

# Pilot: Probabilistic Lightweight Group Communication System for Ad Hoc Networks

Jun Luo, Patrick Th. Eugster, and Jean-Pierre Hubaux, *Senior Member, IEEE*

**Abstract**—Providing reliable group communication is an ever recurring topic in distributed settings. In mobile ad hoc networks, this problem is even more significant since all nodes act as peers, while it becomes more challenging due to highly dynamic and unpredictable topology changes. In order to overcome these difficulties, we deviate from the conventional point of view, i.e., we “fight fire with fire,”<sup>1</sup> by exploiting the nondeterministic nature of ad hoc networks. Inspired by the principles of gossip mechanisms and probabilistic quorum systems, we present in this paper PILOT (Probabilistic Lightweight grOUp communication sysTEM) for ad hoc networks, a two-layer system consisting of a set of protocols for reliable multicasting and data sharing in mobile ad hoc networks. The performance of PILOT is predictable and controllable in terms of both reliability (fault tolerance) and efficiency (overhead). We present an analysis of PILOT’s performance, which is used to fine-tune protocol parameters to obtain the desired trade off between reliability and efficiency. We confirm the predictability and tunability of PILOT through simulations with *ns-2*.

**Index Terms**—Mobile ad hoc networks, group communication systems, quorum systems, reliable multicast, gossiping, data sharing, replication.

## 1 INTRODUCTION

A *Group Communication System* (GCS) [2] is a useful infrastructure on which various reliable distributed computing functions can be built. The need for such a system arises not only in wired networks but also in mobile ad hoc networks. Even some mechanisms traditionally relying on a centralized service have to be implemented in a distributed way in ad hoc networks since the service provided by a single node is not dependable enough. Mobility management [3], [4], for instance, relies on a special group of nodes to continuously track locations of mobile nodes and to serve requests to these location data. The distributed management of cryptographic keys or certificates [5], [6] and group security functions like access control or key agreement [7], [8] represent another class of applications. Last but not least, distributed dynamic host configuration protocols such as naming or addressing services [9], which are essential to build a functional network, need to make agreements within the whole network.

Unfortunately, the complexity of building *reliable* GCSs, which is prohibitively high already in wired networks, is further amplified in ad hoc networks due to highly dynamic and unpredictable topology changes. In fact, even guaranteeing reliability of multicast, a key building block of GCSs, becomes extremely hard. As a result, many distributed computing functions that would depend on reliable GCSs

have to either rely on the fragile “reliability” provided by flooding [9] or make assumptions about such a service while waiting for it to appear [5].

In this paper, we identify two fundamental problems in the context of group communication, namely, 1) multicast and 2) data sharing, and we define notions of probabilistic reliability for these problems, aimed at ad hoc networks. We then present our protocol suite, called Probabilistic Lightweight grOUp communication sysTEM (PILOT) for ad hoc networks, as a solution. Innovating on the principles of gossip mechanisms and probabilistic quorum systems, PILOT provides probabilistic reliability for multicasting and data sharing, based only on a unicast primitive (rather than a multicast primitive like MAODV [10]) in order to improve the adaptability to future technology developments. We present analytical results predicting the performance of PILOT in terms of message overhead and reliability degree. We then compare these results with simulation results obtained with the *ns-2* simulator to show that we can have useful predictions on the performance of PILOT. To the best of our knowledge, the work presented in this paper, as part of the *MICS/Terminodes* project [11], is the first to provide a complete solution to the problems of reliable multicast and data sharing in ad hoc networks, along with both analytical and simulation results. It smoothly integrates, expands, and completes our previous individual results [12], [13] into a compound group communication system.

The remainder of this paper is structured as follows: Section 2 overviews related work. Section 3 details the network model and the problem to be solved. Section 4 presents our PILOT system. Section 5 analyzes PILOT in terms of reliability and efficiency. Section 6 compares those values with simulation results, and also investigates other aspects of PILOT, such as its sensitivity to node failures. Finally, Section 7 concludes the paper.

1. This expression was recently used for another gossip-based protocol [1], albeit one with a different goal from this paper.

• The authors are with the School of Computer and Communication Sciences, Swiss Federal Institute of Technology (EPFL), CH-1015 Lausanne, Switzerland. E-mail: {jun.luo, patrick.eugster, jean-pierre.hubaux}@epfl.ch.

Manuscript received 9 Apr. 2003; revised 23 July 2003; accepted 12 Nov. 2003.

For information on obtaining reprints of this article, please send e-mail to: [tmc@computer.org](mailto:tmc@computer.org), and reference IEEECS Log Number TMC-0049-0403.

## 2 RELATED WORK

The prosperous research on group communication toolkits has led to a multitude of results in wired networks, such as Ensemble [14] and Spread [15]. However, similar systems have not yet appeared in ad hoc networks, although certain supporting mechanisms like token circulation [16], random walk agent [17], reliable broadcast [18], and membership management [19] have been proposed. Our PILOT system is a first step toward building a prototype for a group communication toolkit. Rather than emphasizing the discussion in the framework of GCSs, we will hence focus on the relevant underlying building blocks.

### 2.1 Gossip-Based Probabilistic Reliable Multicast

As opposed to the “perfect” reliability guarantee for multicast, (cf. *reliable broadcast* [20]), approaches to a form of probabilistic reliable multicast (e.g., *probabilistic broadcast* (*pbcast*) [21] and *lightweight probabilistic broadcast* (*lpbcast*) [22]) reduce the protocol overhead by sacrificing safety guarantees such as atomicity through the use of a gossip-based dissemination scheme. These protocols also equally distribute the load over nodes and, thus, outperform the so-called “best effort [21]” reliable multicast (e.g., [23], [24]) by improving the resilience to arbitrary node failures and providing prediction on protocol reliability.

The *Anonymous Gossip* (AG) protocol [25], a descendant of the *pbcast* protocol, pioneered the recent research efforts on gossip-based multicast for ad hoc networks. Through the concept of anonymous gossip, any agreement on membership is avoided during the gossip-based repair phase. This, however, shifts the responsibility for the membership management to the MAODV layer [10], which the AG protocol also relies upon for a preliminary, rough packet dissemination. These prerequisites make the AG protocol more difficult to apply in a broader context than the one offered by MAODV. Furthermore, the property of predictable behavior, an important merit of gossip-based protocols, is lost due to the dependence on MAODV to guide the gossips.

### 2.2 Probabilistic Quorum Systems

*Quorum systems* [26] have been proposed as an alternative to the *state-machine* approach [27] for reliable data sharing. They improve the efficiency of the replication of the stored data by better balancing the overhead between updates and queries. Unfortunately, “original” quorum systems, also termed *strict quorum systems*, do not apply well to highly dynamic environments. This is because the very construction of these quorums is not a trivial task, the outcome of this task being strongly subject to membership changes. By introducing *probabilities* for the intersection of individual quorums, *probabilistic quorum systems* [28] relax the construction rules for quorums and leave more freedom for trading protocol overhead for reliability. While this smoother trade off has constituted the driving force behind probabilistic quorum systems, it turns out that the resulting reduced determinism makes such an approach also more viable for ad hoc networks than a strict approach. The overhead considered in [28] is the charge of computation for individual servers. Our definition of overhead, however, focuses on the consumption of network resources because

computation is much cheaper than communication in wireless networks.

Haas and Liang [29] first introduced probabilistic quorum systems into ad hoc networks for mobility management, under the name of *randomized database groups*. They propose a very interesting way to express both fault tolerance and load as costs of their system, and optimize those costs numerically. Considering the similarity between their system and PILOT, we provide some comparisons between the two solutions in Section 4.6.

### 2.3 Data Management in Ad Hoc Networks

The 7DS system presented in [30] shares certain features of our PILOT system, with respect to the diffusion scheme used for data dissemination. However, since the two systems are designed for different network environments (7DS assumes a rarely connected network, whereas PILOT considers networks of relatively high density), the underlying diffusion mechanisms are quite different. Whereas 7DS passively exploits node mobility to relay data from one node to the other, which can result in a considerable delay for data spreading but has the potential to improve power and bandwidth usages, PILOT more actively “pushes” data to other nodes with a gossip-based protocol. As a result, the analytical models for the two diffusion processes are also different (diffusion controlled process for 7DS and epidemic model for PILOT).

Both [31] and [32] try to guarantee data accessibility upon network partitioning in a replication system by investigating the problem of dynamic replica allocation. While [31] makes assumptions (e.g., data items are not updated) that seem to be too strong to capture the reality of mobile networks and hence has limited application scope, the approach in [32] is more practical in the sense that it takes into consideration topology information (e.g., connection stability) when replicating data; and data replication only happens when necessary, according to certain partition detection schemes. As far as system models are concerned, the problem we solve is somewhat orthogonal to the one of [32]. The mobility model they propose assumes strong correlations between different nodes (e.g., nodes are organized into mobility groups), which might lead to frequent network partitions. We, however, consider a purely random mobility pattern, in which network partitions seldom happen and mobility prediction does not make much sense.

## 3 GOALS AND ASSUMPTIONS

This section models the considered environment and states the problem to be solved.

### 3.1 Model

We consider an ad hoc network consisting of a set  $\mathbf{N}$  of nodes and assume that every node  $i \in \mathbf{N}$  has a unique *id*. Nodes may fail only by crashing, i.e., stopping to function. Failures are not permanent and can be recovered from.<sup>2</sup> All

2. This failure model also captures the case where nodes are deliberately switched off (e.g., for the purpose of battery replacement or operating system rebooting, or because the users do not intend to make use of their devices for a while).

communications between different nodes are assumed to rely on the underlying unicast protocol. We use DSR [33] as an example in this paper but, in practice, our solution can be made to work with any on-demand routing protocol.

### 3.2 Problem Statement

We consider an ad hoc network where reliable group communication primitives are required by mobile nodes. Within the broad scope of group communication, we address two fundamental problems, namely, multicast and data sharing, and associate each of them with a notion of *probabilistic reliability*.

1. *Reliable Multicast Protocol*. The multicast protocol disseminates packets within a multicast group  $G \subset N$ , which, for brevity, will be referred to as *group* hereafter. We define the following two metrics to measure the probabilistic reliability achieved by this protocol:

- *Reliability Degree of Single Packet Dissemination*  $\mathcal{R}_{ds}$ : The fraction of group members that receive the packet sent by a certain member.
- *Reliability Degree of Continuous Packet Dissemination*  $\mathcal{R}_{dc}$ : The fraction of all packets that are received by a certain member, assuming that packets are continuously sent from the same member with rate  $\lambda_o$ .

Both metrics are described by respective *cumulative distribution function* (cdf)  $\mathcal{F}(x) : [0, 1] \rightarrow [0, 1]$ ; it means that  $\mathcal{F}(x)$  is the probability that  $\mathcal{R}_{ds}$  (or  $\mathcal{R}_{dc}$ ) is at most  $x$ .

2. *Reliable Data Sharing Service*. Let  $STS^3 \subset N$  be a storage entity and  $\rho$  be a set of access protocols for STS. The STS holds shared data in a replicated fashion, and the consistency model for data replication is considered to be *shared-private* [36], i.e., the service does not commit itself to any access ordering except FIFO order.<sup>4</sup> Given access rates  $\lambda_u$  and  $\lambda_q$  for updates and queries, respectively, the data sharing service is probabilistically reliable in nature if a query access  $\rho_q(STS, \lambda_q)$  obtains, with a certain probability, the latest version of a data object resulting from an update access  $\rho_u(STS, \lambda_u)$ . The metric for the service is:

- *Reliability Degree of Access*  $\mathcal{R}_{da}$ : The probability that a query operation acquires the most recent update of the corresponding data object, considering both node and channel failures.

