

# Reliable Bursty Convergecast in Wireless Sensor Networks \*

Hongwei Zhang<sup>†</sup>, Anish Arora<sup>‡</sup>, Young-ri Choi<sup>§</sup>, Mohamed G. Gouda<sup>§</sup>

## Abstract

We address the challenges of bursty convergecast in multi-hop wireless sensor networks, where a large burst of packets from different locations needs to be transported reliably and in real-time to a base station. Via experiments on a 49 MICA2 mote sensor network using a realistic traffic trace, we determine the primary issues in bursty convergecast, and accordingly design a protocol, RBC (for Reliable Bursty Convergecast), to address these issues: To improve channel utilization and to reduce ack-loss, we design a window-less block acknowledgment scheme that guarantees continuous packet forwarding and replicates the acknowledgment for a packet; to alleviate retransmission-incurred channel contention, we introduce differentiated contention control. Moreover, we design mechanisms to handle varying ack-delay and to reduce delay in timer-based retransmissions. We evaluate RBC, again via experiments, and show that compared to a commonly used implicit-ack scheme, RBC doubles packet delivery ratio and reduces end-to-end delay by an order of magnitude, as a result of which RBC achieves a close-to-optimal goodput.

**Keywords:** wireless sensor network, bursty convergecast, error control, contention control, experimentation with real networks and testbeds

## 1 Introduction

A typical application of wireless sensor networks is to monitor an environment (be it an agricultural field or a classified area) for events that are of interest to the users. Usually, the events are rare. Yet when an event occurs, a large burst of packets is often generated that needs to be transported reliably and in real-time to a base station. One exemplary event-driven application is demonstrated in the DARPA NEST field experiment “A Line in the Sand” (simply called *Lites* hereafter) [4]. In *Lites*, a typical event generates up to 100 packets within a few seconds and the packets need to be transported from different

---

\*An extended abstract containing some preliminary results of this paper appeared in the 6th ACM International Symposium on Mobile Ad Hoc Networking and Computing (MobiHoc 2005). This work was sponsored by DARPA contract OSU-RF #F33615-01-C-1901.

<sup>†</sup>H. Zhang is with the Department of Computer Science, Wayne State University, Detroit, MI 48202, U.S.A. E-mail: hzhang@cs.wayne.edu.

<sup>‡</sup>A. Arora is with the Department of Computer Science and Engineering, The Ohio State University, Columbus, OH 43210, U.S.A. E-mail: anish@cse.ohio-state.edu.

<sup>§</sup>Y. Choi and M. Gouda are with the Department of Computer Sciences, The University of Texas at Austin. Austin, TX 78712-0233, U.S.A. E-mail: {yrchoi, gouda}@cs.utexas.edu.

network locations to a base station, over multi-hop routes.

The high-volume bursty traffic in event-driven applications poses special challenges for reliable and real-time packet delivery. The large number of packets generated within a short period leads to high degree of channel contention and thus a high probability of packet collision. The situation is further exacerbated by the fact that packets travel over multi-hop routes: first, the total number of packets competing for channel access is increased by a factor of the average hop-count of network routes; second, the probability of packet collision increases in multi-hop networks due to problems such as hidden-terminals. Consequently, packets are lost with high probability in bursty convergecast. For example, with the default radio stack of TinyOS [2], around 50% of packets are lost for most events in Lites.

For real-time packet delivery, hop-by-hop packet recovery is usually preferred over end-to-end recovery [15, 17]; and this is especially the case when 100% packet delivery is not required (for instance, for bursty convergecast in sensor networks). Nevertheless, we find issues with existing hop-by-hop control mechanisms in bursty convergecast. Via experiments with a testbed of 49 MICA2 motes and with traffic traces of Lites, we observe that the commonly used link-layer error control mechanisms do not significantly improve and can even degenerate packet delivery reliability. For example, when packets are retransmitted up to twice at each hop, the overall packet delivery ratio increases by only 6.15%; and when the number of retransmissions increases, the packet delivery ratio actually decreases, by 11.33%.

One issue with existing hop-by-hop control mechanisms is that they do not schedule packet retransmissions appropriately; as a result, retransmitted packets further increase the channel contention and cause more packet loss. Moreover, due to in-order packet delivery and conservative retransmission timers, packet delivery can be significantly delayed in existing hop-by-hop mechanisms, which leads to packet backlogging and reduction in network throughput. (We examine the details in Section 3.)

On the other hand, the new network and application models of bursty convergecast in sensor networks offer unique opportunities for reliable and real-time transport control:

- First, the broadcast nature of wireless channels enables a node to determine, by snooping the channel, whether its packets are received and forwarded by its neighbors.
- Second, time synchronization and the fact that data packets are timestamped relieve transport layer from the constraint of in-order packet delivery, since applications can determine the order of packets by their timestamps.

Therefore, techniques that take advantage of these opportunities and meet the challenges of reliable and real-time bursty convergecast are desired.

**Contributions of the paper.** We study the limitations of two commonly used hop-by-hop packet recov-

ery schemes in bursty convergecast. We discover that the lack of retransmission scheduling in both schemes makes retransmission-based packet recovery ineffective in the case of bursty convergecast. Moreover, in-order packet delivery makes the communication channel under-utilized in the presence of packet- and ack-loss.

To address the challenges, we design protocol RBC (for Reliable Bursty Convergecast). Taking advantage of the unique sensor network models, RBC features the following mechanisms:

- To improve channel utilization, RBC uses a window-less block acknowledgment scheme that enables continuous packet forwarding in the presence of packet- and ack-loss. The block acknowledgment also reduces the probability of ack-loss, by replicating the acknowledgment for a received packet.
- To ameliorate retransmission-incurred channel contention, RBC introduces differentiated contention control, which ranks nodes by their queuing conditions as well as the number of times that the enqueued packets have been transmitted. A node ranked the highest within its neighborhood accesses the channel first.

In addition, we design techniques that address the challenges of timer-based retransmission control in bursty convergecast:

- To deal with continuously changing ack-delay, RBC uses adaptive retransmission timer which adjusts itself as network state changes.
- To reduce delay in timer-based retransmission and to expedite retransmission of lost packets, RBC uses block-NACK, retransmission timer reset, and channel utilization protection.

We evaluate RBC by experimenting with an outdoor testbed of 49 MICA2 motes and with realistic traffic trace from the field sensor network of Lites. Our experimental results show that, compared with a commonly used implicit-ack scheme, RBC increases the packet delivery ratio by a factor of 2.05 and reduces the packet delivery delay by a factor of 10.91. Moreover, RBC achieves a goodput of 6.37 packets/second for the traffic trace of Lites, almost reaching the optimal goodput — 6.66 packets/second — for the trace.

**Organization of the paper.** We describe our testbed and discuss the experiment design in Section 2. In Section 3, we study the limitations of existing hop-by-hop control mechanisms. We present the detailed design of RBC in Section 4, then we present the experimental results in Section 5. We discuss how to extend the basic design to support continuous event convergecast and to deal with queue congestion in Section 6. We discuss related works in Section 7, and we make concluding remarks in Section 8.

## 2 Testbed and experiment design

Towards characterizing the issues in making bursty convergecast both reliable and timely, we conduct an experimental study. We choose experimentation as opposed to simulation in order to gain higher fidelity and confidence in the observations. Before presenting our study, we first describe our testbed and the experiment design.

**Testbed.** We setup our testbed to reflect the field sensor network of Lites, and we use the traffic trace for a typical event in Lites as the basis of our experiments.

The testbed consists of 49 MICA2 motes deployed in a grass field, as shown in Figure 1(a), where the

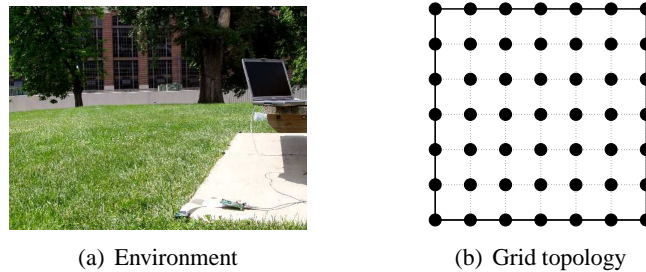


Figure 1: The testbed

grass is 2-4 inches tall. The 49 motes form a  $7 \times 7$  grid with a 5-foot separation between neighboring grid points, as shown in Figure 1(b) where each grid point represents a mote. The mote at the left-bottom corner of the grid is the base station to which all the other motes send packets. The  $7 \times 7$  grid imitates a subgrid in the sensor network of Lites.

The traffic trace (simply called *Lites trace* hereafter) corresponds to the packets generated in a  $7 \times 7$  subgrid of the Lites network when a vehicle passes across the middle of the Lites network. When the vehicle passes by, each mote except for the base station detects the vehicle and generates two packets, which correspond to the start and the end of the event detection respectively and are separated 5-6 seconds on average. Overall, 96 packets are generated each time the vehicle passes by.<sup>1</sup> The cumulative distribution of the number of packets generated during the event is shown in Figure 2. (Interested readers can find the detailed description of the traffic trace in [3].)

