

Packet Reordering is Not Pathological Network Behavior

Jon C. R. Bennett, *Associate Member, IEEE*, Craig Partridge, *Fellow, IEEE*, and Nicholas Shectman

Abstract—It is a widely held belief that packet reordering in the Internet is a pathological behavior, or more precisely, that it is an uncommon behavior caused by incorrect or malfunctioning network components. Some studies of Internet traffic have reported seeing occasional packet reordering events and ascribed these events to “route fluttering”, router “pauses” or simply to broken equipment. We have found, however, that parallelism in Internet components and links is causing packet reordering under normal operation and that the incidence of packet reordering appears to be substantially higher than previously reported. More importantly, we observe that in the presence of massive packet reordering Transmission Control Protocol (TCP) performance can be profoundly effected. Perhaps the most disturbing observation about TCP’s behavior is that large scale and largely random reordering on the part of the network can lead to self-reinforcingly poor performance from TCP.

Index Terms—Communication system traffic, Internet, packet switching

I. INTRODUCTION

THAT THERE *is* packet reordering in the Internet will not come as a great surprise to most people. What may come as a surprise is the frequency and magnitude of the packet reordering that occurs over many routes in the Internet. A further surprise is that the cause of much of this reordering is not due to pathologies such as broken network equipment but is a naturally occurring result of local parallelism.

By local parallelism, we mean that a packet can follow multiple paths within a device or logical link. Examples of local parallelism are link-level striping and switches which allow packets traveling between the same source and destination to take different paths through the switch. One of the challenges of understanding reordering due to local parallelism is that the parallelism is not made externally visible (except through anomalous behavior like reordering). So, unlike Paxson’s experience with multi-path routing [12], where the datagrams’ different paths were clearly indicated by the different addresses of the routers they traversed, the only hint of local parallelism may be that a particular hop exhibits varying delays and reordering.

Manuscript received October 9, 1998; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor S. Pink.

J. C. R. Bennett was with BBN Technologies, Cambridge, MA 02176 USA. He is now with the Division of Applied Sciences, Harvard University, Cambridge, MA 02138 USA and with RiverDelta Networks, Melrose, MA 02176 USA (e-mail: jcrb@riverdelta.com).

C. Partridge and N. Shectman are with BBN Technologies, Cambridge, MA 02176 USA (e-mail: craig@bbn.com; nshectma@bbn.com).

Publisher Item Identifier S 1063-6692(99)09716-2.

Local parallelism is on the increase in the Internet because it reduces equipment and trunk costs. It is often more cost effective to put two components in parallel than to use one component that has twice the speed. It is particularly cost effective when purchasing long-haul serial links because tariffed options often offer link bandwidths that are large multiples of each other. In addition, parallel links are a very useful way to improve reliability. If the parallel links follow different physical paths, the virtual link they implement is less vulnerable to single-point failures (e.g., backhoes). Not surprisingly, large businesses, Internet Service Providers (ISP’s), and their vendors are aggressively pursuing parallel links. In a survey of 38 major ISP’s conducted in mid 1997 [2], only two of the smaller ISP’s did not have parallel paths between nodes.

Compounding this trend is the fact that interface speeds typically trail customer needs, especially ISP’s needs. ISP’s are forced to cobble together the link rates they need by aggregating interfaces. An example of an ISP-like situation where bandwidth demands have led to parallelism is the MAE-East peering point, where we have observed extensive reordering.

In the rest of this paper, we will examine the reordering problem further. Section II reports on tests done to measure the prevalence of reordering. Section III presents how reordering affects Transmission Control Protocol (TCP) performance. Section IV reports in detail on why reordering occurs at MAE-East and then discusses why it is difficult to get rid of reordering. We conclude, in Section V, with a discussion of how this work may affect the construction of networking equipment and the design of networking protocols.

II. OBSERVING AND MEASURING REORDERING

We performed a number of tests between a variety of hosts to determine average frequency and severity of the reordering that any given session might expect to experience. We discovered that reordering is a very complex phenomenon. In particular, we found that reordering is strongly a function of the following properties of the path taken by a session: 1) the existence of parallel links between nodes on the path; 2) the exact configuration of the hardware and software in the nodes; and 3) the load on the nodes.

A. What We Measured

Because we saw signs that reordering was a function of load and configuration, we looked for measurement paths

for which we could get detailed information about load and configuration. In the end, we chose to measure sessions over paths that traversed the MAE-East exchange.

The MAE-East was an excellent example case, for a number of reasons.

- 1) We have full control of, and knowledge about, all the routers between the source of our tests and the MAE-East exchange, allowing us to insure that they are not effecting our results.
- 2) There is a wide diversity of equipment and hosts only one or two hops on the other side of MAE-East from our ISP.
- 3) The topology and makeup of MAE-East are publicly known.
- 4) The traffic statistics at MAE-East are publicly available
- 5) It is a heavily loaded site containing a diversity of links, routing paths, and parallel components likely to encourage reordering.

The reader is strongly cautioned not to infer, however, that reordering at MAE-East is somehow special. We have found many other paths upon which reordering occurs.¹

B. The First Test: How Many Sites Experience Reordering?

We chose 140 Internet hosts which we had reason to believe were topologically close to MAE-East, such as server machines of the ISP's peering at MAE-East. To each host, we sent a burst of five 56-byte ICMP-ping packets to insure that there were no problems caused by route cache misses, followed 5-s later by a back-to-back burst of 50 ping-packets.² We report here the results of two tests spaced a month apart. Each test was run from 6 to 7 pm EST on a Monday evening.

We then measured the number of sessions that experienced out-of-order delivery. To avoid having to define "out-of-order" in the presence of packet loss, we measured the number of probes which received all of the first ten packets, then how many probes that received all the first ten packets received them in order. We also measured the number of probes which received all of the first 20 packets, and how many probes that received all the first 20 packets received them in order.

C. The Second Test: Reordering and Load

Detailed daily traffic load information is available for MAE-East so we decided to conduct a more detailed test of the relationship between traffic load and reordering. We chose a site which demonstrated a particularly high degree of reordering and sent it a 100-packet burst of 512-byte packets every minute for four days.