The overhead is measured by the *Network Load*  $\mathcal{N}_l$ , which is the average number of *unicast packet*  $\times$  *hop* per multicast packet to achieve a certain  $\mathcal{R}_{ds}$  or per unit time to achieve a certain  $\mathcal{R}_{dc}$  or  $\mathcal{R}_{da}$ . This definition is adapted to ad hoc networks by taking into account the number of hops to route a particular packet.  $\mathcal{N}_l$  considers only the load

3. STS is an abbreviation for *Storage Set*, a special group in the network. The algorithm used to initialize the STS will not be discussed here since it is out of the scope of this paper. Refer to [34], [35] for examples of initialization algorithms.

4. All the applications we have mentioned in the introduction comply with this model.

generated by our protocols, which is independent of the various possible implementations of the underlying network functions.

Our goal is to design a set of protocols that achieve a high reliability degree  $\mathcal{R}_d$  (representing  $\mathcal{R}_{ds}$ ,  $\mathcal{R}_{dc}$ , and  $\mathcal{R}_{da}$  hereafter) even under large arrival rates  $\lambda_o$  (the sum of  $\lambda_u$  and  $\lambda_q$  for data sharing), while incurring reasonable overhead  $\mathcal{N}_l$ . We target relatively large-scale networks, i.e., networks with tens or even hundreds of nodes and a random mobility pattern. Under a certain  $\lambda_o$ , the optimal performance with respect to both  $\mathcal{R}_d$  and  $\mathcal{N}_l$  does not exist since one can always be sacrificed to improve the other. Hence, we will study the trade off between the two metrics and show how to fine-tune parameters to trade either for the other.

## 4 PILOT: PROBABILISTIC GROUP COMMUNICATION SYSTEM

In this section, we first present the structure of our PILOT system, then we detail each component of PILOT separately.

### 4.1 Overview: Layered Architecture of PILOT

PILOT is a two-layer system, illustrated by the dark gray part in Fig. 1a. It has a probabilistic multicast protocol, Route Driven Gossip (RDG), as its basis. The protocol is gossip-based [21] in nature: it proceeds round by round while the receivers in each round are randomly chosen and they *relay* packets to the receivers of the later round(s), as shown in Fig. 1c. Upon this layer, two dedicated services are built. R<sup>2</sup>DG (Reliable RDG) is devised for continuous packet dissemination. It exploits the fact that packet losses can be detected by observing gaps in the *pid* (see Section 4.2.1) sequence and, thus, piggybacks a negative acknowledgement with each packet sent (or relayed) to *pull* the lost packet back. The other service, Probabilistic quorum system for Ad hoc Networks (PAN), provides reliable data sharing by assuming the existence of an STS to store the shared data in a replicated manner. Any node  $i \in STS$  is termed *server*, whereas the rest of the nodes are termed *clients* of the STS. Data queries and updates are directed to an arbitrary server in the STS while the message dissemination within the STS is performed by RDG, as shown in Fig. 1b. According to their requirements, applications can either use the upper-layer services or directly call RDG if only single packet dissemination service is required.

### 4.2 RDG: Basic PILOT Multicast Protocol

Our RDG protocol uses a *pure* gossip scheme, as it is not built upon any underlying multicast protocol, in contrast to [25] (the only related approach we are aware of). As opposed to “traditional” gossip protocols that only consider the membership information of a group, RDG adapts to the peculiarity of ad hoc networks by also taking the availability of routing information into account. Although the resulting membership view for each member is just a random subview due to the randomness of routing information that nodes can have, the protocol still works very well in the sense that the reliability is in practice very high and also predictable.

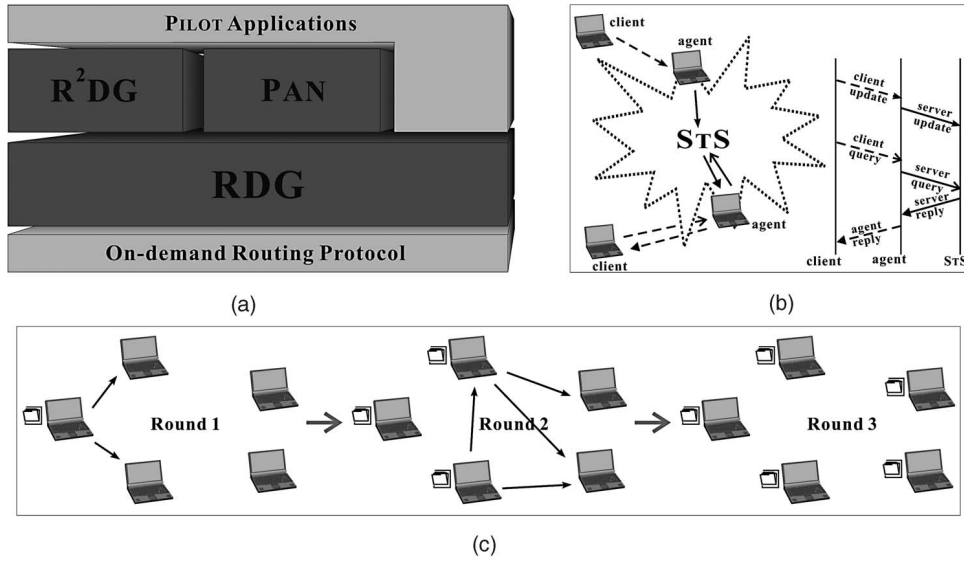


Fig. 1. Principles of PILOT. (a) Architecture of PILOT: The basic probabilistic multicast protocol (RDG) is at the bottom; R<sup>2</sup>DG and PAN are built upon the basic protocol. (b) Message exchanges for updating and querying the STS in PAN. (c) Gossip-based multicasting in RDG.

#### 4.2.1 Protocol Overview

Each packet multicast by RDG is uniquely identified by its identifier  $pid$ , defined as a tuple [group ID ( $gid$ ), source ID ( $sid$ ), pkt seq. no. ( $seq$ )]. The protocol has four data structures. In the data management part,  $pidList$  stores the  $pids$  of the received packets, and  $Buffer$  temporarily stores these packets. The other two are for the membership management of the protocol.  $gidList$  stores the identifiers of all groups that a node belongs to.  $View$  is composed of three fields:

1.  $AView$  stores the  $ids$  of known members, whose corresponding routing or location information is known;
2.  $PView$  stores the  $ids$  of known members, whose corresponding routing or location information is currently unavailable; and
3.  $RView$  stores the  $ids$  of members having indicated their willingness to leave.<sup>5</sup>

All these records are divided into several subsets with each subset being dedicated to a certain group, i.e., each  $node_i$  has four subrecords ( $pidList_i^{gid}$ ,  $Buffer_i^{gid}$ ,  $gidList_i^{gid}$ , and  $View_i^{gid}$ ) for a certain group  $G$  (with identifier  $gid$ ) that it belongs to. In addition, each record is of limited size, noted  $|R|_{max}$  for a given record  $R$ .

RDG offers seven operations, which are grouped into three sessions corresponding to their functionality. The *join* session defines the behavior of the node interested in joining a group and the reactions of other group members. The *leave* session defines the behavior of the node intending to leave the group and the reactions. In the *gossip* session, newly received packets are periodically propagated by a node. Furthermore, nodes respond to the gossip messages received. In relation to the GOSSIP task, three protocol parameters are defined here:

1. the *fanout* ( $F$ ) is the number of gossip destinations randomly selected from the  $AView$  for each gossip emission,
2. the *quiescence threshold* ( $\tau_q$ ) is related to each data packet: a packet will be removed from  $Buffer$  after having been gossiped for  $\tau_q$  rounds by individual nodes, and
3. the *age threshold* ( $\tau_a$ ) limits the propagation range of each packet.

These parameters are set by the upper layer to control the behavior of the protocol (see Section 4.2.3).

#### 4.2.2 Join Session

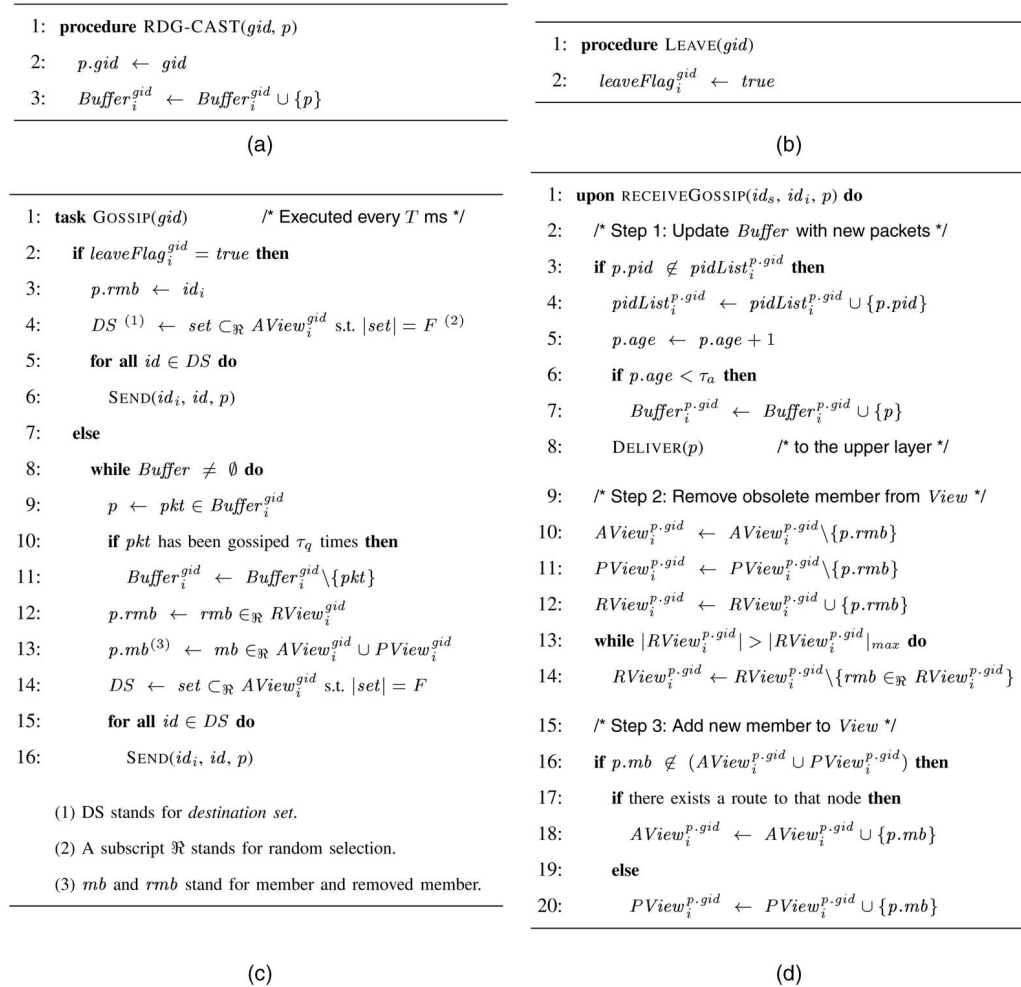
A node intending to join a group floods the network with a GROUPREQUEST message to search for other group members while announcing its existence. Upon receiving such a message from a certain member, all members update their  $AView$  with the new  $id$ . They also return a GROUPREPLY to the request initiator with probability  $P_{reply}$ . The probability is set by each node, according to its own estimation of the group size, in order to avoid GROUPREPLY storms. The initiator of the GROUPREQUEST also updates its  $AView$  after receiving the GROUPREPLY. The detailed description can be found in [12].

By recording the route of each incoming packet, DSR ensures that a new element in  $AView$  has a corresponding route entry in the DSR routing table. The validity of this relationship is periodically checked and the  $AView$  and  $PView$  are updated accordingly. When the size of  $AView$  drops below some threshold, the node has to reinitiate a *join* session.