If we define the burst rate up to a moment in the event as the number of packets generated so far divided by the time since the first packet is generated, the highest burst rate in Lites trace is 14.07 packets/second. Given that the highest one-hop throughput is about 42.93 packets/second for MICA2 motes with B-MAC (the latest MAC component of TinyOS) and that, in multi-hop networks, even an ideal MAC can only

---

<sup>1</sup>We could have chosen a traffic trace where fewer number packets are generated (e.g., when a soldier with a gun passes by), but that would not serve as well in showing the challenges posed by huge event traffic bursts.

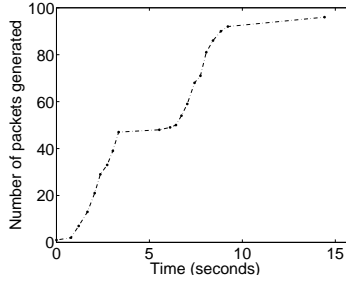


Figure 2: The distribution of packets generated in Lites trace

achieve  $\frac{1}{4}$  of the throughput that a single-hop transmission can achieve [10], the burst rate of Lites trace far exceeds the rate at which the motes can push packets to the base station. Therefore, it is a challenging task to deliver packets reliably and in real time in such a heavy-load bursty traffic scenario.

**Experiment design.** To reflect the multi-hop network of Lites, we let each mote transmit at the minimum power level by which two motes 10 feet apart are able to reliably communicate with each other, and the power level is 9 (out of a range between 1 and 255). We use the routing protocol LGR [6] in our testbed.<sup>2</sup> LGR uses links that are reliable in the presence of bursty traffic, and LGR spreads traffic uniformly across different paths to reduce wireless channel contention. Therefore, LGR provides a reliable and uniform packet delivery service in bursty convergecast [6]. In our testbed, the number of hops in a path is up to 6 and is 3.3 on average.

For each protocol we evaluate, we run the Lites trace 10 times and measure the average performance of the protocol by the following metrics:

- *Event reliability (ER)*: the number of unique packets received at the base station in an event divided by the number packets generated for the event.

Event reliability reflects how well an event is reported to the base station.

- *Packet delivery delay (PD)*: the time taken for a packet to reach the base station from the node that generates it.
- *Event goodput (EG)*: the number of unique packets received at the base station divided by the interval between the moment the first packet is generated and the moment the base station receives any packet the last time.

Event goodput reflects how fast the traffic of an event is pushed from the network to the base station. By definition, the optimal event goodput for Lites trace is 6.66 packets/second, which corresponds to the case where the packet delivery delay is 0 and all the packets are received by the base station.

---

<sup>2</sup>To focus on transport issues, we disable the “base-snooping” in LGR so that the base station does not accept packets snooped over the channel.

- *Node reliability (NR)*: the number of unique packets that are generated by a node and received by the base station divided by the number of packets generated at the node.

(Remark: The study in this paper applies to cases where network topologies other than grid and routing protocols other than LGR are used, since the protocols studied are applicable to other network topologies and routing protocols.)

### 3 Limitations of two hop-by-hop packet recovery mechanisms

Two widely used hop-by-hop packet recovery mechanisms in sensor networks are synchronous explicit ack and stop-and-wait implicit ack. We study their performance in bursty convergecast as follows.

#### 3.1 Synchronous explicit ack (SEA)

In SEA, a receiver switches to transmit-mode and sends back the acknowledgment immediately after receiving a packet; the sender immediately retransmits a packet if the corresponding ack is not received after certain constant time. Using our testbed, we study the performance of SEA when used with B-MAC [12, 2] and S-MAC [19]. B-MAC uses the mechanism of CSMA/CA (carrier sense multiple access with collision avoidance) to control channel access; S-MAC uses CSMA/CA too, but it also employs RTS-CTS handshake to reduce the impact of hidden terminals.

**SEA with B-MAC.** The event reliability, the average packet delivery delay, as well as the event goodput is shown in Table 1, where RT stands for the maximum number of retransmissions for each packet at each

| Metrics          | RT = 0 | RT = 1 | RT = 2 |
|------------------|--------|--------|--------|
| ER (%)           | 51.05  | 54.74  | 54.63  |
| PD (seconds)     | 0.21   | 0.25   | 0.26   |
| EG (packets/sec) | 4.01   | 4.05   | 3.63   |

Table 1: SEA with B-MAC in Lites trace

hop (e.g., RT = 0 means that packets are not retransmitted). The distribution of the number of unique packets received at the base station along time is shown in Figure 3.

Table 1 and Figure 3 show that when packets are retransmitted, the event reliability increases slightly (i.e., by up to 3.69%). Nevertheless, the maximum reliability is still only 54.74%, and, even worse, the event reliability as well as goodput decreases when the maximum number of retransmissions increases from 1 to 2. (The above data is for B-MAC with its default contention window size. We have conducted the experiment with different contention window size of B-MAC, and we found that the performance pattern remains the same.)

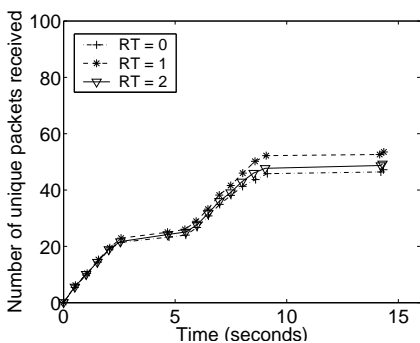


Figure 3: The distribution of packet reception in SEA with B-MAC

**SEA with S-MAC.** Unlike B-MAC, S-MAC uses RTS-CTS handshake for unicast transmissions, which reduces packet collisions. We evaluate SEA when it is used with S-MAC, and the performance data is shown in Table 2 and Figure 4.

| Metrics          | RT = 0 | RT = 1 | RT = 2 |
|------------------|--------|--------|--------|
| ER (%)           | 72.6   | 74.79  | 70.1   |
| PD (seconds)     | 0.17   | 0.183  | 0.182  |
| EG (packets/sec) | 5.01   | 4.68   | 4.37   |

Table 2: SEA with S-MAC in Lites trace

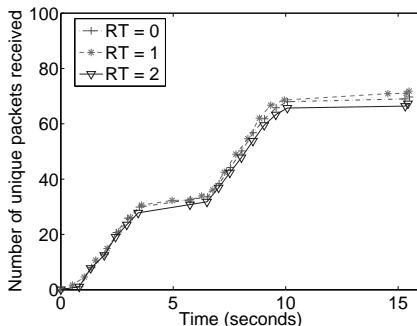


Figure 4: The distribution of packet reception in SEA with S-MAC

Compared with B-MAC, RTS-CTS handshake improves the event reliability by about 20% in S-MAC. Yet packet retransmissions still do not significantly improve the event reliability and can even decrease the reliability.

**Analysis.** We find that the reason why retransmission does not significantly improve — and can even degenerate — communication reliability is that, in SEA, lost packets are retransmitted while new packets are generated and forwarded, thus retransmissions, when not scheduled appropriately, only increase channel contention and cause more packet collision.<sup>3</sup> The situation is further exacerbated by ack-loss

<sup>3</sup>This is not the case in wireline networks and is due to the nature of wireless communications.

(with a probability as high as 10.29%), since ack-loss causes unnecessary retransmission of packets that have been received. To make retransmission effective in improving reliability, therefore, we need a retransmission scheduling mechanism that ameliorates retransmission-incurred channel contention.

### 3.2 Stop-and-wait implicit ack (SWIA)

SWIA takes advantage of the fact that every node, except for the base station, forwards the packet it receives and the forwarded packet can act as the acknowledgment to the sender at the previous hop [11]. In SWIA, the sender of a packet snoops the channel to check whether the packet is forwarded within certain constant threshold time; the sender regards the packet as received if it is forwarded within the threshold time, otherwise the packet is regarded as lost. The advantage of SWIA is that acknowledgment comes for free except for the limited control information piggybacked in data packets.

We evaluate SWIA only with B-MAC, given that the implementation of S-MAC is not readily applicable for packet snooping. The performance results are shown in Table 3 and in Figure 5.

| Metrics          | RT = 0 | RT = 1 | RT = 2 |
|------------------|--------|--------|--------|
| ER (%)           | 43.09  | 31.76  | 46.5   |
| PD (seconds)     | 0.35   | 8.81   | 18.77  |
| EG (packets/sec) | 3.48   | 2.58   | 1.41   |

Table 3: SWIA with B-MAC in Lites trace

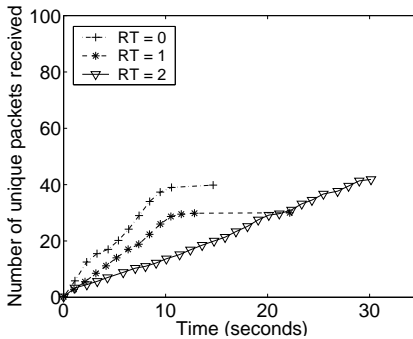


Figure 5: The distribution of packet reception in SWIA with B-MAC

We see that the maximum event reliability in SWIA is only 46.5%, and that the reliability decreases significantly when packets are retransmitted at most once at each hop. When packets are retransmitted up to twice at each hop, the packet delivery delay increases, and the event goodput decreases significantly despite the slightly increased reliability.

**Analysis.** We find that the above phenomena are due to the following reasons. First, the length of data packets is increased by the piggybacked control information in SWIA, thus the ack-loss probability increases (as high as 18.39% in our experiments), which in turn increases unnecessary retransmissions.