In this test, we also developed a metric for how much reordering took place. The metric we chose was the average number of SACK blocks required to cover the out-of-order ping replies received, if the session were a TCP connection. Losses were eliminated from the traces by renumbering of sequence numbers (i.e., if packet 8 was lost, packet 9

TABLE I
REORDERING TESTS

Hosts Receiving Packets	1-10	1-10	1-20	1-20
Date of Test		In order		In order
Early Dec. 1997	117	7	97	5
Early Jan. 1998	60	4	46	1

was renumbered as packet 8 for the purpose of computing the SACK block metric). Thus, all the holes in the SACK scoreboard are due to reordering and not to packet loss.

We chose the SACK block metric after trying several other metrics. Our goal was to find a metric which had the following properties:

- 1) generally matched an intuitive assessment of the amount of scrambling in a given data set;
- 2) had a higher metric if data passed through two scramblers than if data passed through a single scrambler (and ideally, would yield a value equal to the sum of the metrics for the scramblers separately);
- 3) be reasonably independent of the size of the data set;
- 4) be reasonably independent of data loss;
- 5) had a minimal value for in-order data;
- 6) was reasonably easy to calculate.

Among the metrics we tried were displacement metrics, number of monotonic sequences required to include all data, and the number of spurious retransmissions that would have been generated by a standard TCP implementation. Of these, SACK blocks proved to be the best metric, although it is not independent of the size of the data set.

D. Results

Table I shows the results of the first test, pinging 140 sites. We show the number of sessions that received all of the first 10 and first 20 packets, and then the number of these sessions that received all the packets in order. The table suggests the extraordinary result that probability of a session experiencing reordering is over 90%.

The reader may find this result hard to believe. If packet reordering were this common, surely it would have been seen before. We believe there are three reasons why reordering has not been seen before. First, and foremost, the reordering at MAE-East is a function of network load (as shown in the discussion of Fig. 1 below), and is almost nonexistent below a certain load. Prior to middle of 1997, the load at MAE-East almost never exceeded this threshold.

Second, to cope with increasing traffic, MAE East reconfigured its network in November 1997 in a way which increased the likelihood of reordering. Third, the major ISP's all have private peering points which, by and large, do not use parallel links and thus do not suffer reordering.

Fig. 1 shows some of the results from polling one site regularly. The figure shows both the degree of reordering and the traffic load at MAE-East for three days in early 1988. The particular dates were chosen because they vividly illustrate both the relationship of load and reordering and confirm that reordering was occurring at MAE-East. On the third day, between 10–11 pm on February 2nd, a private peering

¹ For example, a quick check at the routing data of the *Traceping* page at Oxford University, U.K., (available HTTP: <http://av1.physics.ox.ac.uk/>) will often show large numbers of parallel links visible at the IP level.

² This test was done using the ping *preload* command.

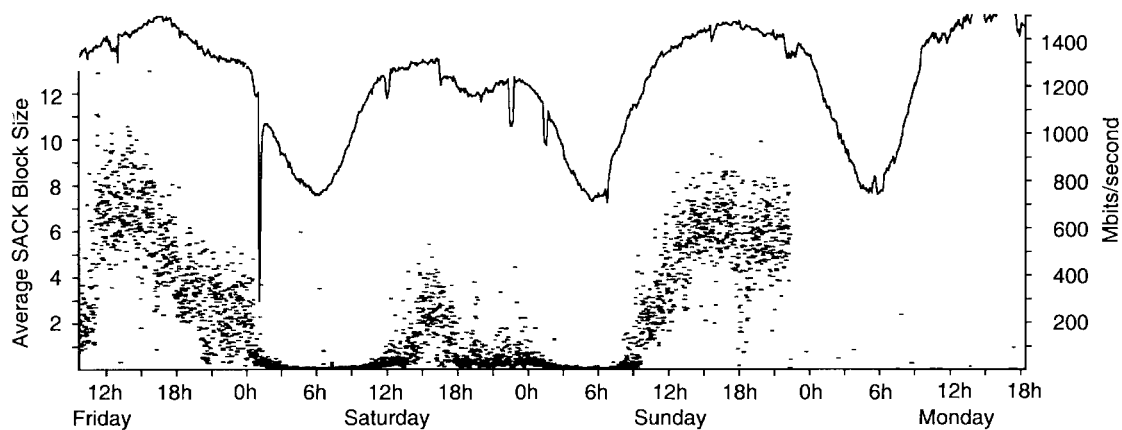


Fig. 1. Polling results.

relationship was established in northern Virginia between the two ISP's on the path, exactly bypassing MAE-East. Observe that the load-curve tracks the SACK curve until the private link was put in place, after which reordering largely (but not entirely) went away.

III. IMPACT OF REORDERING ON TCP

So why does reordering matter? This section provides some motivation for worrying about reordering by examining how packet reordering impacts the performance of the TCP.

TCP is a reliable transport protocol [14], designed to recover from all sorts of misbehavior at the Internet Protocol (IP) layer. As a form of misbehavior, reordering looks pretty tame, since all TCP segments eventually get to the right place. The problem is that large-scale packet reordering tends to drive TCP into two destructive modes. In one mode, it has great difficulty opening its congestion window and makes inefficient use of available link bandwidth. In the other mode, it loses self-clocking and becomes highly bursty.

There are two types of packet reordering that we will examine, data reordering or *forward-path reordering*, and ack reordering, or *reverse-path reordering*. We consider these separately for two reasons. First, the asymmetric nature of Internet routing [12] frequently results in connections which only experience reordering in one direction. Second, the impact of the reordering is different depending on which direction it takes place over. When data is reordered the result is the transmission of duplicate acks, but when acks are reordered, the result is the reception of out-of-order nonduplicate acks.

A. Forward-Path Reordering

In forward-path reordering, TCP segments containing data arrive at the receiver out of order. Reordering of data has five effects: 1) reordering causes unnecessary retransmissions; 2) reordering makes it difficult or impossible to grow TCP's congestion window (*cwnd* and *ssthresh*); 3) reordering can cause actual packet losses to be obscured, so that TCP retransmits lost packets later than it normally would; 4) reordering may cause the round-trip estimator to poorly estimate the round-trip time; and 5) reordering may reduce the efficiency of the

receiving TCP. After briefly presenting the general behavior of a TCP receiver, we discuss each of these effects in turn.