#### 4.2.3 Gossip/Leave Session

When a node wants to multicast a packet  $p$ , it inserts the packet in its  $Buffer$  as shown in Fig. 2a. A node intending to leave a group sets a *leaveFlag* for that group as shown in Fig. 2b. Each member of the group periodically (every

5.  $AView$ ,  $PView$ , and  $RView$  stand for *active view*, *passive view*, and *remove view*, respectively.



(1) DS stands for *destination set*.

(2) A subscript  $\mathfrak{R}$  stands for random selection.

(3)  $mb$  and  $rmb$  stand for member and removed member.

Fig. 2. Gossip/leave session at node  $i$ . (a) Multicast. (b) Node leave. (c) Packet emission. (d) Packet reception.

$T$  ms)<sup>6</sup> gossips packets stored in *Buffer* to  $F$  other nodes randomly chosen from *AView* (see Fig. 2c, lines 8-16). It also piggybacks part of its view of the membership. A data packet is removed from *Buffer* after having been gossiped for  $\tau_q$  times. The SEND primitive is a direct call to the underlying unicast protocol, which will also be used by other parts of PILOT for the same purpose. If the node intends to leave, only the field of  $rmb$  is used (see Fig. 2c, lines 3-6). As illustrated in Fig. 2d, a group member receiving a gossip packet will

1. update the *Buffer* with new packets (lines 3-8),
2. remove the obsolete member from its *View* (lines 10-14), and
3. add the new member to the *View* (lines 16-20).

Note that a packet relayed  $\tau_a$  times will not be gossiped again.

RDG performs message dissemination and membership tracking at the same time. Due to the node mobility and frequent membership changes, it is not practical to have a full membership view for each member. In fact, even if it is possible to have the *ids* of all members, there is no guarantee

6. In order to save bandwidth, we apply the *binary exponential backoff* algorithm to adjust the period when there is no new packet to be sent.

that the corresponding routing or location information is available. Our routing/location-oriented membership management scheme tries to provide each member with a partial view, approximately random in nature, by exchanging membership information between members. The underlying scheme, together with sporadic losses and discoveries of the routing or location information,<sup>7</sup> has a similar effect as the reshuffling of the partial view.

Considering that the locality of network traffic can reduce the network load, we apply a general optimization by raising the awareness of the topology. This optimization is based on the assumption that the underlying routing protocol can provide partial topological information. Our heuristics in the case of DSR work like this: For a given group member, different weights are assigned to the members in *AView* according to the lengths of the routing paths to them, i.e., the longer a path the lower its weight, such that a “near” member is chosen with higher probability to relay a packet. A more detailed protocol description can be found in [12].

7. The information could be lost due to the node mobility or the timeout of route cache timer. On the other hand, a node can also obtain new information by requesting it or tapping it from packets under transmission.

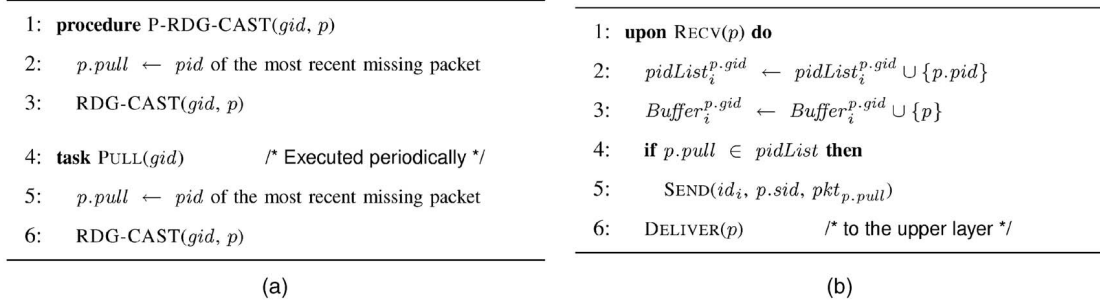


Fig. 3. Multicast and *pull* session at node  $i$ . (a) Multicast and *pull* task. (b) Packet reception and the response to *pull*.

### 4.3 R<sup>2</sup>DG: Continuous Packet Multicasting Service

If a stream of packets is multicast from a source, the *pid* sequence of received packets, at a certain group member, provides important information about packet loss. Based on RDG, our R<sup>2</sup>DG protocol exploits this feature to enhance the reliability of multicasting.

R<sup>2</sup>DG has its own data structures that are the same as for the data management part of RDG, except that the *Buffer* is much larger than that of RDG in order to have enough packets to respond to a negative acknowledgement (or *pull*). Before invoking the RDG primitive, R<sup>2</sup>DG (as shown in Fig. 3a, lines 1-3) inserts the information about a missing packet into the packet header. In addition (task PULL in Fig. 3a, lines 4-6), a packet with an empty payload (*pull*-packet) is periodically sent to the lower layer with similar information attached to it. The period is dynamically adjusted according to the number of missing packets. A group member receiving such a packet will try to respond to the pull with the packets it has, see Fig. 3b.

Considering that R<sup>2</sup>DG passes pull-packets to RDG irregularly, RDG behaves intelligently when gossiping, in the sense that it tries to piggyback the pull information along with a data packet instead of sending the pull-packet directly.

### 4.4 PAN: Reliable Data Sharing Service

Our PAN system relies on the underlying RDG to provide reliable data sharing services. It includes two protocols: a client protocol and a server protocol, as shown in Fig. 1b. In both cases of update and query, a client sends a request to an arbitrary server in the STS.<sup>8</sup> This server, termed *agent* for that client, then performs a corresponding operation of the server protocol. We assume that all messages (updates and queries) for our protocols have relatively small sizes such that they can be fit into single network packets. This requirement is justified by considering the applications we aim at. For example, a public key is only hundreds of bits long and location information might be just a coordinate in a three-dimensional space. We further require that each message be uniquely identified by its identifier *mid*, which is a tuple [source ID (*sid*), object ID (*oid*), version no. (*ver*)],<sup>9</sup> and that there is a way to establish a FIFO order among *mids*.<sup>10</sup> Since the client protocol, a one-to-one connection,

8. A client may have several ways to acquire information (e.g., identity and routing) about servers, depending on certain implementations of the STS initialization algorithm.

9. The elements *oid* and *sid* stand for the ID of the data object to be queried or updated and of the object owner, respectively.

10.  $mid_1 > mid_2$  implies that  $mid_1.sid = mid_2.sid \wedge mid_1.oid = mid_2.oid \wedge mid_1.ver > mid_2.ver$ .

can always implement certain mechanisms (e.g., ARQ [37]) to ensure reliability, we will not consider this protocol in our analysis and simulations. In the rest of this section, we focus on the server protocol.

The server protocol maintains a quorum system building upon the STS with the support from the underlying RDG protocol. We distinguish two types of quorums within the quorum system. A quorum can be a *write quorum*, accessed by an update, or a *read quorum* in the case of an access by query. Throughout the presentation, as well as in the analysis and simulations of the server protocol, we will use two symbols  $\xi_w$  and  $\xi_r$  to represent the *nominal quorum size* and the *real quorum size*, where “?” can be “W” for a write quorum or “R” for a read quorum. The nominal size is the number of servers that a certain update or query attempts to access, while the real size is the number of servers effectively accessed.

#### 4.4.1 Server Update Protocol

The agent diffuses an update message  $m_u$  within the STS by invoking the RDG protocol, as shown in Fig. 4a. Two parameters  $F$  and  $\tau_a$  (see Section 4.2.1 for the definition of these parameters) are set in order to control the size of the resulting quorum. In this paper, the value of  $\tau_a$  is always set to  $\infty$  for the server update protocol to simplify the analysis and simulations.

In order to keep track of the data access, each server keeps a record *midList*. It stores the *mids* of the most recent updates. Each server receiving a new update, including the agent, substitutes the old *mid* with the *mid* of the new update message in its *midList* before delivering the message to the upper layer, as shown in Fig. 4b. At last, all servers that effectively receive the update form a write quorum. The size of the quorum,  $\hat{\xi}_w$ , is predictable thanks to the epidemic nature of the underlying gossip-based protocol, as we will see in Section 5.3.

#### 4.4.2 Server Query Protocol

In the case of a query, the agent again uses RDG to disseminate the query message to other servers. The value of  $\tau_a$  is set to 1 to simplify the protocol evaluation later on. Also, since we consider that the arrival rate of queries is higher than that of updates in most cases, it is justifiable to have a relatively small read quorum.<sup>11</sup>

11. By setting  $\tau_a = 1$ , the nominal read quorum size  $\xi_r$  is directly determined by  $F$  (Fig. 4c, line 6), since a server receiving the query will not relay it further.

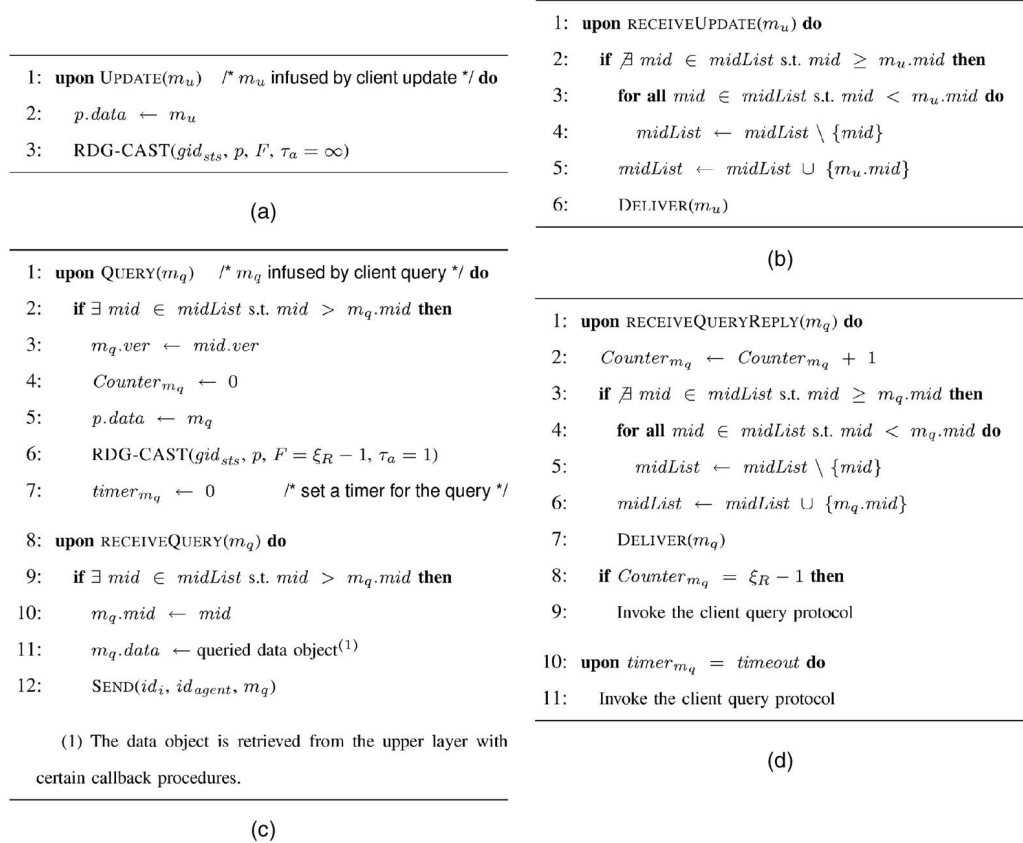


Fig. 4. UPDATE/QUERY operation at node  $i$ . (a) UPDATE emission. (b) UPDATE reception. (c) QUERY emission and reception. (d) REPLY reception at an agent.

