

Internet Engineering Task Force
INTERNET-DRAFT
Intended status: Informational
Expires: 4 January 2010

S. Floyd
ICIR
A. Arcia
D. Ros
TELECOM Bretagne
J. Iyengar
Franklin & Marshall College
4 July 2009

Adding Acknowledgement Congestion Control to TCP
draft-floyd-tcpm-ackcc-06.txt

Status of this Memo

This Internet-Draft is submitted to IETF in full conformance with the provisions of BCP 78 and BCP 79.

This document may contain material from IETF Documents or IETF Contributions published or made publicly available before November 10, 2008. The person(s) controlling the copyright in some of this material may not have granted the IETF Trust the right to allow modifications of such material outside the IETF Standards Process. Without obtaining an adequate license from the person(s) controlling the copyright in such materials, this document may not be modified outside the IETF Standards Process, and derivative works of it may not be created outside the IETF Standards Process, except to format it for publication as an RFC or to translate it into languages other than English.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

The list of current Internet-Drafts can be accessed at
<http://www.ietf.org/ietf/lid-abstracts.txt>.

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>.

This Internet-Draft will expire on 4 January 2010.

Copyright Notice

Copyright (c) 2009 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents in effect on the date of publication of this document (<http://trustee.ietf.org/license-info>). Please review these documents carefully, as they describe your rights and restrictions with respect to this document.

Abstract

This document describes a possible congestion control mechanism for acknowledgement traffic (ACKs) in TCP. The document specifies an end-to-end acknowledgement congestion control mechanism for TCP that uses participation from both TCP hosts, the TCP data sender and the TCP data receiver. The TCP data sender detects lost or ECN-marked ACK packets, and tells the TCP data receiver the ACK Ratio R to use to respond to the congestion on the reverse path from the data receiver to the data sender. The TCP data receiver sends roughly one ACK packet for every R data packets received. This mechanism is based on the acknowledgement congestion control in DCCP's CCID 2. This acknowledgement congestion control mechanism is being specified for further evaluation by the network community.

Table of Contents

1. Introduction	7
2. Conventions and Terminology.	8
3. Overview	9
4. Acknowledgement Congestion Control	10
4.1. The ACK Congestion Control Permitted Option	10
4.2. The TCP ACK Ratio Option.	11
4.3. The Receiver: Implementing the ACK Ratio.	11
4.4. The Sender: Determining Lost or Marked ACK Packets.	12
4.4.1. The Sender: Detecting Lost ACK Packets After a Congestion Event.	14
4.5. The Sender: Adjusting the ACK Ratio	14
4.5.1. Possible Addition: Decreasing the ACK Ratio after a Congestion Window Decrease.	15
4.6. The Receiver: Sending ACKs for Out-of-Order Data Segments	16
4.7. The Sender: Response to ACK Packets	17
4.8. Possible Addition: Receiver Bounds on the ACK Ratio.	18
5. Possible Complications	19
5.1. Possible Complication: Delayed Acknowledgements	19
5.2. Possible Complication: Duplicate Acknowledgements.	19
5.3. Possible Complication: Two-Way Traffic.	19
5.4. Possible Complication: Reordering of ACK Packets.	20
5.5. Possible Complication: Abrupt Changes in the ACK Path.	20
5.6. Possible Complication: Corruption.	20
5.7. Possible Complication: ACKs That Don't Contribute to Congestion.	21
5.8. Possible Complication: TCP Implementations that Skip ACK Packets	24
5.9. Possible Complication: Router or Middlebox-based ACK Mechanisms	24
5.10. Possible Complication: Data-Limited Senders.	25
5.11. Other Issues	25
6. Evaluating ACK Congestion Control.	25
6.1. Contention in Wireless Links or in non-switched Ethernet	26
6.2. Keep-alive and Other Special ACK Packets.	26
7. Measurements of ACK Traffic and Congestion	26
8. Acknowledgement Congestion Control in DCCP's CCID 2.	27
9. Security Considerations.	27

10. IANA Considerations	28
11. Conclusions	29
12. Acknowledgements.	29
A. Appendix: Related Work	29
A.1. ECN-only Mechanisms	30
A.2. Receiver-only Mechanisms.	30
A.3. Middlebox-based Mechanisms.	31
B. Appendix: Design Considerations.	31
B.1. The TCP ACK Ratio Option, or an AckNow bit in data packets?.	31
Normative References (if standardized).	32
Informative References.	33
Acknowledgement	36

TO BE DELETED BY THE RFC EDITOR UPON PUBLICATION:

Changes from draft-floyd-tcpm-ackcc-05.txt:

- * Added a paragraph about Experimental TCP Option numbers. IESG Feedback from Lars Eggert.
- * Added a subsection on "Possible Complication: Data-Limited Senders". Feedback from Ilpo Jarvinen.
- * Added a sentence clarifying the ACK Congestion Control Permitted option.
- * Added a subsection on "The Sender: Detecting Lost ACK Packets After a Congestion Event"

Changes from draft-floyd-tcpm-ackcc-05.txt:

- * Minor editing.

