

The Dynamic Behavior of a Data Dissemination Protocol for Network Programming at Scale

Jonathan W. Hui
Computer Science Division
University of California at Berkeley
Berkeley, CA 94720
jwhui@cs.berkeley.edu

David Culler
Computer Science Division
University of California at Berkeley
Berkeley, CA 94720
culler@cs.berkeley.edu

ABSTRACT

To support network programming, we present *Deluge*, a reliable data dissemination protocol for propagating large data objects from one or more source nodes to many other nodes over a multihop, wireless sensor network. *Deluge* builds from prior work in density-aware, epidemic maintenance protocols. Using both a real-world deployment and simulation, we show that *Deluge* can reliably disseminate data to all nodes and characterize its overall performance. On Mica2-dot nodes, *Deluge* can push nearly 90 bytes/second, one-ninth the maximum transmission rate of the radio supported under TinyOS. Control messages are limited to 18% of all transmissions. At scale, the protocol exposes interesting propagation dynamics only hinted at by previous dissemination work. A simple model is also derived which describes the limits of data propagation in wireless networks. Finally, we argue that the rates obtained for dissemination are inherently lower than that for single path propagation. It appears very hard to significantly improve upon the rate obtained by *Deluge* and we identify establishing a tight lower bound as an open problem.

Categories and Subject Descriptors

C.2.1 [Computer-Communication Networks]: Network Architecture and Design—*Wireless Communication*; C.2.2 [Computer-Communication Networks]: Network Protocols

General Terms

Design, Experimentation, Measurement, Performance, Reliability

Keywords

Wireless Sensor Networks, Dissemination Protocols, Network Programming, Wireless Networks

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

SenSys'04, November 3–5, 2004, Baltimore, Maryland, USA.
Copyright 2004 ACM 1-58113-879-2/04/0011 ...\$5.00.

1. INTRODUCTION

Wireless sensor networks (WSNs) represent a new class of computing with large numbers of resource-constrained computing nodes cooperating on essentially a single application. WSNs must often operate for extended periods of time unattended, where evolving analysis and environments can change application requirements, creating the need alter the network's behavior by introducing new code. Unlike the traditional method of programming a node over a dedicated link, the embedded nature of these systems requires the propagation of new code over the network. As WSN research matures, growing testbeds sized at tens of thousands of nodes are now on the horizon, making code propagation over the network a necessity in the debugging and testing cycle. These factors suggest that *network programming* (the programming of nodes by disseminating code over the network) is required for the success of WSNs. Specifically, we consider the propagation of complete binary images. While virtual machines can provide low-cost retasking with the use of virtual programs, it is still necessary to have the option of reprogramming nodes with a new binary image since the virtual machine itself may need changes [6].

The core service required to enable network programming is the dissemination of a program image over a multihop WSN and presents several problems. First, program images are much larger than the data objects that previous dissemination protocols consider. This is an issue because a sensor node operates on a constrained storage hierarchy where a communication packet (36 bytes) \ll RAM (4K) \ll program size (128K) $<$ external flash (512K). Second, the dissemination must tolerate node densities which can vary by factors of a thousand or more. For example, collectively programming all nodes before deployment could put thousands of nodes within communication range, making suppression of protocol messages essential. Third, complete reliability is required since every byte must be correctly received by all nodes that need reprogramming, even in the presence of high loss rates and evolving link qualities common to WSNs. Fourth, propagation must be a continuous effort to ensure that all nodes receive the newest code since network membership is not static: nodes come and go due to temporary disconnections, failure, and network repopulation. Finally, the dissemination process should require a minimal amount of time, reducing any service interruptions to a deployed application and the debugging and testing cycle.

This paper provides four main contributions. First, we present *Deluge*, a reliable data dissemination protocol for

disseminating a large data object (i.e. larger than can fit into RAM) from one or more source nodes to many other nodes over a multihop wireless network. Deluge’s density-aware, epidemic properties help achieve reliability in unpredictable wireless environments and robustness when node densities can vary by factors of a thousand or more. Representing the data object as a set of fixed-size pages provides a manageable unit of transfer which allows for spatial multiplexing and supports efficient incremental upgrades. Second, we characterize the propagation dynamics of Deluge using a real-world deployment and simulation. In a real deployment, Deluge can disseminate data with 100% reliability at nearly 90 bytes/second, one-ninth the maximum transmission rate of the radio supported under TinyOS. Under simulation, propagation of large data objects in large networks exposes interesting propagation behavior only hinted at in previous dissemination work. Third, we develop a simple model of Deluge’s propagation behavior and use it to identify different factors which limit the overall bandwidth of any multihop communication protocol. Finally, we argue that dissemination is inherently slower than single path propagation and identify establishing a tight lower bound as an open problem. It appears very hard to significantly improve upon the rate obtained by Deluge.

Deluge builds off Trickle, a protocol for maintaining code updates in WSNs [8]. In Trickle, nodes stay up-to-date by occasionally broadcasting a code summary to their neighbors. Trickle’s key contribution is the use of suppression and dynamic adjustment of the broadcast rate to limit transmissions among neighboring nodes. A node suppresses its own broadcast if it recently overhears a similar code summary. When nodes are not up-to-date, the broadcast rate is reduced, but is otherwise increased up to a specified limit. These techniques allow Trickle to scale to thousand-fold changes in network density, disseminate updates quickly, and consume minimal resources in the quiescent state. While Trickle addresses single packet dissemination, Deluge extends it to support large data objects.

In Section 2 of this paper, we review related work. We describe the protocol formally in Section 3 and discuss our evaluation methodology in Section 4. In Section 5, we present Deluge’s overall performance and characterize Deluge’s propagation dynamics. Section 6 presents the current status of Deluge. Section 7 discusses future directions, and this paper concludes with Section 8.

2. RELATED WORK

The problem we address is an important special case of reliable data dissemination. The main differences include the dissemination of data over lossy local links, a constrained storage hierarchy, and the need to retain the most recent copy in order to disseminate data to additional nodes as they become connected over time.

For data dissemination in wireless networks, naive retransmission of broadcasts can lead to the *broadcast storm problem*, where redundancy, contention, and collisions impair performance and reliability [9]. The authors discuss the need to have a controlled retransmission scheme and propose several schemes, such as probabilistic and location-based methods. The experiments conducted by Ganesan et al. identify several interesting effects at the link-layer, notably highly irregular packet reception contours, the likeliness of asymmetric links, and the complex propagation dynamics of sim-

ple protocols [2]. *Scalable Reliable Multicast (SRM)* is a reliable multicast mechanism built for wired networks [4], using communication suppression techniques to minimize network congestion and request implosion at the server.

Demers et al. propose an epidemic algorithm based on randomly chosen point-to-point interactions for managing replicated databases that is robust to unpredictable communication failures [1]. *SPIN-RL* is an epidemic algorithm designed for broadcast networks that makes use of a three-phase (advertisement-request-data) handshaking protocol between nodes to disseminate data [5]. The epidemic property is important since WSNs experience high loss rates, asymmetric connectivity, and transient links due to node failures and repopulation. However, their results show control message redundancy at over 95% as it only considers the suppression of redundant request messages, and SPIN-RL does not perform as well as naive flooding for lossy network models. *Trickle* builds upon this approach by proposing SRM-like suppression mechanisms to minimize redundant transmission of control messages and occasional advertisements to increase reliability, allow for quick propagation, and consume few resources in the steady state [8]. However, Trickle only provides a mechanism for determining when nodes should propagate code. Deluge builds directly off Trickle, adding support for the dissemination of large data objects with a three-phase protocol similar to SPIN-RL.

Because reliability is of top priority, Deluge also borrows ideas from prior work in reliable data transfer protocols. *Pump Slowly, Fetch Quickly (PSFQ)* [16] and *Reliable Multi-Segment Transport (RMST)* [10] are selective NACK-based reliable transport protocols designed for WSNs. Because the cost of end-to-end repair is exponential with the path length, both protocols emphasize hop-by-hop error recovery where loss detection and recovery is limited to a small number of hops (ideally one). Like PSFQ and RMST, Deluge uses a selective NACK-based approach and error recovery is limited to a single hop. However, these approaches do not consider methods for adapting to spatial node density.

Research activity directed at network programming for WSNs has been limited. TinyOS [15] has included limited support for network programming via *XNP* [3]. However, XNP only provides a single-hop solution, requiring all nodes to be within bidirectional communication range of the source. Additionally, repairs are done on a whole file basis, requiring expensive scans through external flash to discover missing data.

Multihop Over-the-Air Programming (MOAP) presents a more comprehensive approach to network programming, supporting multihop networks [11]. Deluge shares many ideas with MOAP, including the use of NACKs, unicast requests, broadcast data transmission, and windowing to efficiently manage which segments are required. However, Deluge considers additional key design options. For example, MOAP does not fragment the image, requiring nodes to receive the entire code image before re-broadcasting, preventing the use of spatial multiplexing to leverage the full capabilities of the network. Additionally, MOAP does not deal with the adverse effects of asymmetric links.

3. DELUGE

Deluge is an epidemic protocol and operates as a state-machine where each node follows a set of *strictly local* rules to achieve a desired global behavior: the quick, reliable dis-

semination of large data objects to many nodes. In its most basic form, each node occasionally advertises the most recent version of the data object it has available to whatever nodes that can hear its local broadcast. If S receives an advertisement from an older node, R , S responds with its object profile. From the object profile, R determines which portions of the data need updating and requests them from any neighbor that advertises the availability of the needed data, including S . Nodes receiving requests then broadcast any requested data. Nodes then advertise newly received data in order to propagate it further.