Second, most packets are queued upon reception and thus their forwarding is delayed. As a result, the piggybacked acknowledgments are delayed and the corresponding packets are retransmitted unnecessarily. Third, once a packet is waiting to be acknowledged, all the packets arriving later cannot be forwarded even if the communication channel is free. Therefore, channel utilization as well as system throughput decreases, and network queuing as well as packet delivery delay increases. Fourth, as in SEA, lack of retransmission scheduling allows retransmissions, be it necessary or unnecessary, to cause more channel contention and packet loss.

## 4 Protocol RBC

To address the limitations of SEA and SWIA in bursty convergecast, we design protocol RBC. In RBC, we design a window-less block acknowledgment scheme to increase channel utilization and to reduce the probability of ack-loss. We also design a distributed contention control scheme that schedules packet retransmissions and reduces the contention between newly generated and retransmitted packets. Moreover, we design mechanisms to address the challenges of bursty convergecast on timer-based retransmission (such as varying ack-delay and timer-incurred delay).

Given that the number of packets competing for channel access is less in implicit-ack based schemes than in explicit-ack based schemes, we design RBC based on the paradigm of implicit-ack (i.e., piggybacking control information in data packets). We elaborate on RBC as follows. (Even though the mechanisms used in RBC can be applied in the explicit-ack paradigm, we relegate the detailed study as our future work.)

### 4.1 Window-less block acknowledgment

In traditional block acknowledgment [5], a sliding-window is used for both duplicate detection and in-order packet delivery.<sup>4</sup> The sliding-window reduces network throughput once a packet is sent but remains unacknowledged (since the sender can only send up to its window size once a packet is unacknowledged), and in-order delivery increases packet delivery delay once a packet is lost (since the lost packet delays the delivery of every packet behind it). Therefore, the sliding-window based block acknowledgment scheme does not apply to bursty convergecast, given the real-time requirement of the latter.

To address the constraints of traditional block acknowledgment in the presence of unreliable links, we take advantage of the fact that in-order delivery is not required in bursty convergecast. Without considering the order of packet delivery, by which we only need to detect whether a sequence of packets

---

<sup>4</sup>Note that SWIA is a special type of block acknowledgment where the window size is 1.

are received without loss in the middle and whether a received packet is a duplicate of a previously received one. To this end, we design, as follows, a *window-less block acknowledgment* scheme which guarantees continuous packet forwarding irrespective of the underlying link unreliability as well as the resulting packet- and ack-loss. For clarity of presentation, we consider an arbitrary pair of nodes  $S$  and  $R$  where  $S$  is the sender and  $R$  is the receiver.

**Window-less queue management.** The sender  $S$  organizes its packet queue as  $(M + 2)$  linked lists, as shown in Figure 6, where  $M$  is the maximum number of retransmissions at each hop. For convenience,

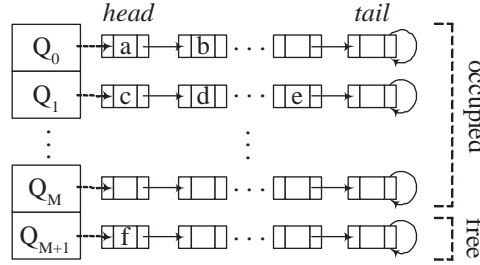


Figure 6: Virtual queues at a node

we call the linked lists *virtual queues*, denoted as  $Q_0, \dots, Q_{M+1}$ . The virtual queues are ranked such that a virtual queue  $Q_k$  ranks higher than  $Q_j$  if  $k < j$ .

Virtual queues  $Q_0, Q_1, \dots$ , and  $Q_M$  buffer packets waiting to be sent or to be acknowledged, and  $Q_{M+1}$  collects the list of free queue buffers. The virtual queues are maintained as follows:

- When a new packet arrives at  $S$  to be sent,  $S$  detaches the head buffer of  $Q_{M+1}$ , if any, stores the packet into the queue buffer, and attaches the queue buffer to the tail of  $Q_0$ .
- Packets stored in a virtual queue  $Q_k$  ( $k > 0$ ) will not be sent unless  $Q_{k-1}$  is empty; packets in the same virtual queue are sent in FIFO order.
- After a packet in a virtual queue  $Q_k$  ( $k \geq 0$ ) is sent, the corresponding queue buffer is moved to the tail of  $Q_{k+1}$ , unless the packet has been retransmitted  $M$  times<sup>5</sup> in which case the queue buffer is moved to the tail of  $Q_{M+1}$ .
- When a packet is acknowledged to have been received, the buffer holding the packet is released and moved to the tail of  $Q_{M+1}$ .

The above rules help identify the relative freshness of packets at a node (which is used in the differentiated contention control in Section 4.2); they also help maintain without using sliding windows the order

<sup>5</sup>Due to block-NACK, to be discussed in Section 4.3.2, a packet having been retransmitted  $M$  times may be in a virtual queue other than  $Q_M$ .

in which unacknowledged packets have been sent, providing the basis for window-less block acknowledgment. Moreover, newly arrived packets can be sent immediately without waiting for the previously sent packets to be acknowledged, which enables continuous packet forwarding in the presence of packet- and ack-loss.

**Block acknowledgment & reduced ack-loss.** Each queue buffer at  $S$  has an ID that is unique at  $S$ . When  $S$  sends a packet to the receiver  $R$ ,  $S$  attaches the ID of the buffer holding the packet as well as the ID of the buffer holding the packet to be sent next. In Figure 6, for example, when  $S$  sends the packet in buffer  $a$ ,  $S$  attaches the values  $a$  and  $b$ . Given the queue maintenance procedure, if the buffer holding the packet being sent is the tail of  $Q_0$  or the head of a virtual queue other than  $Q_0$ ,  $S$  also attaches the ID of the head buffer of  $Q_{M+1}$ , if any, since one or more new packets may arrive before the next enqueued packet is sent in which case the newly arrived packet(s) will be sent first. For example, when the packet in buffer  $c$  of Figure 6 is sent,  $S$  attaches the values  $c$ ,  $d$ , and  $f$ .

When the receiver  $R$  receives a packet  $p_0$  from  $S$ ,  $R$  learns the ID  $n'$  of the buffer holding the next packet to be sent by  $S$ . When  $R$  receives a packet  $p_n$  from  $S$  next time,  $R$  checks whether  $p_n$  is from buffer  $n'$  at  $S$ : if  $p_n$  is from buffer  $n'$ ,  $R$  knows that there is no packet loss between receiving  $p_0$  and  $p_n$  from  $S$ ; otherwise,  $R$  detects that some packets are lost between  $p_0$  and  $p_n$ .

For each maximal sequence of packets  $p_k, \dots, p_{k'}$  from  $S$  that are received at  $R$  without any loss in the middle,  $R$  attaches to packet  $p_{k'}$  the 2-tuple  $\langle q_k, q_{k'} \rangle$ , where  $q_k$  and  $q_{k'}$  are the IDs of the buffers storing  $p_k$  and  $p_{k'}$  at  $S$ . We call  $\langle q_k, q_{k'} \rangle$  the *block acknowledgment* for packets  $p_k, \dots, p_{k'}$ . When  $S$  snoops the forwarded packet  $p_{k'}$  later,  $S$  learns that all the packets sent between  $p_k$  and  $p_{k'}$  have been received by  $R$ . Then  $S$  releases the buffers holding these packets. For example, if  $S$  snoops a block acknowledgment  $\langle c, e \rangle$  when its queue state is as shown in Figure 6,  $S$  knows that all the packets in buffers between  $c$  and  $e$  in  $Q_1$  have been received, and  $S$  releases buffers between  $c$  and  $e$ , including  $c$  and  $e$ .

One delicate detail in processing the block acknowledgment  $\langle q_k, q_{k'} \rangle$  is that after releasing buffer  $q_k$ ,  $S$  will maintain a mapping  $q_k \leftrightarrow q_{k''}$ , where  $q_{k''}$  is the buffer holding the packet sent (or to be sent next) after that in  $q_{k'}$ . When  $S$  snoops another block acknowledgment  $\langle q_k, q_n \rangle$  later,  $S$  knows, by  $q_k \leftrightarrow q_{k''}$ , that packets sent between those in buffers  $q_{k''}$  and  $q_n$  have been received by  $R$ ; then  $S$  releases the buffers holding these packets, and  $S$  resets the mapping to  $q_k \leftrightarrow q_{n''}$ , where  $q_{n''}$  is the buffer holding the packet sent (or to be sent next) after that in  $q_n$ .  $S$  maintains the mapping for  $q_k$  until  $S$  receives a block-NACK  $[n', n]$  (to be discussed in Section 4.3.2) or a block acknowledgment  $\langle q, q' \rangle$  where  $q \neq q_k$ , in which case  $S$  maintains the mapping for  $n$  or  $q$  respectively. Via the buffer pointer mapped as above, the node  $S$  can process the incoming block acknowledgments and block-NACKs. For convenience, we call the buffer being mapped to the *anchor* of block acknowledgments. In the examples discussed above, buffers  $q_{k''}$

and  $q_n$  have been anchors once. We also call the packet in an anchor buffer an *anchor packet*. (The concepts of anchor and anchor packet will be used in Section 6.)