Changes from draft-floyd-tcpm-ackcc-04.txt:

- * Changed desired status of document from Experimental to Informational, with associated editing changes.
- * Specified that ACK packets are only sent as ECN-Capable if ECN-capability has been negotiated as specified in RFC 3168.
- * Added a section on "Possible Addition: Decreasing the ACK Ratio after a Congestion Window Decrease".
- * Minor editing. Feedback from Alfred Hoenes.

Changes from draft-floyd-tcpm-ackcc-03.txt:

- * General editing. Feedback from Alfred Hoenes.
- * General editing. Feedback from Gorry Fairhurst.
- * Added a subsection on "Contention in Wireless Links or in non-switched Ethernet".
- * Modified Section 4.4 to say that the sender doesn't try to detect lost ACK packets during loss events even if the ACK Ratio is two.
- * Added a paragraph to Section 4.4 about "Detecting lost ACK packets after changes in the ACK Ratio".

Changes from draft-floyd-tcpm-ackcc-02.txt:

- * Cited RFC 3449 (PILC), RFC 3135 (PILC), and RFC 2760 (TCPSAT). From December 2007 Working Group meeting.
- * Added a note about the problem of effective ACK Congestion Control for environments with systematic reordering in the data path.
- * General editing. Feedback from Alfred Hoenes.
- * Added more about keep-alive packets and window update packets. Feedback from Anantha Ramaiah.
- * Clarified that SACK "SHOULD" be used. We don't know enough to say "MUST".

Changes from draft-floyd-tcpm-ackcc-01.txt:

- * Added a section on "Keep-alive Packets". Feedback from Anantha Ramaiah.
- * Added a section on "Possible Complication: TCP Implementations that Skip ACK Packets". Motivated by reports at IETF that many high-bandwidth TCPs don't follow the MUST of sending an ACK for every other packet, if they don't have time.
- * Added that receivers might have buffer limitations that require that they ack at least every K packets, for some K. Feedback from Sara Landstrom.
- * Added to the discussion of "Possible Complication: Two-Way Traffic". Feedback from Sara Landstrom.
- * Added a section on "Possible Complication: Router or Middlebox-based ACK Mechanisms". Feedback from Sara Landstrom.
- * Added that SACK is required with ACK congestion control. Feedback from Sara Landstrom.
- * Added a discussion of "Reducing the TCP Acknowledgment Frequency" to the related work section.
- * Moved the Related Work section to the appendix. Feedback from Alfred Hoenes.
- * General editing from feedback from Alfred Hoenes.

- * Added an appendix on "Design Considerations", with a subsection on "The TCP ACK Ratio Option, or an AckNow bit in data packets?".

Changes from draft-floyd-tcpm-ackcc-00.txt:

- * Added a discussion of environments where the reverse path is congested, but the TCP ACK traffic does not significantly contribute to that congestion. In this case, the goal is to minimize the negative impact of AckCC on TCP performance. Feedback from Armando Caro.
- * In Section 5.7, added that when ABC is used with Aggregate Congestion Control, and rate-based pacing is also used, the sender MAY increase cwnd by more than 2 MSS. Feedback from Armando Caro.
- * Added a section about measurements of ACK traffic and congestion. Feedback from Armando Caro.
- * Added a section on the possibility of a TCP receiver-imposed lower bound on the ACK Ratio. Suggested by Mark Allman.
- * Added to the discussion of the minimum ACK sending rate. Suggested by Mark Allman.
- * Added a note that if the TCP receiver doesn't send an ACK for every duplicate data packet, the sender's Fast Recovery procedure will have to be modified to take this into account. Feedback from Mark Allman.
- * Added a discussion of evaluating ACK Congestion Control. From feedback from Mark Allman.
- * Some general editing in response to feedback from Mark Allman.

END OF SECTION TO BE DELETED.

1. Introduction

This document describes a congestion control mechanism for acknowledgements (ACKs) to TCP. This mechanism is based on the acknowledgement congestion control in DCCP's CCID 2 [RFC4340], [RFC4341], which is a successor to the TCP acknowledgement congestion control mechanism proposed by Balakrishnan et al. in [BPK97].

In this document we use the terminology of senders and receivers, with the sender sending data traffic, and the receiver sending acknowledgement traffic in response. In CCID 2's acknowledgement

congestion control, specified in Section 6.1 of [RFC4341], the receiver uses an ACK Ratio R reported to it by the sender, sending roughly one ACK packet for every R data packets received. The CCID 2 sender keeps the acknowledgement rate roughly TCP friendly by monitoring the acknowledgement stream for lost and marked ACK packets and modifying the ACK Ratio accordingly. For every RTT containing an ACK congestion event (that is, a lost or marked ACK packet), the sender halves the acknowledgement rate by doubling the ACK Ratio; for every RTT containing no ACK congestion event, the sender additively increases the acknowledgement rate through gradual decreases in the ACK Ratio.

The goal of this document is to explore a similar congestion control mechanism for acknowledgement traffic for TCP. The assumption is that in some environments with congestion on the reverse path, reducing the sending rate for ACK traffic traversing the congested path can help to reduce the congestion itself. For those environments where the reverse path is congested but where TCP ACK traffic does not appreciably contribute to that aggregate congestion, the goal is for TCP's ACK congestion control to have a minimal negative effect on the performance of the TCP connection.