After receiving a query message from a client, the agent sends it to other servers immediately, along with the version number of the corresponding local data object. The agent also sets a counter and a timer in order to guarantee proper termination of the query session (Fig. 4c, lines 1-7). Each server belonging to the read quorum, upon receiving the message, responds with its own copy of the data object, if its version is more recent than the one of the agent (Fig. 4c, lines 8-12). The agent always delivers a new update returned from other servers. It invokes the corresponding client protocol, after every request either yields a reply or times out, as illustrated in Fig. 4d.

#### 4.5 Examples of Protocol Operations

Fig. 5a gives a visual illustration of the behavior of our RDG protocol with respect to the dissemination of one packet, assuming a single group  $G$  of size  $|G| = 10$  within a 20 nodes network. Another example in Fig. 5b illustrates a simple execution of our PAN system in a network of 50 nodes, assuming an STS consisting of 25 nodes.

#### 4.6 Comparing PILOT with Randomized Database Group

In this section, we compare the work of Haas and Liang [29] with PAN. The comparisons are qualitative rather than quantitative because [29] does not provide simulation results to evaluate the system performance and to confirm the precision of their numerical analysis. In a nutshell, PAN outperforms randomized database group in two aspects. On one hand, the protocol used to access the database

group in [29] consists of multiple unicasts, based on the assumption of perfect routing information. Obviously, such an approach fails under the more realistic assumption of incomplete routing information, while PAN can cope with such incompleteness. On the other hand, the symmetric construction of quorum systems in [29], i.e., the same size for all quorums, is not suitable for all replication systems (for instance, when the arrival rates of queries and updates, respectively, diverge strongly). PAN, on the contrary, can adapt to a given situation by appropriately adjusting parameters. As far as the analytical methodology is concerned, the model in [29] is more application oriented than the one of PAN. It provides an insight into the probabilistic quorum systems from a different perspective.

### 5 ANALYSIS

In this section, we show that the two metrics,  $\mathcal{R}_d$  and  $\mathcal{N}_l$  (defined in Section 3.2), are predictable given certain protocol parameters and information about the network. These analytical results are confirmed by simulations in the next section. Since the behavior of R<sup>2</sup>DG pull depends on far more factors than that of RDG gossip, we will not consider this part of the protocol in the analysis. However, we will show the enhanced reliability by simulations.

#### 5.1 Model

For the multicast protocol, we consider a single group  $G$  composed of  $|G| = n$  members and observe its behavior in terms of the dissemination of a *single* packet (“one run”), but also a *continuous* stream of packets (which is more

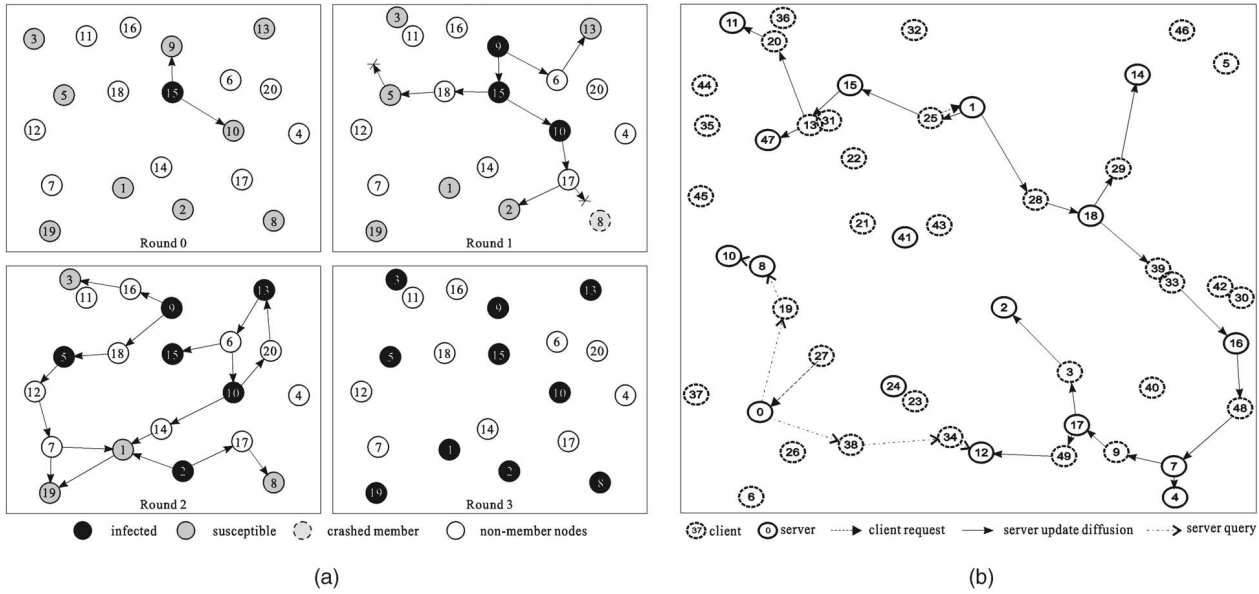


Fig. 5. Illustrations of PILOT. (a) An example of one “run” of RDG with  $F = 2$  and  $\tau_q = 2$  within a group of size 10. The packet initiated by member 15 infects the whole group in only three rounds in spite of the fact that nodes move and even fail. A member may receive duplicates of the same packet (e.g., member 1 at round 2). On the other hand, the packet can get lost at a certain round due to nodes crashing or moving (e.g., members 8 and 3 in round 1), but these losses will be compensated with high probability at a later round. (b) An illustration of an update/query pair in PAN within a network of 50 nodes located in a square area of  $1\text{km}^2$ . When node 25 wants to perform an update, it sends a request to its agent, node 1. The request of this update is diffused to other servers by node 1, using the gossip-based scheme (Only the valid transmissions are shown here. Duplicated transmissions are omitted to simplify the visualization.). If node 27 wants to access the data, it also requests its agent, node 0. Node 0 in turn requests other servers, node 8, 10, and 12. In this case, node 12 is the intersection of the read and write quorum. It is able to reply the requested data of node 25 to node 27. The query reply is omitted here for simplicity.

realistic than related research proposals considering only the “one run” part). Each gossiping operation is modeled as a uniform random selection of  $F$  members out of  $n$ , i.e., without considering the topology-awareness, in order to simplify the tractability. According to the terminology of epidemiology [38], a member that has received a certain packet is termed *infected*, otherwise *susceptible*. An infected member attempting to share the packet with others (i.e., a member who keeps gossiping the packet) is called *infectious*. We analyze our protocol in a network composed of a static set of nodes running closely “synchronized.” More precisely, nodes gossip in synchronous rounds ( $T$  ms, identical for all nodes), and there is an upper bound on the network latency which is smaller than  $T$ .

The probability of packet loss is closely related to the movement and traffic pattern, as well as to the length of the considered routing path. By assuming an identical and independent probability of failure  $p_f$  for each hop along a routing path in a certain network environment, the probability of losing a certain gossip message can be expressed as a function of the number of hops,  $H$ , of that routing path. We further assume that the lengths  $H$  of all routing paths between any two members follow the same distribution  $f(h)$ . On the other hand,  $p_f$  can be split into two parts: 1)  $p_{fc}$  represents the probability of packet loss due to node crash and 2)  $p_{fmo}$  reflects the effects of node mobility and buffer overflow. Since  $p_{fc} \ll p_{fmo}$  in general for mobile wireless networks, we directly use  $p_{fmo}$  to approximate  $p_f$ .

As for the data sharing service, we consider only the server protocol (including both update and query protocols) for analysis. The STS is assumed to consist of  $n$  servers. We also assume that query and update accesses arrive

randomly at an arbitrary server, following Poisson processes with intensities of  $\lambda_q$  and  $\lambda_u$ , respectively. By further assuming that these two processes are independent, the overall access rate is given by  $\lambda_o = \lambda_q + \lambda_u$ .

The dissemination process of the server update is performed by RDG. As this process finishes, all infected servers form a write quorum with real size  $\xi_W$  following a certain probability distribution. We consider only the second query to a data object that was modified by the most recent update, while considering the first query as happening before the update.<sup>12</sup> For example, as shown in Fig. 6a, only the pairs of (update  $\beta$ , query  $\beta_2$ ) and (update  $\gamma$ , query  $\gamma_2$ ) are considered, whereas queries  $\beta_1$  or  $\gamma_1$  are supposed to request previous updates (i.e., updates  $\alpha$  and  $\beta$ , respectively). This assumption makes sense when we consider the time with respect to a server where updates and queries arrive, and also the property of a Poisson process shown in Fig. 6b. Since there is always some delay for the message dissemination, the probability that the actual occurrence of events will follow the order of our assumption at that server is very high, according to different distributions of the time interval between two events within a Poisson process (see Fig. 6b). This makes the present analysis a “viable” lower bound.

We continue using  $p_f$  to represent the network condition, but an empirical value  $p_e$  is also used in the case of queries to represent the server *unavailability* due to failure, at any time instant. One might argue that the server failure should be treated as a Poisson process [29], but this is not justifiable

12. The time of an event is when it happens at an agent.



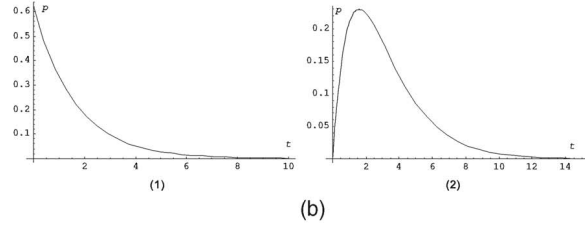
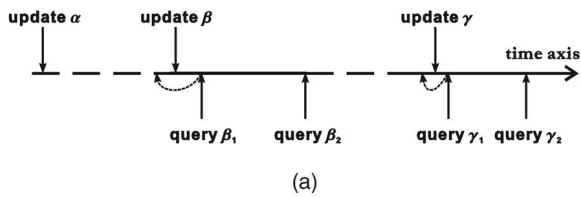


Fig. 6. Time intervals between events and their distributions. (a) The occurrence of events in terms of absolute time. (b) The distributions of time interval between two events: 1) exponential distribution for consecutive events and 2) *Erlang* distribution for nonconsecutive events.

with a failure recovery model, which is usually the case in ad hoc networks (e.g., nodes switching off for the purpose of battery replacement or operating system rebooting).