In the above block acknowledgment scheme, the acknowledgment for a received packet is piggybacked onto the packet itself as well as the packets that are received consecutively after the packet without any loss in the middle. Therefore, the acknowledgment is replicated and the probability for it to be lost decreases significantly, by a factor of 2.07 in Lites trace as analyzed in Appendix A1.

**Duplicate detection & obsolete-ack filtering.** Since it is impossible to completely prevent ack-loss in lossy communication channels, packets whose acknowledgments are lost will be retransmitted unnecessarily. Therefore, it is necessary that duplicate packets be detected and dropped.

To enable duplicate detection, the sender  $S$  maintains a counter for each queue buffer, whose value is incremented by one each time a new packet is stored in the buffer. When  $S$  sends a packet, it attaches the current value of the corresponding buffer counter. For each buffer  $q$  at  $S$ , the receiver  $R$  maintains the counter value  $c_q$  piggybacked in the last packet from the buffer. When  $R$  receives another packet from the buffer  $q$  later,  $R$  checks whether the counter value piggybacked in the packet equals to  $c_q$ : if they are equal,  $R$  knows that the packet is a duplicate and drops it; otherwise  $R$  regards the packet as a new one and accepts it. The duplicate detection is local in the sense that it only requires information local to each queue buffer instead of imposing any rule involving different buffers (such as in sliding-window) that can degenerate system performance.

For the correctness of the above duplicate detection mechanism, we only need to choose the domain size  $C$  for the counter value such that the probability of losing  $C$  packets in succession is negligible. For example, for the high per-hop packet loss probability 22.7% in the case of Lites trace,  $C$  could still be as small as 7, since the probability of losing 7 packets in succession is only 0.003%. (Given the small domain size for the counter value as well as the usually small queue size at each node, the duplicate detection mechanism does not consume much memory. For example, it only takes 36 bytes in the case of Lites.)

In addition to duplicate detection, we also use buffer counter to filter out obsolete acknowledgment. Despite the low probability, packet forwarding at  $R$  may be severely delayed, such that the queue buffers signified in a block acknowledgment have been reused by  $S$  to hold packets arriving later. To deal with this,  $R$  attaches to each forwarded packet the ID as well as the counter value of the buffer holding the packet at  $S$  originally; when  $S$  snoops a packet forwarded by  $R$ ,  $S$  checks whether the piggybacked counter value equals to the current value of the corresponding buffer: if they are equal,  $S$  regards as valid the piggybacked block acknowledgment; otherwise,  $S$  regards the block acknowledgment as obsolete and ignores it.

**Aggregated-ack at the base station.** In sensor networks, the base station usually forwards all the packets it receives to an external network. As a result, the children of the base station (i.e., the nodes that forward packets directly to the base station) are unable to snoop the packets the base station forwards, and the base station has to explicitly acknowledge the packets it receives. To reduce channel contention, the base station aggregates several acknowledgments, for packets received consecutively in a short period of time, into a single packet and broadcasts the packet to its children. Accordingly, the children of the base station adapt their control parameters to the way the base station handles acknowledgments.

## 4.2 Differentiated contention control

In wireless sensor networks where per-hop connectivity is reliable, most packet losses are due to collision in the presence of severe channel contention. To enable reliable packet delivery, lost packets need to be retransmitted. Nevertheless, packet retransmission may cause more channel contention and packet loss, thus degenerating communication reliability. Also, there exist unnecessary retransmissions due to ack-loss, which only increase channel contention and reduce communication reliability. Therefore, it is desirable to schedule packet retransmissions such that they do not interfere with transmissions of other packets.

The way the virtual queues are maintained in our window-less block acknowledgment scheme facilitates the retransmission scheduling, since packets are automatically grouped together by different virtual queues. Packets in higher-ranked virtual queues have been transmitted less number of times, and the probability that the receiver has already received the packets in higher-ranked virtual queues is lower (e.g., 0 for packets in  $Q_0$ ). Therefore, we rank packets by the rank of the virtual queues holding the packets, and higher-ranked packets have higher-priority in accessing the communication channel. By this rule, packets that have been transmitted less number of times will be (re)transmitted earlier than those that have been transmitted more, and interference between packets of different ranks is reduced.

Window-less block acknowledgment already handles packet differentiation and scheduling within a node, thus we only need a mechanism that schedules packet transmission across different nodes. To reduce interference between packets of the same rank and to balance network queuing as well as channel contention across nodes, inter-node packet scheduling also takes into account the number of packets of a certain rank so that nodes having more such packets transmit earlier.

To implement the above concepts, we define the rank  $rank(j)$  of a node  $j$  as  $\langle M - k, |Q_k|, ID(j) \rangle$ , where  $Q_k$  is the highest-ranked non-empty virtual queue at  $j$ ,  $|Q_k|$  is the number of packets in  $Q_k$ , and  $ID(j)$  is the ID of  $j$ .  $rank(j)$  is defined such that 1) the first field guarantees that packets having been transmitted fewer number of times will be (re)transmitted earlier, 2) the second field ensures that nodes

having more packets enqueued get chances to transmit earlier, and 3) the third field is to break ties in the first two fields. A node with a larger rank value ranks higher. Then, the distributed transmission scheduling works as follows:

- Each node piggybacks its rank to the data packets it sends out.
- Upon snooping or receiving a packet, a node  $j$  compares its rank with that of the packet sender  $k$ .  $j$  will change its behavior only if  $k$  ranks higher than  $j$ , in which case  $j$  will not send any packet in the following  $w(j, k) \times T_{pkt}$  time.  $T_{pkt}$  is the time taken to transmit a packet at the MAC layer, and  $w(j, k) = 4 - i$ , when  $rank(j)$  and  $rank(k)$  differ at the  $i$ -th element of the 3-tuple ranks.  $w(j, k)$  is defined such that the probability of all waiting nodes starting their transmissions simultaneously is reduced, and that higher-ranked nodes tend to wait for shorter time.  $T_{pkt}$  is estimated by the method of *Exponentially Weighted Moving Average (EWMA)*.
- If a sending node  $j$  detects that it will not send its next packet within  $T_{pkt}$  time (i.e., when  $j$  knows that, after the current packet transmission, it will rank lower than another node),  $j$  signifies this by marking the packet being sent, so that the nodes overhearing the packet will skip  $j$  in the contention control. (This mechanism reduces the probability of idle waiting, where the channel is free but no packet is sent.)

### 4.3 Timer management in window-less block acknowledgment

In window-less block acknowledgment, a sender  $S$  starts a retransmission timer after sending a packet, and  $S$  retransmits the packet if  $S$  has not received the corresponding ack when the timer times out. Retransmission timers directly affect the reliability and the delay in packet delivery: large timeout values of timers tend to increase packet delivery delay, whereas small timeout values tend to cause unnecessary retransmissions and thus decrease packet delivery reliability. To provide reliable and real-time packet delivery, we design mechanisms to manage the timers in window-less block acknowledgment as follows. (Again, we consider a sender  $S$  and a receiver  $R$ .)

#### 4.3.1 Dealing with varying ack-delay

When the receiver  $R$  receives a packet  $m$  from the sender  $S$ ,  $R$  first buffers  $m$  in  $Q_0$ . The delay in  $R$  forwarding  $m$  depends on the number of packets in front of  $m$  in  $Q_0$ . Since the number of packets enqueued in  $Q_0$  keeps changing, the delay in forwarding a received packet by  $R$  keeps changing, which leads to varying delay in packet acknowledgment. Therefore, the retransmission timer at the sender  $S$  should adapt to the queuing condition at  $R$ ; otherwise, either lost packets are unnecessarily delayed in

retransmission (when the retransmission timer is too large) or packets are unnecessarily retransmitted even if they are received (when the retransmission timer is too small).

To adaptively setting the retransmission timer for a packet, the sender  $S$  keeps track of, by snooping packets forwarded by  $R$ , the length  $s_r$  of  $Q_0$  at  $R$ , the average delay  $d_r$  in  $R$  forwarding a packet after the packet becomes the head of  $Q_0$ , and the deviation  $d'_r$  of  $d_r$ . When  $S$  sends a packet to  $R$ ,  $S$  sets the retransmission timer of the packet as

$$(s_r + C_0)(d_r + 4d'_r)$$

where  $C_0$  is a constant denoting the number of new packets that  $R$  may have received since  $S$  learned  $s_r$  the last time ( $C_0$  depends on the application as well as the link reliability, and  $C_0$  is 3 in our experiments). The reason why we use the deviation  $d'_r$  in the above formula is that  $d_r$  varies a lot in wireless networks in the presence of bursty traffic, in which case the deviation improves estimation quality [9].

At a node, each local parameter  $\alpha$  (such as  $d_r$  for node  $R$ ) and its deviation  $\alpha'$  are estimated by the method of EWMA as follows:

$$\begin{aligned}\alpha &\leftarrow (1 - \gamma)\alpha + \gamma\alpha'' \\ \alpha' &\leftarrow (1 - \gamma')\alpha' + \gamma'|\alpha'' - \alpha|\end{aligned}$$

where  $\gamma$  and  $\gamma'$  are weight-factors, and  $\alpha''$  is the latest observation of  $\alpha$ . Empirically, we set  $\gamma = \frac{1}{8}$  and  $\gamma' = \frac{1}{4}$  in RBC.