Adding acknowledgement congestion control as an option in TCP would require the following:

- * An agreement from the TCP hosts on the use of ACK congestion control. For the mechanism specified in this document, the TCP hosts would use a new TCP option, the ACK Congestion Control Permitted Option.
- * A mechanism for the TCP sender to detect lost and ECN-marked pure acknowledgement packets.
- * A mechanism for adjusting the ACK Ratio. The TCP sender would adjust the ACK Ratio as specified in Section 6.1.2 of [RFC4341].
- * A method for the TCP sender to inform the TCP receiver of a new value for the ACK Ratio. For the mechanism specified in this document, the TCP sender would use a new TCP option, the ACK Ratio Option.

2. Conventions and Terminology

MSS refers to the Maximum Segment Size.

3. Overview

This section gives an overview of acknowledgement congestion control for TCP.

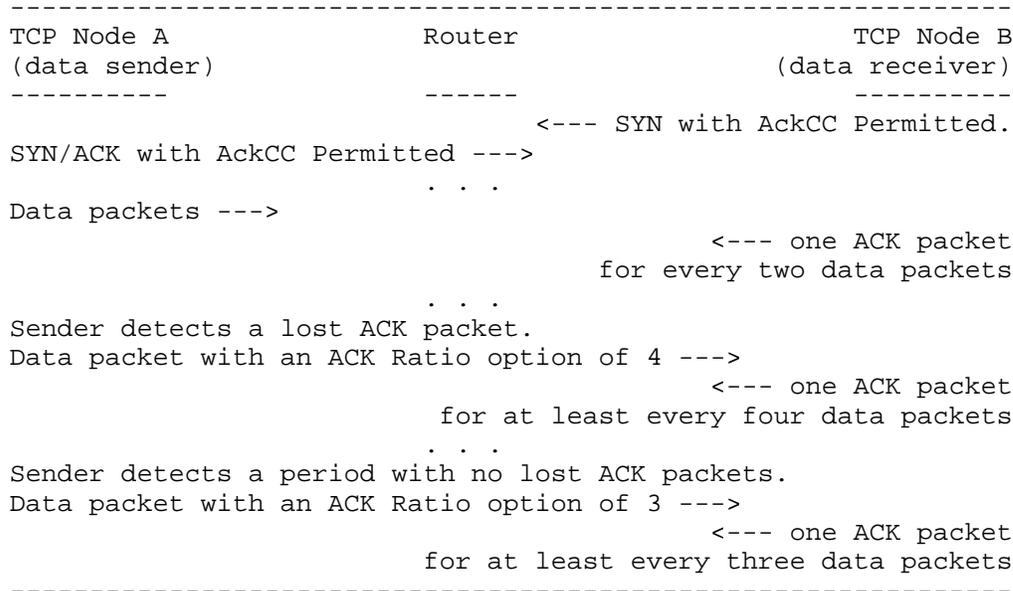


Figure 1: Acknowledgement Congestion Control, Node B as the connection initiator, for a connection without ECN.

Figure 1 gives an example of Acknowledgement Congestion Control (AckCC) with TCP Node B as the connection initiator.

During connection initiation, TCP host B sends an ACK Congestion Control Permitted option on its SYN or SYN/ACK packet. This allows TCP host A (now called the sender) to send instructions to TCP host B (now called the receiver) about the ACK Ratio to use in responding to data packets.

Also during connection initiation, TCP host A sends an ACK Congestion Control Permitted option on its SYN or SYN/ACK packet. In combination with TCP host B's sending of an ACK Congestion Control Permitted option, and with the negotiation of ECN-Capability as specified in RFC 3168, this would allow TCP host B to send its ACK packets as ECN-Capable.

The TCP receiver starts with an ACK Ratio of two, generally sending one ACK packet for every two data packets received.

The TCP sender detects a lost or ECN-marked ACK packet from the TCP receiver, and sends an ACK Ratio option of four to the receiver. The TCP receiver changes to an ACK Ratio of four, sending one ACK packet for at most four data packets. The TCP sender uses Appropriate Byte Counting and rate-based pacing in responding to these ACK packets.

The TCP sender detects a period with no lost ACK packets, and sends an ACK Ratio option of three to the TCP receiver. The TCP receiver changes back to an ACK Ratio of three, sending one ACK packet for at most three data packets.

4. Acknowledgement Congestion Control

The goal of the mechanism proposed in this document is to control pure ACK traffic on the path from the TCP data receiver to the TCP data sender. Note that the approach outlined here is an end-to-end one (as is the approach followed by DCCP's CCID 2 [RFC4341]), but it may also take advantage of explicit congestion information from the network conveyed by ECN [RFC3168], if available. The ECN specification [RFC3168, section 6.1.4] prohibits a TCP receiver from setting the ECT(0) or ECT(1) codepoints in IP packets carrying pure ACKs, but *only* as long as the receiver does *not* implement any form of ACK congestion control. Unlike some of the related work cited in the appendix, in this document we are proposing an end-to-end ACK congestion control mechanism that controls congestion on the reverse path (the path followed by the ACK traffic) by detecting and responding either marked or dropped ACK packets.

4.1. The ACK Congestion Control Permitted Option

The TCP end-points would negotiate the use of ACK Congestion Control (ACKCC) with a TCP option, the ACK Congestion Control Permitted Option. The option number would be allocated by IANA.