While the basic form is quite simple, Deluge considers many subtle issues to achieve high performance. The first is its density-aware capability, where redundant advertisement and request messages are suppressed to minimize contention. Note that while suppression can increase performance by avoiding congestion collapse, its necessary back-offs introduce latency. Second, protocols for WSNs must be robust to asymmetric links, where a link in one direction can have a significantly different loss rate in the other direction. Deluge’s three-phase handshaking protocol helps ensure that a bi-directional link exists before transferring data. Additionally, if a node has not completely received its data after making k requests, it searches for a new neighbor to request data, rather than hanging on to a bad link. However, it may continue making requests if sufficient progress is being made. Third, Deluge dynamically adjusts the rate of advertisements to allow quick propagation when needed while consuming few resources in the steady state. Fourth, Deluge attempts to minimize the set of nodes concurrently broadcasting data within a given cell. Finally, Deluge emphasizes the use of spatial multiplexing to allow parallel transfers of data.

In the remainder of this section, we describe in detail how Deluge represents the large data object and present a formal description of the Deluge protocol.

3.1 Data Representation

To manage the large size, Deluge divides the data object into fixed-size *pages*. The page is the basic unit of transfer and provides three advantages: (i) it limits the amount of state a receiver must maintain while receiving data, (ii) it enables *efficient incremental upgrades* from prior versions and (iii) allows for *spatial multiplexing*. In this section, we discuss the first two in detail and save our discussion of the last for Section 3.2.

The data object of size S_{obj} is divided into *packets* of a fixed size S_{pkt} . To ensure receipt of all packets, the node must keep track of which packets are needed to complete the object. However, because the packet size is generally much smaller than the object, simply maintaining a bit-vector consumes an unacceptably large amount of RAM. Instead, Deluge fragments the data object into P pages each of size $S_{page} = N \cdot S_{pkt}$, where N is a fixed number of packets, as shown in Figure 1. By requiring a node to dedicate itself to receiving a single page at a time, the bit-vector need only be N bits in length.

Both packets and pages include 16-bit cyclic redundancy checks (CRCs). If a packet or a page fails the CRC, all data represented by the CRC is discarded and must be received again. Redundant data integrity checks at both the packet and page level help ensure that data is correctly received by all nodes and is especially important due to Deluge’s epi-

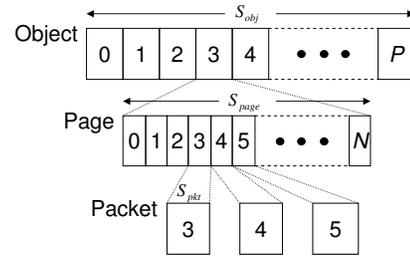


Figure 1: Data Management Hierarchy.

demio nature: a single node with corrupt data can propagate erroneous data to other nodes.

In many cases, an upgrade may only contain a few minor changes to localized areas within the data object. Thus, nodes need only those pages which have changed in order to match the newer version. To support this, we assume that updates to the overall data object are serialized. A *version number* is used to distinguish between different updates and must be monotonically increasing to produce a total order for all updates. A node compares version numbers to determine whether it should request new data.

Upgrading from a prior version requires knowing which pages have changed. Because data from recent deployments show that nodes can fade out of connectivity for a significant amount of time [12], it is necessary to allow efficient upgrades from a version which is more than one version behind. This implies that knowledge about *when* pages were last changed is needed. Deluge represents the set of p pages for a given version, v , by an *age vector*, $\mathbf{a} = \langle a_0, a_1, \dots, a_{p-1} \rangle$, which describes how “old” each page is. More specifically, the contents of page i at version v last changed at version $v - a_i$. A complete description of the object is defined by the tuple (v, \mathbf{a}) , called the *object profile*, and is stored in non-volatile storage along with the data it represents. A node receiving an object profile for a newer version uses the age-vector to determine which pages need updating. In the current implementation, the page age is specified by a nibble, limiting the communication of an object profile to several packets. However, this tradeoff forces nodes sixteen or more versions behind require the transfer of all pages in the object regardless of their age.

3.2 The Protocol

A node operates in one of three states at any time: *MAINTAIN*, *RX*, or *TX*. The set of local rules an individual node follows is a function of its current state and specify what actions and state transitions to take in response to events. We formally describe the set of local rules and discuss how each of these rules contribute to the desired global behavior.

3.2.1 Maintenance

A node in the *MAINTAIN* state is responsible for ensuring that all nodes within communication range have (i) the newest version of the object profile and (ii) all available data for the newest version. To maintain this property, each node occasionally advertises a summary representing the current version of its object profile and the set of pages from the object which are available for transmission. We define a page i as *complete* if every packet for that page has been correctly received. Page i is *available* only if it is complete and all

pages in the range $[0, i)$ are also complete. Thus, the summary need only contain two integers $\{v, \gamma\}$, where v is the version number and γ is the largest numbered page available for transfer.

Ideally, the transmit rate of advertisements in a given cell should be independent of the spatial node density, allow quick propagation, and consume few resources in the steady state. This should hold even when thousands of nodes are within communication range, a scenario representative of programming nodes prior to deployment. Deluge uses Trickle to control the transmission of potentially redundant messages. Trickle divides time into a series of rounds and nodes choose whether or not to broadcast an advertisement in each round. The duration of round i is specified by $\tau_{m,i}$ and is bounded by τ_l and τ_h . In each round, a node maintains a random value r_i in the range $[\frac{\tau_{m,i}}{2}, \tau_{m,i}]$. The local rules for a node with summary ϕ are as follows:

- M.1 During round i with start time $t_i = t_{i-1} + \tau_{m,i-1}$, broadcast an advertisement with summary ϕ at time $t_i + r_i$, only if less than k advertisements with summary $\phi' = \phi$ have been received since time t_i .
- M.2 If any overheard packet indicates an inconsistency among neighboring nodes (i.e. advertisements with $\phi' \neq \phi$, any requests, or any data packets) were overheard during round i , set $\tau_{m,i}$ to τ_l and begin a new round.
- M.3 At the beginning of a round i , if no overheard packet indicates an inconsistency among neighbors during the previous round, set $\tau_{m,i}$ to $\min(2 \cdot \tau_{m,i-1}, \tau_h)$.

Trickle is density-aware in that the threshold k bounds the number of advertisements made in a given cell by suppressing the transmission of redundant advertisements. In a lossless, single-cell network model, an advertisement with a summary $\phi' = \phi$ is only transmitted at most k times in a period $\frac{\tau_{m,i}}{2}$, independent of the density. In a lossy, multi-cell model, the number of transmissions is bounded by $O(\log(n))$, where n is the number of nodes in the cell [8]. The lower bound of $\frac{\tau_{m,i}}{2}$ for r_i is required due to the *short-listen problem*. In an unsynchronized network, a small value of r_i reduces the likelihood of overhearing similar advertisements and results in redundant transmissions. The dynamic adjustment of $\tau_{m,i}$ in the range $[\tau_l, \tau_h]$ allows quick propagation during an upgrade and low resource consumption in the steady state by decreasing and increasing the advertisement rate respectively.

With the advertisement service, a node can determine if any of its neighbors have an old object profile. If so, a node must broadcast the new object profile. The local rule governing this process for a node with version v are as follows:

- M.4 During round i , transmit the object profile for version v at time $t_i + r_i$ only if an advertisement with version $v' < v$ was received at or after time t_i and less than k attempts to update the object profile to version v have been overheard.

This method for updating object profiles is a form of controlled flood, providing a reliable approach which is density-aware. By following the rounds defined by the advertisement service and upgrading an object profile only if less than k redundant upgrade attempts have been transmitted,

we keep Trickle’s beneficial properties of low redundancy and quick propagation while its epidemic property helps to ensure eventual propagation to all nodes.

The local rules defined so far ensure eventual consistency of object profiles, allowing nodes to learn about a newer version and determine which pages need updating in order to match the newer version. We now discuss the method used to initiate the reception or transmission of new pages. We first present the remaining local rules and then discuss their contribution to the overall behavior of Deluge. In Deluge, nodes request data from a single node S at a time. The local rules for a node with with summary $\phi = \{v, \gamma\}$:

- M.5 On receiving an advertisement with $v' = v$ and $\gamma' > \gamma$, transition to RX unless (i) a request for a page $p \leq \gamma$ was previously received within time $t = 2 \cdot \tau_{m,i}$ or (ii) a data packet for page $p \leq \gamma + 1$ was previously received within time $t = \tau_{m,i}$.
- M.6 On receiving a request for data from a page $p \leq \gamma$ from version v , transition to TX .

We leverage the Trickle suppression mechanism to help minimize the set of senders and simplify the decision making process of nodes at any given time. The only trigger that causes R to request data from S is the receipt of an advertisement stating the availability of a needed page. Because Trickle bounds the number of transmitted advertisements in a given cell, the set of nodes which may become senders during any round is also bounded. Deluge simply requests data from the node which most recently advertised the needed page. If R overhears a data packet of the needed page, it suppresses any requests for a full round and attempts to snoop as much as possible. This also suppresses the initiation of any additional senders which can interfere with the current set of senders.