### 4.3.2 Alleviating timer-incurred delay

The packet retransmission timer calculated as above is conservative in the sense that it is usually greater than the actual ack-delay [9]. This is important for reducing the probability of unnecessary retransmissions, but it introduces extra delay and makes network resources under-utilized [20].

To alleviate timer-incurred delay, we design the following mechanisms to expedite necessary packet retransmissions and to improve channel utilization:

- Whenever the receiver  $R$  receives a packet  $m$  from buffer  $n$  of the sender  $S$  while  $R$  is expecting (in the absence of packet loss) to receive a packet from buffer  $n'$  of  $S$ ,  $R$  learns that packets sent between those in buffers  $n'$  and  $n$  at  $S$ , including the one in  $n'$ , are lost. In this case,  $R$  piggybacks a block-NACK  $[n', n]$  onto the next packet it forwards, by which the block-NACK can be snooped by  $S$  immediately.

When  $S$  learns the block-NACK  $[n', n]$  from  $R$ ,  $S$  resets the retransmission timers to 0 for the packets sent between those in  $n'$  and  $n$  (including the one in  $n'$ ), and for each of these packets,

$S$  moves the corresponding buffer to the tail of  $Q_{k-1}$  if the buffer is currently at  $Q_k$ . Therefore, packets that need to be retransmitted are put into higher-ranked virtual queues and are retransmitted quickly.<sup>6</sup>

- Whenever  $S$  learns that the virtual queue  $Q_0$  of  $R$  becomes empty,  $S$  knows that  $R$  has forwarded all the packets it has received. In this case,  $S$  resets the retransmission timers to 0 for those packets still waiting to be acknowledged, since they will not be (due to either packet-loss or ack-loss).

Similarly, when  $S$  snoops the acknowledgment for a packet  $m$ ,  $S$  resets the retransmission timer to 0 for those packets that are sent before  $m$  but are still waiting to be acknowledged.

- When a network channel is fully utilized, it should be busy all the time. Therefore, if the sender  $S$  has packets to send, and if  $S$  notices that no packet is sent by any neighboring node in a period of  $C_1 \times T_{pkt}$  time,  $S$  sends out the packet at the head of its highest-ranked non-empty virtual queue, without considering the retransmission timer even if the packet is to be acknowledged.  $C_1$  is a constant reflecting the desired degree of channel utilization and  $T_{pkt}$  is the time taken to transmit a packet at the MAC layer. (This mechanism improves channel utilization without introducing unnecessary retransmissions because of the “differentiated contention control” in RBC.)

## 5 Experimental results

We have implemented protocol RBC in TinyOS using B-MAC. In the implementation, the control logic takes 185 bytes of RAM when each node maintains a buffer capable of holding 16 packets, and the control information piggybacked in data packets takes 14 bytes.<sup>7</sup> RBC has successfully provided reliable and real-time data transport in the sensor network field experiment ExScal [1] where around 1,200 Mica2/XSM motes were deployed to detect, track, and classify intruders.

We have also evaluated RBC in our testbed. In what follows, we discuss the performance of RBC, compare RBC with SEA and SWIA, discuss the impact of the individual components (i.e., window-less block acknowledgment and differentiated contention control) of RBC, and analyze the timing-shift of packet delivery in RBC.

**Performance of RBC.** Table 4 shows the performance results of RBC, and Figure 7 shows the distri-

---

<sup>6</sup>The movement of NACKed packets affect the buffering order required by block-acknowledgment. Therefore, NACKed packets are specially marked so that they are not mistakenly regarded as having been received; the mark for a NACKed packet is reset after the packet is transmitted once more.

<sup>7</sup>Comparatively, the control logic of SEA uses 150 bytes of RAM when each node maintains a buffer capable of holding 16 packets, and each explicit ack packet takes 16 bytes for the MICA2 radio, including the preamble, the synchronization-code, and the ack-code; the control logic of SWIA uses 68 bytes of RAM when the packet buffer size is 16, and the control information piggybacked in data packets takes 8 bytes.

| Metrics          | RT = 0 | RT = 1 | RT = 2 |
|------------------|--------|--------|--------|
| ER (%)           | 56.21  | 83.16  | 95.26  |
| PD (seconds)     | 0.21   | 1.18   | 1.72   |
| EG (packets/sec) | 4.28   | 5.72   | 6.37   |

Table 4: RBC in Lites trace

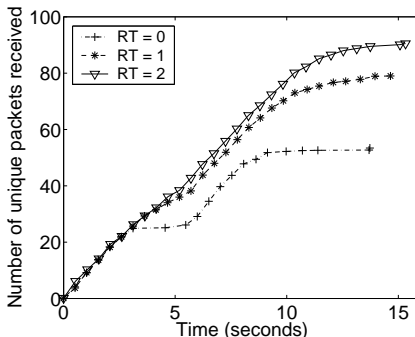


Figure 7: The distribution of packet reception in RBC

bution of packet reception in RBC. From Table 4 and Figure 7, we observe the following properties of RBC:

- The event reliability keeps increasing, in a significant manner, as the number of retransmissions increases. The increased reliability mainly attributes to reduced unnecessary retransmissions (by reduced ack loss and adaptive retransmission timer) and retransmission scheduling.
- Compared with SWIA which is also based on implicit-ack, RBC reduces packet delivery delay significantly. This mainly attributes to the ability of continuous packet forwarding in the presence of packet- and ack-loss and the reduction in timer-incurred delay.
- The rate of packet reception at the base station and the event goodput keep increasing as the number of retransmissions increases. When packets are retransmitted up to twice at each hop, the event goodput reaches 6.37 packets/second, quite close to the optimal goodput — 6.66 packets/second — for Lites trace.

Compared with SWIA, RBC improves reliability by a factor of 2.05 and reduces average packet delivery delay by a factor of 10.91. Compared to SEA with B-MAC (simply referred to as SEA hereafter), RBC improves reliability by a factor of 1.74, but the average packet delivery delay increases by a factor of 6.61 in RBC. Interestingly, however, RBC still improves the event goodput by a factor of 1.75 when compared with SEA. The reason is that, in RBC, lost packets are retransmitted and delivered after those packets that are generated later but transmitted less number of times. Therefore, the delivery delay for lost packets increases, which increases the average packet delivery delay, without degenerating the

system goodput. The observation shows that, due to the unique application models in sensor networks, metrics evaluating aggregate system behaviors (such as the event goodput) tend to be of more relevance than metrics evaluating unit behaviors (such as the delay in delivering each individual packet).

**RBC compared with SEA and SWIA.** To further understand protocol behaviors in the presence of packet retransmissions, we conduct, as follows, a comparative study of RBC, SEA, and SWIA for the case where packets are retransmitted up to twice at each hop.

Figure 8 compares the distribution of packet generation in Lites trace with the distributions of packet

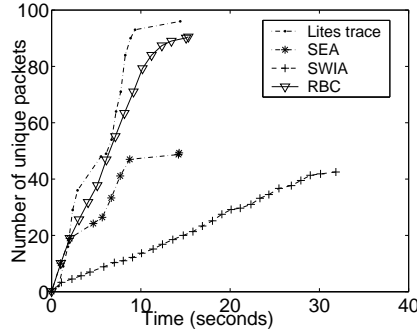


Figure 8: The distributions of packet generation and reception

reception in SEA, SWIA, and RBC. We see that the curve for packet reception in RBC smooths out and almost matches that of packet generation. In contrast, many packets are lost in SEA despite the fact that the rate of packet reception in SEA is close to that in RBC; packet delivery is significantly delayed in SWIA, in addition to the high degree of packet loss.

Based on the grid topology as shown in Figure 1(b), Figures 9(a)-(c) show the node reliability in SEA,

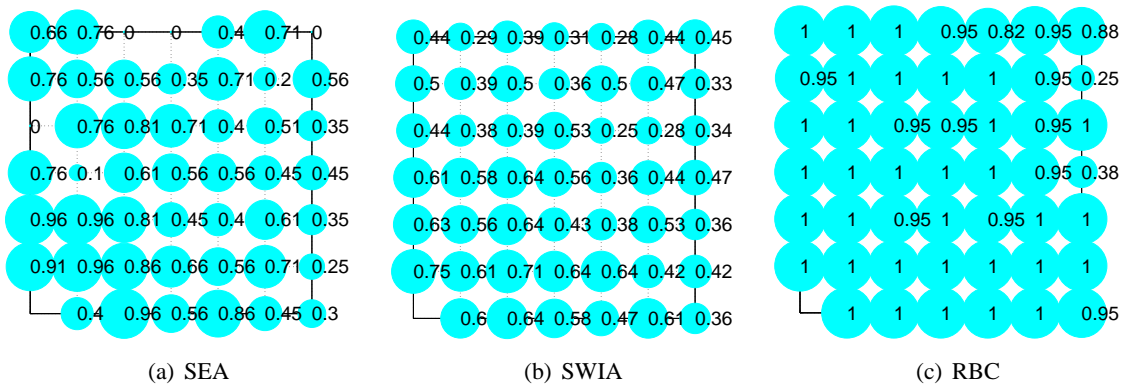


Figure 9: Node reliability

SWIA, and RBC respectively. Figure 10 shows the cumulative distribution of node reliability in SEA, SWIA, and RBC. We see that node reliability improves significantly in RBC: only 4.17% of nodes have a node reliability less than 80% in RBC; yet in SEA and SWIA, above 80% of nodes have a node reliability

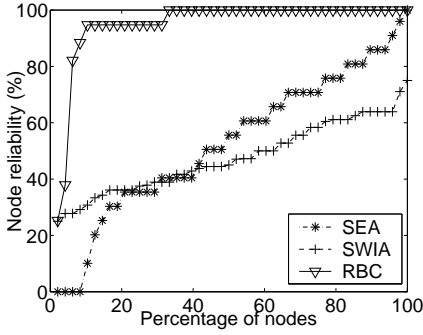


Figure 10: Distribution of node reliability

less than 80%.