1) *What a TCP Receiver Does*: When TCP is receiving data segments, its behavior is very simple. It acks the data it has received. If the data is in order, the acknowledgment is a simple ack which indicates the receiving TCP has received all data through the last byte of the most recent segment.

If the data is out-of-order, things get more complicated. At minimum, the receiving TCP sends a duplicate acknowledgment, which repeats the acknowledgment of the last in-order byte received. If the TCP implements selective acknowledgments (SACK [9]) the ack will both acknowledge the last in-order byte, and also acknowledge the new out-of-order data received.

The effects of reordering in the forward path are determined by how the TCP sender responds to these various acks.

2) *Unnecessary Retransmissions*: Most TCP's implement an algorithm known as fast retransmit [15] to try to recover quickly from occasional packet loss. The fast retransmit algorithm is triggered by a series of duplicate acks. If the TCP sender sees three duplicate acks, it assumes that the data immediately after the byte being acked has been lost, and retransmits that data. The idea is that the duplicate acks indicate a packet has been lost, and by quickly retransmitting the data, the sending TCP can repair the loss before it causes the current pipeline of data to drain.

If the duplicate acks are caused, however, by reordering, then the fast retransmission is unnecessary and wastes bandwidth by sending the same data over the wire twice.

Note that TCP implementations that use forward acknowledgments (FACK) behave even less well because FACK does a fast retransmit if the sender receives an ack that suggests that three segments worth of data are outstanding.³

3) *Difficulty Growing TCP's Congestion Window*: A further problem with reordering is that when a TCP does a fast retransmit, it reduces its congestion window *cwnd*, in the belief that it has seen packet loss. Unfortunately, if there are multiple fast retransmits in a single window (which reordering can trigger), the congestion window is quickly driven down to

³FACK is a sender-side algorithm that uses the information in SACK's to more fully model the receiver's window [10].

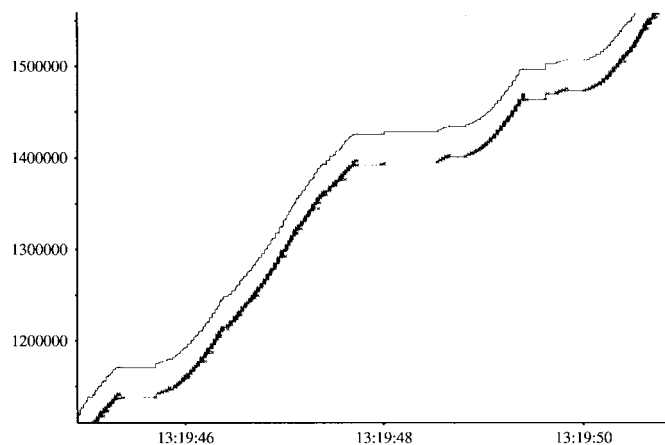


Fig. 2. src1->dst1.4885.

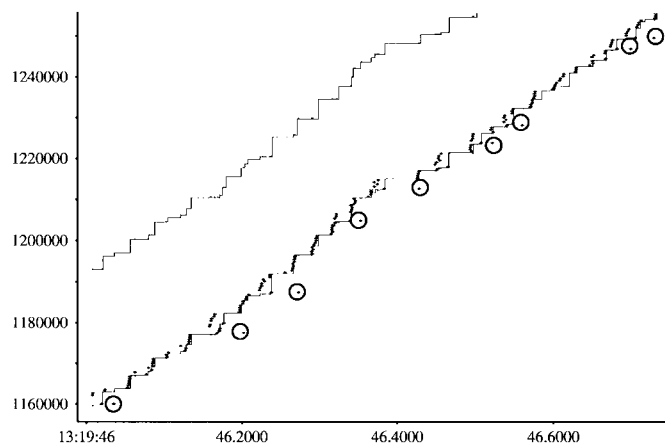


Fig. 3. src1->dst1.4885.

a very low value. The result is that TCP has difficulty opening its congestion window and properly utilizing the link.

This behavior is illustrated in Figs. 2 and 3. Both are plots of part of a 2.1-MB file transfer passing through MAE-East. Observe that in Fig. 2 the usual neat exponential curves of TCP congestion control are interrupted by an extended period of slower growth. Fig. 3 shows part of this period in detail. The TCP is suffering from forward-path reordering, and the duplicate acks are causing regular fast retransmits that reduce the congestion window. The particular behavior is that any time the sender opens its window to about eight segments, segments are reordered by three or more positions. The resulting spurious fast retransmissions are circled. Almost immediately after most of them, the original transmission is acked, and TCP emits new segments but into a reduced congestion window. It is sobering to observe that, in the half-second this graph covers, all but one (that at time 46.475) of the ten fast retransmissions observed are unnecessary.

4) *Obscured Packet Loss*: What happens when, along with reordering, a segment actually gets lost? The unfortunate answer is that TCP often does not initially see the loss because the reordering conceals it. The sending TCP is forced to wait for a timeout to retransmit the loss.

An example may be useful. Suppose the sender transmits segments 1–6, that segment 3 is discarded, and that segment

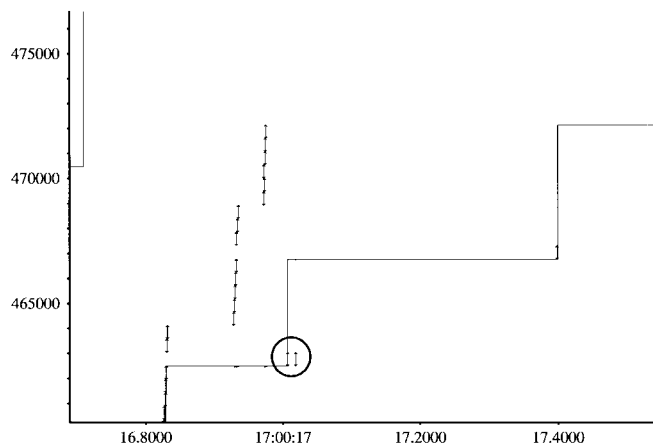


Fig. 4. src1->dst1.4969.