A significant contribution of Deluge is its emphasis on *spatial multiplexing*. Deluge advertises the availability of complete pages even before all pages in the object are complete, allowing the further propagation of newly received pages. Overall throughput is increased by pipelining the transfer of pages across the network. Without spatial multiplexing, propagation across a network of d hops requires d complete object transfers, requiring a time of $o(d \cdot S_{obj})$. Instead, pipelining the transfer approaches a time of $o(d + S_{obj})$. Intuitively, it is the time for the first bit of data to traverse the network in addition to the time required to flush the pipeline. Since Deluge targets both large scale networks and object sizes, spatial multiplexing can significantly enhance performance. One drawback is the entire network must remain powered-on to achieve the full benefits of spatial multiplexing. However, it is not required and transfers can be localized by only powering on a subset of the nodes, trading energy-consumption with completion time.

In order to realize the full benefit of spatial multiplexing, Deluge takes special care to ensure that transfers of different pages do not interfere with each other. First, Deluge constrains nodes by requesting pages in sequential order, that is, a request for page i cannot be made unless data for all pages in the range $[0, i)$ are also up-to-date and complete. This allows neighboring nodes take advantage of the broadcast medium by working together in receiving the same page rather than contending with each other in requesting different pages. An added advantage is that a node need not

decide whether to give up an attempt to receive a specific page p and focus its efforts on a different page. Because all nodes complete pages in sequential order, any node advertising $\gamma > p$ is also able to supply page p .

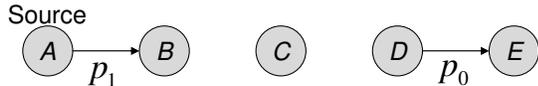


Figure 2: Pipelining. Example four-hop network. A three-hop spacing is required for effective spatial multiplexing. In this example, a simultaneous broadcast from A and C would collide at B.

Second, a transfer of page p will always take higher priority than a transfer of a page $p' > p$. From local rule M.5, nodes may not issue a request for a page p if a transfer of page $p' < p$ is in progress. As shown in Figure 2, this constraint reduces interference, including the hidden terminal problem, caused by messages generated from transfers of different pages.

3.2.2 Request

A node in the RX state is responsible for actively requesting the remaining packets required to complete page $p = \gamma + 1$. Each request operates as a selective negative acknowledgment (SNACK) where a bit-vector specifies which data packets in the page are needed. To achieve density-awareness, requests are made with a random backoff to help minimize collisions with requests from other nodes and allows for suppression if any requests or data packets are overheard during the backoff period. Utilizing the broadcast medium, responses to requests are shared by all receivers.

After a node makes a request, it waits for a response and makes subsequent requests if some of the data was lost in the process. Nodes delay subsequent requests until they detect a period of silence equal to ω packet transmit times to help ensure the completion of data transmissions before requests are made. The value of ω is chosen such that simply dropping a couple packets is not detected as a period of silence.

Deluge is also designed to tolerate asymmetric links. First, a node R receiving advertisements from S may not be able to communicate requests to S . Deluge limits a node to λ requests before returning to $MAINTAIN$. However, if progress (measured as reception rate of data) is above some threshold $0 < \alpha < 1$, Deluge allows the node to continue making requests. Additionally, the rate of requests from R will decrease with the decreasing advertisement rate in the steady state since S will not know that R is not up-to-date. Second, if node S cannot communicate advertisements to R , then R will not transmit any requests to S since requests are only transmitted in response to advertisements.

We now present the local rules for a node in RX . For a node n with summary $\phi = \{v, \gamma\}$, we define S as the node that caused n 's transition into RX .

- R.1 After not receiving a request or a data packet for time $t = \omega \cdot T_{tx} + r$, where r is a random value in the range τ_r , transmit a request to node S .
- R.2 After λ requests with a packet reception rate of $\alpha' < \alpha$, transition to $MAINTAIN$ even if page $\gamma + 1$ is incomplete.

- R.3 If all packets for page p are received and data for page $\gamma + 1$ passes the CRC, mark page as complete and transition to $MAINTAIN$.

Note that a node need not be in RX in order to receive data packets. Deluge takes advantage of the broadcast medium and any packets needed to complete page $\gamma + 1$ are saved regardless of which state the node is in.

3.2.3 Transmit

A node in TX is responsible for broadcasting all requested packets for a given page (continuing to service any subsequent requests for data from the same page) until all requested packets have been broadcast and then transition back to $MAINTAIN$. Deluge services requests by taking the union of any new requests with previous requests not yet serviced. Packets are sent in round-robin order to provide fairness among requesters. We define Π as the set of packets from page p which have been requested and is initialized to the request which caused the transition to this state. The set of local rules are as follows:

- T.1 On receiving a request for packets Π'_{rx} from page p , set Π_{tx} to $\Pi_{tx} \cup \Pi'_{rx}$.
- T.2 Continue broadcasting packets in Π_{tx} in round-robin order and remove each broadcast packet from Π_{tx} until $\Pi = \{\}$, then transition to $MAINTAIN$.

3.3 Design Space

The design space for data dissemination protocols is large and includes: methods for suppressing redundant control and data messages, selection of nodes to transmit data, use of forward error-correction (FEC), the fragmentation of data for spatial multiplexing, use of link quality estimates or other metrics to improve decisions, among others. In the early stages of designing Deluge, we experimented with several of these options in simulation. Due to space limitations, we briefly mention some of our findings.

The suppression mechanisms make up most of the complexity in Deluge. To confirm their importance, we tested Deluge without any suppression and it performed so poorly that the simulations did not complete after days of execution. With just request message suppression, Deluge is similar to SPIN-RL. We again tested Deluge by keeping only the request suppression and it also performed poorly, confirming the results presented by SPIN-RL's authors that show it performing worse than naive flooding in lossy network models.

We experimented with suppressing the transmission of data packets if k redundant data packets were overheard while in TX , where lower values of k represented more aggressive suppression. Lower values of k tended to decrease performance. Recall that Deluge already attempts to limit the number of senders. While too many senders leads to high contention and greater opportunities for the hidden terminal problem, a small number (two or three) of neighboring senders is able to cover more area without much loss in performance. Additionally, the TinyOS MAC layer implements CSMA, minimizing collisions caused by neighboring senders. This result led to the decision of not employing additional techniques for data packet suppression.

Deluge takes a very simple approach to selecting a sender by using the most recent advertisement. We tested more

sophisticated methods based on a hop count metric (i.e. requesting data from nodes nearest to the source, nodes furthest from the source, and nearest neighbor) and on the number of nodes requesting data, but found no significant difference in performance for any of these approaches over Deluge’s current, simple approach.

We also tested the use of FEC which allows receivers to reconstruct the original data from any k -size subset of the encoded data at the expense of transmitting redundant data. It was interesting to see that FEC improved performance in sparse networks while harming performance in dense networks. The decreased performance in dense networks is due to the existence of highly variable link qualities where nearby neighbors had high link qualities and nodes further away had poor link qualities. FEC works best in environments where loss rates are predictable and have low variance, allowing the amount of redundant data transmitted to be tuned to match the link qualities of the network.

4. EVALUATION METHODOLOGY

The metrics we use to evaluate Deluge are driven by the primary motivation for this work: network programming. We list the metrics we consider, ordered from highest to lowest priority.

1. **Complete Reliability.** Every byte of the data must be correctly received by all nodes.
2. **Completion Time.** In deployments, any interruption to their primary service caused by network programming should be minimized. In the development process, network programming should quickly install updates to shorten debugging and testing cycle.
3. **RAM Usage.** Sensor nodes are severely constrained in the amount of memory available to running application and is compounded by the lack of dynamic memory allocation in TinyOS.
4. **Energy Consumption.** Wireless sensor nodes are strictly limited in their energy capacity and a minimal amount of energy should be used in order to lengthen network lifetimes.

The full Deluge protocol as described in this paper is implemented in the nesC programming language on the TinyOS platform. By executing Deluge, we first observe the overall performance of Deluge under different network diameters, densities, and object sizes. We then investigate Deluge’s reaction to these parameters by examining the detailed propagation dynamics which occur. Finally, we develop a simple model to help identify factors which limit overall performance.

To evaluate and investigate the behavior of Deluge, we use two separate mechanisms. The first is a TinyOS hardware testbed of modest size, composed of Mica2-dot nodes [14]. The second is TOSSIM, a bit-level node simulator designed specifically for the TinyOS platform [7]. We use these two methods to provide varying degrees of realism and scale: the former providing the most realistic environment while the latter allows us to scale to hundreds of nodes and different node topologies.

4.1 TinyOS Hardware

To gather empirical data, we use a testbed composed Mica2-dots, a TinyOS supported hardware platform. Each node contains a 7MHz, 8-bit microcontroller as a CPU, which offers 128KB of program memory and 4KB of RAM; a 512KB external flash chip used for storing application generated data; and communicate via a 433MHz radio transceiver which transmits 19.2Kbit/s [14]. In ideal conditions, the Mica2-dot can transmit about 37 packets per second at a size of 36 bytes each (including headers) after encoding and media access.

We fully implemented Deluge in nesC and ran experiments using the Mica2-dot hardware platform with 75 nodes deployed non-uniformly in a 150’ by 100’ office environment [13]. With the source placed at one corner, the diameter of the network is about five hops. To instrument a back-channel, we used specialized hardware to create a UART to TCP bridge, allowing nodes to transmit and receive messages over TCP through an Ethernet adapter. By opening up sockets to each node from a desktop computer, we timestamp each UART message with precision on the order of milliseconds and track the propagation of each page. Note that the UART to TCP bridge is only used as a mechanism for gathering timing information from the network and does not represent any significant source of noise. At the end of each experiment, we verify the integrity of the data object via the TCP connection.

