i} .
- 7: **if** (S_{Q_i} is invalid) **then**
- 8: p sends signature shares to each peer in Q_i .

Final Step:

- 9: p sends $S_{\ell-1}$ to $D \subseteq Q_\ell$ along with m .

RCP-I: RECEIVING PEER $q \in Q_i$

Initial Step:

- 1: **if** ($q \in Q_1$ receives a request by p) **then**
- 2: q checks that a request by p does not violate the rule set. If the request is legitimate, q sends its signature share to p .

Intermediate Steps:

- 3: **if** (q receives $S_{Q_{i-1}}$ and ts_i from p) **then**
- 4: q verifies a $S_{Q_{i-1}}$ using $K_{Q_{i-1}}$ and validates ts_i ; if successful, q sends its signature share, $K_{Q_{i+1}}$ and routing information for Q_{i+1} to p .
- 5: **if** (q receives signature shares from p) **then**
- 6: q verifies all shares using public key shares and informs p of invalid shares.

Figure 3. Pseudocode for RCP-I

individual public/private key shares play a crucial role. To obtain S_i , p sends the received shares to each peer in Q_i using one message per peer. For a share sent to p by a peer in Q_i , each correct peer in Q_i verifies the share using \widehat{K}_{Q_i} . All valid shares are then sent back to p who creates S_i . While members of Q_i may identify those peers which p alleges sent an incorrect share, punitive action is limited since p could be Byzantine. Note that the shares are not recomputed; hence, the adversary can only perform this attack once per step.

Lemma 1. *RCP I guarantees that m is transmitted to a target set of peers $D \subseteq Q_i$ for some quorum Q_i over a path of length ℓ with the following properties:*

- Both the total message complexity and the message complexity of the sending peer is each at

most $2 \cdot s + 4 \cdot s \cdot (\ell - 2) + |D|$.

- Each non-sending peer has message complexity at most 4 messages.
- The latency is at most $2 \cdot (\ell - 2) + 2$.

Proof: First we prove the correctness of our protocol. We show that if p is correct and has not violated the rule set, at each step i of the protocol p either (1) receives a valid signature and routing information for the next step or (2) terminates the protocol by delivering m to all members of D ; correctness follows directly. Our proof is by induction on i :

Base Case: Consider the initial step $i = 1$ where p communicates with the peers in its quorum $Q_p = Q_1$ about sending the message m . If p is correct and has not violated the rule set, upon receiving $[p_{\text{ID}}|p_{\text{addr}}|\text{key}|ts_1]$ all correct peers will send their shares to p . Therefore, p is guaranteed to form S by the goodness invariant. Peer p can then check whether S is valid and, if so, sets S to be S_1 . Otherwise, p must overcome the share corruption attack. Since p belongs to Q_1 , peer p knows the individual public key shares of each peer in Q_1 and can therefore detect which shares are invalid and construct S_1 . Finally, p already has the routing information for Q_2 ; therefore, the base case holds.

Inductive Hypothesis: Assume that at each step up to step $i < \ell$, p has obtained the correct signatures and routing information.

Inductive Step: At step $i + 1$, peer p sends S_i to Q_{i+1} . By the inductive hypothesis, this signature is valid and p possesses the routing information for Q_{i+1} . If $i < \ell - 1$, and no corrupted share attack occurs, then p 's request for S_{i+1} and the routing information for Q_{i+2} will be satisfied due to the goodness invariant. Otherwise, p must overcome the corrupted share attack by sending all signed shares to all other peers in Q_{i+1} . Each correct peer in Q_{i+1} can detect and inform p which peers sent an invalid share. Due to the goodness invariant, peer p can majority filter on these responses to determine the invalid shares and then construct S_{i+1} . If $i = \ell - 1$, p possesses the routing information for Q_ℓ to deliver m to all members of $D \subseteq Q_i$ and the protocol terminates successfully. In either case, the induction holds.

We now analyze the costs of our protocol. In the first step, even in the event that a share corruption attack occurs at most one round-trip round of communication

occurs (between p and Q_1). For steps $i = 2, \dots, \ell - 1$, if a share corruption attack occurs, at most two round-trip rounds of message exchange occur: (1) p sends to Q_i and Q_i sends back to p and (2) p transmits shares to Q_i who then send the correct shares back to p . Adding the last step, the latency is $2 \cdot (\ell - 2) + 2$. In terms of message complexity, in the first round, peer p must send a request to and receive a response from each peer in Q_1 ; this totals at most $2s$ messages. For steps $i = 2, \dots, \ell - 1$ peer p must both send a request to and receive a response from each peer in a quorum; if a corruption attack occurs, p must send another message to each peer in a quorum (with all signed shares collected together) and receive back a response. Therefore, this incurs at most $4 \cdot s$ messages. In the last step p sends to all members of D . Hence, the message complexity is at most $4 \cdot s \cdot (\ell - 2) + |D| + 2s$. For every other involved peer $q \notin D$, the message complexity for q is simply at most 4; clearly, peers in D receive one message. \square

Spamming Attacks: The sending peer p experiences more cost than other participating peers. In part, this is due to the iterative nature of the protocol; however, largely this is because p must send and receive $O(s)$ messages per step. In contrast, other participating peers need only send and receive a constant number of messages over the execution of the protocol.

Peer p may misbehave in other ways. For instance, p may repeatedly contact its quorum to initiate robust communication; however, eventually all correct peers will ignore p . Similarly, using a correct signature, p may repeatedly ask q in another quorum for proof and/or routing information; however, time stamps limit such replay attacks. In conclusion, such actions cannot cause correct peers to perform expensive operations.

4.2. Robust Communication Protocol II

We present RCP-II which is randomized yielding. It is sufficient for each node to possess its own internal random number generator. a small expected message complexity for both the sending peer and non-sending peers. In exchange, join and leave operations incur additional cost in comparison to RCP-I; we discuss this in Section 4.3.

