

# Araneola: A Scalable Reliable Multicast System for Dynamic Environments

Roie Melamed

CS Department, Technion  
mroi@cs.technion.ac.il

Idit Keidar

EE Department, Technion  
idish@ee.technion.ac.il

## Abstract

We present Araneola\*, a scalable reliable application-level multicast system for highly dynamic wide-area environments. Araneola supports multi-point to multi-point reliable communication in a fully distributed manner while incurring constant load on each node. For a tunable parameter  $k \geq 3$ , Araneola constructs and dynamically maintains an overlay structure in which each node's degree is either  $k$  or  $k + 1$ , and roughly 90% of the nodes have degree  $k$ . Empirical evaluation shows that Araneola's overlay structure achieves three important mathematical properties of  $k$ -regular random graphs (i.e., random graphs in which each node has exactly  $k$  neighbors) with  $N$  nodes: (i) its diameter grows logarithmically with  $N$ ; (ii) it is generally  $k$ -connected; and (iii) it remains highly connected following random removal of linear-size subsets of edges or nodes. The overlay is constructed at a very low cost: each join, leave, or failure is handled locally, and entails the sending of only about  $3k$  messages in total.

Given this overlay, Araneola disseminates multicast messages by gossiping over the overlay's links. We show that compared to a standard gossip-based multicast protocol, Araneola achieves substantial improvements in load, reliability, and latency. Finally, we present an extension to Araneola in which the basic overlay is enhanced with additional links chosen according to geographic proximity and available bandwidth. We show that this approach reduces the number of physical hops messages traverse without hurting the overlay's robustness.

---

\* Araneola means "little spider" in Latin.

# 1 Introduction

Our goal is to provide a scalable multi-point to multi-point reliable multicast service for very large groups in wide-area networks. A protocol deployed in such a setting must be able to withstand frequent node failures as well as non-negligible message loss rates [21]. Moreover, studies have shown that users typically frequently join and leave multicast sessions [1]; such behavior is called churn. A major design goal for our work is therefore coping efficiently with churn. Specifically, we address the following challenges: (i) providing high reliability despite considerable message loss and failure rates while incurring constant load on each node; (ii) incorporating joining nodes and removing leaving (or failing) ones with a low *constant* overhead; and (iii) providing an undisrupted service to nodes that are up despite high churn rates.

We present Araneola, a scalable reliable application level multicast (ALM) system for dynamic wide-area environments. Reliability is achieved by constructing a richly-connected overlay and disseminating pertinent information on multiple paths in this overlay. The number of paths in the overlay can be tuned according to the expected failure and loss rates. Araneola is designed to incur small constant load on each node. To this end, it builds an overlay in which each node’s degree is bounded by a small constant. This approach has three advantages: (i) all nodes, including low bandwidth ones, are capable of participating in the overlay; (ii) the load on all nodes is similar, so no user is required to contribute more bandwidth than its fair share; and (iii) nodes have ample remaining bandwidth for connecting to nearby nodes, as we explain below.

Our search for a constant degree overlay leads us to consider  $k$ -regular random graphs. In a  $k$ -regular graph, each node’s degree is  $k$ . A  $k$ -regular random graph with  $N$  nodes is a graph chosen uniformly at random from the set of  $k$ -regular graphs with  $N$  nodes. For  $k \geq 3$ , a  $k$ -regular random graph is almost always a good expander [10], which implies that (i) its diameter grows logarithmically with  $N$  [26]; and (ii) it remains connected after random failures of a linear subset of its nodes and/or edges [11]. In addition, such a graph is generally  $k$ -connected, i.e., at least  $k$  nodes need to be removed in order to cause a partition [26]<sup>2</sup>. We strive to construct and maintain an overlay that resembles a  $k$ -regular random graph: Araneola’s basic overlay converges to a graph in which each node has a degree of either  $k$  or  $k + 1$  and no two neighboring nodes have a degree of  $k + 1$ . Empirically, we show that Araneola’s overlay achieves the desired properties above, namely logarithmic diameter,  $k$ -connectivity, and high robustness.

We construct and maintain the overlay at a very low constant overhead; each join or leave (or failure) incurs sending roughly about  $3k$  messages in a  $k$ -degree overlay, regardless of the number of nodes. Remarkably, in dynamic settings, the cost of handling a single join or leave operation *decreases* as the churn rate increases. This is in contrast to virtually all existing structured peer-to-peer overlays, with which the overhead for handling joins grows logarithmically with  $N$ .

The low maintenance cost is achieved due to the facts that: (i) each join, leave, or failure is handled locally; and (ii) the selection of random neighbors uses partial membership views maintained by a distributed low cost membership service similar to the ones in [9, 24]. The overhead of the membership service is independent of the number of nodes and of the churn rate.

Having built a basic  $k/k+1$ -degree overlay, we next extend it by adding links between geographically-close nodes. The low degree of Araneola’s basic structure allows for allocating plenty of bandwidth for communication with proximate nodes. We show that with this approach, the links in Araneola’s overlay traverse substantially fewer physical hops on average. Moreover, the overlay’s robustness does not deteriorate.

Given Araneola’s overlay, multicast messages are disseminated through gossip between each pair of neighbors. Gossiping in Araneola differs from a standard gossip protocol (e.g., [7, 9, 18]) in that with a standard gossip protocol, each node chooses *different* random nodes to gossip with in each round, whereas in Araneola, each node always gossips with its neighbors in the overlay. We show that this difference leads

---

<sup>2</sup>The probability that a  $k$ -regular random graph is not  $k$ -connected is bounded by  $O(N^{2-k})$ .

to substantial improvements in load, reliability, and latency.

In summary, our contributions include:

- the first algorithm for constructing and maintaining a richly-connected low degree overlay structure in which each join or leave operation incurs a constant overhead;
- the first overlay-based multicast system to provide an uninterrupted multicast service in highly dynamic settings while incurring constant load on each node;
- a complete implementation and a thorough evaluation of Araneola running up to 10,000 nodes on up to 125 machines, in both LAN and WAN, including the first extensive evaluation of the impact of churn on an ALM system; and
- an overlay that allocates ample bandwidth for each node for communication with proximate nodes.

This paper proceeds as follows: Section 2 discusses related work. In Section 3, we summarize our design goals. Section 4 presents Araneola’s design and pseudo code, and Section 5 empirically evaluates Araneola. Section 6 presents and evaluates the extension that exploits network proximity. Section 7 concludes.

## 2 Related Work

In recent years, two leading approaches for supporting scalable ALM in dynamic failure-prone networks have emerged: gossip-based (or epidemic) multicast protocols (e.g., [3, 9, 7]) and dynamic overlay networks (e.g., [12, 5, 24, 8]). With gossip-based protocols, each node periodically chooses other random nodes to propagate the information to. Gossip-based multicast generally achieves good load balancing, high reliability, and uninterrupted service in the presence of message loss and node joins and leaves, but it also induces a high load, as many duplicate messages are sent [9].

Overlay-based ALM systems usually disseminate messages on a tree structure [12, 5, 24, 8]. With tree-based multicast, no duplicate messages are sent. However, in the presence of churn, the tree structure will frequently become partitioned, causing a significant portion of the multicast messages to be lost. Therefore, in order to achieve reliability, such protocols need to detect message loss and recover from it. This can cause recovered messages to be significantly delayed; can induce substantial overhead, especially if failures are frequent; and can inhibit scalability. A second problem with tree-based multicast is uneven load distribution: as recently argued in [4], inner nodes in the tree carry the burden for the multicast, whereas leaf nodes do not share the load. Two recent projects, SplitStream [4] and Bullet [16], address this issue and build a balanced multicast infrastructure; however these two systems are intended for single-source multimedia transfer and do not strive to provide multi-point to multi-point communication or full reliability of message delivery as we do.

Recently, several peer-to-peer overlays with logarithmic diameters and a bounded node degrees have been suggested, e.g., emulating the Butterfly [19], de Bruijn graphs [13], Small Worlds graphs [15], or random expander graphs with degrees  $\geq 8$  [17]. However, none of these systems can guarantee, with high probability, a lower cost than  $O(\log N)$  messages and time for handling joins, since a joining node must search and locate its (random or hashed) joining location prior to joining the system. Chawathe et al. [6] have argued that this logarithmic cost inhibits the scalability of such systems assuming the churn rates measured in Gnutella and Napster [23].

Like our extension of Araneola, several overlay structures, e.g., [20, 24], reduce message delivery latency and communication costs by incorporating links between nearby nodes in addition to the random links required for achieving a good overlay. In comparison with [20, 24], Araneola achieves a smaller average degree than [20, 24] and better load balancing than [24].