## 5.2 Stochastic Behavior of RDG

Considering a packet multicast by a member, we use  $S_r \in \{0, \dots, n\}$  to denote the number of members infected with the packet **after** round  $r$ . With the convention that  $\Pr\{S_r = 0\} = 1$  for  $r < 0$ , it is easy to show that the sequence of random vectors  $\mathbf{S}_r = [S_r, S_{r-1}, \dots, S_{r-\tau_q}]_{r \geq 0}^T$  forms a Markov chain with values taken from the *state space*

$$\mathcal{E} = \overbrace{\{0, \dots, n\} \times \dots \times \{0, \dots, n\}}^{\tau_q + 1}.$$

1. *Recurrence Relation.* Given the probability  $p$  that a certain member is infected by a specific gossip message,  $q = 1 - p$  represents the probability of noninfection. Let  $S_r = i$  (the number of infected members) and  $S_r - S_{r-\tau_q} = k$  (the number of infectious members) in the current round; we introduce a binary random variable,  $X_l$ , for each of the remaining  $n - i$  susceptible members, where  $\Pr\{X_l = 0\} = q^k$ , i.e., the probability that a certain susceptible member is not infected in the next round is the probability that it is not infected by any of the  $k$  infectious members. It is clear that  $S_{r+1} - S_r = \sum_{l=1}^{n-i} X_l$  follows a binomial distribution. Let  $j$  be the number of infected members in the next round; the transition probability is expressed as:

$$\begin{aligned} & \Pr\{S_{r+1} = j | S_r = i, S_r - S_{r-\tau_q} = k\} \\ &= \Pr\left\{ \sum_{l=1}^{n-i} X_l = j - i | S_r - S_{r-\tau_q} = k \right\} \\ &= \begin{cases} \binom{n-i}{j-i} (1-q^k)^{j-i} q^{k(n-j)} & j \geq i \\ 0 & j < i \end{cases} \end{aligned} \quad (1)$$

which leads to the following global balance equation of the chain:

$$\begin{aligned} & \Pr\{\mathbf{S}_{r+1} = \mathbf{s}_{r+1}\} = \\ & \sum_{i_{\tau_q}=0}^{i_{\tau_q}-1} \binom{n-i}{j-i} (1-q^{i_{\tau_q}})^{j-i} q^{(i-i_{\tau_q})(n-j)} \Pr\{\mathbf{S}_r = \mathbf{s}_r\}, \end{aligned} \quad (2)$$

where  $\mathbf{s}_r = [i, i_1, \dots, i_{\tau_q}]^T$ ,  $\mathbf{s}_{r+1} = [j, i, i_1, \dots, i_{\tau_q-1}]^T$ , and  $i = i_0$ . Let the column vector  $\nu_r$ , with  $\nu_r(i) = \Pr\{S_r = i\}$  as its  $i$ th element, be the marginal distribution of  $S_r$ . Given the initial distribution  $\nu_0 = [0, 1, 0, \dots, 0]^T$  and (2),  $\nu_r$  is then computed as:

$$\nu_r(i) = \sum_{i_1=0}^i \sum_{i_2=0}^{i_1} \dots \sum_{i_{\tau_q}=0}^{i_{\tau_q-1}} \Pr\{\mathbf{S}_r = \mathbf{s}_r\}. \quad (3)$$

2. *Computation of  $p$ .* According to our assumptions, the probability of infection  $p$  can be estimated by taking two conditions into account: 1) the considered node is chosen as the gossip destination and 2) the gossip message is successfully received. This results in the following expression (remember that  $F$  is the protocol parameter fanout):

$$p = \overbrace{P_{gossip}}^{(i)} \overbrace{P_{succ}}^{(ii)} = \left( \frac{F}{n-1} \right) P_{succ}. \quad (4)$$

Given a certain length (in hops)  $h$  of a routing path, the probability of a successful delivery is expressed as  $P_{succ} = (1 - p_f)^h$ , i.e., there is no failure in each of the  $h$  hops. So, we have:

$$P_{succ} = \sum_h (1 - p_f)^h \Pr\{H = h\} = \mathbf{E}_H[(1 - p_f)^H]. \quad (5)$$

Therefore,  $p$  is expressed as:

$$p = \left( \frac{F}{n-1} \right) \mathbf{E}_H[(1 - p_f)^H]. \quad (6)$$

The distribution of  $H$  and the value of  $p_f$  are the network information we need. We refer to [12] for discussions about their estimations.

3. *Reliability Degree  $\mathcal{R}_{ds}$  and  $\mathcal{R}_{dc}$ .* With the recurrence relation (3) of the single packet dissemination, the reliability degree can be expressed<sup>13</sup> in terms of  $\nu(i)$  as follows. Note that the distribution of  $\mathcal{R}_{ds}$  is always related to the group size  $n$ , while the distribution of  $\mathcal{R}_{dc}$  is related to the number of packets in a stream, denoted by  $M$  in the formula:

13. The subscript  $r$  is omitted hereafter, because we always consider the final distribution (i.e., after the last round).

$$\text{cdf of } \mathcal{R}_{ds} : \quad \mathcal{F}_n(x) = \sum_{i=1}^{\lfloor nx \rfloor} \nu(i) \quad (7)$$

$$\text{cdf of } \mathcal{R}_{dc} : \quad \mathcal{F}_M(x) = \sum_{i=0}^{\lfloor Mx \rfloor} \binom{M}{i} p_1^i (1-p_1)^{M-i}, \quad (8)$$

where  $p_1 = \sum i \cdot \nu(i)/n$  is the probability that a certain group member receives a single packet in a stream. Here, we assume that the receptions of two distinct packets are independent events.

4. *Network Load*  $\mathcal{N}_l$ . The  $\mathcal{N}_l$  for single packet dissemination is estimated straightforwardly by counting the number of unicast packets sent and the number of hops traveled by each of them:

$$\mathcal{N}_l = \mathbf{E}[S_{\tau_a}] \cdot F \cdot \tau_q \cdot \mathbf{E}[H]. \quad (9)$$

Recall that  $\tau_a$  limits the number of gossip rounds and  $\tau_q$  defines how many times a packet is repeatedly relayed by a certain group member. The expression for  $\mathcal{N}_l$  in the case of continuous packet dissemination is omitted as it becomes trivial with (9) and a given  $\lambda_o$ . This prediction is relatively rough because it is hard to find a way to precisely estimate the distribution of  $H$  as the distribution depends on several factors.

### 5.3 Stochastic Behavior of PAN

Since PAN directly uses RDG to diffuse an update, the distribution of  $\hat{\xi}_W$  can be estimated with (3). The distribution of  $\xi_R$  can be expressed in a similar but more precise way, since  $\tau_a$  is set to 1 (see Section 4.4.2).  $\mathcal{N}_l$  is computable given the two distributions, but information about the time interval between a query and an update is necessary to compute  $\mathcal{R}_{da}$ .

5. *Reliability Degree*  $\mathcal{R}_{da}$ . According to the definition and the protocol description, this value is in fact the probability that a read quorum intersects the most recent corresponding write quorum. More precisely, we are looking for the probability that two subsets with sizes  $\hat{\xi}_W$  and  $\hat{\xi}_R$ , taken from a set of  $n$  servers, intersect. Note that  $\hat{\xi}_R$  is defined as the number of servers that effectively reply to the query back to its forwarding agent.

There exists an  $\tilde{r}$  for which the dissemination process is finished, i.e., no new server is infected when  $r \geq \tilde{r}$ . Based on the assumption of synchronization, we divide the time axis after a given update event  $\beta$  into  $\tilde{r} + 1$  intervals, as shown in Fig. 7. A read quorum, resulting from a query happening in-between two consecutive gossip rounds  $r$  and  $r + 1$ , would have to intersect a write quorum of size  $\hat{\xi}_W$  with a distribution  $\nu_r$ . In order to find the probability of intersection, we need to calculate the read quorum size  $\hat{\xi}_R$  (with a distribution  $\mu$ ) and  $p_r$ , the probability that the query event occurs in-between rounds  $r$  and  $r + 1$ .

The distribution of  $\hat{\xi}_R$ , conditioned on  $\xi_R = s$ , is calculated as follows, with an initial value<sup>14</sup>

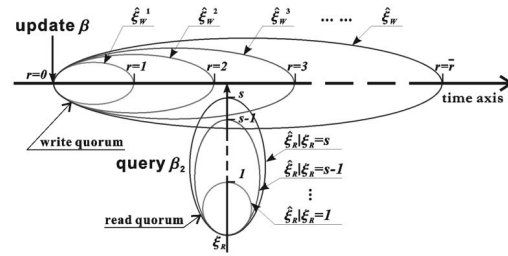


Fig. 7. Incremental processes of read and write quorum size:  $\hat{\xi}_W$  increases round by round, while  $\hat{\xi}_R$  increases with the amount of queries sent by an agent.

$\Pr\{\hat{\xi}_R = 1 | \xi_R = 1\} = 1$  and the convention  $\Pr\{\hat{\xi}_R = k | \xi_R = s\} = 0$  if  $s < 1$ ,  $k < 1$ , or  $k > s$ :

$$\begin{aligned} \mu_s(k) &= \Pr\{\hat{\xi}_R = k | \xi_R = s\} \\ &= \mu_{s-1}(k-1)p + \mu_{s-1}(k)(1-p) \\ &= \Pr\{\hat{\xi}_R = k-1 | \xi_R = s-1\} \cdot p \\ &\quad + \Pr\{\hat{\xi}_R = k | \xi_R = s-1\} \cdot (1-p) \\ &\quad k = 1, \dots, s \text{ and } s \geq 2, \end{aligned} \quad (10)$$

where  $p = \mathbf{E}_H[(1-p_f)^{2H}](1-p_e)$  is the probability that the agent forwarding a query receives the reply from a server belonging to the corresponding read quorum. The estimation of  $\mu$  is somewhat conservative because servers with a relatively old data version do not reply to a query.

The time interval between an update and the second query to it is characterized by an Erlang distribution  $\lambda_q^2 t e^{-\lambda_q t}$ , with the assumption of a Poisson arrival process. Therefore, we have

$$p_r = \begin{cases} \int_{t_r}^{t_{r+1}} \lambda_q^2 t e^{-\lambda_q t} dt & r < \tilde{r} \\ \int_{t_r}^{\infty} \lambda_q^2 t e^{-\lambda_q t} dt & r = \tilde{r}. \end{cases} \quad (11)$$

Now, the probability of intersection, i.e.,  $\mathcal{R}_{da}$ , is expressed by taking an average over all possible cases:

$$\mathcal{R}_{da} = \sum_{r=0}^{\tilde{r}} \sum_{i=1}^n \sum_{j=1}^s \left( 1 - \frac{\binom{n - \hat{\xi}_W}{\hat{\xi}_R}}{\binom{n}{\hat{\xi}_R}} \right) \mu_s(j) \nu_r(i) p_r. \quad (12)$$

6. *Network Load*  $\mathcal{N}_l$ . For a certain  $\mathcal{R}_{da}$  with its parameter pair  $F$  and  $\xi_R$ , we evaluate the corresponding  $\mathcal{N}_l$  by averaging the load over a certain time unit (e.g., 1s), taking into account the arrival rate of updates and queries. Therefore, the loads generated by a single update and query are calculated separately and then  $\mathcal{N}_l$  is obtained by summing the products of the loads of the individual operations and their corresponding arrival rates:

$$\mathcal{L}_W = \mathbf{E}[\hat{\xi}_W] \cdot F \cdot \tau_q \cdot \mathbf{E}[H], \quad (13)$$

14. Because the agent, one of the servers, has already received the query, it is sure that  $\xi_R = 1$ , if  $\xi_R$  is set to 1.

$$\mathcal{L}_R = 2 \cdot \xi_R \cdot \mathbf{E}[H], \quad (14)$$

$$\mathcal{N}_l = \lambda_u \mathcal{L}_W + \lambda_q \mathcal{L}_R. \quad (15)$$

This estimation is conservative in the same sense as we mentioned before. Again, it is relatively rough compared with the one for  $\mathcal{R}_{da}$ , because we do not take into account the following two facts: 1) many packets get dropped before reaching their destinations, and 2) packets, especially those eventually dropped, may travel quite a long way due to stale routing information. We will show with simulations that the former fact has a dominating effect in most cases, but these facts tend to offset each other in some cases.

## 6 SIMULATIONS

This section presents the simulation results of our PILOT system in five parts. Sections 6.2 and 6.3 are dedicated to RDG/R<sup>2</sup>DG and Sections 6.4, 6.5, and 6.6 are devoted to PAN. The main goal is to confirm our claim that both the reliability degree  $\mathcal{R}_d$  and the network load  $\mathcal{N}_l$  of PILOT are predictable. The impact of the message arrival rate  $\lambda_o$  and of the server failures  $p_e$  on PAN is also investigated by simulations in different settings. We refer to [12] for comparisons of simulation results between RDG/R<sup>2</sup>DG and AG [25].

### 6.1 Model and Parameters

The simulator we use is *ns-2* [39] with the Monarch Project wireless and mobile extensions. It provides both implementations of DSR and wireless MAC, based on the Lucent WaveLAN IEEE 802.11 product, with a 2Mbps transmission rate and a nominal range of 250m. We adopt the two-ray ground reflection model [40] as the radio propagation model.