RCP-II utilizes signed routing table information. Each entry of a routing table has the form $[u_{\text{ID}}, u'_{\text{ID}}, u_{\text{addr}}, K_{Q_i}, ts_r]$, where u is a peer, u' is the peer with the next largest identifier, K_{Q_i} is the quorum public key of Q_i to which u belongs, and ts_r is a time stamp for when this entry was created. \mathcal{RT}_{Q_j} denotes the routing table information for all peers in Q_j . $[K_{Q_j}]_{k_{Q_i}}$ is the quorum public key of Q_j signed

using the private quorum key of Q_i ; recall, neighbors in the quorum topology know each others' public key. $[\mathcal{RT}_{Q_j}]_{k_{Q_i}}$ is the routing information signed with the private key of Q_i ; entries of the routing table are signed separately. Routing table information is time stamped and re-signed periodically when DKG is executed.

Overview: We sketch RCP-II here. For simplicity, we temporarily assume that peers act correctly; our pseudocode in Figure 4 is complete for when peers fail to respond to requests by p . Initially, each correct peer in Q_1 receives $[p_{\text{ID}}|p_{\text{addr}}|k_{\text{eY}}|ts]$ from p . The time stamp ts is chosen by p and peers in Q_1 will acquiesce to the value if it agrees with the rule set to within some bound to compensate for clock drift. If the request does not violate the rule set, then the information is signed allowing p to form $M_1 = [p_{\text{ID}}|p_{\text{addr}}|k_{\text{eY}}|ts]_{k_{Q_1}}$.

In the second step of the protocol, p knows the membership of Q_2 and selects a peer $q_2 \in Q_2$ uniformly at random (u.a.r.) without replacement. Peer p then sends M_1 to q_2 . Assuming q_2 is correct, it verifies M_1 using K_{Q_1} and checks that the ts is valid; the duration for which a time stamp is valid would be specified by the rule set. Once verified q_2 sends p the information $[K_{Q_1}]_{k_{Q_2}}$, $[\mathcal{RT}_{Q_3}]_{k_{Q_2}}$ and $[K_{Q_3}]_{k_{Q_2}}$. Peer p knows K_{Q_2} since Q_1 links to Q_2 and verifies $[K_{Q_1}]_{k_{Q_2}}$, $[\mathcal{RT}_{Q_3}]_{k_{Q_2}}$ and $[K_{Q_3}]_{k_{Q_2}}$, and checks that the time stamp on the routing information is valid. If so, p constructs $M_2 = [M_1|[K_{Q_1}]_{k_{Q_2}}]$. Here $[K_{Q_1}]_{k_{Q_2}}$ will allow some peer in Q_3 to verify K_{Q_1} and M_1 , while the signed verified K_{Q_3} will allow p to check the response from that peer in Q_3 .

This process repeats with minor changes for the remaining steps. Using \mathcal{RT}_{Q_3} from the previous step, p selects a peer q_3 randomly from Q_3 and sends M_2 . Since Q_3 is linked with Q_2 in the quorum topology, q_3 knows K_{Q_2} , which it uses to verify $[K_{Q_1}]_{k_{Q_2}}$; this allows q_3 to verify M_1 signed with k_{Q_1} . Peer q_3 then confirms that ts is valid and sends $[K_{Q_2}]_{k_{Q_3}}$, $[\mathcal{RT}_{Q_4}]_{k_{Q_3}}$ and $[K_{Q_4}]_{k_{Q_3}}$ to p . Peer p has a verified public key K_{Q_3} from the previous step and uses it to verify $[K_{Q_2}]_{k_{Q_3}}$, $[\mathcal{RT}_{Q_4}]_{k_{Q_3}}$, and $[K_{Q_4}]_{k_{Q_3}}$. Then p constructs $M_3 = [M_2|[K_{Q_2}]_{k_{Q_3}}] = [M_1|[K_{Q_1}]_{k_{Q_2}}|[K_{Q_2}]_{k_{Q_3}}]$. This process continues until m is delivered. Figure 4 gives the pseudocode for RCP-II. Every peer contacted by p verifies a chain of certificates, which can be converted into a single signature using the concept of aggregate signatures [8].

It is possible that p chooses a Byzantine peer that may not respond. In that case, after some appropriate time interval, p will select an additional peer in the quorum. Let X be a random variable denoting the time required for a correct peer to respond. We make a weak

RCP-II: SENDING PEER p

Initial Step:

- 1: p sends the following to each peer $q \in Q_1$:
 $[p_{\text{ID}}|p_{\text{addr}}|k_{\text{eY}}|ts]$
- 2: p gathers all responses and constructs:

$$M_1 \leftarrow [p_{\text{ID}}|p_{\text{addr}}|k_{\text{eY}}|ts]_{k_{Q_1}}$$

Intermediate Steps:

- 3: **for** $i = 2$ to $\ell - 1$ **do**
- 4: **while** (p does not have M_i and has waited time Δ since previous selection) **do**
- 5: p sends M_{i-1} to $q \in Q_i$ selected u.a.r. without replacement.
- 6: **if** ($[K_{Q_{i-1}}]_{k_{Q_i}}$, $[\mathcal{RT}_{Q_{i+1}}]_{k_{Q_i}}$ and $[K_{Q_{i+1}}]_{k_{Q_i}}$ are received from any peer in Q_i previously selected) **then**
- 7: p uses K_{Q_i} to verify $K_{Q_{i+1}}$, $\mathcal{RT}_{Q_{i+1}}$ and $K_{Q_{i-1}}$.
- 8: **if** ($K_{Q_{i+1}}$, $\mathcal{RT}_{Q_{i+1}}$ and $K_{Q_{i-1}}$ are all verified) **then**
- 9: $M_i \leftarrow [M_{i-1}][K_{Q_{i-1}}]_{k_{Q_i}}$

Final Step:

- 10: p sends $M_{\ell-1}$ to $D \subseteq Q_\ell$ along with m .

RCP-II: RECEIVING PEER q

Initial Step:

- 1: **if** ($q \in Q_1$ receives $[p_{\text{ID}}|p_{\text{addr}}|k_{\text{eY}}|ts]$ from $p \in Q_1$) **then**
- 2: q checks that p 's request is legitimate and, if so, sends its signature share.

Intermediate Steps:

- 3: **if** ($q \in Q_i$ receives M_{i-1} from p) **then**
- 4: **for** $j = i - 1$ downto 1 **do**
- 5: q uses K_{Q_j} to verify $K_{Q_{j-1}}$.
- 6: Peer q uses K_{Q_1} to verify M_1 .
- 7: **if** verification is successful **then**
- 8: q sends $[K_{Q_{i-1}}]_{k_{Q_i}}$, $[\mathcal{RT}_{Q_{i+1}}]_{k_{Q_i}}$ and $[K_{Q_{i+1}}]_{k_{Q_i}}$ to p .