### 3 Design Goals

The purpose of Araneola is to support scalable reliable multi-point to multi-point communication in dynamic wide-area settings where nodes frequently join and leave (or fail). We have set the following requirements for our service: (i) high reliability and graceful degradation in the face of increasing failure rates; (ii) low latency; (iii) low constant load on each node; (iv) low constant cost for handling joins and failures; and (v) quick failure recovery and prompt incorporation of joining nodes. Araneola is designed to achieve these goals without using any infrastructure, servers, or any elaborate communication mechanism beyond point-to-point UDP communication. We assume that every pair of nodes can communicate with each other.

Araneola thus strives to build an overlay structure with the following characteristics: (i) multiple disjoint paths between every pair of nodes, where the number of paths is a configurable parameter  $k$ ; (ii) robustness to random removal of a certain percentage of the nodes or edges; (iii) low diameter and average distance; (iv) bounded degree (of  $k + 1$ ); and (v) support for local addition and removal of nodes at a constant cost. These characteristics are naturally achieved by  $k$ -regular random graphs. Empirically, we show that these desirable characteristics are also achieved by Araneola’s overlay structure.

### 4 Araneola’s Design

Araneola builds an overlay structure for each multicast group. Since each group is handled independently, we present the protocol for a single group, and omit the group’s name. All Araneola nodes run the same code. The code has two main components: one constructs and maintains the overlay, as described in Section 4.1, and the second implements the multicast service, as described in Section 4.2.

When joining the overlay, a node randomly selects several other nodes to connect to. This requires each node to know some other nodes’ identities. To this end, we implement a scalable randomized membership protocol similar to [9, 24], where membership information is gossiped over the overlay’s links. Each node maintains a small set of node identities, called a *membership view*, which evolves over time. Periodically, the membership protocol piggybacks a small amount of information on gossip messages. Views have been shown to become uniformly distributed over time [9, 24]. Experimentally, we observe that it suffices to piggyback membership information infrequently, e.g., once a minute. When a new node joins for the first time, it can ask one other for its membership view, and use that as its initial view. We do not detail the membership protocol in this paper; the interested reader is referred to [9].

Araneola’s data structures are presented in Fig. 1. The set *neighbors* holds the node’s current neighbors in the overlay, with their respective degrees and heartbeats (for failure detection purposes). The degree of a node is the size of its neighbors set, i.e.,  $|neighbors|$ . The set *missing\_msg* holds identifiers of messages that the node heard of but did not receive. A function *heard\_from* maps each identifier in this set to nodes from which it was heard. *recent\_mids* holds identifiers of messages received in the latest gossip round. *last\_connect\_to\_time* records the latest time a CONNECT\_TO message was handled, as explained below. The parameter  $L$  determines the graph’s target degree ( $k$ ), and  $H$  defines the maximum allowed node degree. A number of timeout values are defined in order to control the frequency at which different events occur.

#### 4.1 Building and Maintaining the Overlay

Three tasks participate in the construction and maintenance of the overlay: (i) the *connect task* (see Fig. 2) adds new connections when a node’s degree is below  $L$ ; (ii) the *disconnect task* (see Fig. 3) tries to reduce the node’s degree if it is above  $L$ , without causing any node’s degree to drop below  $L$ ; and (iii) the *failure detector* (see Fig. 2) detects failures and recovers from them.

**Data structures:**

*id* – this node’s identifier.  
*neighbors* – set of triples (id,degree,heartbeat), initially  $\emptyset$ .  
*messages* – queue of messages tagged with *m.id*, initially  $\emptyset$ .  
*missing\_msg* – set of messages identifiers, initially  $\emptyset$ .  
*heard\_from:missing\_msgs*  $\rightarrow$  list of nodes.  
*recent\_mids* – set of messages identifiers, initially  $\emptyset$ .  
*last\_connect\_to\_time* – a time.

**Constants:**

*L* – target number of neighbors.  
*H* – upper bound on the number of neighbors.  
 MIN\_HEARTBEAT – failure detection parameter.  
 Timeouts: *connect\_timeout*, *failure\_detection\_timeout*,  
*disconnect\_timeout*, *connect\_to\_timeout*,  
*gossip\_round\_timeout*.

Figure 1: Araneola’s data structures and constant definitions.

When a node’s degree is below *L*, the connect task periodically attempts to set up as many new connections as it is missing to randomly chosen nodes (lines 1–7). The target nodes are chosen at random from the local membership view. For each attempted connection, the node sends a CONNECT request (line 5). At bootstrap time, the node issues CONNECT requests to *L* nodes, and then sleeps for *connect\_timeout*. It is expected that during this period enough new connections will be formed, although since some of the chosen nodes may be faulty or overloaded, there may be a need to attempt more connections after the timer expires. The connect task can be awakened by other tasks before the timer expires (line 36).

A node that receives a CONNECT request (line 8) *accepts* it, by calling *add\_connection*, provided that its degree is smaller than *H*, and otherwise it *redirects* the request, as will be explained shortly. *add\_connection* adds the sender to *neighbors* (line 30) and responds with a CONNECT\_OK. Upon receiving the CONNECT\_OK (line 15), the requester registers the new connection, unless its degree has already reached *H*, in which case it sends a LEAVE message (line 19). A LEAVE message causes its receiver to remove its connection with the sender. Redirecting is done by sending a REDIRECT message to the requester, naming the sender’s lowest degree neighbor (line 12). This causes the requester to send a new CONNECT request to the named neighbor (line 14). CONNECT and CONNECT\_OK messages carry the sender’s current degree for initializing the *degree* in the *neighbors* data structure. In addition, every node periodically sends its degree to its neighbors (this is not shown in the code).

A node that voluntarily leaves the system sends a LEAVE message to all its neighbors. In case a node fails, the failure detector task (Fig. 2, lines 22–29) eventually detects the failure. Every node keeps a heartbeat counter for each one of its neighbors. If a node receives less than MIN\_HEARTBEAT messages from a neighbor *n* during an interval of length *failure\_detection\_timeout*, it removes *n* and sends a LEAVE message to it. The heartbeat counter of *n* is increased whenever a message from it is received (this not shown in the code).

There are two rules for removing connections: *Rule 1* and *Rule 2*. Rule 1 removes the connection between a pair nodes that both have degrees higher than *L*. Specifically, if a node *n*’s degree is  $L + i$ , then *n* attempts to remove *i* of its neighbors. The *i* neighbors with the highest degrees are candidates for removal; they are inserted into the set *cands* (line 5). Nodes with degrees  $\leq L$  are then deleted from *cands* (line 8). If *n* has a higher id than a node *c* in *cands*, then *n* sends a DISCONNECT message to *c* (line 10). Upon receiving this message (line 16), if *c*’s degree is still higher than *L*, it removes the connection with *n*, and sends a DISCONNECT\_OK message. Upon receipt of a DISCONNECT\_OK (line 20), *n* removes the connection with *c*. Note that Rule 1 never reduces a node’s degree to be below *L*.

With Rule 1, it is still possible for a node to have degree *H* while all of its neighbors have degree *L*. Rule 2 is only invoked at a node *n* when all of *n*’s neighbors’ degrees are  $\leq L$ . With Rule 2, node *n* chooses its two neighbors with the highest and lowest degrees, *h* and *l*, resp. (lines 12–13). If *n*’s degree is at least *l.degree* + 2, then *n* tries to cause *h* to shift one of its connections from *n* to *l*. But before removing *h*’s connection with *n*, we ensure that *l* is willing to accept *h*’s connection. Therefore, *n* contacts *l* (rather than *h*) and asks it to try to connect to *h*, and to ask *h* to remove its connection with *n*. To this

**Connect task:**

```

1. loop forever
2. gap ← L - |neighbors|
3. while gap > 0
4.   n ← random node
5.   send ⟨CONNECT, |neighbors|⟩ to n
6.   gap ← gap - 1
7.   sleep (connect_timeout)

Event handles:
8. upon receive ⟨CONNECT, d⟩ from n do
9.   if |neighbors| < H then
10.    add_connection(n, d, true)
11.  else
12.    send ⟨REDIRECT, lowest degree neighbor⟩ to n

13. upon receive ⟨REDIRECT, n'⟩ from n do
14.  send ⟨CONNECT, |neighbors|⟩ to n'

15. upon receive ⟨CONNECT_OK, d⟩ from n do
16.  if |neighbors| < H then
17.    add_connection(n, d, false)
18.  else
19.    send ⟨LEAVE⟩ to n

20. upon receive ⟨LEAVE⟩ from n do
21.  remove_connection(n)