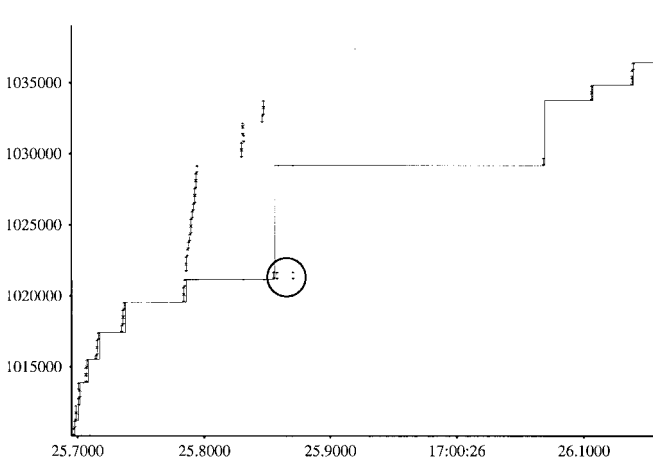


Fig. 5. src1->dst1.4969.

2 is reordered so it arrives after segment 6. The receiver will generate an in-order ack for segment 1, and then three duplicate acks for segment 1 as it receives 4, 5, and 6. These duplicate acks will cause the sender to perform a spurious fast retransmit of segment 2. When segment 2 finally arrives at the receiver, the receiver generates an ack for segment 2, but since there are no more segments in flight, no duplicate acks will be generated, and the sender will not be aware, until it times out, that segment 3 is missing. If the segments had not been reordered, then the duplicate acks would have been for segment 2 and spurred a fast retransmit of segment 3. Figs. 4 and 5 show examples of obscured packet losses. In each figure the reordered packet and its false retransmission are circled. In both cases, the lost packet must be recovered by timeout. One of the modifications proposed for FACK TCP, moving the duplicate ack retransmission to the SACK scoreboard so that fast retransmit can be performed for each hole in the scoreboard, would enable TCP to recover from this scenario, but in cases with more severe reordering it would suffer from many more false retransmissions.

5) *Poor Round-Trip Time Estimation*: Another impact of reordering is that it can, at least in theory, cause some TCP round-trip time estimators to stop working due to lack of round-trip time samples.

Most TCP implementations do round-trip time estimation in one of two ways. If the TCP implements the TCP Timestamps Option and the option has been negotiated at connection setup, then the sending TCP estimates the round-trip times by subtracting the timestamps echoed in the acks from the present time to get an accurate round-trip time sample [5]. This algorithm works correctly even in the face of reordering.

However, if TCP timestamps are not in use, most TCP's use an algorithm that samples one round-trip time per window of data. The usual implementation is to note the sequence number of the first byte of the segment currently being sent, start a timer, and then sample the timer when an ack that includes the sequence number of that byte is received. However, if the byte being sampled is retransmitted, the sample must be discarded, to avoid a recurrence relationship [6], [7].

Unfortunately, in the presence of spurious fast retransmits, TCP is likely to have to discard most of its potential samples. As a result, the round-trip time estimator will not sample the round-trip time very frequently and will not keep a good estimate of the round-trip time. A poor round-trip estimate will often increase the impact of obscured packet loss, because the lost packet will be transmitted later than it needs to be.

6) *TCP Receiver Efficiency*: TCP requires that the receiving application receive data in order. As a result, network reordering of data segments places serious burdens on the TCP receiver.

First, the receiver must buffer out-of-order data until data is received to fill gaps in the sequence space. At some stage this process requires sorting data into proper order, an expensive operation that is never required if data arrives in order.

Second, the data being buffered is being withheld from the receiving application. Rather than getting its data in a smooth series of modest chunks, the application gets the data in bursts. The overall efficiency of the system suffers because it is not work conserving.

Third, many TCP implementations use the *header prediction* [4] algorithm to reduce the costs of TCP processing. Header prediction, however, only works for in-order TCP segments. So, if segments are reordered, most TCP implementations do far more processing than they would for in-order delivery.

B. Reverse-Path Reordering

In reverse-path reordering, the acks traveling back to the sender are reordered. For the purposes of this section, we will assume that no reordering affects the forward path, since that simplifies the discussion.

It is important to observe that the ack patterns are very different for forward and reverse-path reordering. In forward-path reordering, the reordering of data segments causes the receiver to send duplicate acks. In reverse-path reordering, the receiver is sending cumulative acks (because data is arriving in order), but the acks are being reordered on their way back to the sender.

Reverse-path reordering has one major effect: a loss of self-clocking leading to highly bursty transmission patterns.

1) *Loss of Self-Clocking*: In the late 1980's, TCP was found to have a nice self-clocking property: the acks roughly

reflected the rate at which data could transit the network and be removed from the network by the receiver [3]. During in-order delivery, the sender sees an even sequence of acks, each acking one or two segments worth of data and notifying the sender it is free to inject one or two more segments into the network.

Reordering acks scrambles self-clocking. The problem is that, because all the acks are cumulative, reordering acks means that an ack for a later byte arrives early. What the sender sees is an occasional ack, acknowledging a large amount of data, followed by a number of apparently redundant acks which acknowledge data that the first ack already acknowledged. This ack pattern has at least two harmful effects.

First, it affects the rate at which the sending TCP opens its congestion window. The sender opens its congestion window for each ack that acknowledges previously unacknowledged data. The window is opened by one maximum segment size per ack during slow start, and by a fraction of the maximum segment size during congestion avoidance. The goal is to roughly double the window size per congestion window's worth of segments during slow start, and to grow the congestion window by one segment per window's worth of data after a slow start.

However, if acks are reordered, there may be only one or two acks that acknowledge new data, even if tens or hundreds of segments have been sent. The unfortunate result is that the congestion window will grow far too slowly (e.g., one-tenth or -hundredth of its intended rate). Particularly severely hurt by this behavior are long delay bandwidth paths like satellite paths and high-speed transcontinental links.

The second effect of ack reordering is to encourage bursts of segments. An in-order ack frees up space within the window and permits the sender to transmit an amount of data equal to the amount acknowledged. So, if acks are reordered and one ack arrives early and acknowledges a large amount of data, the sender is likely to launch a volley of new segments all at once. Some spacing of the segments will occur at bottleneck links, but bursts typically lead to closely spaced acks and, if ack compression occurs on the reverse path, makes more ack reordering likely. In short, this behavior is self-reinforcing, at least for modest periods. If the sender gets an out-of-order ack, the sender will inject a burst of data that makes receiving another out-of-order ack more likely. Note that this effect has almost nothing to do with the congestion window (which the first effect modifies).