4.2 TOSSIM

In addition to hardware experiments, we use TOSSIM, a discrete-event network simulator, to investigate and evaluate Deluge at networks of much greater scale and differing structures. TOSSIM compiles directly from unmodified TinyOS application code and simulates communication between nodes at the bit level. Because TOSSIM models a previous generation of TinyOS hardware, the absolute values in time may not match closely with the Mica2-dot hardware. However, it is very useful for exposing overall behavior. In the rest of this section, we focus our discussion on TOSSIM’s radio model since the radio is the most significant shared resource of a WSN and is the most important component when simulating Deluge.

Capturing sufficient detail when simulating the communication between nodes is essential since the behavior of dissemination protocols can be highly sensitive to low level factors. We have experimented with using high-level simulators, but they were unable to capture unique behaviors which appear when simulating low-level details. TOSSIM captures data-link level network interactions with high fidelity by simulating communication between nodes at the bit level, allowing TOSSIM to capture the entire TinyOS network stack and all of its complex behaviors, including the CSMA MAC layer packet transmission delays and backoffs, link-level acknowledgments, packet CRC checks, SECDED packet encoding scheme, sender-receiver synchronization, and hardware-specific timing. The main advantage with simulating at the bit-level is that the transmission and reception of bits govern the actions of each layer, rather than modeling each layer with its own set of parameters.

The network itself is modeled by a weighted directed graph $G = (V, E)$, which is static for the entire duration of the simulation. Each vertex $v' \in V$ represents a node and each edge $(u, v) \in E$ specifies a bit-error rate in the direction

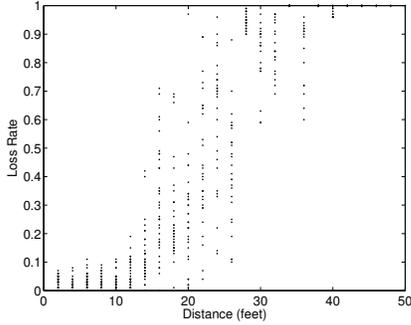


Figure 3: TOSSIM Packet Loss Rates vs. Distance

from u to v . The bit-error rate represents the probability a given bit is flipped while in transmission from u to v . Each bit-error rate is independently chosen based on the distance between u and v and a random distribution derived from empirical data collected on Mica nodes. Figure 3 shows this model’s packet loss rate over distance from an experiment in TOSSIM. Sampling the bit-error rate for each edge independently allows for asymmetric links where the a bit-error rate for (u, v) is significantly different than (v, u) . This is important since asymmetric links are a real problem in WSNs.

TOSSIM treats the transmission of each individual bit as an event. If a node receives bits from multiple senders at the same time, the union is taken as the result. In doing so, real-world wireless problems dealing with interference and the hidden terminal problem can be observed in simulation. However, TOSSIM does make some significant simplifications from the real world. For example, the transmission strength for each node is uniform within a 50 foot radius. This implies that while a close node may have a lower bit-error rate, it is unable to overpower the signal of a further node. This aspect of TOSSIM is overly pessimistic when compared to the real world where signal strength fades at a polynomial rate with distance.

5. PROPAGATION

5.1 Empirical Results

The procedure of each experiment follows that of a normal deployment scenario. In its initial state, all but one of the nodes are deployed and operating in the steady-state, meaning that they have the same version v of the object profile and the same set of completed pages for version v . The remaining node acts as the source node. Initially disconnected from the deployment, a new object with version $v' > v$ and all page ages set to 0 is loaded onto the source via a physical connection from a desktop computer. The experiment begins by introducing the source into the 75 node deployment, beginning the dissemination process. We measure the time to propagate each page to individual nodes, relative to the first advertisement made by the source. No concurrent services, other than those required by Deluge, execute for the duration of each experiment. In the current implementation, each page is 1104 bytes with 48 data packets per page and each data packet has a data payload of 23 bytes. For the maintenance service, we set $\tau_l = 2$ seconds, $\tau_h = 60$ seconds, and $k = 1$. For requests, we set $\tau_r = 0.5$, $\lambda = 2$, and $\omega = 8$.

In every test, the data was correctly and completely received by all nodes in the network. The empirical results are shown in Figure 4. Figure 4(a) shows the average completion time of 10 propagations when varying the object size with 95% confidence intervals. The completion time is linear with the size of the object. With this deployment, Deluge pushes 88.4 bytes/second on average. To see the beneficial effects of spatial multiplexing, the dotted line is a projection of the completion time *without* spatial multiplexing and is represented by multiplying the average completion time for one page by the number of pages in the object. On average, the completion time would increase by 31% without spatial multiplexing. We expect a more significant improvement in larger networks as there is greater opportunity for parallel transfers.

Figure 4(b) shows the distribution of completion times for individual nodes across the same 10 experiments. All nodes have a completion time greater than 220 seconds, but less than 275 seconds. We note that the average range of completion times for an individual run is 18.6 ± 2.60 seconds. The narrow range of completion times is a direct result of spatial multiplexing, where the range of completion times represents the amount of time to flush the pipeline and is essentially the time to disseminate a single page.

To get a better understanding of what packets are transmitted during propagation, we counted the number of transmissions for each packet type on each node when propagating a 20 page object to all 75 nodes. As shown in Figure 4(c), 18% of the transmitted packets are control packets while 82% are data packets.

Using the transmission count information on advertisements, we plot a histogram of the average advertisement rate for each node by dividing the count by the completion time. As shown in Figure 4(d), the average advertisement rate is 0.048 packets/second and the median is 0.042 packets/second. This low rate shows the effectiveness of advertisement suppression. Without suppression, each node would broadcast 0.5 packets/second, more than ten times the average rate than with suppression.

Figure 4(e) displays a histogram based on the number of request transmissions for each node. The average number for all nodes is 10.52 packets while the median is 8. Again, these low numbers show the effectiveness of suppression on request transmissions. Considering that there are 20 pages in the object, many nodes did not transmit requests in order to receive entire pages. In fact, more than half the nodes transmitted less than 10 requests, meaning that those nodes took advantage of the broadcast channel and did not transmit requests to receive a majority of their pages. A few nodes transmitted many more requests and was due to their poor link qualities with neighboring nodes.

Finally, we counted the number of data packet receptions to calculate the ratio of total data packets received to the minimal number of data packets required. Figure 4(f) displays a histogram of the ratio for each node. The minimal number of data packets required is the number of packets per page times the number of pages. On average, a node receives about 3.35 times the minimum number of required data packets, while half receive fewer than 3.34 times the minimum number. This redundancy is mainly due to the single-channel, broadcast network that is used by all nodes. In a simple linear setup, as shown in Figure 2, we would expect a node to receive twice the minimum number: once by

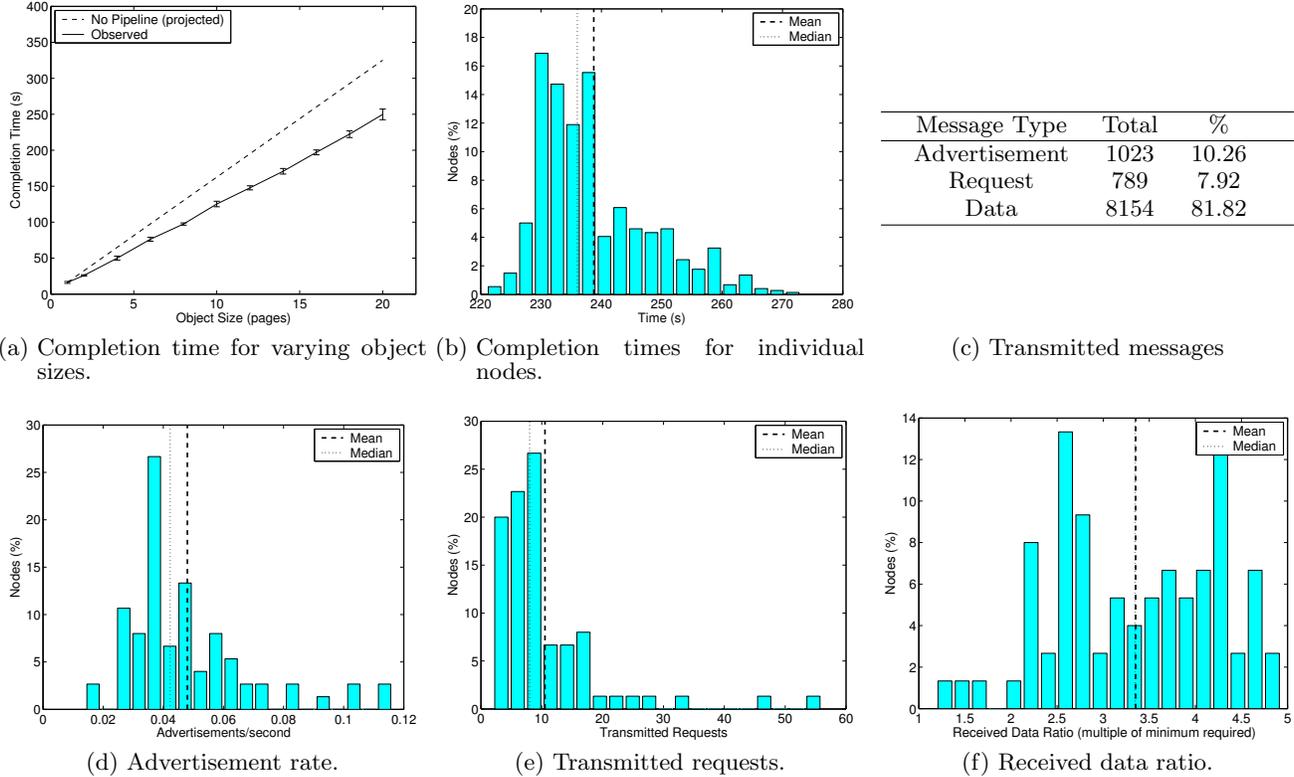


Figure 4: Empirical Data. These results are from a network of 75 nodes. Figure (a) shows the completion time when varying object size. Figure (b) shows the completion time for individual nodes. Figures (c)-(f) are derived from transmitting 20 pages over 75 nodes.

first receiving the data from its parent and once by overhearing its child retransmit the data. In a dense linear topology, we would expect a node to receive nearly three times the minimum number, as a neighboring node may retransmit the data on its behalf since it is in the same radio cell. In a non-linear topology, we would expect a node to receive three to five times the minimum number as a node may have additional neighbors which need to retransmit data. In our experiments, none of the nodes received more than five times the amount of required data packets. Some nodes experienced significantly less redundancy since they lie on the edge of the deployment.