```

**Failure detector task:**

```

22. loop forever
23.  sleep (failure_detection_timeout)
24.  foreach n ∈ neighbors
25.    if n.heartbeat < MIN_HEARTBEAT then
26.      send LEAVE to n
27.      remove_connection(n)
28.    else
29.      n.heartbeat ← 0

Procedures:
Procedure add_connection (node_id n, int d, boolean ack)
30. neighbors ← neighbors ∪ {n, d, MIN_HEARTBEAT}
31. if ack = true then
32.  send ⟨CONNECT_OK, |neighbors|⟩ to n

Procedure remove_connection (node n)
33. remove n from neighbors
34. remove n from all heard_from lists
35. if |neighbors| < L then
36.  wake up connect task

```

Figure 2: Overlay construction: connect and failure detector tasks.

end,  $n$  sends a  $\langle \text{CONNECT\_TO}, h \rangle$  message to  $l$ . If upon receiving this message  $l$ 's degree is still  $\leq L$ , and  $l$  has not handled another  $\text{CONNECT\_TO}$  request in the last  $\text{connect\_to\_timeout}$  (lines 23–24), then  $l$  sends a  $\text{CHANGE\_CONNECTION}$  message to  $h$ . The recipient,  $h$ , connects to  $l$  (line 29) and sends a  $\text{DISCONNECT}$  message to  $n$  (line 30). This can increase  $l$ 's degree, but not to become higher than  $L+1$ , since  $l$  handles at most one  $\text{CONNECT\_TO}$  request at a time, and only if its degree  $\leq L$ . Moreover, note that if  $l$ 's degree will become higher than  $L$ , and  $n$ 's degree will remain above  $L$ , then Rule 1 will eventually reduce  $l$ 's degree back to  $L$ .

**Proposition 1.** *If from some point onward no nodes join, leave, or are detected as faulty, then each node's degree is eventually either  $L$  or  $L+1$ , and at most 50% of the nodes have degree  $L+1$ .*

**Proof (sketch):** The connect task ensures that eventually, each node's degree is between  $L$  and  $H$ . Rule 1 removes the connection between every two neighbors with degrees  $> L$ , without reducing any node's degree below  $L$ . Thus, barring additional joins, leaves or detected failures, Rule 1 ensures that the overlay converges to a state in which each node's degree is between  $L$  and  $H$  and no two neighboring nodes have a degree  $> L$ . This implies that at least 50% of the nodes have a degree of  $L$ . However, with Rule 1, it is still possible for a node  $n$  to have a degree  $> L+1$  when all of  $n$ 's neighbors have a degree of  $L$ . In this case, Rule 2 is invoked at node  $n$ , reducing  $n$ 's degree by one and increasing the degree of  $n$ 's lowest degree neighbor  $l$  by one without changing the rest of  $n$ 's neighbors' degrees. Now,  $l$ 's degree becomes  $L+1$  and Rule 1 becomes enabled again, removing the connection between  $n$  and  $l$ . Thus, after activating Rule 2 and Rule 1 consecutively,  $n$ 's degree is reduced by 2 while the degrees of the rest of  $n$ 's neighbors remain  $L$ . If  $n$ 's degree is still above  $L+1$ , subsequent activations of the two reduction rules continue to reduce  $n$ 's degree until it becomes either  $L$  or  $L+1$ .

**Disconnect task:**

```

1. loop forever
2. sleep (disconnect_timeout)
3.  $i \leftarrow |\text{neighbors}| - L$ 
4. if  $i > 0$  then
    /* Rule 1 */
5.    $\text{cands} \leftarrow$  set of  $i$  neighbors with highest degrees
6.   foreach  $c \in \text{cands}$ 
7.     if  $c.\text{degree} \leq L$  then
8.        $\text{cands} \leftarrow \text{cands} \setminus \{c\}$ 
9.     else if  $c.\text{id} < id$  then
10.      send (DISCONNECT) to  $c$ 
    /* Rule 2 */
11.  if  $\text{cands} = \emptyset$  then
12.     $h \leftarrow$  neighbor with highest degree
13.     $l \leftarrow$  neighbor with lowest degree
14.    if  $|\text{neighbors}| \geq l.\text{degree} + 2$  then
15.      send (CONNECT_TO,  $h$ ) to  $l$ 

```

## Event handlers:

```

16. upon receive (DISCONNECT) from  $n$  do
17.   if  $|\text{neighbors}| > L$  then
18.     send (DISCONNECT_OK) to  $n$ 
19.     remove_connection( $n$ )
20. upon receive (DISCONNECT_OK) from  $n$  do
21.   remove_connection( $n$ )
22. upon receive (CONNECT_TO,  $n'$ ) from  $n$  do
23.   if  $|\text{neighbors}| \leq L \wedge$ 
24.      $\text{clock} - \text{last\_connect\_to\_time} > \text{connect\_to\_timeout}$  then
25.     send (CHANGE_CONNECTION,  $|\text{neighbors}|, n$ ) to  $n'$ 
26.      $\text{last\_connect\_to\_time} \leftarrow \text{clock}$ 
27. upon receive (CHANGE_CONNECTION,  $d, n'$ ) from  $n$  do
28.   if  $|\text{neighbors}| < H$  then
29.     add_connection( $n, d, \text{true}$ )
30.     send (DISCONNECT) to  $n'$ 

```

Figure 3: Reducing node degrees.

## 4.2 Gossip-Based Multicast

**Gossip task:**

```

1. loop forever
2. sleep (gossip_round_timeout)
   /* Send gossip messages to neighbors */
3. foreach  $n \in \text{neighbors}$ 
4.   create new gossip message  $m$ , with new  $m.\text{id}$ 
5.    $m.\text{degree} \leftarrow |\text{neighbors}|$ 
6.    $m.\text{ids} \leftarrow \{\text{recent\_mids} \mid \text{source} \neq n\}$ 
7.    $m.\text{reqs} \leftarrow \emptyset$ 
8.   foreach  $mid \in \text{missing\_msgs}$ 
9.     if  $\text{heard\_from}(mid).\text{first} = n$  then
10.       $m.\text{reqs} \leftarrow \text{reqs} \cup \{mid\}$ 
11.   send (GOSSIP,  $m$ ) to  $n$ 
   /* Update data structures */
12. move 1st element of each  $\text{heard\_from}(mid)$  list to end
13.  $\text{recent\_mids} \leftarrow \emptyset$ 

```

## Event handlers:

```

14. upon receive (GOSSIP,  $m$ ) from  $n$  do
15.   foreach  $id \in m.\text{ids} \wedge id \notin \text{messages}$ 
16.      $\text{missing\_msgs} \leftarrow \text{missing\_msgs} \cup \{id\}$ 
17.     append  $n$  to  $\text{heard\_from}(id)$ 
   /* Send requested messages to  $n$  */
18.   foreach  $r \in \text{reqs}$ 
19.     send (DATA, message with identifier =  $r.\text{id}$ ) to  $n$ 
20. upon receive (DATA,  $m$ ) from  $n$  do
21.    $\text{messages.enqueue}(m)$ 
22.    $\text{missing\_msgs.remove}(m.\text{id})$ 
23.    $\text{recent\_mids} \leftarrow \text{recent\_mids} \cup \{m.\text{id}\}$ 

```

Figure 4: Gossip-based multicast.

Each Araneola node gossips about recent messages identifiers with its neighbors and requests missing messages from them. The *gossip task* is presented in Fig. 4. Every *gossip\_round\_timeout*, a node sends a gossip message to each of its neighbors. A gossip message  $m$  sent by a node  $a$  to its neighbor  $n$  is identified by a message identifier,  $m.\text{id}$ , which includes  $a$ 's identifier (e.g., IP address and port) and a one byte serial number (cyclic counter). The field  $m.\text{degree}$  holds  $a$ 's current degree (line 5). The set  $m.\text{ids}$  includes message identifiers that  $a$  has received in the last gossip round and has not heard about from  $n$  (line 6). Finally,  $m.\text{reqs}$  are message identifiers that  $a$  is requesting from  $n$ . After sending the gossip messages, the first element in each *heard\_from* list is moved to the end of that list (line 12), in order to vary the node from which the message is requested.