Figs. 6 and 7 show this self-reinforcing behavior in different levels of detail. The trick is to look at the bottom line, which shows the sequence number acknowledged by the latest ack. Cases where the line goes up and then down again reflect acks which have been reordered, such that they arrive early. Fig. 6 shows a long interval of reordering. Fig. 7 illustrates a particular period of clearly self-reinforcing behavior, where four successive bursts were triggered by the reordering of acks from a prior burst.

C. Forward- and Reverse-Path Reordering

What happens when a TCP experiences both forward- and reverse-path reordering? The traces we have so far suggests

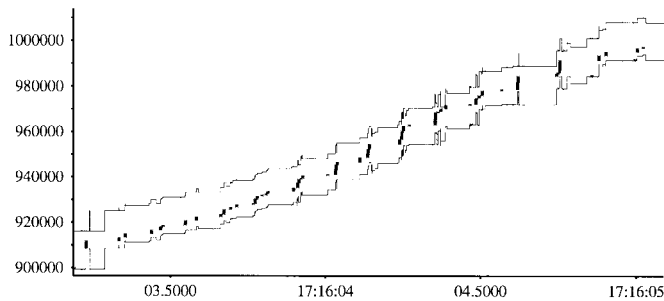


Fig. 6 SRC-A->Dst-B.

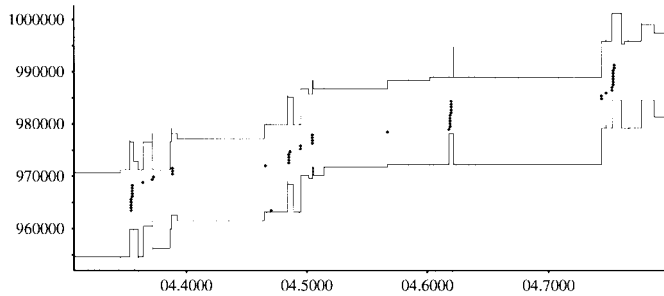


Fig. 7 SRC-A->Dst-B.

that TCP tends to oscillate between forms of misbehavior. In some situations the effects of forward-path reordering dominate, and the TCP has trouble doing slow-start. In other situations reverse-path reordering dominates and we see bursty behavior. One can see some of this style of behavior in Fig. 6, where on the left side of the figure, we see the TCP having difficulty opening its congestion window due to forward-path reordering, while on the right side, we see bursts due to reverse-path reordering.

IV. UNDERSTANDING THE CAUSES OF REORDERING

We believe that reordering will continue to be a fact of life on the Internet. It is a natural result of the increasing use of parallelism in network devices and surprisingly hard to prevent.

To explain our views, this section starts with a practical example: why reordering occurs at MAE-East. The section then proceeds to discuss the challenges in eliminating reordering.

A. Hunt Groups at MAE-East

To the best of our knowledge MAE-East is currently the largest Internet exchange point in the world in terms of number of bits per second of traffic⁴ and handling perhaps as much as a third of *all* long-haul Internet traffic [2]. The primary network device used to switch this volume of traffic is the DEC Gigaswitch/FDDI which is a multiport FDDI switch with a crossbar based switch fabric. A particular feature of the Gigaswitch is its support of what DEC calls *hunt groups*. A hunt group is a collection of ports which operate as a virtual

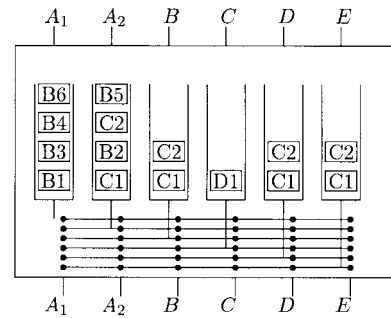


Fig. 8 Hunt groups.

link; that is, they allow two layer-2 bridges to be connected by multiple parallel links for increased bandwidth capacity.

Fig. 8 is a simplified diagram of a Gigaswitch with six links labeled A_1 , A_2 , B , C , D , E . The first two links have been labeled with subscripts to indicate that they are members of the same hunt group. Any packet which wants to go over virtual link “A” can take either link A_1 or link A_2 . The input buffers on the Gigaswitch/FDDI are single FIFO’s, which results in *head-of-line blocking*. While this blocking increases reordering, it is not the fundamental cause of reordering. In the figure each buffer contains a sequence of packets labeled with the letter of the output port for which they are destined and a transmission sequence number. In this example, assume that all packets between a given input port/output port pair belong to the same “flow.” Thus, one can see that packets belonging to the flow between the input hunt group ports A_1 , A_2 and the output port B has been split between the two hunt group ports, but that the packets were still effectively received in order.

Because the input buffers are FIFO, they always try to send the packet at the front of their queue to its output port. Switch allocation is done through a bidding system, in which the input requests to connect to the output and waits for the output port to *grant* the connection. The output port may receive requests from multiple input ports but may only grant one because the crossbar can only connect one output port to one input port at a time. While an input port is waiting for a grant it can not send any packets due to the FIFO nature of the buffer.

Now consider the behavior of the switch shown in Fig. 8. For simplicity, assume that all the packets are the same size so that we can speak in terms of *time slots*.⁵ In time slot 1, only port A_1 has a packet for port B , so it will certainly get a grant to send packet $B1$ to port B . In time slot 1, ports A_2 , B , D and E all have packets bound for port C . Assuming a fair arbitration scheme, there is only a one-in-four chance that port A_2 will get to send its packet to port C . If port A_2 does not get to send a packet to port C , then in time-slot 2, port A_1 will send packet $B3$ to port B , while packet $B2$ is still

⁴ As of this writing, the average load at MAE-East is *continuously* greater than 1 Gb/s for the better part of the working day, as reported on the MAE-East statistics web page, available at <http://www.mfst.com/MAE/east.html>.