5.2 Simulation Results

While the empirical results are promising, we were unable to experiment with networks of large scale since the testbed configuration does not scale easily. Instead, we used TOSSIM to evaluate and investigate the behavior of Deluge with network sizes on the order of hundreds of nodes and tens of hops. Recall that TOSSIM models a previous generation of TinyOS hardware and we use it mainly for exposing overall behavior, rather than absolute time measurements. This section starts by evaluating the propagation performance with different network diameters, densities, and object sizes for a square topology. We then examine the propagation dynamics in detail to understand how Deluge reacts to the different parameters. Finally, the overall performance of Deluge in a linear network is provided and a simple model

is derived to help identify the different factors which limit overall performance.

5.2.1 Overall Performance

We briefly discuss the overall performance of Deluge for a square topology. Due to execution times on the order of tens of hours for each simulation, we do not provide confidence intervals. We begin by executing application-level code identical to that used in the empirical experiments, including all parameter values except that each page has 24 packets rather than 48 packets. The smaller page size reduces the simulation time while showing the effect of propagating multiple pages. Because Deluge requires nodes to keep their radios on, the vast majority of energy is spent in the listening state where the radio continuously listens to the channel. Thus, we focus on completion time as it closely resembles energy consumption.

Figure 5(a) shows the total propagation time for different object sizes in square topologies of different diameters. Node density is kept constant with nodes spaced 15 feet apart. At network diameters less than 8, the propagation time is roughly a function of the *product* of the network diameter and object size. With network diameters greater than 8, the slope of the propagation times for the multi-page objects approaches that of a single page, indicating the effectiveness of spatial multiplexing (i.e. the propagation time is the *sum* of the network diameter and object size). Figure 5(b) plots the propagation times for various object sizes in a 20×20 grid topology, showing that propagation time is

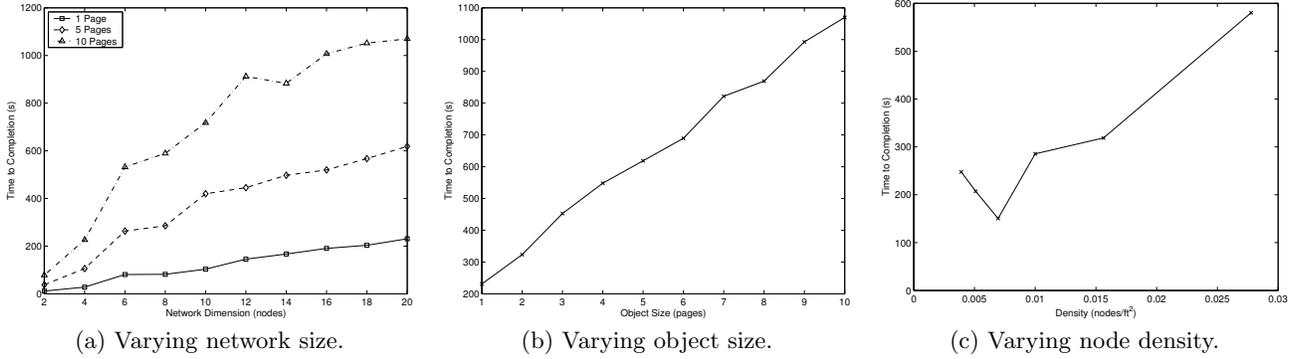


Figure 5: Simulated Propagation Time for Square Structures. These results are from a $N \times N$ network where $N = 2, 4, 8, \dots, 20$ with 15' spacing. Figure (c) is for a $180' \times 180'$ field for different densities.

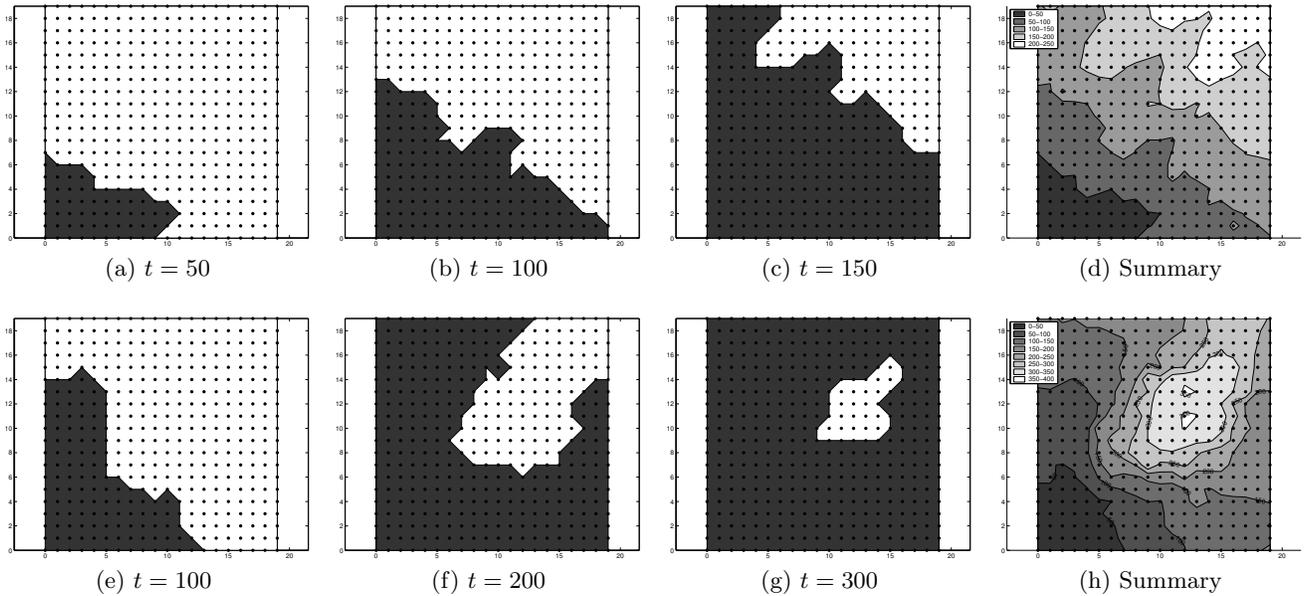


Figure 6: Simulated Propagation Time for 1 Page in a 20×20 Grid Topology. Figures (a) - (d) and (e) - (h) are from a network with 15' and 10' node spacing, respectively.

linear with object size. Figure 5(c) shows that an increase in density harms overall performance. In the next section, we investigate the propagation dynamics of Deluge to see why density affects performance.

5.2.2 Dynamic Behavior

In this section, we investigate the propagation behavior of Deluge for a square topology. Figure 6 shows the propagation of a single page from a corner node through a 20×20 network with nodes spaced 15 feet apart for a single experiment and is representative of all other experiments of similar topologies. Figures 6(a)-6(c) show a time series of the propagation while Figure 6(d) summarizes the completion time for all nodes. With this topology, the propagation behaves as expected: the propagation progresses at a fairly constant rate in a nice wavefront pattern from corner to corner. The irregularities are due to the non-uniform loss rates and contention.

An interesting behavior emerges as we increase density. To show this behavior, we repeat the experiment with a 20×20 network with adjacent nodes spaced 10 feet apart, increasing the density by $2\frac{1}{4}$ times. As shown in Figure 6(e), the propagation began as it did in the sparse case. But soon after, the propagation along the diagonal slowed significantly while quick propagation along the edge continued and completely wrapped around the edge before filling in the middle.