When a node  $a$  receives a gossip message  $m$  from neighbor  $n$ , for each  $id$  in  $m.\text{ids}$  that is not in the

*messages* buffer, *id* is inserted into *missing\_msgs* (line 16) and *n* is appended to *heard\_from(id)* (line 17). Then, *a* sends to *n* all the messages requested in *m.reqs*. When a data message arrives, it is enqueued in *messages*, removed from *missing\_msgs*, and its identifier is inserted into *recent\_mids* (lines 21–23). Periodically, old messages are purged from *messages* and *missing\_msgs*. This garbage collection mechanism is straightforward, and is omitted from the pseudo code.

#### 4.2.1 Eliminating partitions

The probability that the overlay will become partitioned is negligible; in our dozens of experiments with thousands of nodes, the overlay never became partitioned. Nevertheless, should a partition occur, we implement a simple mechanism for detecting this situation and recover from it. This mechanism is not included in the pseudo code. It works as follows: at the end of each multicast session, every node *a* saves its local view to a log file. Occasionally, during a multicast session, *a* chooses a random non-neighbor node *n* from its log file and sends *n* a gossip message *m*. Upon receiving a message not from a neighbor, *n* waits a number of rounds that is larger than the maximum hop-count with which it receives multicast messages (the hop-count is piggybacked on every multicast message), and then checks if it has received the identifiers in *m.ids* from one of its neighbors. If *n* did not receive any of these identifiers, then it deduces that the overlay is partitioned. In this case, *n* connects to *a* in order to merge the two connected components.

## 5 Evaluation

We have implemented Araneola in Java using UDP/IP. In this section, we evaluate Araneola on a single LAN in Netbed [25]. In the next section, where we consider exploiting network proximity, we will evaluate Araneola also on a WAN. We run multiple Araneola nodes per machine, and therefore need to space the gossip rounds sufficiently as to allow all the nodes running on the same machine to complete their gossip operation during a round. Thus, we chose a fairly large round duration of 5 seconds. The *disconnect\_timeout* is set to 30 seconds, and the *connect\_timeout* is 20 seconds. We begin our study, in Section 5.1, by evaluating Araneola’s scalability and performance in a static setting. In Section 5.2 we study Araneola’s fault-tolerance. In Section 5.3, we consider high churn.

### 5.1 Static Evaluation

In our static evaluation, all the nodes are created simultaneously, and remain up throughout the experiment. In each round, a single data message is injected into the system, each time from a different machine. A total of 200 data messages are sent in each experiment. Each experiment (with a given number of nodes and choice of parameter settings) was run at least 4 times, for a total of several dozens.

#### 5.1.1 The impact of L

Araneola’s parameter L is very significant: it affects Araneola’s load, latency, and robustness. Fig. 5 shows the impact of this parameter on message propagation rates in runs with 8000 nodes (on 100 Netbed machines) and values of L ranging from 3 to 10. Each curve in the figure depicts the CDF of the average number of nodes that receive a message by each hop count for a given value of L. As expected, the message latency decreases as L increases.

In most of the experiments below, we set L to 5. We chose this value because it provides a good balance between the desired properties: each node sends a given message identifier only 4 times, and the latency to reach all nodes is reasonable: 7 rounds with 1000 nodes and 9 with 10,000. Moreover, as we shall see below,



it yields a highly robust overlay and achieves 100% reliability at churn rates exceeding those measured on the MBone [1]. We choose  $H$  to be  $L+5$  because we have observed that this choice achieves low overhead.

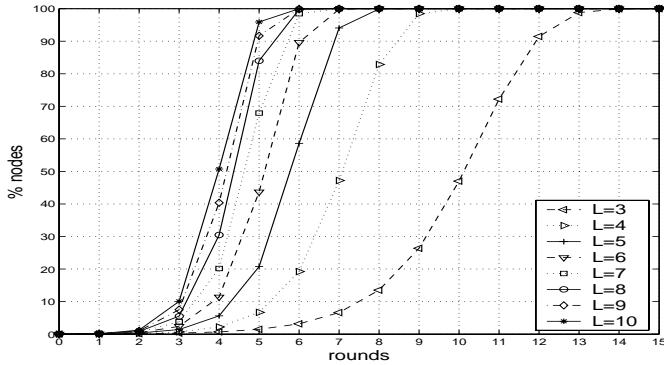


Figure 5: CDF: average % of nodes that receive a message for different degree Araneola overlays, 8000 nodes.

N	% nodes degree= 5	diameter	avg distance	avg # paths
500	91.8	6–7	4.18	5.01
1000	91.4	7	4.69	5.01
2000	92	7–8	5.16	5.01
4000	91.45	8	5.63	5.01
6000	90.42	8–9	5.93	5.01
8000	90.33	9	6.12	5.01
10000	90.36	9	—	—

Figure 6: Scalability of Araneola’s overlay.

### 5.1.2 Overlay properties and scalability

In order to understand Araneola’s scalability, we vary  $N$ , the group size, from 500 nodes (on 10 Netbed machines) to 10,000 nodes (on 125 Netbed machines).  $L$  and  $H$  are set to 5 and 10, resp. At the end of each experiment, we take a snapshot of the overlay structure, and then analyze its properties offline. We measure node degrees as well as the overlay’s diameter, average distance, and connectivity. The results are summarized in Table 6.<sup>3</sup> The first column shows the percentage of nodes whose degree is  $L$  (i.e., 5). The remaining nodes’ degrees are  $L+1$ . In all of our experiments, Araneola converges to a state in which over 90% of the nodes have degree  $L$ . The exact percentage of nodes with degree  $L$  does not seem related to  $N$ . The next column presents the smallest and largest measured diameters for every value of  $N$ . In cases where the same diameter was measured in all experiments, we present only one value. The diameter gives a measure for the *worst case* latency (in the absence of failures and message loss), whereas the average latency depends on the average distance between two nodes in the overlay. This average is presented in the following column, and it increases gradually with  $N$ . Finally, we measure the overlay’s connectivity. In over 90% of our experiments, the overlay is 5-connected, i.e., there are at least 5 disjoint paths between every pair of nodes. In the few cases where the connectivity was less than 5, there were at most 4 nodes with a connectivity of 4, whereas the rest of the nodes had a connectivity of 5. The average number of node-disjoint paths between every pair of nodes is presented in the last column. It does not vary with  $N$ .

Fig. 7 depicts the message propagation rates measured for values of  $N$  ranging from 500 to 10,000. As  $N$  increases, messages take longer to propagate, but the slow-down is gradual. The average hop-count in each experiment is very close to the average distance in that experiment’s overlay (see Table 6, 3rd column).

### 5.1.3 Comparison with gossip protocol

We now compare Araneola to a standard gossip protocol [18] implemented using Araneola’s gossip-based multicast module. The gossip protocol takes a parameter  $F$ , its fan-out. Where an Araneola node sends gossip messages to its neighbors, the gossip protocol sends gossip messages to  $F$  randomly selected nodes from its membership view. Whereas Araneola sends each message identifier downstream only, the gossip

<sup>3</sup>We did not analyze the average distance and connectivity for the experiments with  $N = 10,000$ .

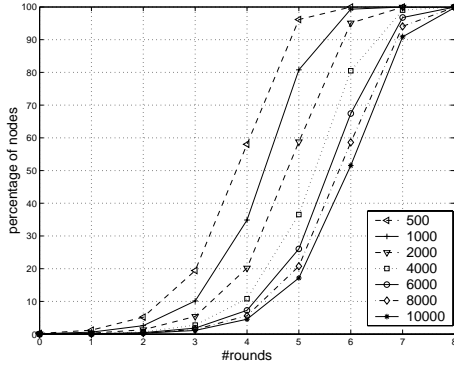


Figure 7: Message propagation rates for different group sizes,  $L=5$ .

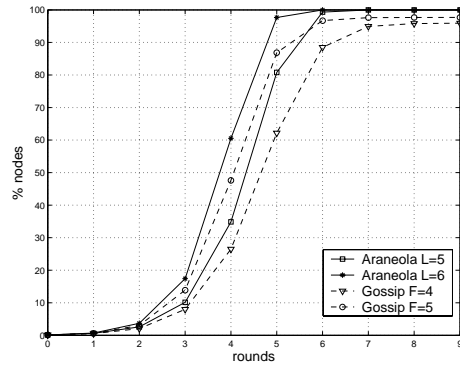


Figure 8: Araneola versus gossip, 1000 nodes.

protocol sends all its *recent\_mids* to all the chosen targets. Thus, gossip protocol instantiated with a fan-out of  $F$  sends information as many times as Araneola with  $L = F + 1$ .