⁵ In fact, if we consider the general case of different-sized packets, we see another case where, even without FIFO queuing, we can see very large degrees of reordering. Suppose an input A has packets bound for two outputs C and D and is sending a long packet to C when D finishes sending a packet it received from a different input, B , which has a long stream of packets bound for output D . Then, the only way that A can send its next packet to D is if D waits for A to finish sending its packet to C . If D starts accepting a new packet from B while A is busy sending to C , then in fact, A may never send a packet to D . In such a switch, the only way to fairly arbitrate amongst the different ports is for the outputs to be *nonwork conserving*.

sitting in the queue at port A_2 ! Even if port A_2 sent packet C_1 during the first time slot, there is only a one-in-two chance that port A_2 will get to send packet B_2 during time-slot 2. Thus, reordering is a consequence of the bid scheme.

The reordering can be quite dramatic. For instance, the DEC Gigaswitch Operation Manual version 3.2 notes that hunt groups may cause packet reordering which may cause TCP to incorrectly retransmit packets and that to prevent this users may wish to set TCP's duplicate ack retransmission threshold *tcp_rmxthresh* to 100! (Compare this value with the more common value of three that is widely used).

At this point, one may be tempted to simply declare that the Gigaswitch has a bug. That would be a mistake. Reordering is, in fact, a natural consequence of parallelism. To explain why this might be so, the next several subsections discuss various approaches to solving the reordering problem, and why they generally fail.

B. Fixing the Switch

Suppose we try to fix the Gigaswitch's switch to maintain packet ordering. In particular, we want to find a way to allow traffic to be split across two input ports of the switch and yet stay in the order they were received. (We leave for later the problem of ensuring that the order in which packets are received is the same as the order in which they were transmitted). Furthermore, we look to fix the Gigaswitch in a work conserving way. *Work conserving* means that no part of the switch is idle when there is data that could be sent through that part of the switch [16]. We seek to have the switch be work conserving because a switch which is nonwork conserving typically has to run faster to handle the same load of traffic.

Packet switch architectures can be generally classified into one of four types: 1) shared buffer; 2) output buffered; 3) multistage; and 4) input buffered. Multistage switches look particularly unsuited to this problem. They have buffering and queues inside the switch and, as a result, packets from different inputs can experience wildly different delays through the switch. But the other three switches, at least at first glance, appear amenable to retaining order.

1) *Shared-Buffer and Output-Buffered Switches*: In a shared buffer switch, the switch fabric is comprised of a central memory system. In the case of a *nonblocking*⁶ switch, it must have an I/O bandwidth equal to the sum of the *bi-directional* link capacities.

In an output-buffered switch, the switch fabric is a trivial copy network, which delivers a copy of every packet received to all the output ports. Each output port has a device called a *concentrator* which is responsible for filtering out packets not destined for the output. The concentrator's input must have an I/O bandwidth equal to the sum of the inbound link capacities. Because several packets for the output may arrive at the same time, the concentrator's output must have an

I/O bandwidth which is much greater than the output link bandwidth.

Clearly, both of these switches can maintain packet order. Packets will enter the shared memory in the order received on a shared buffer switch, and with appropriate management, can be encouraged to leave in order. Similarly, packets arrive in the output concentrator in the order received and presumably the concentrator can maintain the order. Unfortunately, neither switch design scales well.

The limiting factor in building shared memory switches is the speed of memory chips and the physical/electrical difficulties of multiplexing all the inputs into the memory. For example consider a shared memory ATM switch using SRAM with a 10-ns cycle time. Even if the memory is 424 bits wide so that an entire cell can be written into the memory in one cycle, and an entire cell can be read out from the memory during the next cycle, the switch would only be able to process 50-million cells/s for a total capacity of about 20 Gb/s. While this could be doubled with dual-ported memories, and an IP router might be able to use somewhat wider memories, these are only constant factor improvements.

While the challenge of building a concentrator is substantially simpler, we are still bound, on the outbound side of the concentrator, by the limits of memory speeds.

In short, even though they can maintain ordering, these switch designs are not expected to scale up to the speeds the Internet requires.

2) *Input-Buffered Switches*: Input-buffered switches have the majority of their buffer capacity located on the input side of the switch fabric. The primary reason for a switch to be input buffered is so that the switch fabric and buffers only need to run at, or slightly above the speed of the fastest link as opposed to output buffered switches where the fabric and the buffers must run at n times the link speed for an n port switch. Currently, the most popular switch fabric for implementing input-buffered switches is the *crossbar*.

The challenge in constructing a fast (or large) crossbar is designing the switch *arbiter*. The arbiter controls the connection of input ports to output ports. Using the best algorithms, the complexity of the arbiter is a roughly linear function of the product of the number of ports and the number of bids each port may submit. To eliminate head-of-line blocking, most input-buffered switches require each input port to bid in parallel for all output ports for which it holds packets. A good arbiter also has the property that the way it services competing bids appears, for practical purposes, to be near random. In particular, no one port should be predictably favored over another.

The GigaSwitch is an input-buffered switch and, as Sections II and IV-A showed, it suffers from reordering. The reason (beyond head-of-line blocking effects) is that the arbiter will randomly mix traffic from A_1 and A_2 , and there isn't an obvious way to solve this problem. (The one apparently obvious solution, namely having the transmitter sequentially number the packets it sends on A_1 and A_2 , does not work very well. It reintroduces head-of-line blocking, in a new form, and makes the arbiter hideously complex to boot).

⁶It is important to understand the definition of the term *nonblocking*, since outside of academic circles (read "in marketing literature") this term has been greatly abused in recent years. We use the classical circuit-switch definition of *nonblocking*, which is that at a given instant, a nonblocking switch is capable of connecting all inputs to all outputs in an arbitrary permutation.

Observe that the problem is cost and complexity. We might be able to build a input-buffered switch that maintained ordering, but it would be more expensive.

3) *Make the Switch Ports Faster*: Another approach to fixing reordering is to make the switch ports faster, and attach multiple interfaces through one port. The links for A_1 and A_2 now come into the same interface board. (Note that speeding up the ports does not imply that the overall must be faster overall; we can have fewer ports that are clocked at a higher bandwidth).

By keeping all the links in a hunt group on the same board, the problem of reordering becomes a purely local problem between the transmitting and receiving interface boards. It is still not a simple problem (see Appendix I), but good algorithms to preserve order exist [1], including one that is work-conserving [11].