To see the behavior with spatial multiplexing, Figure 7 shows the propagation of a five page object in both the sparse and dense case. With multiple pages, the complex rate of progress becomes apparent even in the sparse case. A closer look at the sparse case when propagating a single page shows that the propagation along the diagonal is actually slightly slower than along the edge, which accounts for the linear wavefront shape. The difference in propagation times is not enough to clearly show the behavior as

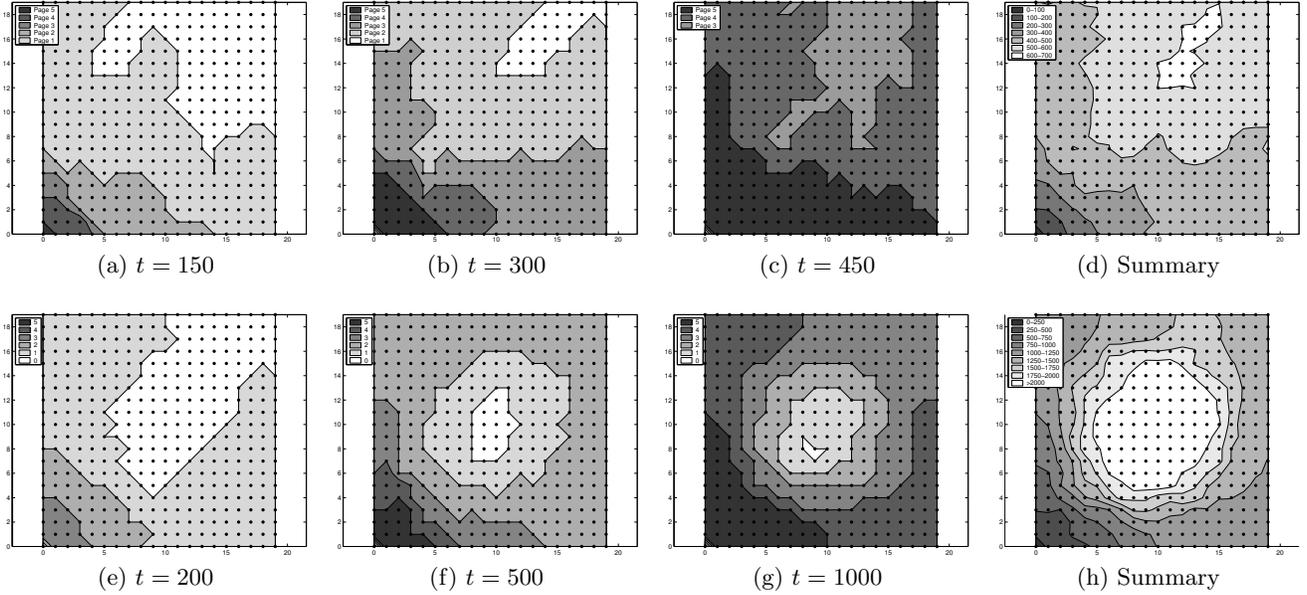


Figure 7: Simulated Propagation Time for 5 Pages in a 20×20 Grid Topology. Figures (a) - (d) and (e) - (h) are from a network with 15' and 10' node spacing, respectively.

experienced by the dense case. However, when propagating multiple pages, the delay in reaching the middle nodes accumulates over time and thus forms the same pattern as in the dense case. In the dense case, the behavior becomes even more pronounced.

To examine this behavior in greater detail, we plot the rate of propagation for a single page along the edge and diagonal in Figure 8(a). We use the propagation times from only those nodes on the bottom edge (i.e. $(i, 0) \forall i \in [0, 19]$) and along the diagonal (i.e. $(i, i) \forall i \in [0, 19]$) to represent the propagation along the edge and diagonal respectively. The slope of the plotted times represents the propagation rate. Early in the process, the propagation rate along both the edge and diagonal are identical. However, the propagation rate along the diagonal drops to nearly 20% of the original rate once it reaches node (5,5). The propagation rate remains fairly constant after the drop.

The root cause of this behavior is the *hidden terminal problem*, which occurs when two nodes A and C communicating within the range of node B are unable to coordinate their transmissions, thus causing collisions when transmitting to B . Nodes in the center of the network have more neighboring nodes and experience a greater probability of collisions than those on the edge of the network. Considering that nodes are spaced 10' apart with TOSSIM's 50' interference range, a significant slow down at node (5,5) is understandable since it, along with all other nodes in the center, experiences the greatest number of nodes within its interference range.

To see the effects of interference, we count the number of transmissions within interference range and the number of messages overheard which pass the CRC check. Figures 9(a) and 9(b) show the number of each message type sent in a 10 second window within the interference range of nodes (2,2) and (5,5) respectively. Node (5,5) clearly shows a higher number of sustained transmissions within its interference

range, causing a higher likelihood of collisions. Figures 9(c) and 9(d) show the effects of the interference by plotting the number of messages overheard which pass the CRC check. Even though the rate of transmissions within node (5,5)'s interference range is greater, the rate of correctly received messages is, at best, half of node (2,2).

Deluge's reaction to collisions is compounded by a more subtle problem. Recall that Deluge achieves its density-awareness by overhearing similar packets and suppresses redundancies by employing a linear backoff scheme. Thus, it relies on overhearing packets to estimate node density. When the channel is pushed to saturation, the high number of collisions can cause such a mechanism to underestimate the number neighbors, stimulating transmission of more redundant messages, causing more collisions, and leading to congestion collapse. While Deluge borrows the suppression mechanisms from Trickle, it differs in that Deluge operates near channel capacity to quickly disseminate large amounts of data. The authors of Trickle studied its performance with a relatively low data rate. Figure 9(b) shows that the advertisement rate increases in proportion with other activity, providing evidence that Deluge is underestimating the node density.

To test how Deluge performs when the channel is not pushed to saturation, we increase the backoff time for a request. We chose to vary the request backoff time for several reasons. First, decreasing the rate of requests has the effect of decreasing the rate of invoking nodes to begin transmitting data. Second, in contrast to data messages, request messages can be transmitted by any node needing data. With a large set of requesters, request messages are more likely to increase collisions due to the hidden terminal problem. Because requests are unicast to the node that most recently advertised, it is unlikely for many senders in a region to begin transmitting data. Finally, Figure 9(d) clearly shows that the amount of request traffic can signifi-

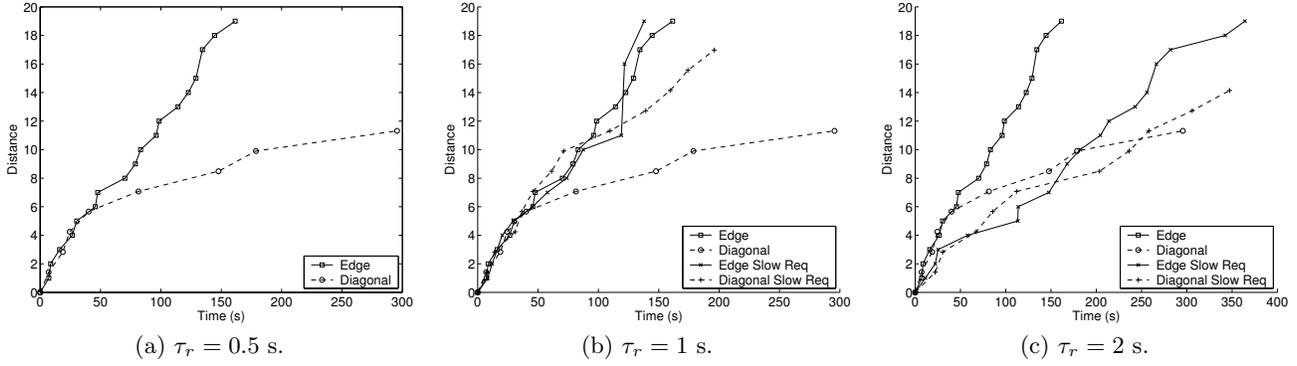


Figure 8: Simulated Propagation Rate of a Single Page. Figure (a) shows the propagation rate through the diagonal and edge while Figures (b) and (c) compare the propagation rate for different values of τ_r .

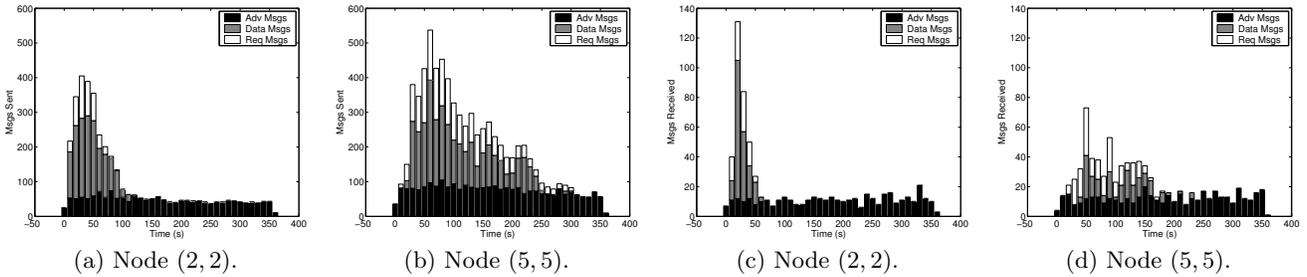


Figure 9: Number of Messages Sent and Received. Figures (a) and (b) show the number of messages sent within the interference range. Figures (c) and (d) show the number of messages overheard by the node.

cantly overcome the amount of data traffic for those nodes which experience slow propagation times. By decreasing the request traffic, useful data traffic can utilize more of the channel.

Figure 8(b) compares the propagation of a single page along the edge and diagonal for a τ_r of 1 second with the original τ_r of 0.5 seconds. By doubling τ_r , the propagation rate along the diagonal improves by about 2.7 times while the propagation rate along the edge remains nearly identical, leading to an improvement in overall propagation performance. With these improved effects, we continue to increase τ_r to 2 seconds. Figure 8(c) compares the propagation of a single page along the edge and diagonal for a τ_r of 2 seconds with the original τ_r of 0.5 seconds. While the propagation rate across the diagonal is slightly improved over the original test, it does not match Deluge’s performance when τ_r is 1 second. The propagation rate along the edge is worse than either of the other tests. This leads to lowered overall performance. Note that the propagation rate across the diagonal only experiences a 28% rate reduction relative to the edge rather than the 80% reduction seen in the original test. This shows that lowering channel utilization is effective in eliminating the hidden terminal problem and minimizing the difference in propagation rates between the edge and diagonal. However, it is also at odds with lowering the overall propagation rate, a primary goal of Deluge.