We simulate ad hoc networks in a square area of 1km<sup>2</sup>. The movement pattern is defined by the “random waypoint” model [41]. The simulation parameters such as network size and maximum node speed are specified for each simulation. The STS for PAN always contains half of the network nodes. We do not justify this number,<sup>15</sup> but only use it as an example. The servers in the STS are assumed to be predefined in order to simplify the simulation.<sup>16</sup> The client protocol is omitted to reduce side effects.

The gossip period is set to 200ms. For RDG/R<sup>2</sup>DG, a CBR traffic generator produces 64 byte packets at regular intervals of 200ms, which gives a  $\lambda_o = 5\text{pkt/s}$ . The effect of the sending rate will be investigated in our future work. The arrival of queries or updates in PAN is emulated by a Poisson traffic source attached to each server, generating packets of 128 bytes with rate  $\lambda_o$ . We first investigate the impact of the overall access rate  $\lambda_o$  on the performance of

15. It is not the goal of this paper to find the optimal size for an STS, but we note that generally, the larger the size, the heavier the network is loaded, whereas the load on individual servers becomes smaller.

16. Although the clustering algorithm is a popular way to elect some representatives of the network, introducing such an algorithm into our simulation may only bring more overhead to this task, without any help to show the essence of our system.

PAN, then we take an appropriate value for all simulations. Due to space limitations, we use  $\lambda_o = 8\lambda_u$  for all simulations. With certain simulation parameters (network size, maximum speed, pause time, and arrival rate), we vary the protocol parameters  $F$  and  $\xi_R$  in order to show the trade off between the two metrics  $\mathcal{R}_{da}$  and  $\mathcal{N}_l$ . As the last simulation parameter,  $p_e$  is first set to 1 percent, and then varied to show the sensitivity of PAN to server failures.

Both RDG/R<sup>2</sup>DG and PAN are operated more than 400 seconds of simulated time. The first 50 seconds of the simulation are used for system initialization. Then, each traffic source continues generating traffic according to the predefined intensity until the end. Each simulation is carried out 10 times with different scenario files created by *ns-2*.

### 6.2 Single Packet Dissemination Reliability $\mathcal{R}_{ds}$

Fig. 8 shows the analytical and simulation results (i.e., the evolution of the infection processes) of the basic RDG protocol. The figure contrasts one process with another instead of providing the value of  $\mathcal{R}_{ds}$  explicitly. These comparisons basically confirm that the theoretical prediction of the relationship between the reliability and the latency is valid.

It is easy to observe that the reliability of the protocol with  $F = 3$  is better than the one with  $F = 2$  because the fanout has a significant effect on the reliability. However, when we further increase the fanout, the reliability decreases instead of increasing (analysis) or only marginally increases (simulation). The reason is that increasing the fanout has the same effect as increasing the number of connections, and  $p_f$  increases dramatically because of the network congestion. A similar reason accounts for what happens when  $\tau_q$  changes from 1 to 2.

In fact, there is always a trade off between certain requirements on reliability and the introduced overhead, characterized by the values of  $F$  and  $\tau_q$ . Considering the network capacity imposes a further limitation not considered in other research proposals (considerably large  $F$  [42] or unbounded  $\tau_q$  [22]). Therefore, for all simulations later in this paper, we always take  $F \leq 3$  and  $\tau_q = 1$  for RDG.

### 6.3 Continuous Packet Dissemination Reliability $\mathcal{R}_{dc}$ and Network Load $\mathcal{N}_l$

Fig. 9 shows  $\mathcal{R}_{dc}$  and  $\mathcal{N}_l$  of both RDG and R<sup>2</sup>DG with different mobility patterns and group sizes. We provide here the mean value of  $\mathcal{R}_{dc}$  and its standard deviation, which characterize the distribution function  $\mathcal{F}$ . The figures again exhibit the similarity between the simulation and analytical results with respect to RDG (see Section 5.2 for the explanation of the rough prediction on  $\mathcal{N}_l$ ). As expected, R<sup>2</sup>DG always performs better than RDG in terms of reliability, while the improvement is significant in high mobility and large group scenarios, thanks to the gossip-pull mechanism. We also note that only a slight reliability degradation is observed (especially in the case of R<sup>2</sup>DG) when the mobility or group size is increased (with a sublinear increment of  $\mathcal{N}_l$  in the case of increasing group size), illustrating the scalability of our protocols.

Note that two simulation parameters are paired to represent the mobility pattern such that each node has a

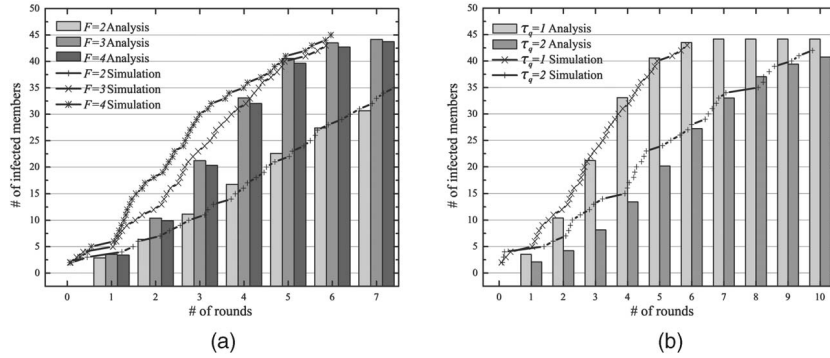


Fig. 8. Average number of infected members (simulation) and expected number of infected members (analysis) in time (expressed in rounds) with  $n = 50$  in a network of 100 nodes. Each node has a maximum speed of 2m/s and an average pause time of 40s. (a)  $\tau_q = 1$  with different values for  $F$ . (b)  $F = 3$  with different values for  $\tau_q$ .

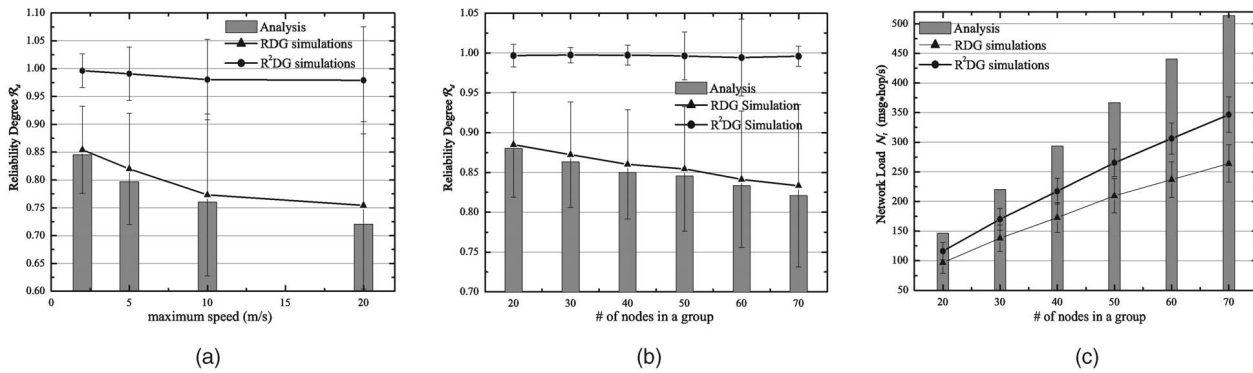


Fig. 9. Reliability degree  $\mathcal{R}_{dc}$  and network load  $\mathcal{N}_l$  versus mobility and group size in 100 nodes networks. (a) 50 nodes in a group. (b)  $Speed_{max} = 2m/s$ ;  $Time_{pause} = 10s$ . (c)  $Speed_{max} = 2m/s$ ;  $Time_{pause} = 10s$ .

maximum speed of 2m/s, 5m/s, 10m/s, and 20m/s, and a corresponding average pause time of 10s, 20s, 40s, and 80s, respectively (maximum speed is used as a symbol of the mobility pattern in this case). This concept of mobility pattern will be used throughout the rest of this section.

#### 6.4 Impact of $\lambda_o$ on PAN Performance

Fig. 10 shows the performance of PAN (assuming  $F = 2$  and  $\xi_R = 4$ ) with respect to  $\lambda_o$ , the overall access rate. We observe that PAN performs in a relatively stable way for  $1.5s^{-1} \leq \lambda_o < 3s^{-1}$ , and  $\mathcal{R}_{da}$  begins to degrade if we further increase  $\lambda_o$ , since the request arrival rate becomes larger than the service rate that PAN can provide. It is also natural to see that  $\mathcal{N}_l$  increases linearly with  $\lambda_o$  by (15). However, it may seem somewhat odd to observe that  $\mathcal{R}_{da}$  is very low in high mobility scenarios, when  $\lambda_o < 1.5s^{-1}$ . The main reason for this is the increased amount of stale routing information. In practice, this effect does not appear in the presence of background traffic. This problem can also be solved actively by requiring each STS server to send control packets during idle time in order to keep routing information fresh. Based on these observations, we apply  $\lambda_o = 2s^{-1}$  for all other simulations.

The evaluations of  $\mathcal{R}_{da}$  are presented in two ways. The “pessimistic”  $\mathcal{R}_{da}$  refers to the probability that a query reaches the most recent update (with the same assumption as in Section 5.1 about the event order), whereas for the “optimistic” one, we consider a query to be successful even if it only retrieves the result of an update that occurred right

before the most recent update. This second evaluation makes sense because, in practice, there are different data objects stored in an STS, and the probability that a queried data object has been modified by the most recent update is quite low. We will use these notations for all graph illustrations in the rest of this section.

#### 6.5 Access Reliability $\mathcal{R}_{da}$ and Network Load $\mathcal{N}_l$

Fig. 11 shows comparisons between simulation and analytical results for networks of “normal” density, i.e., 50 nodes in an area of  $1km^2$ , and “high” density, i.e., 100 nodes in an area of  $1km^2$ . We vary the maximum speed and pause time to test the impact of mobility on the performance of PAN. The protocol parameters  $F$  and  $\xi_R$  are adjusted to cope with the increased network size. We note that a real number  $x.y$  for the value of  $F$  means that each server, when gossiping an update, takes  $F = x$  with probability  $1 - y/10$  and  $F = x + 1$  with probability  $y/10$ .

We make the following observations:

1. The simulation and analytical results of  $\mathcal{R}_{da}$  match very well; this confirms the predictability on  $\mathcal{R}_{da}$ .
2. The analytical results of  $\mathcal{N}_l$  provide certain information about the system overhead, such as the trend of its changes in different situations.
3. The optimistic  $\mathcal{R}_{da}$  is always much higher than the pessimistic one; this basically means that the potential of PAN is much higher than what could be expected from the analytical results.

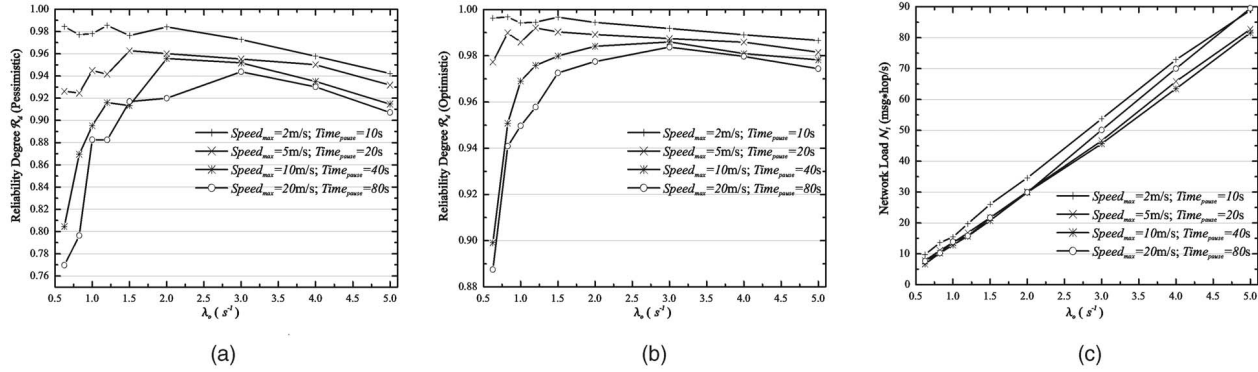


Fig. 10. Reliability degree  $\mathcal{R}_{da}$  and network load  $\mathcal{N}_l$  versus overall access rate  $\lambda_o$  for 50 nodes networks. (a) Pessimistic reliability degree. (b) Optimistic reliability degree. (c) Network load.