The ACK Congestion Control Permitted option can only be sent on packets that have the SYN bit set. If TCP end-point A receives an ACK Congestion Control Permitted option from TCP end-point B, then the TCP end-points may use ACK Congestion Control on the pure acknowledgements sent from B to A. This means that TCP end-point A may send ACK Ratio values to TCP end-point B, for TCP end-point B to use on pure acknowledgement packets. Equivalently, if TCP end-point A *does not* receive an ACK Congestion Control Permitted option from TCP end-point B, then TCP end-point A knows not to waste its time detecting lost ACK packets and adjusting and sending the ACK Ratio values.

If TCP end-point B receives an ACK Congestion Control Permitted option from TCP end-point A, then the TCP end-points may use ACK Congestion Control on the pure acknowledgements sent from A to B.

If TCP end-point B receives an ACK Congestion Control Permitted option from TCP end-point A and also sent an ACK Congestion Control Permitted option to TCP end-point A, and ECN-capability has been negotiated, then TCP end-point B can send its pure ACK packets as ECN-Capable.

TCP ACK Congestion Control Permitted Option:

Kind: TBD1

```
+-----+-----+
| Kind=TBD1 | Length=2 |
+-----+-----+
```

When ACK Congestion Control is used, the default initial ACK Ratio is two, with the receiver acknowledging at least every other data packet.

4.2. The TCP ACK Ratio Option

The sender uses a ACK Ratio TCP Option to communicate the ACK Ratio value from the sender to the receiver.

TCP ACK Ratio Option:

Kind: TBD2

```
+-----+-----+-----+
| Kind=TBD2 | Length=3 | ACK Ratio |
+-----+-----+-----+
```

The ACK Ratio Option is only sent on data packets. Because TCP uses reliable delivery for data packets, the TCP sender can tell if the TCP receiver has received an ACK Ratio Option.

4.3. The Receiver: Implementing the ACK Ratio

With an ACK Ratio of R, the receiver should send one pure ACK for every R newly received data packets unless the delayed ACK timer expires first. A receiver could simply maintain a counter that increments by one for each new data packet received, and send an ACK packet when the counter reaches R. The receiver would reset the

counter to zero whenever a pure or piggybacked ACK is sent,

If the receiver has buffer limitations, the receiver might have to acknowledge K packets, for some K less than the current ACK Ratio R . In this case, the sender could observe from the acknowledgements that the receiver is acknowledging less than R packets.

It is possible for there to be lost or marked ACK packets when there haven't yet been any lost or marked data packets. Thus, the sender could increase the ACK Ratio R even during the initial slow-start.

[RFC2581] recommends that the receiver SHOULD acknowledge out-of-order data packets immediately, sending an immediate duplicate ACK when it receives a data segment above a gap in the sequence space, and sending an immediate ACK when it receives a data segment that fills in all or part of a gap in the sequence space.

When ACK Congestion Control is being used and the ACK Ratio is at most two, the TCP receiver acknowledges each out-of-order data packet immediately. For an ACK Ratio greater than two, Section 4.6 specifies in detail the receiver's behavior for sending ACKs for out-of-order data packets.

4.4. The Sender: Determining Lost or Marked ACK Packets

The TCP data sender uses its knowledge of the ACK Ratio in use by the receiver to infer when an ACK packet has been lost.

Because the TCP sender knows the ACK Ratio R in use by the receiver, the TCP sender knows that in the absence of dropped or reordered acknowledgement packets, each new acknowledgement received will acknowledge at most R additional data packets. Thus, if the sender receives an acknowledgement acknowledging more than R data packets, and does not receive a subsequent acknowledgement acknowledging a strict subset (with a smaller cumulative acknowledgement, or with the same cumulative acknowledgement but a strict subset of data acknowledged in SACK blocks), then the sender can infer that an ACK packet has been dropped. The use of SACK options in ACK packets would help the sender in detecting lost ACK packets.

Similarly, the TCP sender knows that in the absence of dropped or delayed data packets from the sender, and in the absence of delayed acknowledgements due to a timer expiring at the receiver, each new pure acknowledgement received will acknowledge at least R additional data packets. In terms of ACK congestion control, the TCP sender does not have to take any actions when it receives an acknowledgement acknowledging less than R additional packets.

Out-of-order data packets: If the ACK Ratio is at most two, then the TCP receiver sends a DupACK for every out-of-order data packet. In this case, the TCP sender should be able to detect lost DupACK packets by counting the number of DupACKs that arrive between the beginning of the loss event and the arrival of the first full or partial ACK, and comparing this number with the number of DupACKs that should have arrived (based on the number of packets being ACKed by the full or partial ACK). Simulations and/or experiments will be needed to determine whether, in practice, it works for the TCP sender to assess lost ACK packets during loss events, for an ACK Ratio of at most two.

If the ACK Ratio is greater than two, the TCP receiver does not send a DupACK for every out-of-order data packet, as specified in Section 4.6. For simplicity, the TCP sender does not attempt to detect lost ACK packets during loss events involving forward-path data traffic. That is, as soon as the sender infers a packet loss for a forward-path data packet, it stops detection of ACK loss on the reverse path. The sender waits until a new cumulative acknowledgement is received that covers the retransmitted data, and then restarts detection of ACK loss for reverse-path traffic.