In each experiment thus far, the propagation began at a corner node. We now examine the behavior when the source node is placed at the center of the network. One might suggest that starting the propagation in the center

might help to eliminate the behavior of following the edge and also decrease the propagation time by about half. We repeat the simulations with τ_r at 0.5 seconds except with a 21×21 grid and the source node at the center of the network rather than at the corner.

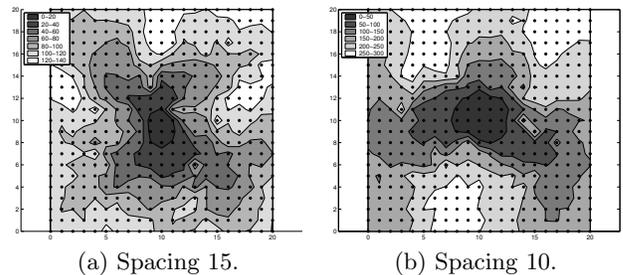


Figure 10: Simulated Propagation for 1 Page from the Center in a 21×21 Grid Topology.

Figure 10(a) shows the propagation behavior for the sparse case (spacing 15). The time to reach all nodes is reduced by approximately 40%. Additionally, the propagation still does not behave in a nice circular manner. One cause is Deluge’s depth-first tendency, where propagation of a single page along good links is not blocked by delays caused by poor links. Notice that the propagation speeds up as it approaches the edge of the network. Figure 10(b) shows the propagation behavior for the dense case (spacing 10).

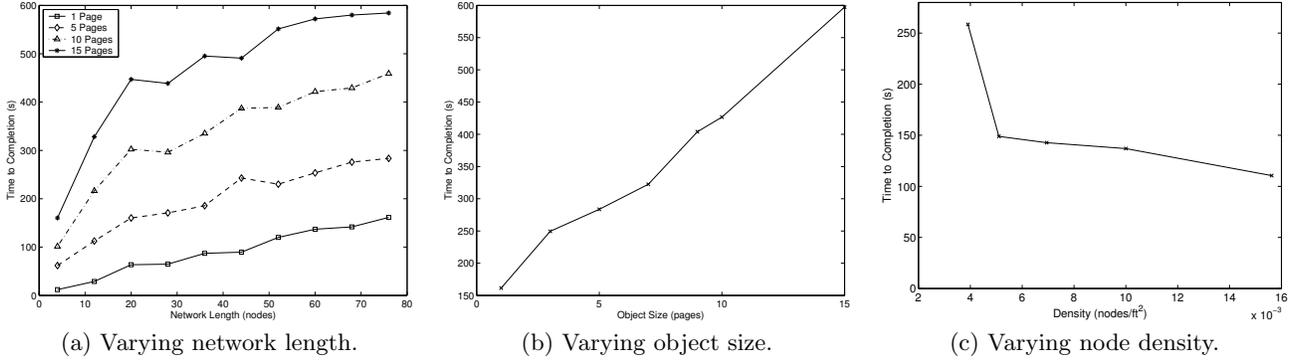


Figure 11: Simulated Propagation Time for Linear Structures. These results are from a $2 \times N$ network where $N = 4, 12, 20, \dots, 76$ and $10'$ spacing.

The time to reach all nodes is reduced by approximately 25%, significantly less than the expected 50%. Even though placing the source at the center effectively reduces the diameter by about half, Deluge is unable to take advantage of the quick edges since nodes in the center experience a greater number of collisions. However, once the propagation reaches an edge, it begins to speed around the edge as before.

5.2.3 Linear Structures

We briefly discuss the Deluge's overall performance in linear networks. Figure 11 shows similar tests as in the square topology case but for a $2 \times N$ topology, where N ranges from 4 to 76 nodes. The linear topology can be considered a special case of the grid by representing just the edge. Figure 11(a) shows the effectiveness of the spatial multiplexing and Figure 11(b) shows that propagation time is linear with object size. Unlike the square topology, Figure 11(c) shows that an increase in density actually improves performance.

5.2.4 Model

For the linear case, the simulations show that Deluge takes about 40 seconds to disseminate each page to 152 nodes across 15 hops. Considering that the link bandwidth between two nodes in TOSSIM can reach 32 packets/second, this seems suboptimal, only achieving 4% of this rate. To help understand this difference and the various factors which contribute to it, we develop a simple model of the Deluge process for the linear case.

Unlike wired networks, the broadcast nature of single-channel, wireless networks prevents any multihop communication protocol from achieving the maximum transmission rate of the radio. Recall from Figure 2 that a three hop spacing is required to prevent collisions from simultaneous transmissions, limiting any multihop communication scheme to 33% of the maximum transmission rate.

Assume that the expected time to transmit a page is $E[T_{tPage}]$ and the delay between completing a page and requesting a new page is D_{newReq} , the expected time to transmit an object of n pages across a d hop network is

$$E[T_{obj}] = \min(d \cdot n, d + 3(n - 1)) \cdot (E[T_{tPage}] + D_{newReq}). \quad (1)$$

Given a packet loss rate $E[r_l]$, we define the expected number of transmissions for a given packet as $E[N_{tPkt}] =$

$\frac{1}{1 - r_l}$. The expected time required to transmit just the data packets is

$$E[T_{tx}] = E[N_{tPkt}] \cdot T_{tPkt} \cdot N, \quad (2)$$

where T_{tPkt} is the transmission time for a single packet. Given that $T_{tPkt} = 32.2$ ms and $E[r_l] = 0.1$, as derived from the simulation data, $E[T_{tx}]$ is relatively small.

Much of the remaining difference from the maximum transmission rate is due to the various delays and backoffs within Deluge. The first source of delay is the random backoffs before transmitting requests, where the expected backoff delay of $E[\tau_r] = \frac{\tau_r}{2}$. When receiving a page, the expected time spent backing off while making requests is

$$E[T_{req}] = E[N_{tPkt}] \cdot E[N_{reqs}] \cdot E[\tau_r], \quad (3)$$

where $E[N_{reqs}]$ is the expected number of requests a node must make to complete a given page.

The second source of delay comes from the relatively low rate of advertisements. The expected time for receiving advertisements from nodes closer to the source is given by

$$E[T_{rAdv}] = E[N_{tPkt}] \cdot \frac{\tau_l}{2} \cdot (1 + E[N_{Supp}]) \quad (4)$$

where $E[N_{Supp}]$ specifies the expected number of times that an advertisement from hop h is suppressed by hop $h - 1$ before transmitting an advertisement (in the linear case, $E[N_{Supp}] = 1$). Assume node R is at hop $h + 1$ while S is at hop h . After S acquires new data to transmit, the expected time for R to learn that S has new data and transition to RX is $E[T_{rAdv}]$.

Additionally, when a node exceeds its limit of λ requests, it must transition to *MAINTAIN* and wait for another advertisement before making additional requests. Thus, the expected amount of time spent waiting for additional advertisements when the request limit of λ is exceeded is

$$E[T_{tGiveUp}] = E[N_{tPkt}] \cdot \left(\frac{E[N_{reqs}]}{\lambda} - 1 \right) \cdot E[T_{rAdv}]. \quad (5)$$

We now define $E[T_{tPage}]$ (the expected time required to transmit a page across a single hop) used in (Equation 1) as

$$E[T_{tPage}] = E[T_{rAdv}] + E[T_{tx}] + E[T_{req}] + E[T_{tGiveUp}]. \quad (6)$$

We test the accuracy of the model by fitting it with the simulated data from the linear case. The parameters $N_{hops} =$

15, $T_{pkt} = 0.0322$ seconds, $r_l = 0.1$, and $E[N_{reqs}] = 5.4$ represent averages derived from the simulated data while the remaining parameters are identical to the ones used in the experiments. Figure 12 shows that the model fits fairly well the simulated data.

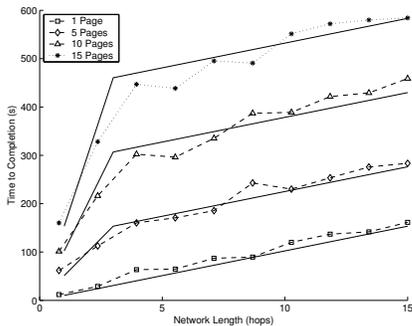


Figure 12: Comparison of Model with Simulation. The solid lines represent the predicted times from the model.

Pipelining	$E[T_{req}]$	$E[T_{tGiveUp}]$	$E[T_{rAdv}]$	$E[T_{tx}]$
66%	14%	11%	5%	4%

Table 1: Contributions to Reduced Bandwidth from the Maximum Transmission Rate.

From the model, the items that contribute to the difference from the maximum transmission rate are shown in Table 1. As mentioned, spatial multiplexing alone accounts for 66% of the difference and is a fundamental limit of any multihop communication protocol in single-channel, wireless networks. The remaining 30% of the difference is due to the delays and backoffs in Deluge and represent a trade-off: while we would like to minimize these delays to increase bandwidth, they are also necessary for suppression and without suppression, contention would slow propagation. Thus, there is an inherent limit to the dissemination rate which is much less than that of simply routing a single message. With a greater density, a larger delay period reduces any collisions caused by the hidden terminal problem. We have tested Deluge in very sparse cases with low delay settings and it is able to achieve higher transfer rates.