We experiment with 1000 nodes on 20 Netbed machines. In each experiment, 400 messages are sent. Fig. 8 compares the average message propagation rates of Araneola with  $L=5$  and 6 to those of the gossip protocol with the corresponding fan-outs  $F=4$  and 5. Evidently, Araneola propagates information much more effectively than the gossip protocol. Initially, the propagation rates are similar, but after about 6 rounds, Araneola continues to effectively propagate the message, while the gossip protocol tapers off. Araneola succeeds in disseminating all the messages to 100% of the nodes in 7 rounds with  $L=5$ , and in 6 rounds with  $L=6$ . In contrast, the gossip protocol only reaches 95.91% of the nodes on average with  $F=4$ , and 97.69% with  $F=5$ . Indeed, according to previous studies [14], a fan-out of 14 is required for a gossip protocol with 1000 nodes. This is due to the fact that with a gossip protocol only the out-degree (fan-out) is balanced, while the in-degree (fan-in) may be highly unbalanced. In contrast, Araneola’s in-degrees and out-degrees are balanced as all links in the overlay are bi-directional. As more nodes have a given message, the gossip protocol is more likely to “waste” its gossip on nodes that already have the message than Araneola, and therefore is less effective at spreading the information to additional nodes.

## 5.2 Fault-tolerance and graceful degradation

We now study the fault-tolerance and robustness of the Araneola overlay. We consider two kinds of failures: communication link failures and node failures. We study the overlay’s robustness with an offline analysis of the overlay snapshot obtained at the end of static experiments with 1000 and 2000 nodes. To study communication failures, we remove random subsets of edges from the overlay graph and analyze the resulting graphs. This allows us to predict Araneola’s reliability and latency in the presence of message loss. Similarly, we study Araneola’s resistance to node failures by removing random subsets of nodes.

We model node and edge failures as *independent and identically distributed (IID)*. For node failures the IID assumption has no significance since the overlay structure is random. Moreover, Bhagwan et al. have found that host failures are indeed independent [2]. For edge failures, the IID assumption fails to capture a situation in which some nodes have poorer links than others. The analysis of non-IID edge failure patterns is an interesting subject for future work.

We first analyze the impact of edge removals on the overlay with  $L=5$  and  $N=1000$ . This overlay has 2547 edges. For each percentage  $p \leq 50$  of the edges, we remove 10 different random subsets consisting of  $p\%$  of the edges from the overlay graph. The overlay becomes partitioned for the first time in one of the

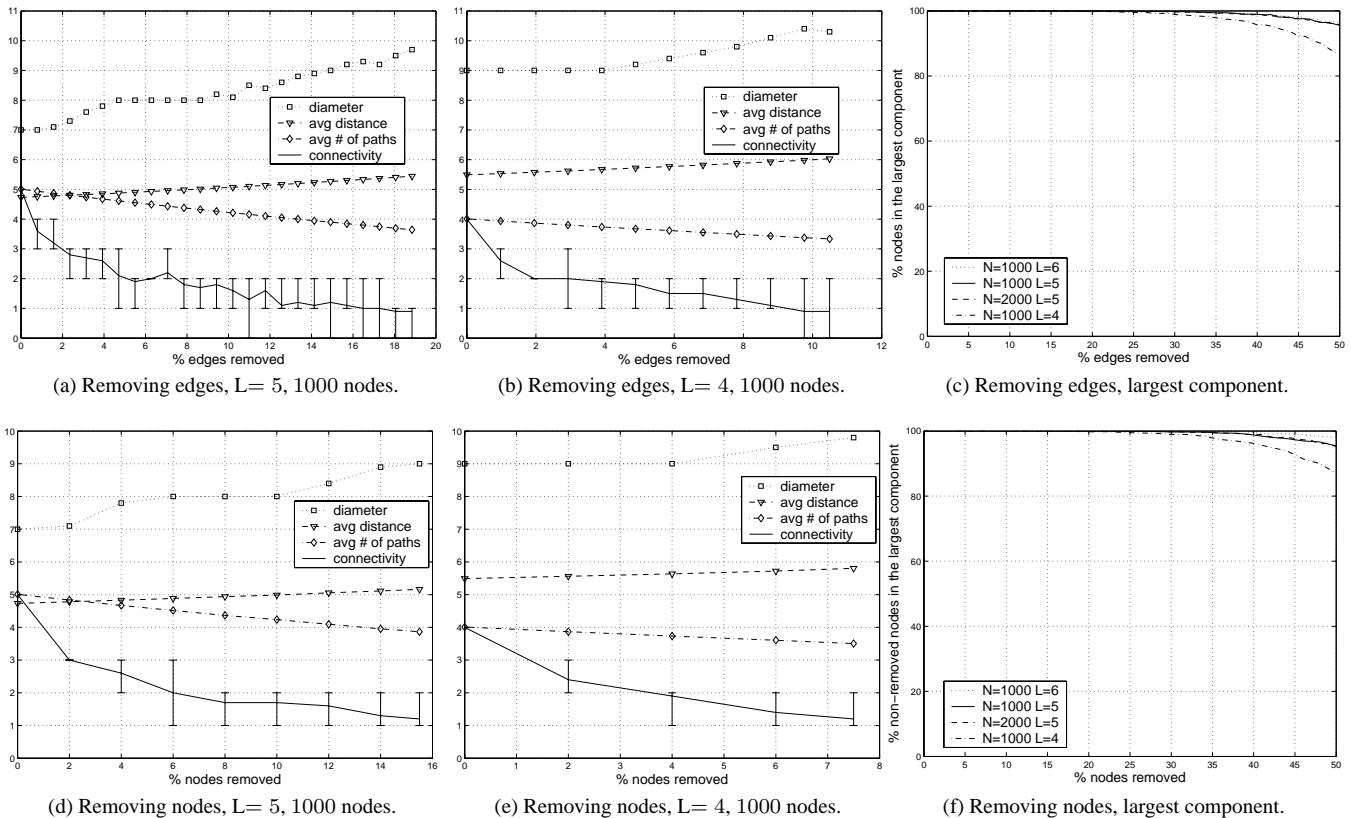


Figure 9: Resilience and graceful degradation of Araneola’s overlay.

ten experiments removing 11% (280) of the edges, and then in one of the experiments removing 15% (380). In both cases, a single node became disconnected from the rest. Fig. 9(a) shows how the removal of up to 19% of the edges affects the overlay’s characteristics. For each  $p$  in this range, the overlay is partitioned in at most one out of ten experiments in which  $p\%$  of the edges are removed. We observe that the average diameter increases from 7 to about 8 when 5–10% of the edges are removed, and to 9 when 15% of the edges are removed. The average distance increases more gradually, suggesting that message loss has a very moderate effect on the average latency. The average number of disjoint paths also decreases gradually with the failure rate. The bottom curve illustrates the average connectivity. The bars around each data point show the maximum and minimum connectivity observed in experiments with this  $p$ ; when the minimum goes to 0, there was a partition in one of the 10 experiments. We next experiment with  $L=4$  and  $N=1000$ . The overlay is less robust in this case— it partitions in more than 10% of the cases whenever  $p > 11\%$ . Fig. 9(b) shows the overlay’s degradation when up to 11% of the edges are removed.

We next examine how many of the nodes are still connected to each other, i.e., what is the size of the largest connected component in the graph. Fig. 9(c) depicts the average size of the largest connected component after random edge removals for  $L=4/5/6$  with  $N=1000$  and for  $L=5$  with  $N=2000$ . We can clearly see that the overlay’s resilience to the removal of a given percentage of its edges is *completely independent of  $N$* , as is expected in  $k$ -regular random graphs [11]: the curves for  $N=2000$  and  $N=1000$  (both with  $L=5$ ) are barely distinguishable. As expected, the value of  $L$  does impact the overlay’s robustness, but the difference between  $L=5$  and  $L=6$  is negligible. Remarkably, for  $L=5$ , after the removal of up to 38% of the edges, 99% of the nodes are still connected to each other, and only 1% of the nodes are partitioned from the rest.

We now turn our attention to node failures. Fig. 9(d) shows how node removals affect the properties of an overlay with 1000 nodes and  $L=5$  when up to 15% of the nodes are removed. None of the experiments with up to 15% removed nodes resulted in partitions. The overlay becomes partitioned in two of the ten experiments in which 16% (160) of the nodes are removed. This suggests that even if 15% of the nodes running Araneola fail during the brief time interval that it takes to detect and recover from failures (e.g., one minute), Araneola can continue to deliver messages reliably to surviving nodes. As with edge removals, the overlay exhibits graceful degradation: the diameter and average path length increase very moderately, while the average number of disjoint paths moderately decreases. When  $L=4$ , the overlay is half as robust to node failures as with  $L=5$ . It becomes partitioned in two of the ten runs with 8% of the nodes removed. Fig. 9(e) shows the overlay's degradation when up to 7.5% nodes are removed.