Detecting lost ACK packets after changes in the ACK Ratio: In detecting lost ACK packets, the sender relies on its knowledge of the ACK Ratio used by the receiver. But when the sender makes a change in the ACK Ratio, and then receives ACK packets, how does the sender know whether the receiver was using the new or the old ACK Ratio when it sent those ACK packets? As specified in the next section, the sender can make only one of two possible changes to the ACK Ratio within one round-trip time. The sender can decrease the ACK Ratio by one, from R to $R-1$, or the sender can double the ACK Ratio, increasing it from R to $2R$. But, in detecting lost ACK packets after an increase in the ACK Ratio, the sender needs to know whether the receiver was using the old ACK Ratio R or the new ACK Ratio $2R$.

The sender sends ACK Ratio Options only on data packets, and these data packets are acknowledged by the receiver. One possibility would be for the sender to save the sequence number of the last data packet that contained an ACK Ratio Option, and to remember whether that ACK Ratio Option was for an increase or a decrease in the ACK Ratio. Then, if the sender receives an ACK packet acknowledging the saved sequence number, then the sender knows that the receiver has begun using the new ACK Ratio.

It *might* be sufficient for the sender just to save the information of whether the last change in the ACK Ratio was an increase or a decrease, without saving the sequence number associated with the last ACK Ratio Option. In this way, if the sender recently increased the

ACK Ratio from R to $2R$, the sender could be more cautious in detecting lost ACK packets. Another possibility would be that after sending an ACK Ratio Option, the sender waits until that data has been ACKed, with the new ACK Ratio in use by the receiver, before resuming the detection of lost ACK packets. However, we do not explore either of these approaches in more detail in this document.

4.4.1. The Sender: Detecting Lost ACK Packets After a Congestion Event

After a sender's retransmit timeout or fast retransmit, the sender might retransmit a number of data packets dropped from a single window of data. In particular, during a loss recovery period (from the sender's detection of the congestion event up until the sender receives an acknowledgement of all data packets transmitted before the loss recovery period began) retransmitted data packets can fill holes in the receiver's sequence space, resulting in irregular jumps in the cumulative acknowledgement field in ACK packets from the receiver. These jumps in the cumulative acknowledgement field make it difficult for the sender to reliably detect lost ACK packets during a loss recovery period.

Because of this uneven progress of the cumulative acknowledgement field during a loss recovery period, the sender should not attempt to detect lost ACK packets during a loss recovery period. As a consequence, the sender will not increase the ACK Ratio in response to ACK packets that are lost during a loss recovery period.

4.5. The Sender: Adjusting the ACK Ratio

The TCP sender will adjust the ACK Ratio as specified in Section 6.1.2 of [RFC4341], as follows.

The ACK Ratio always meets the following three constraints.

- (1) The ACK Ratio is an integer.
- (2) The minimum ACK sending rate: The ACK Ratio does not exceed $\max(2, \text{cwnd}/(K \cdot \text{MSS}))$, rounded up, for $K=2$. As a result, the TCP receiver generally sends at least two ACKs in response to a window of at least four full-sized segments.
- (3) If the congestion window is at least as large as four full-sized segments, then the ACK Ratio is at least two. In other words, an ACK Ratio of one is only allowed when the congestion window is at most three full-sized segments.

The sender changes the ACK Ratio within those constraints as follows. For each congestion window of data with lost or marked ACK packets,

the ACK Ratio R is doubled; and for each $cwnd/(MSS*(R^2 - R))$ consecutive congestion windows of data with no lost or marked ACK packets, the ACK Ratio is decreased by 1. (See Appendix A of RFC 4341 for the derivation. Note that Appendix A of RFC 4341 assumes a congestion window W in packets, while we use $cwnd$ in bytes.) As stated in the previous section, when the ACK Ratio is greater than two the sender does not attempt to detect lost ACK packets during loss events for forward-path traffic.

For a constant congestion window, these modifications to the ACK ratio give an ACK sending rate that is roughly TCP friendly. Of course, $cwnd$ usually varies over time; the dynamics will be rather complex, but roughly TCP friendly. We recommend that the sender determines when to decrease the ACK Ratio by one (i.e., by calculating the number of in-order data packets to count) right after an ACK loss event.

The frequency of ACK Ratio negotiations: The sender need not keep the ACK Ratio completely up to date. For instance, it may rate-limit ACK Ratio renegotiations to once every four or five round-trip times, or to once every second or two. The sender should not attempt to change the ACK Ratio more than once per round-trip time. In particular, before sending a packet with a new value for the ACK Ratio, the sender should verify that the receiver has acknowledged a data packet containing an ACK Ratio Option for the old value of the ACK Ratio. Additionally, the sender may enforce a minimum ACK Ratio of two, or it may set the ACK Ratio to one for half-connections with persistent congestion windows of 1 or 2 packets.