Unfortunately, this solution is also somewhat limiting. The maximum bandwidth that can be bundled together is the amount that can be fit on one interface card. In situations like MAE-East, where a few core switches may need a large amount of bandwidth between them, this solution may not serve.

C. What Can IP or TCP Do?

Yet another approach is to ask whether the IP or TCP layers can do anything to help solve this problem.

1) *IP*: It turns out that IP can help the problem slightly, but also that current trends in forwarding engine design may actually cause more problems.

If IP is made aware of the presence of the parallel links and can direct which link each packet is sent over, then IP can ensure that traffic between a particular source and destination always goes over the same link. This trick is easy to do (the common approach is to use a simple hash of the source and destination, modulo the number of parallel links) and ensures that no conversation sees reordered packets. But, it has a number of pitfalls. First, it means that no one conversation can use the entire bandwidth of the parallel links. Second, depending on the mix of conversations and the quality of the hash function, there may be uneven assignments of conversations to links. In an extreme case, one might find all the conversations sharing one link while the others are idle.

Looking forward a bit, we can also see a potential problem for the IP layer. In the past year, several vendors have developed inexpensive ASIC implementations of IP forwarding, each capable of forwarding 1–8 million packets/s. As OC-192c interfaces have become available, some vendors may need to use two or more forwarding ASIC's to get enough forwarding power to keep pace with OC-192c's 10-Gb/s data rate. That situation creates a new opportunity for reordering, if two successive packets are handed to different ASIC's. There are plenty of scenarios in which, depending on the state of the two ASIC's, one ASIC may process the same header faster than another, especially if the ASIC's keep a recently-used route cache. The hashing trick works here too: We can hash the headers to allocate them to ASIC's. But the same problems exist too: we can easily poorly utilize the ASIC's.

2) *TCP*: What can TCP do to better handle reordering?

First, TCP can detect when it mistakenly does a fast retransmit and undo the congestion window changes. The obvious marker of a fast retransmit is that the ack comes back far too soon (typically just after the fast retransmit is done).

Second, the TCP can support SACK. SACK's give the sending TCP a fighting chance to figure out that reordering is taking place by more completely reflecting what is happening in the window (even with reordering). For instance, the SACK specification [9] indicates SACK's should reflect the order in which segments are received. As a result, a sender can judge, based on the ordering of SACK's, whether forward-path reordering is taking place and can adjust the three segment threshold for fast retransmit accordingly. (Indeed, if no reordering is taking place, one duplicate ack might suffice, while if reordering is ongoing, it may be 30 or 40 duplicates before we can safely assume a segment is lost).

These two patches, however, are far from a complete answer to the problem.

D. Numbering Packets

So far, we have demonstrated that getting a parallel system to naturally keep packet order, without introducing inefficiencies and higher costs, is very hard and there are no obvious correctives that can be applied in IP and TCP.

We conclude this section by considering the obvious question: why not just number the packets on a parallel path and then sort them back into order? (And do the analogous thing for IP forwarding ASIC's).

This approach works, and indeed, good algorithms even exist which do the numbering implicitly, at some risk of reordering [1], but there is an annoying limitation. As Appendix I proves, for a broad range of situations, a scheduling system that keeps packets in order over parallel links must be nonwork conserving. That means there is some inefficiency in the system; we are wasting some system capacity to retain ordering.

All is not completely lost. Work-conserving algorithms exist to keep the packets in order, but they require us to divide the packet processing into uniform units of work. To illustrate this idea, suppose we striped packets bit-by-bit onto two parallel links and removed them bit-by-bit at the opposite end. Ignoring the problems of errors and the inevitable slight variations in link delays, this scheme delivers packets in-order and is work conserving. An approach based on this idea exists for multiplexing parallel links [11]. Applying this idea to forwarding ASIC's, or switch schedulers, remains an unsolved problem.

V. CONCLUSION

We hope we have convinced the reader that reordering is now a part of the Internet's behavior and can cause serious performance problems, especially for TCP. Furthermore, we hope the reader agrees that eliminating reordering is a hard problem, both economically and theoretically. Given those two observations, how should we go forward?

First, we should try to make reordering rare. Good algorithms exist [1], [11] to minimize or eliminate reordering in some situations. In other cases, algorithms like numbering packets can be used to at least minimize reordering. These algorithms should be used to reduce the frequency of reordering.

Second, we need to begin to accept that some reordering is still going to occur, and probably often enough to affect performance. In short, we need to start designing protocols to recover more gracefully from reordering. We think observing the SACK patterns to learn if reordering is occurring is a nice idea, but it is only a starting point. This topic desperately needs more attention.

APPENDIX

REORDERING FROM A THEORETICAL PERSPECTIVE

It may be useful to look at reordering from a theoretical perspective. Many switch and router architects design their equipment to be *work conserving*. Informally, a work conserving system is one where a packet is always sent as soon as the link on which it is to be sent becomes idle. (For a formal definition, see [16], [8]). A nice feature of work conserving systems is that they have the minimum achievable average queuing delay if we assume the system does not discard packets. Furthermore, a nonwork conserving system is clearly wasting some amount of capacity (it is leaving a line idle while work for that line builds up). While the wastage can be small, in many nonwork conserving schemes, the wastage (and thus reduced effective bandwidth) is substantial.

Unfortunately, it is very easy to prove that, for a fairly simple and common set of properties, a system cannot be work conserving and maintain packet ordering while doing link-level load balancing. In particular, given the following four conditions, we can prove that no system can maintain packet ordering and be work conserving:

- 1) unit of transmission is the packet (we do not try to break the packet up across links);
- 2) links across which we balance operate at the same rate;
- 3) packet sizes vary by more than a factor of two;
- 4) transmission delay over all links is the same.

The trick is to observe when the last bit of each packet leaves the sender. If the last bit of the second packet leaves before the last bit of the first packet, then the receiving system will either reorder the packets or be nonwork conserving, because the receiving system will receive the last bit of the second packet first (as a consequence of conditions 2 and 4), and be forced to either delay the second packet awaiting the first packet, which is nonwork conserving, or send the second packet on, causing reordering in relation to the first packet.