Figure 4. Pseudocode for RCP-II

assumption that $Pr[X \leq \Delta] \geq c$ where Δ is any duration of time and $c > 0$ is any constant probability. This does not circumscribe a particular distribution for response times; in fact, *any* distribution suffices, including the Poisson, exponential, and gamma distributions previously used to characterize round trip time (RTT) over the Internet. In practice, a peer p would set its own Δ value by sampling the network using methods for estimating RTT [27]. Since there are only a constant fraction of Byzantine peers, taking the median from a sufficiently large sample will allow p to determine Δ . As p receives a response from any of the

previously selected peers in Q_i , this is in accordance with the weak synchrony assumption in Section 3.

Lemma 2. *RCP- II guarantees that m is transmitted to a target set of peers $D \subseteq Q_i$ for some quorum Q_i over a path of length ℓ with the following properties:*

- Both the total message complexity and the message complexity of the sending peer is each at most $2 \cdot s + \frac{(\ell-2)}{(1-\epsilon) \cdot c} + (\ell - 2) + |D|$.
- Each non-sending peer has expected message complexity at most $\frac{2}{(1-\epsilon) \cdot c \cdot s}$.
- The expected latency is at most $\frac{(\ell-2)}{(1-\epsilon) \cdot c} + 2$.

Proof: First we prove the correctness of our protocol and, as before, we show that if p is correct and has not violated the rule set, at each step i of the protocol p either (1) establishes a valid M_i and receives the routing information for the next hop or (2) terminates the protocol by delivering m to all members of D . Our proof is by induction on i .

Base Case: Consider the initial step $i = 1$ where p communicates with the peers in its quorum $Q_p = Q_1$ about sending the message m . If p is correct and has not violated the rule set, upon receiving $[p_{\text{id}} | p_{\text{addr}} | \text{key} | ts]$ all correct peers will send their shares to p . Therefore, p is guaranteed to obtain M_1 by the majority invariant. Peer p already has the routing information for Q_2 ; therefore, the base case holds.

Inductive Hypothesis: Assume that at step $i < \ell$, p has obtained a correct M_i and routing information for Q_{i+1} .

Inductive Step: First assume that $i = \ell - 1$. Then, by the induction hypothesis, peer p possesses $M_{\ell-1}$ and the necessary routing information to send this signature and message m to $D \subseteq Q_{i+1}$; thus the protocol terminates correctly. Otherwise, assume $i < \ell - 1$; we consider step $i + 1$. Peer p sends M_i to a peer $q \in Q_{i+1}$ selected uniformly at random without replacement. By the inductive hypothesis, the contents of M_i are valid and p possesses the necessary routing information. If q is a Byzantine peer, then p 's request can fail and p can detect an invalid response using $K_{Q_{i+1}}$ obtained from the previous step. It is also possible that q is a correct but slow node and do not respond in a predefined time period. In this case, p re-issues its request to another randomly selected peer in Q_{i+1} ; eventually, one of selected correct peers will respond by sending $[K_{Q_i}]_{k_{Q_{i+1}}}$, $[\mathcal{RT}_{Q_{i+2}}]_{k_{Q_{i+1}}}$ and $[K_{Q_{i+2}}]_{k_{Q_{i+1}}}$ to p . Peer p will verify this information and create a valid M_{i+1} . Therefore, at this point p

possesses a correct M_{i+1} and routing information for Q_{i+2} ; therefore, the induction holds.

Since RCP-II is a randomized algorithm, our costs are given in expectation. We assume the following: let X_i be a random variable denoting the time required for the i^{th} correct peer (note we condition on correctness) selected u.a.r without replacement by p to respond to p 's request. We assume that $Pr[X_i \leq \Delta] = c$ where $c > 0$ is some constant probability.

We now calculate loose upper bounds of the expected resource costs. In the first step, in communicating with Q_1 , peer p handles at most $2 \cdot s$ messages and the round-trip latency is 1. Then for each step $i = 2, \dots, \ell - 1$, let Y_i be the random variable with value 1 if the i^{th} peer is both correct and responds within time Δ ; 0 otherwise. Then $Pr[Y_i = 1] \leq (1 - \epsilon) \cdot c$; for simplicity, set $\rho = (1 - \epsilon) \cdot c$ to be this probability of success. Let $Y = \sum_{i=1}^s Y_i$. The expected number of selections $E[Y]$ before p receives a response from a correct peer is at most:

$$\sum_{k=0}^s (1-\rho)^k \cdot \rho \cdot (k+1) = \rho \left(\sum_{k=0}^s (1-\rho)^k \cdot k + \sum_{k=0}^s (1-\rho)^k \right)$$

where the first term is upper-bounded by the well-known telescoping series and the second is simply a geometric series. Therefore $E[Y] \leq \frac{1}{(1-\epsilon) \cdot c}$ and including the last step, the expected latency is at most $\frac{\ell-2}{(1-\epsilon) \cdot c} + 2$. The ℓ^{th} step requires D messages and one hop. In terms of expected message complexity, since each step requires at most 2 messages and the last step requires $|D|$ messages, we can give a crude upper bound of $2s + \frac{2}{(1-\epsilon) \cdot c} \cdot (\ell - 2) + |D|$. However, note that once p hears back from a node, any message from any other previously selected nodes in the current step can be easily ignored/filtered. Therefore, per step, p handles $\frac{1}{(1-\epsilon) \cdot c} + 1$ messages. We can now give a more accurate upper bound of $2s + \frac{\ell-2}{(1-\epsilon) \cdot c} + (\ell - 2) + |D|$. Finally, while latency is measured in the number of communication rounds, we note that the expected duration of time required for each intermediate round is $\frac{\Delta}{(1-\epsilon) \cdot c}$.

In terms of the expected message complexity of a non-sending peer $q \notin D$ in a quorum along the lookup path, a correct peer chosen by p receives one message and sends one message. The probability that q is chosen is at most $1/((1-\epsilon) \cdot s)$; therefore the expected message complexity for q is at most $2/((1-\epsilon) \cdot s)$. \square

While latency is measured in communication rounds, the time for executing RCP-II depends on Δ and we discuss this briefly. Accounting for the

response time incurred in the intermediate steps, p waits for at most time $\frac{\Delta}{(1-\epsilon)\cdot c}$ per step in expectation as shown in Lemma 2. Since peer p will have knowledge of the response time distribution, p may optimize performance by selecting Δ so that $\frac{\Delta}{c}$ is minimized.