The minimum ACK sending rate: From rule (2) above, the TCP receiver always sends at least $K=2$ ACKs for a window of data, even in the face of very heavy congestion on the reverse path. We would note, however, that if congestion is sufficiently heavy, all the ack packets are dropped, and then the sender falls back on an exponentially backed-off timeout. Thus, if congestion is sufficiently heavy on the reverse path, then the sender reduces its sending rate on the forward path, which reduces the rate on the reverse path as well. One possibility would be to use a higher minimum ACK sending rate, adding a constant upper bound on the ACK Ratio. That is, if the ACK Ratio also had an upper bound of J , independent of $cwnd$, then the receiver would always send at least one ACK for every J data packets, regardless of the level of congestion on the reverse path.

4.5.1. Possible Addition: Decreasing the ACK Ratio after a Congestion Window Decrease

After a lost or ECN-marked data packet, the data sender halves the

congestion window, thus halving the sending rate for data packets, while making no change to the ACK Ratio R. As a result, after a congestion event involving a data packet, the sending rate for ACK packets on the return path is also halved. If the congestion event was a lost or ECN-marked data packet, this was due to congestion on the forward path, which may have been unrelated to conditions on the reverse path. Thus, it has been suggested that the sender could decrease the ACK Ratio R when it halves the congestion window; in this case, the halving of the sending rate for data packets would not be accompanied by a halving of the sending rate for ACK packets also.

However, there are a few cases where a congestion event involving data packets could in fact have been caused by congestion on the reverse path. As one example, the path could include a congested multiaccess link where forward-path and reverse-path traffic can interfere with each other. Thus, in this case it might be desirable if a congestion event resulted in a reduction in the sending rate of ACK packets as well as of data packets.

As a second example of a congestion event involving congestion of the reverse path, a congestion event could be caused not by a dropped or ECN-marked data packet, but by a window of dropped ACK packets, resulting in a retransmit timeout at the data sender. After a retransmit timeout, the TCP sender will slow-start, reducing the congestion window to the initial window, and setting the ACK Ratio to at most two.

Until further investigation, the sender will not decrease the ACK Ratio as a result of a congestion event involving a data packet.

4.6. The Receiver: Sending ACKs for Out-of-Order Data Segments

RFC 2581 says that "a TCP receiver SHOULD send an immediate duplicate ACK when an out-of-order segment arrives." After three duplicate ACKs are received, the TCP sender infers a packet loss and implements Fast Retransmit and Fast Recovery, retransmitting the missing packet. When the ACK Ratio is at most two, the TCP receiver should still send an immediate duplicate ACK when an out-of-order segment arrives.

In general, when the ACK Ratio is greater than two, the TCP receiver still should send an immediate duplicate ACK for each of the first three out-of-order segments that arrive in a reordering event. (We define a reordering event at the receiver as beginning when an out-of-order segment arrives, and ending when the receiver holds no more out-of-order segments.) However, when the ACK Ratio is greater than two, after the first three duplicate ACKs have been sent, the TCP receiver should perform ACK congestion control on the remaining ACKs

to be sent during the current reordering event. That is, after the first three duplicate ACKs have been sent, the TCP receiver should return to sending an ACK for every R segments, instead of sending an ACK for every out-of-order segment in that reordering event. [We note that the Fast Recovery procedure of the TCP sender might have to be modified to take this change into account.] In addition, a receiver must not withhold an ACK for more than 500 ms.

We note that in an environment with systematic reordering in the data path (e.g., every set of K data packets arrives in inverted order, for some value of K), the guideline above could result in the receiver sending an ACK for every data packet, regardless of the ACK Ratio. In such an environment with persistent reordering, the receiver may decide not to send an immediate duplicate ACK for each of the first three out-of-order segments that arrive in a reordering event. We leave the investigation of mechanisms for effective ACK Congestion Control in environments with systematic reordering for future work.

4.7. The Sender: Response to ACK Packets

The use of a large ACK Ratio can generate line rate data bursts at a TCP sender. When the ACK Ratio is greater than two, the TCP sender should use some form of burst mitigation, or rate-based pacing for sending data packets in response to a single acknowledgement. The use of rate-based pacing will be limited by the timer granularity at the TCP sender.

We note that the interaction of ACK congestion control and burst mitigation schemes needs further study.

Byte counting at the sender: In addition to the impact of a large ACK Ratio on the burstiness of the TCP sender's sending rate, a large ACK Ratio can also affect the data sending rate by slowing down the increase of the congestion window *cwnd*. As specified in RFC 2581, in slow-start the TCP sender increases *cwnd* by one full-sized segment for each new ACK received (in this context, a "new ACK" is an ACK that acknowledges new data). RFC 2581 also specifies that in congestion avoidance, the TCP sender increases *cwnd* by roughly $1/cwnd$ full-sized segments for each ACK received, resulting in an increase in *cwnd* of roughly one full-sized segment per round-trip time. In this case, the use of a large ACK Ratio would slow down the increase of the sender's congestion window.

RFC 2581 notes that during congestion avoidance it is also acceptable to count the number of bytes acknowledged by new ACKs, and to increase *cwnd* based on the number of bytes acknowledged, rather than on the number of new ACKs received. Thus, the sender should use this

form of byte counting with Acknowledgement Congestion Control, so that the Acknowledgement Congestion Control doesn't slow down the window increases for the data traffic sent by the sender. Because rate-based pacing should be used with Acknowledgement Congestion Control, as recommended earlier in this section, the TCP sender may increase the congestion window by more than two MSS for each ACK.