In fig. 9(f), we examine the size of the largest connected component that survives following node failures, for  $L=4/5/6$  with  $N=1000$  and for  $L=5$  with  $N=2000$ . Again, the overlay's resilience shows exactly the same trend with  $N=1000$  as it does with  $N=2000$ . This suggests that Araneola's resilience to simultaneous failures of a certain percentage of its nodes is also independent of  $N$ . When  $L=5$ , the largest component still includes 99% of the nodes following the failure of up to 38% of the nodes. When  $L=4$ , 99% of the nodes are still connected following the failure of 28% of the nodes. When 50% of the nodes fail, the largest component with  $L=5$  still includes over 95% of the nodes, (giving 90% reliability), and with  $L=4$ , it includes 87% (76% reliability). As with edge removals, increasing  $L$  from 5 to 6 achieves only slightly better robustness to node removals when there is an unrealistically high failure percentage.

## 5.3 Dynamic Evaluation

### 5.3.1 Methodology

Our model for this evaluation is based on studies of user behavior in multicast groups on the MBone [1], and in file sharing applications [23]. Both of these studies model the join and leave rates of most of the nodes using an exponential distribution. Moreover, both studies observe that a small portion of the nodes have substantially longer life times than others. However, these studies greatly differ in the mean life times they measure: the mean life time measured on the MBone is generally very short, e.g., 7 minutes in a typical multicast session, whereas the average measured life time in a file sharing application is roughly one hour.

We designate a small subset (roughly 7%) of the nodes as *perseverant*. Perseverant nodes created at the beginning of the experiment and remain active throughout the experiment. Subsequently, every minute, 50 additional (non-perseverant) nodes are awakened, until all nodes (1000 or 2000) are up. Each non-perseverant awakened node join the multicast group (becomes *active*) with probability 0.5. Otherwise, the node remains *inactive*. This gradual joining is modeled after the Berkeley session in [1]. Throughout the experiment, each non-perseverant node once a minute flips a coin with probability  $\lambda$  in order to decide whether to change its state from active to inactive and vice versa. We experiment with values of  $\lambda$  ranging from 0.01 (yielding a mean life time of 100 minutes) to 0.15 (giving a mean life time of 6.7 minutes). As a baseline, we also experiment with  $\lambda=0$ , in which case nodes do not change their states. There are roughly 1000 nodes alive at the end of each experiment with  $N=2000$ , (and resp., 500 when  $N=1000$ ), regardless of  $\lambda$ , since the join rate is equal to the leave rate.

### 5.3.2 Join/Leave overhead

We now measure the cost of constructing and maintaining the overlay in terms of the average number of control messages received by each node, where a control message is any message other than DATA or GOSSIP. In the appendix, we also analyze the expected number of control messages incurred by a single join or leave operation when the system is stable. Fig. 10 shows the overhead measured for different values

of  $\lambda$  with  $N = 1000$  and  $N = 2000$ . Remarkably, the overhead *decreases* as the rate of such events increases, the only exception occurring when  $\lambda$  increases from 0 to 0.01. Note that when  $\lambda = 0$ , no leave events occur. The measured average cost per join operation in this case is 15.6, which is very similar to the expected overhead calculated at the appendix. The overhead decreases as the churn rate rises because when many join and leave events occur concurrently, their costs can be amortized. E.g., a join event may increase a node’s degree while a leave event is reducing it, eliminating the need for correcting the overlay. Furthermore, we observe that the overhead does not increase with  $N$ . This is especially impressive given that the overhead for handling joins in structured overlays based on DHTs increases logarithmically with the number of nodes.

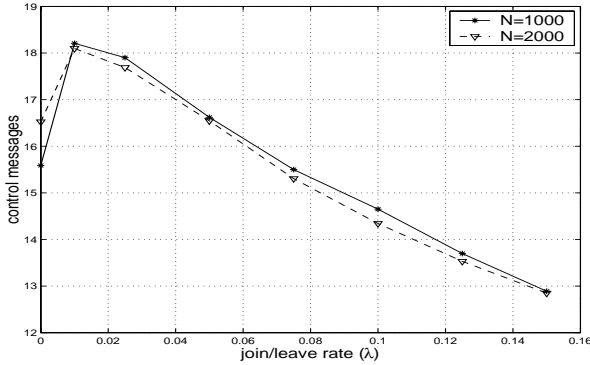


Figure 10: Average cost per join/leave with increasing churn rates for different group sizes,  $L = 5$ .

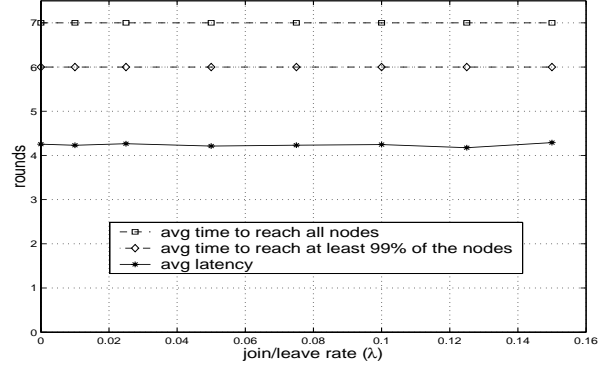


Figure 11: Average latencies in the presence of increasing churn rates, 1000 nodes,  $L = 5$ .

### 5.3.3 Undisrupted service

The second challenge we address is providing an undisrupted service in the presence of churn. For each message  $m$ , we define *nodes that are up during  $m$ 's transmission* to be nodes that have joined at least 12 rounds before  $m$ 's transmission, and did not leave at least 12 rounds after the transmission. We chose 12 as a very gross over-estimate. In fact, nodes can normally begin to receive messages reliably immediately upon requesting to join. In all of our dynamic experiments, each message is received by 100% of the nodes that were up during its transmission. Moreover, messages are delivered with *the same latency as in static runs*. We illustrate this in Fig. 11 for  $N = 1000$ ; similar results were obtained with  $N = 2000$ .

## 6 Exploiting Network Proximity

We now present an extension to Araneola that exploits network proximity by incorporating additional links between nearby nodes. This extension runs in parallel with and independently of the basic overlay construction and maintenance code presented in Section 4.1. The extension code has two components: (i) a mechanism for locating nearby nodes; and (ii) a `connect_nearby` task. The first component discovers nearby nodes and stores them in a set named *nearby\_cand*. The second component uses this set.

Generally speaking, Araneola can use a variety of mechanisms for locating nearby nodes. Our implementation does this as follows: at bootstrap time, each node  $n$  measures the network-level hop-count distances to the nodes in its local view using the UNIX `tracpath` utility, and inserts them to the *nearby\_cand* set in an ascending order of their network-level hop-count distances from  $n$ .

The `connect_nearby` task closely resembles the `connect` task presented in Section 4.1, except that no reduction rules are applied and no `REDIRECT` messages are sent. Specifically, there are three control

messages: CONNECT\_NEARBY, CONNECT\_OK\_NEARBY, and LEAVE\_NEARBY, which correspond to CONNECT, CONNECT\_OK, and LEAVE. In addition, both L and H are replaced by the parameter NB, which is the maximum number of nearby neighbors the node is willing to be connected to, and the *neighbors* set is replaced by the *nearby* set, which holds the node’s current nearby neighbors. Note that every node can set its own NB parameter to reflect its available bandwidth. Each CONNECT\_NEARBY request is issued to the closest node in *nearby\_cand*, rather than to a random node from the local view.

We evaluate this mechanism over the Internet, running 500 nodes over 25 Planet Lab [22] physical machines, with no two machines at the same site. In all the experiments presented in this section, all the nodes are created simultaneously, and remain up throughout the experiment. Although in principal, each node can choose its own NB parameter, in our experiments, we use the same value of NB for all nodes. We denote an experiment in which each node chooses L random neighbors and NB nearby neighbors as  $\langle L, NB \rangle$ .