Spamming Attacks: Due to the iterative nature of RCP-II, p sends more messages than other participating peers, but not to the degree seen in RCP-I. Instead of making it expensive for p to perform robust communication, RCP-II uses the following two properties to deter spamming: (1) it is inexpensive for a correct peer to detect spam and (2) the congestion suffered by a correct peer is low since the number of messages is not magnified by the use of quorums.

To address our first point, p may launch as many robust communication operations as the rule set allows; p may even try to circumvent the rule set by directly sending to a correct peer q ; however, it is inexpensive for q to verify that the proof being sent is invalid. The operation terminates at that point since q will not reply. In contrast to the passive scenario of Section 3.5, q need not keep a history to judge the legitimacy of a request; it simply verifies the accompanying certificate.

Our second point, and a key difference between RCP-I and RCP-II, is that with RCP-II an operation incurs only expected $O(\ell)$ messages which compares favourably to a system *without a quorum topology*. Therefore, the congestion caused by such requests is not significantly magnified by the use of quorums which was a key concern regarding spamming.

Adversarial peers may misbehave in other ways with many of the same consequences and remedies as discussed in RCP-I. Even with a generous upper bound on the expiration of ts , the congestion p can cause with a replay attack is again limited since only p can use the certificate. A notable attack, unique to RCP-II, occurs when a faulty peer gives p stale routing table information. Since entries are signed and time stamped, we are guaranteed that in the fairly recent past, the location indicated by the stale information was indeed correct. This fact, coupled with the standard assumption that ID collisions do not occur, guarantees that the adversary *cannot* engineer a situation where requests are forwarded to a faulty peer. Consequently, the impact of this attack is limited. The search path may be slightly lengthened by forwarding to an older location. Alternatively, stale information may point to a peer that no longer exists or is not the correct recipient, which forces p to backtrack one hop. These cases are easy to handle, but for ease of exposition, they are not treated in our pseudocode in Figure 4. In short, routing integrity is not compromised. The fact that

routing tables can be signed periodically every several minutes without significant CPU cost (see Section 5) implies that the impact of such an attack is negligible.

4.3. The Join Protocol and Membership Updates

For the sake of being self-contained, we describe how a peer would join our system. This first involves a discussion of a result by Awerbuch and Scheideler [5] which allows a DHT to be robust even if the number of join and leave events is polynomial in the size of the network n . More precisely, in each *round* the adversary may opt to insert a Byzantine peer (assuming the total number of Byzantine peers in the system does not exceed the allotted amount) or remove a Byzantine peer from the system. The key protocol in defending against is the *cuckoo-rule*. We assume the identifier space of the DHT is normalized to be $[0,1)$. For any interval $I \subset [0,1)$, the cuckoo rule maintains two invariants. The first is the *balancing invariant* which guarantees that I contains $\Theta(|I| \cdot n)$ peers. The second is the *majority invariant* which guarantees the majority of peers in I are correct. The authors show that for $|I| = \Theta(\log(n)/n)$ both invariants can be maintained with high probability over n^c join and leave operations, where c is a constant that can be tuned according to the parameters of the protocol. It follows that the peers in I can be used to form a quorum.

It is important to understand the resource costs of the cuckoo rule which functions as follows. The ring $[0,1)$ is assumed to be broken into disjoint segments of constant length k/n for some constant k . Each segment is called a *k-region* and $R_k(x)$ denotes the unique *k-region* containing x . When a peer p joins the network, it is assigned a random identifier $x \in [0,1)$ and placed in this location on the ring. All nodes in $R_k(x)$ are evicted from their locations and placed into new locations chosen uniformly and independently at random from $[0,1)$. Figure 5 illustrates the operations performed by the cuckoo rule.

The node placements required by the cuckoo rule can be executed by having quorums use robust communication in order to inform each other about the arrival of the evicted nodes at their new locations. Once a quorum Q_i knows about the presence of a recently evicted node q , all correct members of Q_i update their membership lists, share IP addresses, and aid q in setting up any required links (i.e. such as finger links in Chord). A detailed discussion of how this could be done is presented in [17] and random numbers can be generated using the protocol of [4]. We finish our discussion of a join protocol by discussing the steps

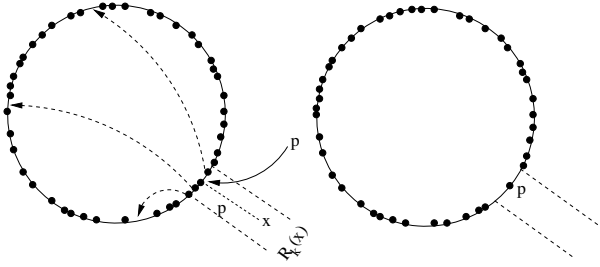


Figure 5. An illustration of the cuckoo rule. (Left) Peer p is placed in its random location x . All peers in the k -region, $R_k(x)$ denoted by dashed lines, are assigned random locations in $[0, 1)$. (Right) After the cuckoo rule is executed, peer p is the only peer in $R_k(x)$.

necessary for maintaining DKG and the consequent cost of membership changes under both protocols:

RCP-I: Consider a quorum Q_i to which a new peer is added. The membership update protocol of DKG [31] is executed to redistribute the shares of the public/private quorum key pair over all members of Q_i . In the process, the individual public/private key shares are also updated. Notably, *no other quorums are affected by this process* as the quorum key pair remains the same and the individual key shares need only be known to members of Q_i . When a peer leaves Q_i , the departure can be treated as a crash and so long as the number of crashes does not exceed the crash-limit f , the DKG (share renewal) protocol need not be executed. We use this to associate the system churn rate to DKG session time. Note that the adversary may crash some of its t nodes, and in principle, the system can handle $t + f$ node leaves. However, we cannot associate these additional t crashes with the system churn due to the inherent arbitrary nature of Byzantine peers.