Figure 11 shows the average node reliability in SEA, SWIA, and RBC as the number of routing hops

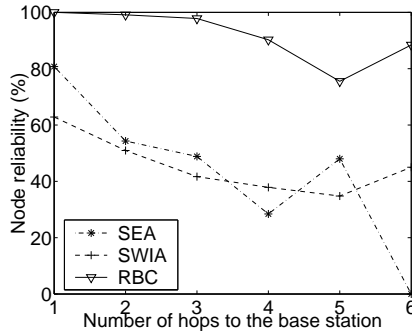


Figure 11: Node reliability as a function of routing hops

(to the base station) increases. We see that the node reliability in RBC is much higher than that in SEA and SWIA at every routing hop, and that the reliability at the farthest hop in RBC is even greater than that at the closest hop in SEA and SWIA. (Note that, in RBC, the reason why nodes 5 hops away from the base station have lower average delivery rate than nodes 6 hops away can be due to the specific traffic pattern and the difference among nodes' hardware.)

**Breakdown of RBC.** To understand the individual impact of window-less block acknowledgment and differentiated contention control in RBC, we evaluate the performance of RBC without using differentiated contention control (i.e., RBC with window-less block acknowledgment only). Table 5 shows the

| Metrics          | RT = 0 | RT = 1 | RT = 2 |
|------------------|--------|--------|--------|
| ER (%)           | 54.90  | 77.19  | 82.29  |
| PD (seconds)     | 0.22   | 1.12   | 1.52   |
| EG (packets/sec) | 4.04   | 4.13   | 4.12   |

Table 5: RBC without differentiated contention control

performance results. Comparing Table 5 with Table 4, we observe the following:

- Differentiated contention control improves packet delivery performance even when there is no retransmission (i.e.,  $RT = 0$ ). This is because the contention control reduces channel contention by prioritizing channel access according to the degree of queue accumulation at different nodes.
- Without differentiated contention control, packet delivery reliability also improves significantly when  $RT$  (maximum number of per-hop retransmissions) increases from 0 to 1, but the improvement becomes far less when  $RT$  increases from 1 to 2. This is because differentiated contention control plays an increasingly important role when  $RT$  (thus channel contention) increases.

Comparing Table 5 with Tables 1 and 3, we see that, with window-less block acknowledgment alone, RBC significantly improves the packet delivery performance of SEA and SWIA. The reasons are as follows:

- Compared with SEA, the channel contention is less in window-less block acknowledgment because no explicit acknowledgment packet is generated (thus reducing the number of packets in the network). Moreover, the intra-node packet prioritization (via the queue management) in window-less block acknowledgment also improves the packet delivery reliability.
- Compared with SWIA, window-less block acknowledgment improves packet delivery reliability by reducing ack-loss probability (and thus reducing unnecessary packet retransmissions) and employing intra-node packet prioritization. Window-less block acknowledgment also significantly reduces packet delivery delay by careful timer management and by enabling continuous packet transmission in the presence of packet- and ack-loss.

**Timing-shift of packet delivery.** In this paper, we focus on scenarios where packets are timestamped and thus we do not need to precisely preserve the relative timing between packets as it is when they are generated. Nevertheless, to characterize how RBC affects the relative timing of packets, we measure the *timing-shift* of packet delivery as follows:

Given a packet  $P1$  received at the base station, the timing-shift for  $P1$  is calculated as

$$|(R_1 - R_0) - (S_1 - S_0)|$$

where  $R_1$  denotes the time when  $P1$  is received at the base station,  $R_0$  denotes the time when a packet  $P0$  is received at the base station immediately before  $P1$ ,  $S_1$  denotes the time when  $P1$  is generated at some node in the network, and  $S_0$  denotes the time when  $P0$  is generated at some node

in the network.<sup>8</sup> For convenience, we set the timing-shift to 0 for the first packet received at the base station.

Based on the above definition, we find that the average timing-shift is 1.0035 second for the packets received in the case of RBC, and that the average timing-shift is around 0.1698 second in both SEA and SWIA. Even though the timing-shift in RBC is predictably greater than those in SEA and SWIA, it is still small enough for real-time event-driven applications such as Lites and ExScal [1] where high-level decisions are based on data in the order of seconds. (Note that, if need be, the timing-shift in packet delivery can be reduced by tuning the queue management policy in window-less block acknowledgment. But the detailed study is beyond the scope of this paper.)

## 6 Discussion

In this section, we discuss how to extend the basic design of RBC to support continuous event convergecast, and how to avoid queue overflow via flow control.

### 6.1 Continuous event convergecast

In event-driven applications, events are mostly rare and well-separated in time. But it may happen (though with relatively low probability) that several events happen continuously during a period of time. In what follows, we first analyze the potential issues of the basic RBC in supporting continuous event convergecast, then we extend RBC to solve the challenges.

**Starvation of orphan packets.** In RBC, window-less block acknowledgment and differentiated contention control ensures that packets having been transmitted fewer number of times will be (re)transmitted earlier. This works without any problem in the case of single event convergecast. In continuous event convergecast, however, fresh packets can be generated continuously as long as events keep occurring. Therefore, if RBC was applied without any adaptation, packets that need to be retransmitted might be delayed in transmission or never get the chance for channel access, since there might exist fresh packets at the node itself or its neighbors most of the time; that is, the packets requiring retransmission may starve in channel access.

On the other hand, not all queued packets will starve. To understand this, we divide the queued packets into two categories (according to the order of their transmission relative to the *anchor packet*<sup>9</sup>) as follows:

---

<sup>8</sup>In this definition, we do not consider the packets that are lost, only calculating the relative timing-shift between packets that are received at the base station.

<sup>9</sup>Defined in Section 4.1.

- *Orphan packets*: the packets sent earlier than the anchor packet.

We regard these packets as orphans in the sense that their acknowledgments or NACKs have been lost. In the window-less block acknowledgment, orphan packets will not be acknowledged by the receiver any more, and they may starve if fresh packets keep arriving.

- *Awaiting packets*: the packets sent later than the anchor packet.

We regard these packets as awaiting in the sense that they are yet to be ACKed or NACKed, but the corresponding acknowledgment packet has not been received. Given that nodes having fresh packets tend to have higher priority in channel access, and that the probability of losing several ACK or NACK packets in succession is low (e.g., 5.15% for losing 2 packets in Lites), the number of awaiting packets tend to be small. Moreover, awaiting packets will either be quickly acknowledged or quickly become orphan. If an awaiting packet is acknowledged, it will be released or moved to  $Q_0$  for retransmission. Therefore, awaiting packets will not starve even if fresh packets keep arriving.

Given that only orphan packets may starve in the case of continuous event convergecast, we only need to adapt the packet scheduling of RBC to avoid starving orphan packets.

**Probabilistic scheduling of orphan packets.** One simple way to avoid starving orphan packets is to always put them in  $Q_0$  once they become orphan. Nevertheless, one shortcoming of this approach is that the detection of new events may be delayed, because the delivery of new packets is delayed. This becomes undesirable especially in the case of incremental event detection where the first few packets related to an event should be delivered as soon as possible. Moreover, an orphan packet may have already been received (if the packet becomes orphan because of ack loss).

Therefore, we first analyze the importance of orphan packets in terms of the new information they may carry (since there is no need to transmit an orphan packet if it has already been received). To this end, we calculate the *probability  $P_{loss}$  that an orphan packet has not been received* as follows: (Interested readers can check the detailed derivation in Appendix A2.)

$$P_{loss} = \begin{cases} (1 - (P'_{rbc})^k) \frac{p}{p + P'_{rbc}} & \text{if } k > 0 \\ \frac{p}{p + P'_{rbc}} & \text{if } k = 0 \end{cases}$$

where  $k$  is the number of times that an orphan packet has been retransmitted due to timeout of retransmission timers,  $p$  is the packet loss rate, and  $P'_{rbc} = p - \frac{p(1-3p+4p^2-2p^3)}{1-p+p^2}$ . In the case of Lites [3],  $p = 22.7\%$ . Then,  $P_{loss}$  is 71.86% and 65.47% for  $k = 0$  and  $k = 1$  respectively. Therefore, the probability that an orphan packet has been lost in previous transmissions is high (e.g., up to 71.86% in Lites).