The proof is very simple. Assume our system is work conserving and seeks to retain packet ordering. Further assume the system is currently idle, and the sender is given packets in sequence order at a rate which is greater than or equal to the sum of the individual rates r of the n outbound links. Now assume the system receives two packets to transmit. The first packet has length $2l + 1$ bits and is one bit longer than twice the length of the second packet which is of length l .

We can quickly observe that we must send the first packet first to retain ordering. (Since the first packet is longer and transmission delays are the same on all links, we cannot send the first packet later and have it arrive first). So, the question is: can we send the two packets in order, in a work-conserving fashion, and retain their ordering at the receiver?

The equations vary slightly, depending on whether the system starts transmitting each packet when it gets the first bit of the packet, or waits until it has the last bit of the packet. If we begin transmitting on the first bit, then when the first bit of the second packet starts to be transmitted, we will have transmitted $(2l + 1)r/nr$ bits of the first packet. The second packet takes l bit times to transmit, during which we will also send l more bits of the first packet. Thus, when the last bit of the second packet departs, we will have sent $(2l + 1/n) + l$ bits of the first packet. We know that n is at least 2, so this value is less than $2l + 1$ and thus we will not have completed sending the first packet before the last bit of the second packet departs. The analysis is the same if the system transmits after receiving the last bit of the packet, except that we transmit even less of the first packet (just l/n bits) before the last bit of the second packet departs. In either case, the system completes sending the second packet before completing the first packet, and since the link delays are identical, the receiver will receive the complete second packet first.

From this proof we can conclude that if we take the straightforward approach to link-level load balancing, we must either be nonwork conserving or reorder packets. Neither is an appealing idea.

ACKNOWLEDGMENT

The authors would like to thank D. Li for her very constructive comments.

REFERENCES

- [1] H. Adishesu, G. Parulkar, and G. Varghese, "A reliable and scalable striping protocol," in *Proc. ACM SIGCOMM'96*, pp. 131–142.
- [2] R. Gareiss, Is the internet in trouble. *Data Communications Magazine*, Sept. 1997.
- [3] V. Jacobson, "Congestion avoidance and control," *Proc. ACM SIGCOMM'88*, pp. 314–329.
- [4] ———, "4bsd header prediction," *ACM SIGCOMM Computer Communication Rev.*, Apr. 1990, pp. 13–15.
- [5] V. Jacobson, R. T. Braden, and D. Borman, *TCP Extensions for High Performance*, rfc-1323, May 1992.
- [6] R. Jain, "Divergence of timeout algorithms for packet retransmissions," in *Proc. 5th Phoenix Conf. Computers and Communications*, Mar. 1986, pp. 174–179.
- [7] P. Karn and C. Partridge, "Estimating round-trip times in reliable transport protocols," *ACM Trans. Comput. Syst.*, vol. 17, pp. 2–7, Nov. 1991.
- [8] L. Kleinrock, *Queueing Systems; Volume 2: Computer Applications*. New York: Wiley, 1976.
- [9] M. Mathis, J. Mahdavi, S. Floyd, and A. Romanow, *TCP Selective Acknowledgment Options*, rfc-2018, Oct. 1996.
- [10] M. Mathis, J. Semke, and J. Mahdavi, "The macroscopic behavior of the TCP congestion control algorithm," *Comput. Commun. Rev.*, vol. 27, no. 3, pp. 67–82, July 1997.
- [11] C. Partridge and W. Milliken, "Method and apparatus for byte-by-byte multiplexing of data over parallel communications links," patent application, 1998.
- [12] V. Paxson, "End-to-end routing behavior in the internet," in *Proc. ACM SIGCOMM'96*, Stanford, CA, Oct. 1996, pp. 25–39.
- [13] ———, "End-to-end internet packet dynamics," in *Proc. ACM SIGCOMM'97*, Cannes, France, Sept. 1997, pp. 139–154.

- [14] J. Postel, Transmission control protocol; rfc-793, Sept. 1979.
- [15] W. R. Stevens, *TCP Slow Start, Congestion Avoidance, Fast Retransmit, and Fast Recovery Algorithms*, rfc 2001, Jan. 1997.
- [16] R. W. Wolff, *Stochastic Modeling and the Theory of Queues*. Englewood Cliffs, NJ: Prentice Hall, 1989.

Jon C. R. Bennett (M'95–A'96) received the B.S. degree (with research honors) in applied mathematics/computer science from Carnegie Mellon University, Pittsburgh, PA, in 1992, and is working toward the Ph.D. degree at Harvard University, Cambridge, MA.

He is currently Chief Engineer at RiverDelta Networks, Melrose, MA, responsible for overall QoS design and implementation issues. Prior to joining RiverDelta Networks, he was a Senior Systems Architect with Xylan Corporation, where he was an Architecture Board member, responsible for advanced switch and router design, as well as overall architectural vision. Prior to that, he was a Senior Network Scientist at BBN Technologies, Cambridge, MA, engaged in research on advanced internetworking technologies. He came to BBN from FORE Systems, where was as a Systems Architect, involved with the design and implementation of advanced ATM switches. In addition to his design roles, he served as FORE Systems Representative to the traffic management working groups of both the ATM Forum and the IETF. He holds seven patents in the areas of ATM switch design, high speed scheduling, and hardware-sorting algorithms.

Dr. Bennett is a member of ACM SIGCOMM, the IEEE Communications Society, and is on the editorial board of *IEEE Network Magazine*.



Craig Partridge (M'88–SM'91–F'99) received the A.B., M.Sc., and Ph.D. degrees from Harvard University, Cambridge, MA.

He is a Chief Scientist with BBN Technologies, Cambridge, MA, where he conducts research on high-speed and high-performance networking. He is also a Consulting Associate Professor at Stanford University, Stanford, CA.

Dr. Partridge is the Chair of the ACM Special Interest Group on Data Communication (ACM SIGCOMM).

Nicholas Shectman received the A.B. degree from Harvard College, Cambridge, MA, in 1992.

He is a Scientist in the Internetwork Research Department at BBN Technologies, Cambridge, MA, where his research interests include high-speed internet-router design, transport protocol enhancements, and theoretical research on quality-of-service. Prior to joining BBN, he was an Engineer with EVI in Columbia, MD, involved with advanced wireless communication.

Mr. Shectman is a member of ACM.