RCP-II: When a peer q joins Q_i , the DKG protocol needs to be executed as in the case of RCP-I; however, there are additional costs due to the need to update and re-sign the routing table information. In particular, not only do the peers in Q_i need to update and have signed their routing table information to reflect the addition of q , all quorums to which Q_i is linked under the quorum topology also need to update and re-sign their routing table information; note that this *does not require any revocation* since the public key does not change. Therefore, a join event under this scheme *does affect other quorums*. When a peer leaves Q_i , DKG may be required as in the case of RCP-I. However,

routing table information for Q_i and the quorums to which it links must again update and re-sign their routing table information. Therefore, while RCP-II reduces message complexity, the cost of join/leave operations is higher in comparison to RCP-I.

5. Experimental Results

We examine the performance of DKG and our two protocols on the PlanetLab platform [40]. Based on our experimental results and known churn rates, we propose parameters for DHTs using our protocols.

5.1. Implementation and Microbenchmarks

The DKG protocol is a crucial component of our protocols. It is required to initiate a threshold signature system in a quorum and to securely manage membership changes. We use a C++ implementation [30] to measure the performance of DKG. We incorporate threshold BLS signatures into this implementation and realize our two protocols using this setup on PlanetLab.

Distributed Key Generation: We test the DKG implementation for quorum sizes $s = 10, 15, 20, 25, 30$ and present median completion times and median CPU usage for our experiments in Table 1. The median completion periods vary from 6 seconds for $s = 10$ to more than 5 minutes for $s = 30$. The bulk of this latency is due to network delays; in contrast, the required CPU time is much smaller than the completion periods.

In the next subsection, we examine the feasibility of these completion periods. Our experiments with DKG assume that 30% of the peers may crash and 10% of the peers may be Byzantine. While we can tolerate any fraction of Byzantine peers less than $1/3$, we use these numbers since in many practical scenarios we expect the fraction of Byzantine faults to be less than 10% and modest compared to the fraction of crash failures.

RCP-I and RCP-II: For our RCP-I and RCP-II experiments, we set $s = 30$, $t = 3$, and $f = 10$. In RCP-I, a node requires 0.14 seconds on average to obtain a threshold signature from a quorum, if all of the obtained signature shares are correct. The average execution time increases to 0.23 seconds in case of a share corruption attack. Extrapolating to a path length ℓ , an operation should take $0.14 \cdot \ell$ to $0.23 \cdot \ell$ seconds on average. For a DHT with 10^5 nodes, the average total time for RCP-I is then 3 to 4 seconds with $\ell = 20$.

In RCP-II, a node takes 0.04 seconds on average to obtain the required signed public keys and the signed routing information from a correct peer. An individual signature verification takes 0.004 seconds

Table 1. Median values of DKG completion time and CPU time per node for various s values.

s	t	f	Completion Time (sec)	CPU Seconds/Node
10	1	3	5.73	0.76
15	2	4	18.0	1.94
20	2	6	68.0	2.55
25	3	7	290.9	6.13
30	3	10	336.7	7.27

on average. The median latency value over Planetlab is approximately equal to 0.08 seconds. [14] That is, $\Delta = 0.08$ seconds for $c = 0.5$. With a chain of signed public keys of length ℓ , the total communication time is approximately $0.14 + 0.04 \cdot (\ell - 1) + \frac{\Delta \cdot (\ell - 2)}{c \cdot (1 - \epsilon)} + 0.004 \cdot \frac{\ell(\ell - 1)}{2}$ which for 10% Byzantine peers, is 4.68 seconds in expectation. To a first-approximation, the execution times of our protocols seem quite reasonable.

System Load: We address the issue of system load under the assumption that signature verification is the most significant computational operation. We make back-of-the-envelope calculations to obtain the expected order of magnitude for our performance figures. For RCP-I, from the above discussion, each signature verification takes 0.004 seconds; thus, the total CPU time required per execution is $0.004 \cdot \ell \cdot (1 + s + s^2)$; this includes the costs due to share corruption attacks. For $\ell = 20$ and $s = 30$, this value is 75 CPU seconds, spread out over 600 nodes. Therefore, the number of executions of RCP-I that can be started per second on average is $n/75 \approx 10^3$ when $n = 10^5$. Note this rate value is for *the entire system*. Now, if no share corruption attacks occur, the total CPU time required per execution becomes $0.004 \cdot \ell \cdot (1 + s)$ which, for the same parameter values, is 2.5 CPU seconds. This implies that $4 \cdot 10^4$ executions can be started per second on average in the entire system. For RCP-II, the total CPU time required for execution is given by $0.004 \cdot \left(\ell + \frac{(\ell - 1) \cdot \ell}{2 \cdot (1 - \epsilon)} \right)$ which, for the same parameters and $\epsilon = 1/10$ is roughly 1 CPU second on average. Therefore, approximately 10^5 executions can be started per second on average in the entire system.

5.2. Analysis and Discussion

As mentioned in Section 3.1, important questions remain with regards to translating theoretical results to a practical setting. In particular, two quantities of interest are the size of quorums, s , and the number of quorums to which each peer belongs, n_Q . Unfortunately, pinning down these quantities is non-trivial. Only asymptotic analysis is present in the literature. Furthermore, it is not a simple case of substituting hard numbers because s depends on a number of parameters: (1) the exact

Table 2. The expected number of seconds before a quorum experiences a membership change (r_Q).

s	10			15			20		
n_Q	1	2	3	1	2	3	1	2	3
r_Q	526	351	175	350	234	117	263	132	88
	25			30					
	1	2	3	1	2	3			
	210	140	70	175	87	58			

guarantees being made (i.e. time until quorum failure), (2) algorithms for quorum maintenance, (3) the tools of analysis (i.e. form of Chernoff bounds used) and many more. Evaluating these parameters is outside the scope of this work. Instead, we assume a range of values for s and n_Q . As our protocols appear to be the most efficient to date, the following results illuminate what currently seems possible in practice.

System Churn and DKG: The performance of our two protocols will likely depend on system churn. A common metric for measuring the degree of churn is *session time*: the time between when a node joins the network and when it leaves [42]. As discussed in Section 3.5, we make the standard assumption that the cost of joining the network can be made large enough so as to prevent the adversary from substantially increasing the rate of churn through rapid rejoin operations.