To take into account the probability that an orphan packet carrying new information, we adapt the intra-node and inter-node packet scheduling of RBC as follows:

- *Intra-node scheduling.* In short, the scheduling is adapted so that an orphan packet is regarded as a fresh packet with the probability that the packet has been lost. More specifically, the adaptation is as follows:
  - If the head packet  $pkt_0$  of the highest ranked non-empty virtual queue  $Q_k$  ( $0 \leq k \leq M$ ) is ready for (re)transmission,<sup>10</sup> the node  $S$  selects the head of the orphan packets<sup>11</sup> to transmit with the probability  $P_{loss}$  that the packet has not been received. Accordingly,  $S$  selects packet  $pkt_0$  to transmit with probability  $(1 - P'_{loss})$ .
  - If the head packet  $pkt_0$  of the highest ranked non-empty virtual queue  $Q_k$  ( $0 \leq k \leq M$ ) is not ready for (re)transmission,  $S$  selects the head orphan packet to transmit with probability 1.
- *Inter-node scheduling.* This relates to the distributed contention control in RBC. Similar to intra-node scheduling, the definition of the node rank is adapted to regard an orphan packet as one in  $Q_0$  with certain probability. More specifically, if a node  $S$  has  $M$  orphan packets, then, in calculating its node rank,  $S$  regards these orphan packets as  $(\sum_{i=1}^M P_i)$  number of packets in  $Q_0$ , where  $P_i$  denotes the probability that the  $i$ -th orphan packet has not been received.

Via the above adaptation to intra-node scheduling, we have

**Theorem 1 (Freedom of packet accumulation)** *The orphan packets at a node do not accumulate indefinitely, as long as the packet loss rate along a link is less than 49.14%.*

(Interested readers can check Appendix A3 for the proof.)

For routing in wireless sensor networks, reliable links are usually chosen (especially for heavy-load bursty convergecast). In Lites, for instance, the average per-hop packet delivery rate is still as high as 77.3 % despite the high channel contention. Therefore, it is reasonable to assume that routing links do have packet delivery rate greater than 50.86% (i.e.,  $(100 - 49.14)\%$ ) in practice. Thus, the adapted intra-node scheduling not only prevents orphan packets from starving but also guarantees that orphan packets do not accumulate indefinitely. (In the worst case when some routing links happen to have reliability lower than 50.86%, the mechanisms to be discussed in the next subsection can deal with the problem of queue accumulation.)

We have implemented the above extensions to RBC, and experimentally evaluated the performance of the extended RBC in our testbed by injecting Lites traffic trace in succession. We observe the following:

---

<sup>10</sup>A packet is ready for transmission if it is in  $Q_0$  or if the retransmission timer associated with the packet is 0.

<sup>11</sup>The orphan packets are organized as an ordered list via the window-less queue management. Moreover, by the way RBC operates, the retransmission timers for orphan packets are 0.

- The performance (e.g., in event reliability and packet delivery delay) of the extended RBC in continuous event convergecast is very similar to that of the basic RBC in single event convergecast. Moreover, there are very few orphan packets (no more than 1 per node on average) and they do not accumulate.
- In the case of single event convergecast, the extension does not degenerate the performance of RBC. For instance, the difference between the extended and the basic version of RBC in average event reliability is within 2%.

## 6.2 Flow control

In the presence of high traffic load, the packet queue at a node may accumulate and overflow if the corresponding senders transmit too many packets in a short time. This issue can be avoided by proper flow control, and has been well studied in [18, 8]. We have implemented a simple hop-by-hop flow control mechanism (to work with RBC) as follows:

- When forwarding packets, a node piggybacks the number of free queue buffers at its place.
- Whenever a sender  $S$  detects that the number  $L_r$  of free queue buffers at the receiver  $R$  is below a threshold  $L$ ,  $S$  will stop sending any packet in the following  $(L - L_r) \times d_{e,R}$  time.  $L$  is a constant chosen such that the probability of losing  $L$  packets in succession is negligible (by which the sender will not fail to detect the congestion state at the receiver), and  $d_{e,R}$  is the average interval between  $R$  releasing one buffer and the next one while there are packets enqueued at  $R$ . ( $R$  estimates  $d_{e,R}$  by the method of EWMA.)
- After learning the number  $L_r$  of free buffers at the receiver  $R$  each time, the sender  $S$  will send at most  $L_r$  packets to  $R$  in the following  $L_r \times d_{e,R}$  time unless  $S$  snoops another packet forwarded by  $R$ .
- To help relieve queue congestion, the nodes having less than  $L$  queue buffers are not subject to the differentiated contention control.

Via testbed-based experiments and outdoor deployment in ExScal [1], the above mechanism has been proved to be highly effective in avoiding queue overflow.

## 7 Related work

The performance of packet delivery in dense sensor networks has been studied in [21]. The results show that, in the presence of heavy channel load, a commonly used loss recovery scheme at link layer (i.e.,

lost packets are retransmitted up to 3 times) does not mask packet loss, and more than 50% of the links observe 50% packet loss. The observation shows the challenge of reliable communication over multi-hop routes, since the reliability decreases exponentially as the number of hops increases.

The limitations of timers in TCP retransmission control have been studied in [20]. The author analyzes the intrinsic difficulties in using timers to achieve optimal performance and argues that additional mechanisms should be used. Despite its focus on TCP, the study also applies to retransmission control in sensor networks.

Block acknowledgment [5] has been proposed for error as well as flow control in the Internet. It considers the problem of in-order packet delivery. Therefore, a lost packet blocks the delivery of all the packets that are behind the lost one but have reached the receiver, as a result of which packet delivery delay is increased. Moreover, a sender can send packets at most up to its window size once a packet is sent but unacknowledged, thus the channel resource may be under-utilized.<sup>12</sup> Block acknowledgment also uses timers without addressing their limitations, which can render additional delay in packet delivery.

For packet-loss detection and retransmission control, DFRF [11] uses stop-and-wait implicit ack (SWIA). Yet DFRF does not address the issue of retransmission-incurred channel contention. Moreover, the retransmission timers in DFRF do not adapt to the varying ack-delay, which can introduce more retransmission or delay than necessary. To reduce the number of packet transmissions, DFRF uses raw data aggregation where multiple short packets are concatenated to form a longer packet. In the type of bursty convergecast as experienced by Lites and ExScal [1], it is more difficult to perform data aggregation at the mote level because motes detecting an event can be multiple hops away from one another and the length of a single sensor data entry is more than half of the packet length. Therefore, the current implementation of RBC is based on the paradigm of implicit-ack to reduce the number of packets competing for channel access. On the other hand, we believe that the methodologies developed in RBC (e.g., window-less block acknowledgment and differentiated contention control) can also be applied when there is data aggregation, in which case we can use explicit-ack packets to send out control information. The detailed study on this is a part of our future work.

RMST [15] and PSFQ [17] have shown the importance of hop-by-hop packet recovery in sensor networks. Yet RMST and PSFQ do not focus on bursty convergecast. Therefore, they do not cope with retransmission-incurred channel contention, they do not design mechanisms to alleviate delay incurred by retransmission timers (whose timeout values are conservatively chosen to reduce unnecessary retransmissions), and they do not design mechanisms to reduce the probability of ack-loss. Our work

---

<sup>12</sup>The situation is worsened by the fact that the window size is less than half of the buffer size in block acknowledgment and the fact that the buffer size is usually small (e.g., 16) in wireless sensor networks due to limited RAM.

complements theirs by identifying issues with existing hop-by-hop mechanisms in bursty convergecast and proposing approaches to address the issues.

CODA [18] and ESRT [13] have studied congestion control in sensor networks. They consider a traffic model where multiple sources are continuously or periodically generating packets. Therefore, they do not focus on real-time packet delivery in bursty convergecast, and they do not consider retransmission-incurred delay as well as channel contention. Recent work in [8] and [7] has studied different techniques for mitigating congestion and guaranteeing fairness in wireless sensor networks. Our work complements theirs by focusing on retransmission-based error control and retransmission scheduling. Several transport protocols, such as ATP [16] and WTCP [14], are also proposed for wireless ad hoc networks. Again, they do not face the challenges of reliable bursty convergecast.

## 8 Concluding remarks

Unlike most existing literature on reliable transport in sensor networks that focuses on periodic traffic, we have focused on bursty convergecast where the key challenges are reliable and real-time error control and the resulting contention control. To address the unique challenges, we have proposed the window-less block acknowledgment scheme which improves channel utilization and reduces ack-loss as well as packet delivery delay; we have also designed mechanisms to schedule packet retransmissions and to reduce timer-incurred delay, which are critical for reliable and real-time transport of bursty traffic. With its well-tested support for reliable and real-time transport of bursty traffic, RBC has been applied in the sensor network field experiment ExScal [1] where about 1,200 Mica2/XSM motes were deployed.

From protocol RBC, we see that bursty convergecast not only poses challenges for reliable and real-time transport control, it also provides unique opportunities for protocol design. Tolerance of out-of-order packet delivery enables the window-less block acknowledgment, which not only guarantees continuous packet delivery in the presence of packet- and ack-loss but also facilitates retransmission scheduling. Overall, the unique network as well as application models in sensor networks offer opportunities for new methodologies in protocol engineering and are interesting areas for further exploration.

In designing RBC, we have focused on reliable bursty convergecast in event-driven applications. Nevertheless, we believe some techniques of RBC (e.g., differentiated contention control) can be applied to the case where data traffic is periodic or continuous. Detailed study of this is beyond the scope of this paper, and we regard it as a part of the future work.

## Acknowledgment

We thank Hui Cao for his help with the initial data analysis.