We note that for Appropriate Byte Counting (ABC) as specified in [RFC3465], during Slow-Start the sender is allowed to increase the congestion window by at most two MSS for each ACK. It has not yet been determined whether, with Acknowledgement Congestion Control, the TCP sender could use ABC during Slow-Start. If ABC is used with Acknowledgement Congestion Control, then when the TCP sender is in slow-start and the ACK Ratio is greater than two, the TCP sender may increase the congestion window by more than two MSS in response to a single ACK. Section 4.2 of [LL07] explores some of the issues with the use of ABC for TCP connections with a fixed ACK Ratio greater than two.

Inferring lost data packets: As cited earlier, RFC 2581 infers that a packet has been lost after it receives three duplicate acknowledgements. Because ACK Congestion Control is only used when there is congestion on the reverse path, after a packet loss one or more of the three duplicate ACKs sent by the receiver could be lost on the reverse path, and the receiver might wait until it has received R more out-of-order segments before sending the next duplicate ACK. All this could slow down Fast Recovery and Fast Retransmit quite a bit. The use of SACK can help reduce the potential delay in detecting a lost packet. With SACK, a TCP sender can use the information in the SACK option to detect when the receiver has received at least three out-of-order data packets, and to initiate Fast Retransmit and Fast Recovery in this case, even if the TCP sender has not yet received three dup ACKs.

4.8. Possible Addition: Receiver Bounds on the ACK Ratio

It has been suggested that in some environments, the TCP receiver might want to set lower bounds on the ACK Ratio. For example, the TCP receiver might know from configuration or from past experience that the bandwidth on the return path is limited, and might want to set a lower bound (greater than two) on the ACK Ratio R. If this is included, this would require a TCP Option from the TCP receiver to the TCP sender reporting the lower bound on the ACK Ratio. Care would also be needed so that the lower bound on the ACK Ratio was only in effect when the TCP sender's congestion window was sufficiently high.

5. Possible Complications

5.1. Possible Complication: Delayed Acknowledgements

The receiver could send a delayed acknowledgement acknowledging a single packet, even when the ACK Ratio is two or more.

This should not cause false positives (when the TCP sender infers a loss when no loss happened). The TCP sender only infers that a pure ACK packet has been lost when no data packet has been lost, and an ACK packet arrives acknowledging more than R new packets.

Delayed acknowledgements could, however, cause false negatives, with the TCP sender unable to detect the loss of an ACK packet sent as a delayed acknowledgement. False negatives seem acceptable; this would result in approximate ACK congestion control, which would be better than no ACK congestion control at all. In particular, when this form of false negative occurs, it is because the receiver is sending acknowledgements at such a low rate that it is sending delayed acknowledgements, rather than acknowledging at least R data packets with each acknowledgement.

5.2. Possible Complication: Duplicate Acknowledgements.

As discussed in Section 4.3, RFC 2581 states that "a TCP receiver SHOULD send an immediate duplicate ACK when an out-of-order segment arrives", and that "a TCP receiver SHOULD send an immediate ACK when the incoming segment fills in all or part of a gap in the sequence space" [RFC2581]. When ACK Congestion Control is used, the TCP receiver instead uses the guidelines from Section 4.6 to govern the sending of duplicate ACKs. More work would be useful to evaluate the advantages and disadvantages of this approach in terms of the potential delay in triggering Fast Retransmit, and to explore alternate possibilities.

5.3. Possible Complication: Two-Way Traffic.

In a TCP connection with two-way traffic, the receiver could send some pure ACK packets, and some acknowledgements piggy-backed on data packets. The receiver would still follow the rule of only sending a pure ACK packet when there is a need for a delayed ACK, or there are R new data packets to acknowledge.

In a connection with two-way traffic, the TCP sender would not always be able to infer when a pure ACK packet had been lost. For example, the receiver could send a pure ACK packet acknowledging packet K, and soon afterwards, the receiver could send a newly-generated data packet for the reverse-path flow also acknowledging packet K. The

pure ACK packet could be dropped in the network, and the sender would not be able to detect this drop.

Fortunately, there are limitations to the potential problems caused by undetected ACK losses in two-way traffic. The sender will only fail to detect the loss of a pure ACK packet if the ACK packet was followed by a data packet with the same acknowledgement number. If the reverse-path traffic for the connection is dominated by data traffic, then the congestion control for the data traffic is more important than the congestion control for the pure ACK traffic. If the reverse-path traffic is dominated by pure ACK traffic, then the sender would detect any losses of pure ACK packets followed by other pure ACK packets, and this would include most of the pure ACK packets for that connection. Thus, the sender's failure to detect the loss of a pure ACK packet followed by a data packet with the same acknowledgement number would not disable acknowledgement congestion control for a TCP connection with two-way traffic.

5.4. Possible Complication: Reordering of ACK Packets.

It is possible for ACK packets to be reordered on the reverse path. The TCP sender could either use a parallel mechanism to the DupACK threshold to infer when an ACK packet has been lost, as with TCP, or, more robustly, the TCP sender could wait an entire round-trip time before inferring that an ACK packet has been lost [RFC4653].

5.5. Possible Complication: Abrupt Changes in the ACK Path.

What happens when there are abrupt changes in the reverse path, such as from vertical handovers? Can there be any problems that would be worse than those experienced by a TCP connection that is not using ACK congestion control?