Part I - An Argument for Batching: Investigations have yielded differing measurements for median session times. The Kazaa system was found to have a median session time of 144 seconds [20]. In the Gnutella and Napster networks, the median session time was measured to be approximately 60 minutes [47]. Measurements of the Skype P2P network yielded a median session time of 5.5 hours for super-peers [26]. Here, we temporarily assume a median session time of 60 minutes and a standard Poisson model of peer arrivals/departures as in [33], [42]. To calculate churn rate, r (number of arrivals/departures per second), based on the median session time t_{med} (in seconds), we use the formula of [42]: $r = (n \cdot \ln 2) / t_{med}$. For $n = 10^5$ and $t_{med} = 3600$ seconds, $r \approx 19$. Assuming that join and leave events occur independently of each other, Table 2 gives the expected number of seconds, r_Q , at which point a quorum will undergo a membership change when each peer belongs to n_Q quorums. Our choice of $n_Q \leq 3$ is based upon the reasonable assumption that overlap occurs only with immediate neighboring quorums in the ID space.

In several cases, the r_Q values are less than the corresponding median DKG completion times in Table 1. Therefore, a quorum may not be able to execute DKG often enough to accommodate each membership change. However, join operations can be queued and

performed in batches. Executing DKG for a batch of joins does not increase the message complexity and message size increases only linearly in the batch size (see [31, Sec. 6]). Therefore, batching can mitigate the effects of churn and it seems plausible that peers would tolerate some delay in joining in exchange for security.

Part II - Batching and the Security Threshold:

Batching join events improves performance; however, many peers might leave a quorum before a new batch is added, thus violating the security threshold. Hence, we are interested in the session time value required such that this is not likely to occur. Based on Table 1 for $s = 20$ and $n_Q = 1$, DKG completes within 68 seconds. The number of leave events a quorum can suffer while not exceeding the crash limit is $f = 6$. If Byzantine peers leave, more crashes are tolerable; however, identifying such events is impossible, so we assume the worst case of $f = 6$. Assuming DKG executes every $r_{DKG} = 1200$ seconds, we seek the median session time such that at most 6 peers leave the system within 1268 seconds. With $n/s = 5000$ quorums in the system, each experiencing 6 leave events within 1268 seconds, the system churn rate is $r = 23.7$. This gives $t_{med} = 2930$ or, equivalently, 49 minutes. Therefore, with this t_{med} , we expect the system to remain secure. Moreover, a quorum only spends $68/1268 = 5.4\%$ of the time executing DKG.

Certain parameters can be tuned to offer performance trade-offs. Decreasing r_{DKG} yields smaller required median session times; however, the percentage of time spent on DKG increases. Such tuning would depend on the desired system performance, the application, and s and n_Q . Table 3 gives session time calculations for other values of s , r_{DKG} and n_Q .

Required session times increase with s . Notably, for $s = 30$ and $n_Q = 1$, t_{med} does not far exceed the 60 minutes in [47]. As n_Q increases, the required session times grow linearly. However, our maximum of 3.7 hours is still less than t_{med} measured for super-peers in the Skype network [26]. We tentatively conclude that our protocols can be deployed in applications where session times range from 10 minutes to a few hours and that such applications currently exist.

5.2.1. Final Comments. The implementation of the DKG protocol used here is an academic version, and more efficient implementations may be possible. In terms of performance improvements, using the aggregate signature scheme [8], messages can be made more compact yielding a savings in RCP-II. For certain applications, it may be possible to restrict membership to those peers that meet certain latency and bandwidth criteria. Our choice of $n_Q \leq 3$ is based upon

Table 3. Median session times (in hours) derived from values for s , n_Q and r_{DKG} (in hours).

s	10			15			20		
r_{DKG}	0.167			0.25			0.33		
n_Q	1	2	3	1	2	3	1	2	3
t_{med}	0.19	0.39	0.58	0.66	1.32	2.00	0.81	1.62	2.44
	25			30					
	0.42			0.5					
	1	2	3	1	2	3			
	1.23	2.46	3.70	1.23	2.47	3.70			

the assumption that quorums only overlap with their immediate neighbors in the ID space. In terms of fault-tolerance, our experiments were performed with 10% of nodes suffering Byzantine faults; however, we generally expect the fraction of Byzantine nodes to be less, thus reducing execution times.

Future Work: The performance of a complete system is an important open question. The quorum topology chosen is crucial and optimizing this in practice is a topic of future work. Another pertinent question is whether our protocols can be modified to achieve recursive routing since this could result in lower message delivery time. While we focus on DHTs, our results may apply to other P2P designs and more general settings where groups of machines, some with untrustworthy members, must communicate; it would be of interest to identify such applications.

References

- [1] M. Abd-El-Malek, G. R. Ganger, G. R. Goodson, M. K. Reiter, and J. J. Wylie. Fault-Scalable Byzantine Fault-Tolerant Services. In *SOSP*, pages 59–74, 2005.
- [2] A. Adya, W. Bolosky, M. Castro, G. Cermak, R. Chaiken, J. Douceur, J. Howell, J. Lorch, M. Theimer, and R. Wattenhofer. FARSITE: Federated, Available, and Reliable Storage for an Incompletely Trusted environment. In *OSDI*, pages 1–14, 2002.
- [3] J. Aspnes, N. Rustagi, and J. Saia. Worm versus alert: Who Wins in a Battle for Control of a Large-Scale Network? In *OPODIS*, pages 443–456, 2007.
- [4] B. Awerbuch and C. Scheideler. Robust Random Number Generation for Peer-to-Peer Systems. In *OPODIS*, pages 275–289, 2006.
- [5] B. Awerbuch and C. Scheideler. Towards a Scalable and Robust DHT. In *SPAA*, pages 318–327, 2006.
- [6] B. Awerbuch and C. Scheideler. Towards Scalable and Robust Overlay Networks. In *IPTPS*, 2007.
- [7] A. Boldyreva. Threshold Signatures, Multisignatures and Blind Signatures Based on the Gap-Diffie-Hellman-Group Signature Scheme. In *PKC*, pages 31–46, 2003.
- [8] D. Boneh, C. Gentry, B. Lynn, and H. Shacham. Aggregate and Verifiably Encrypted Signatures from Bilinear Maps. In *EUROCRYPT*, pages 416–432, 2003.
- [9] E. Bortnikov, M. Gurevich, I. Keidar, G. Kliot, and A. Shraer. Brahms: Byzantine Resilient Random Membership Sampling. In *PODC*, pages 145–154, 2008.