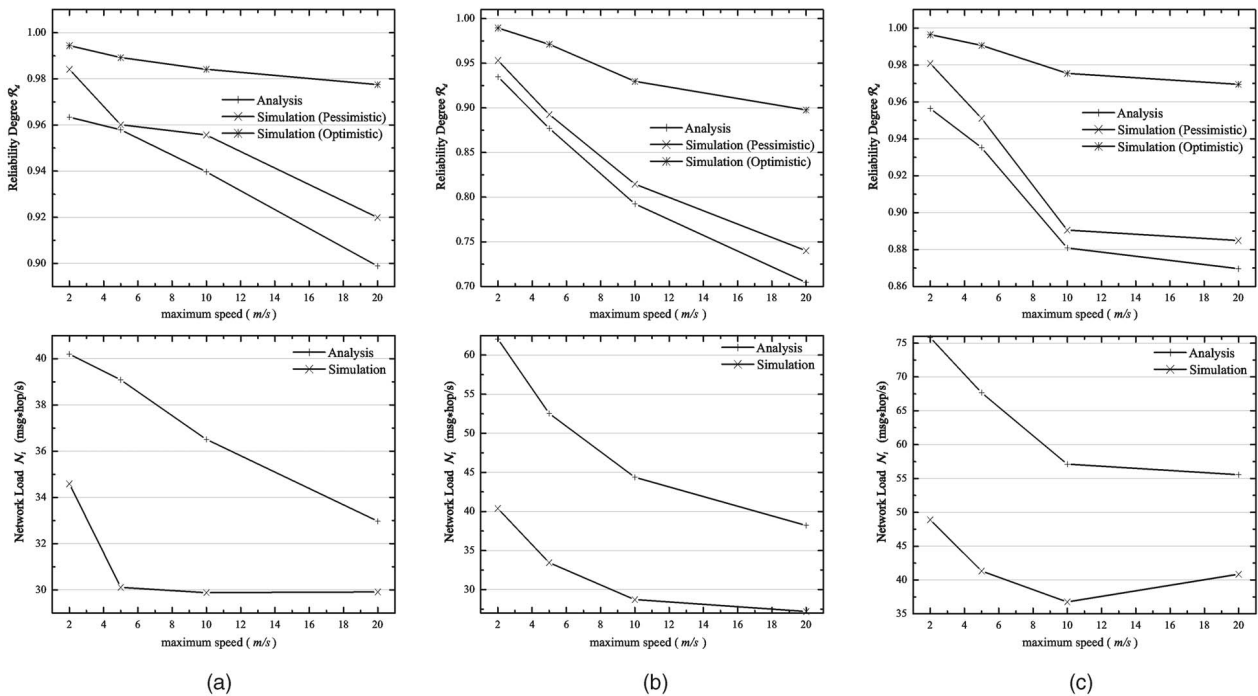


Fig. 11. Analytical and simulation results for reliability degree  $\mathcal{R}_{dl}$  and network load  $\mathcal{N}_l$  versus mobility pattern. (a) Normal density network,  $F = 2$  and  $\xi_R = 4$ . (b) High density network,  $F = 2$  and  $\xi_R = 4$ . (c) High density network,  $F = 2.2$  and  $\xi_R = 5$ .

- As the network size and the maximum node speed grow, protocol parameters have to be adjusted to maintain a good performance of  $\mathcal{R}_{da}$ , at the cost of an increased system overhead.

## 6.6 Sensitivity to Server Unavailability $p_e$

According to the simulation results shown in Fig. 12, the sensitivity of PAN (assuming  $F = 2$  and  $\xi_R = 4$ ) to  $p_e$  increases as the node mobility grows. In addition, the sensitivity of PAN considering optimistic  $\mathcal{R}_{da}$  is lower than the sensitivity considering pessimistic  $\mathcal{R}_{da}$ .

We also observe that the increase of  $p_e$  leads to an improvement of  $\mathcal{R}_{da}$  in some cases. This paradox indeed suggests a way to optimize our system, i.e., a server belonging to a certain read quorum would not always try to reply to a query back to its agent, even if the server is "alive" and has a new version of the queried data object. With such a behavior, PAN could avoid the case where more than one server replies to an agent with the same data object,

thereby reducing the probability of packet collisions and, in turn, improving  $\mathcal{R}_{da}$ . However, we do not actually apply this optimization to PAN because it is not as stable as the topology-awareness optimization in dynamic environments.

## 7 CONCLUSION

In this paper, we are concerned with probabilistic reliable group communication in mobile ad hoc networks. We have studied two fundamental aspects, i.e., multicast and data sharing, within this framework and specified performance metrics that take the peculiarities of mobile ad hoc networks into account. We have proposed our PILOT system as a solution, based on the principle of gossip mechanisms and probabilistic quorum systems. The performance of PILOT has been analyzed by making use of, notably, epidemic theory. The evaluation and investigation of PILOT have also been carried out by simulations in *ns-2*.

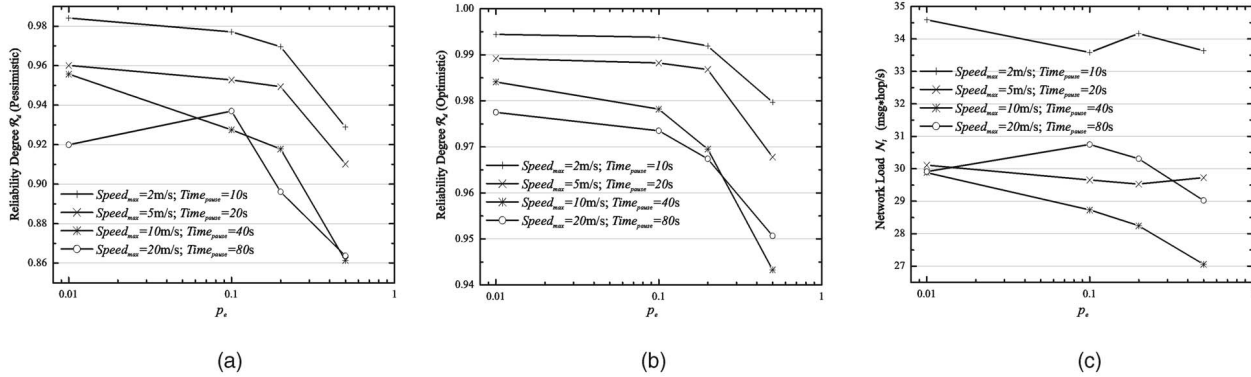


Fig. 12. Reliability degree  $\mathcal{R}_{da}$  and network load  $\mathcal{N}_l$  versus server unavailability  $p_e$  for 50 nodes networks. (a) Pessimistic reliability degree. (b) Optimistic reliability degree. (c) Network load.

As the first step toward building a probabilistic group communication toolkit, our PILOT system consists of two layers: RDG, at the bottom layer, is a gossip-based probabilistic reliable multicast protocol. At the upper layer, two dedicated services,  $R^2DG$  and PAN, provide continuous reliable packet dissemination and reliable data sharing, respectively. Our main contributions are:

1. an ad hoc adapted gossip mechanism,
2. a hybrid gossip including both push and pull,
3. gossip-based quorum access protocols, and
4. an asymmetric quorum construction.

We have proposed an analytical model to predict the performance of both RDG and PAN. The validity of these predictions has been evaluated by simulations. The results show that our analytical model provides predictions that are adequate for tuning the trade off between reliability degree  $\mathcal{R}_d$  and network load  $\mathcal{N}_l$ . Our simulation results also show that, even under frequent topology changes, the reliability degrees of RDG/ $R^2DG$  and PAN are fairly high in practice. Finally, we have investigated also other aspects of PAN with intensive simulations, which confirm its robustness, in the sense that it can sustain a large access rate  $\lambda_o$ , different network sizes, and up to 50 percent server failures.

We are in the process of determining a probabilistic notion of membership consistency and improving the analytical model by taking this notion into account. In addition, we are considering other models [42], [29] in order to further understand the benefits of gossip-based protocols and to provide numerical comparisons between PILOT and similar systems for ad hoc networks, which would better justify the deployment of PILOT. Finally, we intend to take into consideration, in our simulations, the recently recommended modifications to the “random waypoint” model [43].

## APPENDIX

### NOMENCLATURE

1.  $\mathcal{R}_{ds}$ : Reliability degree of single packet dissemination (Sections 3.2, 5, and 6).

2.  $\mathcal{R}_{dc}$ : Reliability degree of continuous packet dissemination (Sections 3.2, 5, and 6).
3.  $\mathcal{R}_{da}$ : Reliability degree of access (Sections 3.2, 5, and 6).
4.  $\mathcal{N}_l$ : Network load (Sections 3.2, 5, and 6).
5.  $\mathcal{F}$ : Cumulative distribution function of the reliability degree (Sections 3.2 and 5).
6.  $\lambda_o$ : Overall access rate to PILOT (Sections 3.2, 5, and 6).
7.  $\lambda_u$ : Update rate to PILOT (Sections 3.2, 5, and 6).
8.  $\lambda_q$ : Query rate to PILOT (Sections 3.2, 5, and 6).
9.  $F$ : Gossip fanout (Sections 4.2.1, 5, and 6).
10.  $\tau_q$ : Gossip quiescence threshold (Sections 4.2.1, 5, and 6).
11.  $\tau_a$ : Gossip age threshold (Sections 4.2.1, 5, and 6).
12.  $\xi_W, \hat{\xi}_W$ : Write quorum nominal size and real size, respectively (Sections 4.4, 5, and 6).
13.  $\xi_R, \hat{\xi}_R$ : Read quorum nominal size and real size, respectively (Sections 4.4, 5, and 6).
14.  $n$ : Group (or, in particular, STS) size (Sections 5 and 6).
15.  $H$ : Random variable representing the length of an arbitrary routing path (Section 5).
16.  $p, q$ : Probability of infection and noninfection, respectively (Section 5).
17.  $p_f$ : Failure probability for each hop (Section 5).
18.  $p_r$ : Probability of a query occurring within rounds  $r + 1$  after an update (Section 5).
19.  $p_e$ : Server unavailability in the case of query (Sections 5 and 6).
20.  $r, \tilde{r}$ : A certain gossip round and the final round, respectively (Section 5).
21.  $S_r$ : Number of infected nodes after round  $r$  (Section 5).
22.  $\nu_r, \nu$ : Distribution of  $S_r$  (also  $\hat{\xi}_W^r$ ) and its eventual value, respectively (Section 5).
23.  $\mu$ : Distribution of  $\hat{\xi}_R$  (Section 5).