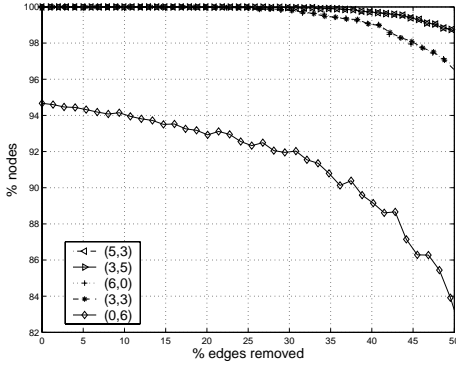
It is known that in order to achieve the good properties of  $k$ -regular graphs, each node should choose at least three random neighbor [26]. Thus, we run experiments in which each node chooses three random neighbors and three nearby neighbors ( $\langle 3, 3 \rangle$ ). We contrast these experiments against experiments in which each node chooses six random neighbors ( $\langle 6, 0 \rangle$ ), and against experiments in which each node chooses six nearby neighbors ( $\langle 0, 6 \rangle$ ). In addition, we run experiments in which the each node’s degree is roughly eight ( $\langle 3, 5 \rangle$ , and  $\langle 5, 3 \rangle$ ). Note that all the overlays we experiment with have a low degree compared to those used in previous systems [4, 16, 24]. For each selection of  $\langle L, NB \rangle$ , we run three experiments. In all our experiments, more than 97% of the nodes end up with NB nearby neighbors, and more than 90% of the nodes have exactly L random neighbors; the overall average node degrees in experiments with  $\langle 3, 3 \rangle$ ,  $\langle 6, 0 \rangle$ , and  $\langle 0, 6 \rangle$  are almost identical as are those of experiments with  $\langle 3, 5 \rangle$  and  $\langle 5, 3 \rangle$ .

We quantify the effectiveness of our approach by measuring the average number of physical hops that links in the extended overlay traverse. This metric is significant because a smaller hop-count distance implies reduced communication latencies as well as less stress on physical links. The results are summarized in Table 12. The first column shows the percentage of links between two nodes running on the same machine. The second column shows the percentage of short links with a hop-count distance of 3. These are Internet2 links between machines deployed at different sites belonging to the same enterprise. Finally, the third column shows the average hop-count in the overlay. Clearly, as NB is increased at the expense of L, there are more local and short links and the average number of physical hops that each link traverses is reduced.

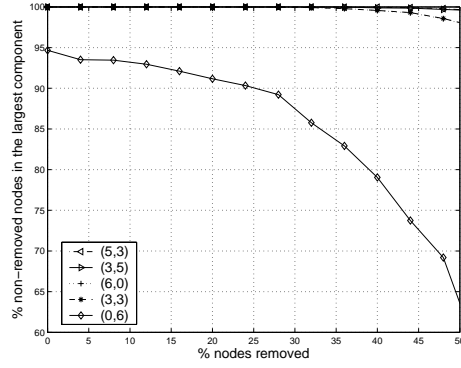
$\langle L, NB \rangle$	% of links on the same machine	% of short links	average hop count
$\langle 3, 3 \rangle$	<b>34.43</b>	<b>15.27</b>	<b>5.21</b>
$\langle 6, 0 \rangle$	4.97	6.93	8.69
$\langle 0, 6 \rangle$	74.23	3.4	1.88
$\langle 3, 5 \rangle$	<b>51.18</b>	<b>12.25</b>	<b>3.82</b>
$\langle 5, 3 \rangle$	35.6	10.46	5.54

Figure 12: Hop-count statistics with different selections of  $\langle L, NB \rangle$ .

Having verified that the mechanism achieves its goal, we next check its impact on the overlay’s robustness. We repeat the experiments of Section 5.2, i.e., we remove random subsets of edges and nodes from the overlay graphs and measure the sizes of the largest remaining components. The top two curves in Fig. 13(a) and Fig. 13(b) are for experiments with  $\langle 5, 3 \rangle$  and  $\langle 3, 5 \rangle$ . These curves are indistinguishable. Slightly below these are the curves for experiments with  $\langle 6, 0 \rangle$  and  $\langle 3, 3 \rangle$ , which are also conjoined. The bottom curve in both figures is for experiments with  $\langle 0, 6 \rangle$ . Remarkably, the robustness of an overlay with  $\langle 5, 3 \rangle$  is almost identical to that with  $\langle 3, 5 \rangle$ , and the robustness of an overlay with  $\langle 6, 0 \rangle$  is virtually identical to



(a) Removing edges, largest component, 500 nodes.



(b) Removing nodes, largest component, 500 nodes.

Figure 13: Graceful degradation with different selection of  $\langle L, NB \rangle$ .

that with  $\langle 3,3 \rangle$ . We believe that this stems from the fact that there is sufficient randomness in the choice of links since: (i) the nodes in *nearb\_cand* are chosen from the randomized local view; and (ii) each node is connected to at least 3 random neighbors.

The curves for experiments with  $\langle 0,6 \rangle$  show why it is important to choose random nodes as neighbors: in all these experiments, the overlay is partitioned even before we remove any edge or node. Moreover, as the percentage of removed edges or nodes increases, the robustness of the overlay deteriorates much quicker than when random edges are used.

We conclude from the experiments in this section that it is preferable for each node to have three random neighbors, and to allocate the rest of its available bandwidth for communication with nearby nodes.

## 7 Conclusions

We have presented Araneola, a scalable reliable multi-point to multi-point application-level multicast system for dynamic environments. We have evaluated Araneola over both a LAN and a WAN, and have shown that Araneola is highly scalable. The only aspect of Araneola that varies with the number of nodes is message latency, which increases logarithmically with the group size, whereas Araneola's load, reliability, resilience to message loss, resilience to simultaneous node failures, and overhead for handling join and leave events are all independent of the group size. Araneola can deliver messages with high reliability and bounded latency in the presence of sizable message loss rates, simultaneous failures of a certain percentage of the nodes, and high churn. The failure rates that Araneola can withstand depend on a tunable parameter. As the failure rate increases beyond its expectation, Araneola's reliability degrades gracefully. We have also shown how to extend Araneola to exploit available bandwidth for communication with nearby nodes. Such an approach substantially reduces the communication costs and message latency without hurting the overlay's robustness.

## Acknowledgements

We thank Ophir Ovadia for implementing part of the code and Dahlia Malkhi for helpful comments. We are grateful to the Flux research group at the University of Utah, and especially Leigh Stoller and Jay Lepreau, for allowing us to use their network emulation testbed and assisting us with our experiments.