- [10] C. Cachin and J. Poritz. Secure Intrusion-tolerant Replication on the Internet. In *DSN*, pages 167–176, 2002.
- [11] M. Castro and B. Liskov. Byzantine Fault Tolerance Can Be Fast. In *DSN*, pages 513–518, 2001.
- [12] M. Castro and B. Liskov. Practical Byzantine Fault Tolerance and Proactive Recovery. *ACM Transactions on Computer Systems*, 20(4):398–461, 2002.
- [13] J. Cowling, D. Myers, B. Liskov, R. Rodrigues, and L. Shrira. HQ Replication: A Hybrid Quorum Protocol for Byzantine Fault Tolerance. In *OSDI*, pages 177–190, 1999.
- [14] F. Dabek, J. Li, E. Sit, J. Robertson, M. F. Kaashoek, and R. Morris. Designing a DHT for Low Latency and High Throughput. In *NSDI*, pages 85–98, 2004.
- [15] J. Douceur. The Sybil Attack. In *IPTPS*, 2002.
- [16] J. Falkner, M. Piatek, J. P. John, A. Krishnamurthy, and T. Anderson. Profiling a Million User DHT. In *IMC*, pages 129 – 134, 2007.
- [17] A. Fiat, J. Saia, and M. Young. Making Chord Robust to Byzantine Attacks. In *ESA*, pages 803–814, 2005.
- [18] R. Geambasu, T. Kohno, A. A. Levy, and H. M. Levy. Vanish: Increasing Data Privacy with Self-Destructing Data. In *USENIX Security Symp.*, pages 299–315, 2009.
- [19] R. Gennaro, S. Jarecki, H. Krawczyk, and T. Rabin. Robust Threshold DSS Signatures. In *Advances in Cryptology - EUROCRYPT*, pages 354–371, 1996.
- [20] P. K. Gummadi, R. J. Dunn, S. Saroiu, S. D. Gribble, H. M. Levy, and J. Zahorjan. Measurement, Modeling, and Analysis of a Peer-to-Peer File-Sharing Workload. In *SOSP*, pages 314–329, 2003.
- [21] A. Gupta, B. Liskov, and R. Rodrigues. One Hop Lookups for Peer-to-Peer Overlays. In *9th Conf. on Hot Topics in Operating Systems*, pages 2–2, 2003.
- [22] K. Hildrum and J. Kubiatowicz. Asymptotically Efficient Approaches to Fault-Tolerance in Peer-to-peer Networks. In *DISC*, pages 321–336, 2004.
- [23] D. Boneh, B. Lynn, and H. Shacham. Short Signatures from the Weil Pairing. In *ASIACRYPT*, pages 514–532, 2001.
- [24] M. Castro, P. Druschel, A. Ganesh, A. Rowstron, and D. Wallach. Secure Routing for Structured Peer-to-Peer Overlay Networks. In *OSDI*, pages 299–314, 2002.
- [25] P. Feldman. A Practical Scheme for Non-Interactive Verifiable Secret Sharing. In *FOCS*, pages 427–437, 1987.
- [26] S. Guha, N. Daswani, and R. Jain. An Experimental Study of the Skype Peer-to-Peer VoIP System. In *IPTPS*, 2006.
- [27] H. Jiang and C. Dovrolis. Passive Estimation of TCP Round-Trip Times. *CCR*, 32:75–88, 2002.
- [28] H. Johansen, A. Allavena, and R. van Renesse. Fireflies: Scalable Support for Intrusion-Tolerant Network Overlays. In *OSR*, pages 3–13, 2006.
- [29] A. Kapadia and N. Triandopoulos. Halo: High-Assurance Locate for Distributed Hash Tables. In *NDSS*, 2008.
- [30] A. Kate and I. Goldberg. Asynchronous Distributed Private-Key Generators for Identity-Based Cryptography. Cryptology ePrint Archive, Report 355, 2009.
- [31] A. Kate and I. Goldberg. Distributed Key Generation for the Internet. In *ICDCS*, pages 119–128, 2009.
- [32] J. Kubiatowicz, D. Bindel, Y. Chen, S. Czerwinski, P. Eaton, D. Geels, R. Gummadi, S. Rhea, W. Weimer, H. Weatherspoon, C. Wells, and B. Zhao. OceanStore: An Architecture for Global-Scale Persistent Storage. In *ASPLOS*, pages 190–201, 2000.
- [33] D. Liben-Nowell, H. Balakrishnan, and D. Karger. Analysis of the Evolution of Peer-to-Peer Systems. In *PODC*, pages 233–242, 2002.
- [34] P. Mittal and N. Borisov. ShadowWalker: Peer-to-peer Anonymous Communication using Redundant Structured Topologies. In *CCS*, 2009.
- [35] Michael Mitzenmacher. *The Power of Two Choices in Randomized Load Balancing*. PhD thesis, University of California at Berkeley, 1996.
- [36] A. Nambiar and M. Wright. Salsa: A Structured Approach to Large-Scale Anonymity. In *CCS*, pages 17–26, 2006.
- [37] M. Naor and U. Wieder. A Simple Fault Tolerant Distributed Hash Table. In *IPTPS*, pages 88–97, 2003.
- [38] R. Ostrovsky and M. Yung. How to Withstand Mobile Virus Attacks (Ext. Abstract). In *PODC’91*, pages 51–59, 1991.
- [39] T. P. Pedersen. Non-Interactive and Information-Theoretic Secure Verifiable Secret Sharing. In *CRYPTO*, pages 129–140, 1991.
- [40] L. Peterson, T. Anderson, D. Culler, and T. Roscoe. A Blueprint for Introducing Disruptive Technology into the Internet. *CCR*, 33(1):59–64, 2003.
- [41] M. K. Reiter. The Rampart Toolkit for Building High-Integrity Services. In *Intl. Workshop on Theory and Practice in Distributed Systems*, pages 99–110, 1995.
- [42] S. Rhea, D. Geels, T. Roscoe, and J. Kubiatowicz. Handling Churn in a DHT. In *USENIX Annual Technical Conf.*, pages 127–140, 2004.
- [43] R. Rodrigues, P. Kouznetsov, and B. Bhattacharjee. Large-Scale Byzantine Fault Tolerance: Safe but Not Always Live. In *HotDep*, 2007.
- [44] R. Rodrigues and B. Liskov. Rosebud: A Scalable Byzantine-Fault-Tolerant Storage Architecture. Technical Report TR/932, MIT LCS, December 2003.
- [45] R. Rodrigues, B. Liskov, and L. Shrira. The Design of a Robust Peer-to-Peer System. In *10th ACM SIGOPS European Workshop*, page 2002, 117-124.
- [46] J. Saia and M. Young. Reducing Communication Costs in Robust Peer-to-Peer Networks. *Information Processing Letters*, 106(4):152–158, 2008.
- [47] S. Saroiu, P. K. Gummadi, and S. D. Gribble. A Measurement Study of Peer-to-Peer File Sharing Systems. In *MMCN*, pages 314–329, 2002.
- [48] F. B. Schneider. Implementing Fault-Tolerant Services Using the State Machine Approach: A Tutorial. *ACM Computing Surveys*, 22(4):299–319, 1990.
- [49] V. Shoup. Practical Threshold Signatures. In *Advances in Cryptology - EUROCRYPT*, pages 207–220, 2000.
- [50] E. Sit and R. Morris. Security Considerations for Peer-to-Peer Distributed Hash Tables. In *1st Intl. Workshop on Peer-to-Peer Systems*, pages 261–269, 2002.
- [51] M. Steiner, T. En-Najjary, and E. W. Biersack. A Global View of KAD. In *IMC*, pages 117 – 122, 2007.
- [52] I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, and H. Balakrishnan. Chord: A Scalable Peer-to-peer Lookup Service for Internet Applications. In *SIGCOMM*, pages 149–160, 2001.
- [53] G. Urdaneta, G. Pierre, and M. van Steen. A Survey of DHT Security Techniques. *ACM Computing Surveys*, 2009. To appear.
- [54] V. Pappas and D. Massey and A. Terzis and L. Zhang.