## References

- [1] ExScal. <http://www.cse.ohio-state.edu/exscal>.
- [2] Wireless embedded systems. <http://webs.cs.berkeley.edu>.
- [3] A Lites event traffic trace. <http://www.cse.ohio-state.edu/~zhangho/publications/Lites-trace.txt>, 2003.
- [4] A. Arora and et al. A Line in the Sand: A wireless sensor network for target detection, classification, and tracking. *Computer Networks (Elsevier)*, 46(5), 2004.
- [5] G. Brown, M. Gouda, and R. Miller. Block acknowledgment: Redesigning the window protocol. In *ACM SIGCOMM*, pages 128–134, 1989.
- [6] Y.-R. Choi, M. Gouda, H. Zhang, and A. Arora. Stabilization of grid routing in sensor networks. *AIAA Journal of Aerospace Computing, Information, and Communication*, to appear.
- [7] C. T. Ee and R. Bajcsy. Congestion control and fairness for many-to-one routing in sensor networks. In *ACM SenSys*, pages 134–147, 2004.
- [8] B. Hull, K. Jamieson, and H. Balakrishnan. Mitigating congestion in wireless sensor networks. In *ACM SenSys*, pages 134–147, 2004.
- [9] V. Jacobson. Congestion avoidance and control. In *ACM SIGCOMM*, pages 314–329, 1988.
- [10] J. Li, C. Blake, D. D. Couto, H. Lee, and R. Morris. Capacity of ad hoc wireless networks. In *ACM MobiCom*, pages 61–69, 2001.
- [11] M. Maroti. The directed flood routing framework. In *Technical report, Vanderbilt University, ISIS-04-502*, 2004.
- [12] J. Polatre, J. Hill, and D. Culler. Versatile low power media access for wireless sensor networks. In *ACM SenSys*, 2004.
- [13] Y. Sankarasubramaniam, O. Akan, and I. Akyildiz. ESRT: Event-to-sink reliable transport in wireless sensor networks. In *ACM MobiHoc*, pages 177–188, 2003.
- [14] P. Sinha, N. Venkitaraman, R. Sivakumar, and V. Bharghavan. WTCP: a reliable transport protocol for wireless wide-area networks. In *ACM MobiCom*, pages 231–241, 1999.
- [15] F. Stann and J. Heidemann. RMST: Reliable data transport in sensor networks. In *IEEE SNPA*, pages 102–112, 2003.
- [16] K. Sundaresan, V. Anantharaman, H. Hsieh, and R. Sivakumar. ATP: A reliable transport protocol for ad-hoc networks. In *ACM MobiHoc*, pages 64–75, 2003.
- [17] C. Wan, A. Campbell, and L. Krishnamurthy. PSFQ: A reliable transport protocol for wireless sensor networks. In *ACM WSNA*, pages 1–11, 2002.
- [18] C. Wan, S. Eisenman, and A. Campbell. CODA: Congestion detection and avoidance in sensor networks. In *ACM SenSys*, pages 266–279, 2003.
- [19] W. Ye, J. Heidemann, and D. Estrin. An energy-efficient MAC protocol for wireless sensor networks. In *IEEE InfoCom*, pages 1567–1576, 2002.
- [20] L. Zhang. Why TCP timers don't work well. In *ACM SIGCOMM*, pages 397–405, 1986.
- [21] J. Zhao and R. Govindan. Understanding packet delivery performance in dense wireless sensor networks. In *ACM SenSys*, pages 1–13, 2003.

# Appendix

## A1: ack-loss probability in RBC

For convenience, we define the following notations:

- $p$  : the probability of losing a single (data) packet;
- $N$  : the number of packets received in succession without any loss in the middle;
- $N'$  : the number of packets lost in succession;
- $B$  : the number of packets received in succession without any loss in the middle, after a packet is already received;
- $A$  : the number of times that the acknowledgment for a packet is received at the sender.

Assuming that packet losses are independent of one another, we have the probability mass functions for random variables  $N$  and  $N'$  as follows.

$$\begin{aligned}P[N = k] &= p(1 - p)^k \\P[N' = k] &= (1 - p)p^k\end{aligned}$$

In RBC, when a packet  $m$  is received at a receiver  $R$ , the acknowledgment for  $m$  can reach back to the sender  $S$  in two ways:  $S$  snoops  $m$  when it is forwarded by  $R$  later, with probability  $P_{self}$ ; or  $S$  does not snoop  $m$  but snoops a packet whose block acknowledgment acknowledges the reception of  $m$ , with probability  $P_{ba}$ . Therefore, the probability  $P_{rbc}$  of  $S$  receiving the acknowledgment for  $m$  can be derived as follows:

$$\begin{aligned}P_{self} &= 1 - p \\P_{ba} &= p \sum_{k=0}^{\infty} P[B = k]P[A \geq 1|B = k] \\&= p \sum_{k=0}^{\infty} P[B = k](1 - P[A = 0|B = k]) \\&= p \sum_{k=0}^{\infty} P[N = k + 1](1 - P[N' = k]) \\&= \frac{p(1-3p+4p^2-2p^3)}{1-p+p^2} \\P_{rbc} &= P_{self} + P_{ba} \\&= 1 - p + \frac{p(1-3p+4p^2-2p^3)}{1-p+p^2}\end{aligned}$$

Then, the probability  $P'_{rbc}$  of losing the acknowledgment for a packet in RBC is  $1 - P_{rbc}$ .

In the case of Lites trace and implicit-ack,  $p = 22.7\%$ . Therefore  $P'_{rbc} = 8.89\%$ , reducing the ack-loss probability of SWIA by a factor of 2.07.

## A2: probability that an orphan packet has not been received

According to the analysis in Appendix A1, if the probability of losing a single data packet is  $p$ , then the probability  $P'_{rbc}$  of losing an ACK for a packet is calculated as follows:

$$P'_{rbc} = p - \frac{p(1 - 3p + 4p^2 - 2p^3)}{1 - p + p^2} \quad (1)$$

Thus, if a packet has been retransmitted  $k$  ( $k > 0$ ) times not due to NACK (but due to retransmission timer timeout), then the probability that the packet has been received by the receiver is  $(P'_{rbc})^k$ .

Given that the NACK for a lost packet is piggybacked onto a data packet only once, the probability  $P_{nack}$  of losing a NACK for a lost packet is simply calculated as:

$$P_{nack} = p \quad (2)$$

Therefore, if an orphan packet has been retransmitted  $k$  ( $k \geq 0$ ) times not due to NACK, the probability  $P_{loss}$  that the packet has not been received can be calculated as follows:

$$P_{loss} = \begin{cases} (1 - (P'_{rbc})^k) \frac{P_{nack}}{P_{nack} + P'_{rbc}} & \text{if } k > 0 \\ \frac{P_{nack}}{P_{nack} + P'_{rbc}} & \text{if } k = 0 \end{cases}$$

Therefore,

$$P_{loss} = \begin{cases} (1 - (P'_{rbc})^k) \frac{p}{p + P'_{rbc}} & \text{if } k > 0 \\ \frac{p}{p + P'_{rbc}} & \text{if } k = 0 \end{cases} \quad (3)$$

## A3: proof of Theorem 1

**Theorem 1 (Freedom of packet accumulation)** The orphan packets at a node do not accumulate indefinitely, as long as the packet loss rate along a link is less than 49.14%.

*Proof:* We first compute the probability  $P_{orph}$  that an enqueued packet becomes an orphan packet. According to the definition of orphan packet, a packet becomes orphan when either of the following hold: it has been received, but the corresponding ACK is lost; it has been lost, but the corresponding NACK is lost. Therefore, we calculate  $P_{orph}$  as follows:

$$P_{orph} = (1 - p) \times P'_{rbc} + p \times P_{nack} \quad (4)$$

where  $p$  is packet loss rate along a link.

To assure that orphan packets do not accumulate indefinitely even if fresh packets keep arriving, we

only need to ensure that  $P_{orph}$  is less than  $P_{loss}$ , the probability that an orphan packet will be transmitted.

From Formula 3, we see that  $P_{loss}$  is minimum when  $k = 1$ , that is,  $\min(P_{loss}) = \frac{(1-P'_{rbc})P_{nack}}{P_{nack}+P'_{rbc}}$ . To guarantee that  $P_{orph}$  is less than  $P_{loss}$ , we only need  $P_{orph} < \min(P_{loss})$ , that is,

$$(1 - p) \times P'_{rbc} + p \times P_{nack} < \frac{(1 - P'_{rbc})P_{nack}}{P_{nack} + P'_{rbc}}$$

Solving this inequality, we get  $p < 49.14\%$ .

More intuitively, Figure 12 shows the relationship between  $P_{orph}$  and  $\min(P_{loss})$ , as  $p$  changes.

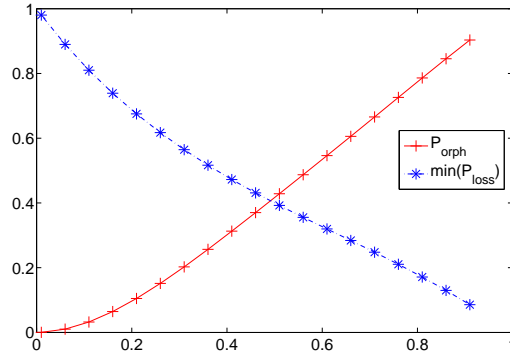


Figure 12:  $P_{orph}$  vs.  $\min(P_{loss})$

□