## References

- [1] K. C. Almeroth and M. H. Ammar. Collecting and modeling the join/leave behavior of multicast group members in the mbone. In *High Performance Distributed Computing (HPDC)*, August 1996.
- [2] R. Bhagwan, S. Savagen, and G. Voelker. Understanding availability. In *2nd Intl. Workshop on Peer-to-Peer Systems (IPTPS)*, 2003.
- [3] K. P. Birman, M. Hayden, O. Ozkasap, Z. Xiao, M. Budiu, and Y. Minsky. Bimodal multicast. *ACM Trans. Comput. Syst.*, 17(2):41–88, 1999.
- [4] M. Castro, P. Druschel, A.-M. Kermarrec, A. Nandi, A. Rowstron, and A. Singh. Splitstream: High-bandwidth multicast in a cooperative environment. In *ACM SIGOPS Symposium on Operating Systems Principles (SOSP)*, October 2003.
- [5] M. Castro, P. Druschel, A.-M. Kermarrec, and A. Rowstron. SCRIBE: a large-scale and decentralized application-level multicast infrastructure. *IEEE J. Selected Areas in Comm. (JSAC)*, 2002.
- [6] Y. Chawathe, S. Ratnasamy, L. Breslau, N. Lanham, and S. Shenker. Making gnutella-like p2p systems scalable. In *ACM SIGCOMM*, August 2003.
- [7] A. Demers, D. Greene, C. Hauser, W. Irish, J. Larson, Shenker, Stuygis, D. Swinehart, and D. Terry. Epidemic algorithms for replicated database maintenance. In *6th ACM Symposium on Principles of Distributed Computing (PODC)*, pages 1–12, 1987.
- [8] S. El-Ansary, L. O. Alima, P. Brand, and S. Haridi. Efficient broadcast in structured P2P networks. In *2nd International Workshop On Peer-To-Peer Systems (IPTPS)*, Feb. 2003.
- [9] P. T. Eugster, R. Guerraoui, S. B. Handurukande, A. M. Kermarrec, and P. Kouznetsov. Lightweight probabilistic broadcast. In *The International Conference on Dependable Systems and Networks (DSN)*, 2001. Full version to appear in *ACM Trans. Comput. Syst.*
- [10] J. Friedman. On the second eigenvalue and random walks in random d-regular graphs. *Combinatorica*, vol. 11, pp. 331–362, 1991.
- [11] A. Goerdts. The giant component threshold for random regular graphs with edge faults. *Theoretical Comput. Sci.*, 259(1-2):307–321, 2001.
- [12] J. Jannotti, D. K. Gifford, K. L. Johnson, M. F. Kaashoek, and H. W. O. Jr. Overcast: reliable multicasting with an overlay network. In *Symp. Operating Systems Design and Implementation (OSDI)*, 2000.
- [13] F. Kaashoek and D. Karger. Koorde: A simple degree-optimal hash table. In *2nd Intl. Workshop on Peer-to-Peer Systems (IPTPS)*, 2003.
- [14] A.-M. Kermarrec, L. Massouli, and A. J. Ganesh. Probabilistic reliable dissemination in large-scale systems. *IEEE Transactions on Parallel and Distributed Systems*, 14(3):248–258, March 2003.
- [15] J. Kleinberg. The small-world phenomenon: An algorithmic perspective. In *32nd ACM Symp. on Theory of Computing (STOC)*, pages 163–170, 2000.
- [16] D. Kostić, A. Rodríguez, J. Albrecht, and A. Vahdat. Bullet: High bandwidth data dissemination using an overlay mesh. In *ACM SIGOPS Symposium on Operating Systems Principles (SOSP)*, October 2003.
- [17] C. Law and K. Siu. Distributed construction of random expander networks. In *IEEE Infocom*, April 2003.
- [18] M. J. Lin, K. Marzullo, and S. Masini. Gossip versus deterministically constrained flooding on small networks. In *14th International Symposium on Distributed Computing (DISC)*, pages 253–267, 2000.
- [19] D. Malkhi, M. Naor, and D. Ratajczak. Viceroy: A scalable and dynamic emulation of the butterfly. In *ACM Symposium on Principles of Distributed Computing (PODC)*, July 2002.
- [20] L. Massoulié, A.-M. Kermarrec, and A. J. Ganesh. Network awareness and failure resilience in self-organising overlay networks. In *22st IEEE International Symposium on Reliable Distributed Systems (SRDS)*, 2003.
- [21] V. Paxson. End-to-end Internet packet dynamics. In *ACM SIGCOMM*, September 1997.
- [22] L. Peterson, T. Anderson, D. Culler, and T. Roscoe. A blueprint for introducing disruptive technology into the internet. In *Proceedings of ACM HotNets-I*, October 2002.
- [23] S. Saroiu, K. Gummadi, and S. Gribble. A measurement study of peer-to-peer file sharing systems. In *Multimedia Computing and Networking*, January 2002.
- [24] K. Shen. Structure management for scalable overlay service construction. In *the First USENIX/ACM Symposium on Networked Systems Design and Implementation (NSDI'04)*, San Francisco CA, March 2004.
- [25] B. White, J. Lepreau, L. Stoller, R. Ricci, S. Guruprasad, M. Newbold, M. Hibler, C. Barb, and A. Joglekar. An integrated experimental environment for distributed systems and networks. In *Symp. Operating Systems Design and Implementation (OSDI)*, pages 255–270, Boston, MA, Dec. 2002.
- [26] N. Wormald. Models of random regular graphs. *Surveys in Combinatorics*, 276:239–298, 1999.



## A Analyzing the Join/Leave Overhead

We calculate the overhead for the simple case where a single join or leave occurs when the system is stable, i.e., each node's degree is either  $L$  or  $L+1$ , and no two neighboring nodes have a degree of  $L+1$ . We assume that the probability that a node has a degree of  $L$  is  $p$ , and the probability that a node has a degree of  $L+1$  is  $1 - p = q$ .

### A.1 The overhead for join

We begin by calculating the expected overhead for a single CONNECT request. Assume that node  $c$  issues a CONNECT request to node  $t$ . There are three possible cases: (i)  $t$  has a degree of  $L$  and all of  $t$ 's neighbors also have a degree of  $L$ ; (ii)  $t$  has a degree of  $L$  and at least one of its neighbors has a degree of  $L+1$ ; or (iii)  $t$  has a degree of  $L+1$ . In the latter case, all of  $t$ 's neighbors have a degree of  $L$ . The probability for case (i) is  $p^{L+1}$ , the probability for case (ii) is  $p(1 - p^L)$ , and for case (iii) it is  $1 - p$ .

In case (i),  $c$  sends one CONNECT message to  $t$  and in return,  $t$  sends one CONNECT\_OK message to  $c$ , total of two messages. In case (ii),  $c$  first sends one CONNECT message to  $t$  and  $t$  responds with a CONNECT\_OK. In addition, two messages (DISCONNECT and DISCONNECT\_OK) are sent in order to reduce the degree of  $t$  and one of its neighbors,  $n$ , from  $L+1$  to  $L$ . Thus, a total of four control messages are sent. In case (iii),  $c$  again first sends one CONNECT message to  $t$  and in return,  $t$  sends one CONNECT\_OK to  $c$ . Now,  $t$ 's degree becomes  $L+2$  while the rest of  $t$ 's neighbors have a degree of  $L$ . Thus,  $t$  activates the second reduction rule (see Section 4.1). First,  $t$  sends to one of its neighbors,  $n'$ , a CONNECT\_TO message with the identity of another neighbor,  $n''$ . Then,  $n'$  sends a CHANGE\_CONNECTION message to  $n''$  with the identity of node  $t$ . In return,  $n''$  sends a CONNECT message to  $n'$  and a DISCONNECT message to  $t$ . Finally,  $n'$  sends a CONNECT\_OK message to  $n''$  and  $t$  sends a DISCONNECT\_OK message to  $n''$ . Now, the degrees of  $t$  and  $n'$  are  $L+1$  and the degree of  $n''$  remains  $L$ . In the next iteration of the reduce algorithm, either  $t$  or  $n'$  sends a DISCONNECT message to the other and the other replies with a DISCONNECT\_OK. The total number of messages sent in this case is ten. The expected number of control messages sent for a single CONNECT request is therefore:  $2p^{L+1} + 4p(1 - p^L) + 10q = 4p + 10q - 2p^{L+1}$ . Thus, the expected overhead associated with a single join operation during a stable period is:  $L(4p + 10q - 2p^{L+1})$ . Recall that when the system is stable, roughly 92% of nodes have a degree of  $L$  (when  $L$  is set to 5), and the rest of the nodes have a degree of  $L+1$ . Substituting 5 for  $L$ , 0.92 for  $p$ , and 0.08 for  $q$ , we get an overhead of roughly 16.326 messages.

### A.2 The overhead for leave

Assume that node  $l$  sends a LEAVE message to node  $t$ . There are two possible cases: (i)  $t$  has a degree of  $L$ ; or (ii)  $t$  has a degree of  $L+1$ . The probability of case (i) is  $p$ , and the probability of case (ii) is  $q$ . In the first case,  $l$  sends a LEAVE message to  $t$ . Subsequently,  $t$  sends a CONNECT request to a random new node. We showed above that the expected overhead for sending a CONNECT request is  $4p + 10q - 2p^{L+1}$ . Thus, the expected number of messages sent in the first case is  $1 + 4p + 10q - 2p^{L+1}$ . In the second case,  $l$  sends a LEAVE message to  $t$ . However, in this case,  $t$  does not send any messages as its degree is  $L$ . Thus, the total expected overhead for sending a LEAVE message is:  $p(1 + 4p + 10q - 2p^{L+1}) + q$ . The expected number of LEAVE messages sent upon a node leaving the system is:  $pL + q(L + 1) = L + q$ . Thus, the expected number of messages sent upon a node leaving the system is:  $(L + q) * [p(1 + 4p + 10q - 2p^{L+1}) + q]$ . Substituting 5 for  $L$ , 0.92 for  $p$ , and 0.08 for  $q$ , we get an overhead of roughly 20.35 messages.

### A.3 Theory versus practice

As we have shown, the overhead for a leave operation is inherently larger than the overhead for a join operation. Recall that experiments in Section 5.3.2 shows that the average overhead for a single join or leave operation is roughly between 18.1 and 18.2 control messages when  $\lambda = 0.01$ . Such experiments are very similar to a stable system, as the churn rate is low. We observe that the measured overhead is closer to the overhead for a join operation than to the overhead for a leave operation. This is due to the fact that in a dynamic experiment there are more joins than leaves, since each dynamic experiment starts with no live nodes and ends with several alive nodes. E.g., in a dynamic experiment with 2000 nodes and  $\lambda = 0.01$  there were 1411 joins and 387 leaves. Thus, according to the formula above, the expected overhead for a join or leave operation is  $(1411 * 16.326 + 387 * 20.35)/1798 \cong 17.19$ . The difference between this expected overhead and the measured overhead in Section 5.3.2 stems from the fact that the system is not fully stable, and hence some CONNECT messages are redirected to other nodes when a node with a degree of H receives a CONNECT message.