## ACKNOWLEDGMENTS

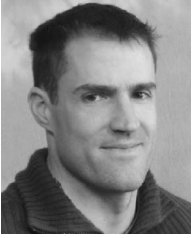
The authors would like to thank Milan Vojnovic, Matthias Grossglauser, and Pierre Brémaud for several instructive discussions. The work presented in this paper was supported (in part) by the National Competence Center in Research on Mobile Information and Communication Systems (NCCR-MICS), a center supported by the Swiss National Science Foundation under grant number 5005-67322. (<http://www.terminodes.org>)

## REFERENCES

- [1] I. Gupta, K.P. Birman, and R. van Renesse, "Fighting Fire with Fire: Using Randomized Gossip to Combat Stochastic Scalability Limits," *J. Quality and Reliability Eng. Int'l*, vol. 18, no. 3, pp. 165-184, 2002.
- [2] D. Powell et al. "Group Communication (special issue)," *Comm. ACM*, vol. 39, no. 4, pp. 50-97, 1996.
- [3] Z.J. Haas and B. Liang, "Ad Hoc Mobility Management with Uniform Quorum Systems," *IEEE/ACM Trans. Networking*, vol. 7, no. 2, pp. 228-240, 1999.
- [4] G. Pei and M. Gerla, "Mobility Management for Hierarchical Wireless Networks," *ACM/Kluwer Mobile Networks and Applications (MONET)*, vol. 6, no. 4, pp. 331-337, 2001.
- [5] L. Zhou and Z. Haas, "Securing Ad Hoc Networks," *IEEE Network*, vol. 13, no. 6, pp. 24-30, 1999.
- [6] S. Capkun, L. Buttyán, and J.-P. Hubaux, "Self-Organized Public-Key Management for Mobile Ad Hoc Networks," *IEEE Trans. Mobile Computing*, vol. 2, no. 1, pp. 52-64, Jan-Mar. 2003.
- [7] L. Buttyán and J.-P. Hubaux, "Report on a Working Session on Security in Wireless Ad Hoc Networks," *ACM SIGMOBILE Mobile Computing and Comm. Rev.*, vol. 6, no. 4, 2002.
- [8] N. Asokan and P. Ginzboorg, "Key-Agreement in Ad-Hoc Networks," *Computer Comm.*, vol. 23, no. 17, pp. 1627-1637, 2000.
- [9] S. Nesargi and R. Prakash, "MANETconf: Configuration of Hosts in a Mobile Ad Hoc Network," *Proc. INFOCOM*, pp. 1059-1068, 2002.
- [10] E.M. Royer and C.E. Perkins, "Multicast Operation of the Ad-Hoc On-Demand Distance Vector Routing Protocol," *Proc. MobiCom*, pp. 207-218, 1999.
- [11] J.-P. Hubaux, T. Gross, J.-Y. Le Boudec, and M. Vetterli, "Toward Self-Organized Mobile Ad Hoc Networks: The Terminodes Project," *IEEE Comm. Magazine*, vol. 39, no. 1, pp. 118-124, 2001.
- [12] J. Luo, P.Th Eugster, and J.-P. Hubaux, "Route Driven Gossip: Probabilistic Reliable Multicast in Ad Hoc Networks," *Proc. INFOCOM*, pp. 2229-2239, 2003.
- [13] J. Luo, J.-P. Hubaux, and P.Th Eugster, "PAN: Providing Reliable Storage in Mobile Ad Hoc Networks with Probabilistic Quorum Systems," *Proc. MobiHoc*, pp. 1-12, 2003.
- [14] M.G. Hayden, "The Ensemble System," PhD thesis, Dept. of Computer Science, Cornell Univ., 1997.
- [15] Y. Amir, C. Danilov, and J. Stanton, "A Low Latency, Loss Tolerant Architecture and Protocol for Wide Area Group Communication," *Proc. Int'l Conf. Dependable Systems and Networks*, pp. 327-336, 2003.
- [16] N. Malpani, N.H. Vaidya, and J.L. Welch, "Distributed Token Circulation in Mobile Ad Hoc Networks," *Proc. 10th IEEE Int'l Conf. Network Protocols*, pp. 4-13, 2002.
- [17] S. Dolev, E. Schiller, and J. L. Welch, "Random Walk for Self-Stabilizing Group Communication in Ad Hoc Networks," *Proc. 21st Symp. Reliable Distributed Systems*, pp. 70-79, 2002.
- [18] E. Vollset, "The Autograph Protocol—a Preliminary Report on a Reliable Broadcast Protocol for Mobile Ad-Hoc Networks," *Proc. Int'l Conf. Dependable Systems and Networks-Student Forum*, 2003.
- [19] G.-C. Roman, Q. Huang, and A. Hazemi, "Consistent Group Membership in Ad Hoc Networks," *Proc. ISCE*, 2001.
- [20] V. Hadzilacos and S. Toueg, "Fault-Tolerant Broadcasts and Related Problems," *Distributed Systems*, pp. 97-145, chap. 5, second ed., Addison-Wesley, 1993.
- [21] K.P. Birman, M. Hayden, O. Ozkasap, Z. Xiao, M. Budiu, and Y. Minsky, "Bimodal Multicast," *ACM Trans. Computer Systems*, vol. 17, no. 2, pp. 41-88, 1999.
- [22] P. Eugster, R. Guerraoui, S. Handurukande, A.M. Kermarrec, and P. Kouznetsov, "Lightweight Probabilistic Broadcast," *ACM Trans. Computer Systems*, 2003.
- [23] S. Floyd, V. Jacobson, C-G. Liu, S. McCanne, and L. Zhang, "A Reliable Multicast Framework for Light-Weight Sessions and Application Level Framing," *IEEE/ACM Trans. Networking*, vol. 5, no. 6, pp. 784-893, 1997.
- [24] S. Paul, K.K. Sabnani, J.C. Lin, and S. Bhattacharyya, "Reliable Multicast Transport Protocol," *IEEE J. Selected Areas Comm.*, vol. 15, no. 3, pp. 784-893, 1997.
- [25] R. Chandra, V. Ramasubramanian, and K. Birman, "Anonymous Gossip: Improving Multicast Reliability in Mobile Ad-Hoc Networks," *Proc. 21st IEEE Int'l Conf. Distributed Computing Systems*, pp. 275-283, 2001.
- [26] D. Barbara and H. Garcia-Molina, "The Reliability of Vote Mechanisms," *IEEE Trans. Computers*, vol. 36, no. 10, pp. 1197-1208, Oct. 1987.
- [27] F.B. Schneider, "Replication Management Using the State-Machine Approach," *Distributed Systems*, pp. 169-197, chapter. 6, second ed., Addison-Wesley, 1993.
- [28] D. Malkhi, M.K. Reiter, and A. Wool, "Probabilistic Quorum Systems," *Information and Computation*, vol. 170, no. 2, pp. 184-206, 2001.
- [29] Z.J. Haas and B. Liang, "Ad Hoc Mobility Management with Randomized Database Groups," *Proc. Int'l Conf. Comm.*, vol. 3, pp. 1756-1762, 1999.
- [30] M. Papadopoulou and H. Schulzrinne, "Effects of Power Conservation, Wireless Coverage and Cooperation on Data Dissemination among Mobile Devices," *Proc. MobiHoc*, pp. 117-127, 2001.
- [31] T. Hara, "Effective Replica Allocation in Ad Hoc Networks for Improving Data Accessibility," *Proc. INFOCOM*, pp. 1568-1576, 2001.
- [32] K.H. Wang and B. Li, "Efficient and Guaranteed Service Coverage in Partitioned Mobile Ad-Hoc Networks," *Proc. INFOCOM*, pp. 1089-1098, 2002.
- [33] D.B. Johnson, D.A. Maltz, and Y-C. Hu, "The Dynamic Source Routing Protocol for Mobile Ad Hoc Networks(DSR)," Internet-Draft, draft-ietf-manet-dsr-08. txt, work in progress, Feb. 2003.
- [34] K. Xu, X. Hong, and M. Gerla, "An Ad Hoc Network with Mobile Backbones," *Proc. Int'l Conf. Comm.*, vol. 5, pp. 3138-3143, 2002.
- [35] R. Sivakumar, P. Sinha, and V. Bharghavan, "CEDAR: A Core-Extraction Distributed Ad Hoc Routing Algorithm," *IEEE J. Selected Areas in Comm.*, special issue on ad-hoc routing, vol. 17, no. 8, pp. 1454-1465, 1999.
- [36] A.W. Fu and D.W. Cheung, "A Transaction Replication Scheme for a Replicated Database with Node Autonomy," *Proc. Very Large Data Bases Conf.*, pp. 214-225, 1994.
- [37] L.-G. Alberto and I. Widjaja, *Communications Networks*. McGraw Hill Higher Education, 2000.
- [38] J.D. Murray, *Mathematical Biology*. second ed. Berlin: Springer, 1993.
- [39] "The ns Manual," *The VINT Project*, K. Fall and K. Varadhan, eds., UC Berkeley, LBL, USC/ISI, and Xerox PARC, Apr. 2002, available from <http://www.isi.edu/nsnam/ns/>.
- [40] T.S. Rappaport, *Wireless Communications: Principles and Practice*, second ed. Upper Saddle River, N.J.: Prentice Hall, 2002.
- [41] D.B. Johnson and D.A. Maltz, "Dynamic Source Routing in Ad Hoc Wireless Networks," *Mobile Computing*, T. Imielinski and H. Korth, eds., chap. 5, pp. 153-181, Kluwer Academic Publishers, 1996.
- [42] A.-M. Kermarrec, L. Massoulié, and A. Ganesh, "Probabilistic Reliable Dissmination in Large-Scale Systems," *IEEE Trans. Parallel and Distributed Systems*, vol. 14, no. 3, pp. 248-258, Mar. 2003.
- [43] J. Yoon, M. Liu, and B. Noble, "Random Waypoint Considered Harmful," *Proc. INFOCOM*, pp. 1312-1321, 2003.



**Jun Luo** received the BS and MS degrees both in electrical engineering from Tsinghua University, Beijing, PRC, in 1997 and 2000, respectively. As a research assistant of Laboratory for Computer communication and Application (LCA), he is now working toward the PhD degree in computer science in the Swiss Federal Institute of Technology (EPFL). His research interests include multicasting, mobile computing (especially in ad hoc networks), reliable group communication, and network security. He is a student member of the ACM. For more information, check <http://lcawww.epfl.ch/luo>.



**Patrick Eugster** received the master's and PhD degrees, both from the Swiss Federal Institute of Technology in Lausanne (EPFL). After having worked for some time at Chalmers University of Technology in Goteborg, Sweden, he returned to EPFL as postdoctoral researcher. His research focuses on algorithms and programming abstractions for reliable distributed systems.



**Jean-Pierre Hubaux** joined the faculty of the Swiss Federal Institute of Technology-Lausanne (EPFL) in 1990; he was promoted to full professor in 1996. His research activity is focused on mobile networking and computing, with a special interest in fully self-organized wireless ad hoc networks. In particular, he has performed research on cooperation aspects, security, power efficiency, and distributed algorithms for ad hoc and sensor networks. During the last few years, he has been strongly involved in the definition and launching phases of a new National Competence Center in Research named "Mobile Information and Communication Systems" (NCCR/MICS), see <http://www.terminodes.org>. He served as the general chair for the Third ACM Symposium on Mobile Ad Hoc Networking and Computing (MobiHoc 2002), held in June on the EPFL campus. He is an associate editor of the *IEEE Transactions on Mobile Computing* and of the *Elsevier Journal on Ad Hoc Networks*. Within his first year at EPFL, he defined the first curriculum in communication systems. From October 1999 until September 2001, he was the first chairman of the Communication Systems Department. He has held visiting positions at the IBM T.J. Watson Research Center and at the University of California at Berkeley. He has published more than 60 papers in the area of networking. In a former life, he spent 10 years in France with Alcatel, where he was involved in R&D activities, mostly in the area of switching systems architecture and software. He is a senior member of the IEEE. For more information, check <http://www.lcawww.epfl.ch/hubaux>.

▷ **For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).**