6. CURRENT STATUS

The fully implemented Deluge protocol is included as a core part of TinyOS 1.1.8 to support network programming on the Mica2, Mica2-dot, MicaZ, and Telos platforms, the latter two using 2.4 GHz, IEEE 802.15.4 radios. The current implementation consumes 84 bytes of RAM, of which 43 bytes is a message buffer. In addition to the features described in this paper, Deluge supports multiple objects on each node by including an object summary of each object in the advertisement. This allows a user to switch between two or more program objects without continuously propagating a new object.

We codeveloped *TOSBoot*, a bootloader that can reprogram a node by transferring a Deluge object from external flash into program flash. TOSBoot is a static piece of code that executes each time the node exits the reset state. Arguments to TOSBoot are passed via the internal flash of the

MCU and indicate whether to program the node and which Deluge object to program the node with. If there is no command to program the MCU, TOSBoot simply invokes the TinyOS application. The use of non-volatile storage to pass arguments provides robust operation in the face of node failures.

In addition to Deluge’s objects, the hardware manufacturer can install a program that contains Deluge and other minimal services in the write-protect region of external flash. TOSBoot programs a node with the factory installed object on request by the user (through a radio command or a physical gesture) or if the node experiences failures. This factory installed object is crucial for node reliability by providing a solid software base that allows user interaction without programming the node using a physical connection.

7. LOOKING FORWARD

While Deluge provides good, robust performance, there is room for potential improvement. From our experiments, additional design options become apparent. One possibility is to employ an exponential backoff rather than a linear backoff to accurately estimate the node density and avoid congestion collapse. Dynamically adjusting the backoffs may have the additional effect of minimizing wasted time caused by backoffs. In areas where nodes are sparse, the backoffs can be relatively small while contention remains low and the suppression mechanisms remain effective.

The increased performance along the edge introduces an interesting concept for disseminating large data objects. Ideally, we would like to duplicate the beneficial edges in the center of the network. One way is to select a linear set of nodes through the center of the network to initially disseminate data. This allows for quick propagation across the center of the network, splitting the network in two. We might recursively continue this process, creating a fractal structure for propagation. This structured approach should improve the propagation time for a single page, but inhibits the use of pipelining since it is more difficult to minimize interference between transfers of different pages. Also, methods for creating these paths in a distributed manner need to be explored.

We have experimented with both adaptive delays and structured approaches using TOSSIM. In some cases, we were able to improve performance. However improvement was limited to a 28% reduction in completion time even when the protocols relied on ideal information. For example, we assumed that nodes knew the location of their one-hop neighbors and had accurate bi-directional link-quality estimates for each neighbor. We were unable to achieve comparable results with on-line estimators. A complete analysis of our results is beyond the scope of this paper.

This raises an important open question: What is the true limit on the rate of dissemination in wireless networks? It is clear that this bound is much lower than simply routing a single message across the network. For large data objects, we have shown that the use of spatial multiplexing provides a lower bound than without it. However, spatial multiplexing limits a node’s broadcast rate to no greater than one-third the maximum rate due to the single-channel, broadcast medium. The limit of one-third assumes a perfectly formed, linear network topology with no interference between nodes other than their direct children and grandchildren. Many networks are not of this form even under a uniform distri-

bution because of the unpredictable nature of wireless communications. In Section 5.1, we argued that nodes receive up to four or five times the minimum number of required packets in realistic deployments with single-channel, broadcast radios. This would reduce the transmission rate limit to one-fourth or one-fifth the maximum rate. This limit does not even consider packet-loss. Additionally, the unpredictable interference causes additional sources of potential collisions between nodes and leads to the use of selective and delayed retransmissions, implemented in Deluge through the use of suppression mechanisms. These delays further limit the maximum transmission rate. From the empirical results, Deluge operates at about one-ninth the maximum rate when using TinyOS defaults.

The real-world deployment shows much better results for Deluge than simulation. The propagation dynamics described are determined by very low level factors and hence may be influenced by the simulator itself. Because TOSSIM models a uniform signal strength within a 50' radius, the impact of interfering transmissions may be overestimated. The same experiments using a 20' interference radius model do not experience the edge effect and show a 70% reduction in completion time, better representing the real-world deployment. However, we are confident that such effects occur at some density in real deployments and that the actual density may depend heavily on the way the network behaves. While previous work has hinted at similar behaviors in real deployments, the effects are too complex to state anything conclusive and only small amounts of data are disseminated.

8. CONCLUSIONS

In this paper, we presented *Deluge*, a reliable data dissemination protocol for propagating large data objects from one or more source nodes to many other nodes over a multi-hop WSN. With its density-aware, epidemic mechanisms, we showed that Deluge can reliably disseminate data to all nodes at a rate of 90 bytes/second in a real-world deployment, one-ninth the maximum transmission rate of the radio supported under TinyOS. Control messages are limited to 18% of all transmissions. At scale, Deluge exposes propagation dynamics only hinted at by previous work, showing the impact of the hidden terminal problem on dissemination. We presented a simple model of Deluge's propagation behavior, describing factors which limit its propagation performance. Finally, we argued that dissemination is inherently slower than single path propagation and identified establishing a tight lower bound as an open problem. Unlike wired networks, dissemination protocols cannot achieve an aggregate bandwidth near the link capacity due to the single-channel radio, spatial multiplexing, and delays necessary for suppression. It appears very hard to significantly improve upon the rate obtained by Deluge.

9. ACKNOWLEDGMENTS

Special thanks to Gilman Tolle for his helpful input and early contributions to Deluge. This work was supported by the Defense Advanced Research Projects Agency (grant F33615-01-C-1895), the National Science Foundation (grants EIA-0122599 and IIS-033017), the Center for Information Technology Research in the Interest of Society (CITRIS), Intel, Sun Microsystems, Hewlett Packard, Microsoft, and the California MICRO program.

10. REFERENCES

- [1] A. Demers, D. Greene, C. Hauser, W. Irish, and J. Larson. Epidemic algorithms for replicated database maintenance. In *Proceedings of the Sixth Annual ACM Symposium on Principles of Distributed Computing*, pages 1–12. ACM Press, 1987.
- [2] D. Ganesan, B. Krishnamachari, A. Woo, D. Culler, D. Estrin, and S. Wicker. Complex behavior at scale: An experimental study of low-power wireless sensor networks. Technical Report UCLA/CSD-TR 02-0013, UCLA, 2002.
- [3] J. Jeong, S. Kim, and A. Broad. *Network Reprogramming*. University of California at Berkeley, Berkeley, CA, USA, August 2003.
- [4] S. K. Kaser, G. Hjálmtýsson, D. F. Towsley, and J. F. Kurose. Scalable reliable multicast using multiple multicast channels. *IEEE/ACM Transactions on Networking*, 8(3):294–310, 2000.
- [5] J. Kulik, W. R. Heinzelman, and H. Balakrishnan. Negotiation-based protocols for disseminating information in wireless sensor networks. *Wireless Networks*, 8(2-3):169–185, 2002.
- [6] P. Levis and D. Culler. Maté: a tiny virtual machine for sensor networks. In *10th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-X)*, pages 85–95. ACM Press, 2002.
- [7] P. Levis, N. Lee, M. Welsh, and D. Culler. TOSSIM: Accurate and scalable simulation of entire tinycos applications. In *Proceedings of the First ACM Conference on Embedded Networked Sensor Systems (SenSys 2003)*. ACM Press, November 2003.
- [8] P. Levis, N. Patel, S. Shenker, and D. Culler. Trickle: A self-regulating algorithm for code propagation and maintenance in wireless sensor networks. Technical report, University of California at Berkeley, 2004.
- [9] S.-Y. Ni, Y.-C. Tseng, Y.-S. Chen, and J.-P. Sheu. The broadcast storm problem in a mobile ad hoc network. In *Proceedings of the Fifth Annual ACM/IEEE International Conference on Mobile Computing and Networking*, pages 151–162. ACM Press, 1999.
- [10] F. Stann and J. Heidemann. RMST: Reliable data transport in sensor networks. In *Proceedings of the First International Workshop on Sensor Net Protocols and Applications*, pages 102–112, Anchorage, Alaska, USA, April 2003. IEEE.
- [11] T. Stathopoulos, J. Heidemann, and D. Estrin. A remote code update mechanism for wireless sensor networks. Technical report, UCLA, Los Angeles, CA, USA, 2003.
- [12] R. Szewczyk, J. Polastre, A. Mainwaring, and D. Culler. Lessons from a sensor network expedition. In *Proceedings of the First European Workshop on Sensor Networks (EWSN)*, Berlin, Germany, Jan. 2004.
- [13] University of California, Berkeley. Smote testbed. <http://smote.cs.berkeley.edu/>.
- [14] University of California, Berkeley. Mica2-dot schematics. http://webs.cs.berkeley.edu/tos/hardware/design/ORCAD_FILES/MICA2/6310-0306-01ACLEAN.pdf, March 2003.
- [15] University of California, Berkeley. Tinyos. <http://www.tinyos.net/>, 2004.
- [16] C.-Y. Wan, A. T. Campbell, and L. Krishnamurthy. PSFQ: A reliable transport protocol for wireless sensor networks. In *Proceedings of the 1st ACM International Workshop on Wireless Sensor Networks and Applications*, pages 1–11. ACM Press, 2002.