- A Comparative Study of the DNS Design with DHT-Based Alternatives. In *INFOCOM*, pages 1–13, 2006.
- [55] D. S. Wallach. A Survey of Peer-to-Peer Security Issues. In *ISSS*, pages 253–258, 2002.
- [56] P. Weng, N. Hopper, I. Osipkov, and Y. Kim. Myrmic: Secure and Robust DHT Routing. Technical Report 2006/20, University of Minnesota, 2006.

6. Appendix

The first protocol of [46] yields a robust communication protocol with expected constant message complexity. Here, we prove that, with probability $1 - o(1)$, an execution of this protocol requires a logarithmic number of peers to exceed this expected cost.

Theorem 2. *Consider the first lookup scheme in [46] over $\Theta(\log n)$ hops. With probability $1 - o(1)$, at least $\Theta(\log n)$ peers will have a message complexity of $\Omega\left(\frac{\log \log n}{\log \log \log n}\right)$ over the course of the protocol.*

Proof: Consider a message being passed from quorum L to quorum R . The scheme of [46] works as follows. Let B denote a set of $\ln n$ bins. Each peer in R is mapped to c bins in B uniformly at random. Then each peer in L is mapped to a single bin in B uniformly at random. A peer in R listens only to those peers that get mapped to the same bin as it does. Our first goal is to show that, with probability $1 - o(1)$, there exists some bin to which $O(\ln \ln n / \ln \ln \ln n)$ peers in R get mapped. We proceed by adapting the balls-and-bins argument of [35]. For the Poisson distribution, let p be the probability that a bin has k balls or more, then:

$$p_k \geq \frac{\left(\frac{C \ln n}{\ln n}\right)^k}{k! \cdot e^{\frac{C \ln n}{\ln n}}} = \frac{C^k}{k! \cdot e^C}$$

The probability that no bin has at least $k = \ln \ln n$ balls is at most $(1 - p_k)^{\ln n} \leq e^{-p_k \cdot \ln n}$. We wish to now show that $e^{-p_k \cdot \ln n} < \frac{1}{\ln^2 n}$. Taking the logarithm of each side twice, this is equivalent to proving:

$$k \cdot \ln C + \ln^{(2)} n > \ln 2 + \ln^{(3)} n + \ln(k!) + C$$

substituting $k = \frac{\ln \ln n}{\ln \ln \ln n}$, the left side of the inequality is:

$$\left(\frac{\ln^{(2)} n}{\ln^{(3)} n}\right) \cdot \ln C + \ln^{(2)} n$$

and by using Stirling's approximation for $x!$, the right side of the inequality is at most:

$$\ln 2 + \ln^{(3)} n + \ln \left(\sqrt{2 \cdot \pi \cdot \frac{\ln^{(2)} n}{\ln^{(3)} n}} \right) + \ln^{(2)} n + C$$

Therefore, for sufficiently large n , the inequality holds and, by the results in [35] relating tail bounds of the Poisson and Binomial distributions, the result holds up to a constant factor.

We will call a bin *overloaded* if it has $\Omega\left(\frac{\log \log n}{\log \log \log n}\right)$ peers mapped to it under the scheme in [46]. We now

consider the number of peers that get mapped from L to an overloaded bin over $\ell = \Theta(\log n)$ steps of the message passing protocol in [46]. Let X_i be the indicator random variable that has value 1 if in step i of the protocol, at least one overloaded bin exists and at least one peer in L is mapped to an overloaded bin; otherwise, X_i is zero. Then $\Pr[X_i = 0] \leq o(1) + (1 - 1/\ln n)^{C \ln n} \leq o(1) + e^{-C}$ for each step $i = 1, \dots, \ell$. Therefore, $\Pr[X_i = 1] \geq 1 - o(1) - e^{-C} = \Theta(1)$. Let $X = \sum_i X_i$, then by linearity of expectation $E[X] = \Theta(\log n)$. Since the X_i s are independent and i.i.d., it follows by applying standard Chernoff bounds that $\Pr[X < (1 - \delta) \cdot E[X]] < e^{-\Theta(\log n)} = n^{-\Theta(1)} = o(1)$. Therefore, with probability $1 - o(1)$, over $\Theta(\log n)$ steps of the protocol, $\Theta(\log n)$ peers will have message complexity $\Omega\left(\frac{\log \log n}{\log \log \log n}\right)$. \square