5.6. Possible Complication: Corruption.

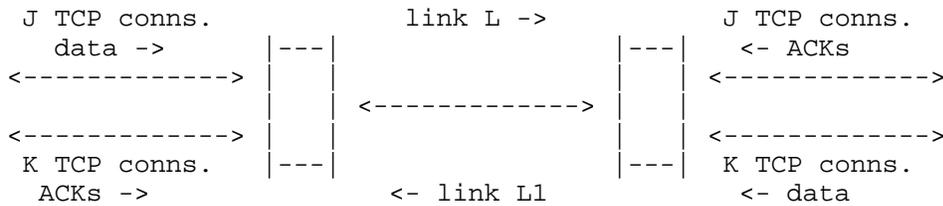
As with data packets, it is possible for ACK packets to be dropped in the network due to corruption rather than congestion. The current assumption of ACK congestion control is that all losses should be taken as indications of congestion. If there is ever some better mechanism for identifying and responding to corrupted TCP data packets, the same solution hopefully would apply to corrupted ACK packets as well.

One problem with the interaction of packet corruption and congestion control, for both data and ACK packets, is that it is not always obvious when the packet corruption is related to congestion, and when the packet corruption is independent of the level of congestion on the corrupting link. In environments where packet corruption exists

and is independent of the level of congestion on the corrupting link, applying ACK congestion control would only make the connection more sensitive to ACK packet corruption, by reducing the number of ACKs that are sent.

5.7. Possible Complication: ACKs That Don't Contribute to Congestion.

It is possible for the ACK packets in a TCP connection to traverse a congested path where ACK packets are dropped, but where the ACK packets themselves don't significantly contribute to the congestion on the path. In scenarios where ACK packets are dropped but where ACK traffic doesn't make a significant contribution of the congestion on the path, the use of ACK Congestion Control would not contribute to reducing the aggregate congestion on the path. In this case, one goal is to minimize the negative impact of ACK Congestion Control on the overall performance of the TCP connection.



A scenario with J forward and K reverse TCP connections.

To explore the relative contribution of ACK traffic on congestion, it is useful to consider a simple scenario with a congested unidirectional link L carrying data traffic from J TCP connections (the forward TCP connections) and ACK traffic from K TCP connections (the reverse TCP connections). We assume that all TCP connections have the same round-trip time R and the same data packet size S of 1500 bytes. We further assume that all of the forward TCP connections have the same data packet drop rate p and the same congestion window W, and that all of the reverse TCP connections have the same congestion window W1 and the same ACK packet drop rate p1. (The packet drop rate for data packets is defined as the fraction of arriving data packets that are dropped; similarly, the packet drop rate for ACK packets is the fraction of arriving ACK packets that are dropped.) The J TCP connections each use a bandwidth on link L of 1500*W/R bytes per second, and the K TCP connections, without ACK Congestion Control, each use a bandwidth on link L of 40*(W1/2)/R bytes per second. This gives a ratio of 75*(J/K)*(W/W1) for TCP data bandwidth to TCP ACK bandwidth on link L. The ratio J/K is the ratio between the number of forward and reverse TCP connections on link L,

and could have a wide range of values (e.g., large for an access link from a web server, and small for an access link to a web server). For this scenario, the ratio W/W_1 is largely a function of the different levels of congestion on the forward and reverse paths.

To explore the possibilities, we will consider some of the range of congestion control mechanisms for the congested link. First, we consider scenarios where the limitation on the congested path is in the link bandwidth in bytes per second.

Cases (1), (2), (3), (5), and (7) below represent the best scenarios for ACK Congestion Control, where the fraction of packet drops for TCP ACK packets roughly matches the TCP ACK packets' contribution to congestion. [In several of these cases this is at best a rough match because the data packets are a factor in the bandwidth and in the queue limitations, while the TCP ACK packets are only a factor in the queue limitations.] Cases (4) and (8) below represent problematic scenarios where the fraction of packet drops for TCP ACK packets is much higher than the TCP ACK packets' contribution to congestion (in terms of taking space in a congested queue, using scarce CPU cycles at the congested router, or using scarce bandwidth). Case (6) below represents scenarios where ACK Congestion Control would not be effective because it would not be invoked. In the scenarios in case (6), the fraction of packet drops for TCP ACK packets would be much smaller than the TCP ACK packets' contribution to congestion.

(1) The Drop-Tail queue for link L is measured in packets. In this case, the congested queue can accommodate N packets, regardless of packet size, there is a limitation of both bandwidth in bytes per second and also in queue space in packets, and large data packets and small TCP ACK packets should see similar packet drop rates. Although TCP ACK packets most likely aren't a major factor in the bandwidth limitation, they can be a significant contribution to the limitation of queue space. So, while the packet drop rate for ACK packets could be high in times of congestion, the ACK packets are contributing to that congestion somewhat by using scarce buffer space.

(2) The Drop-Tail queue is measured in bytes. In this case, the congested queue can accommodate M bytes of packets, and TCP ACK packets don't make a significant contribution to either the bandwidth limitation or to the limitation in queue space. It is also the case that in this scenario, even if there is heavy congestion, the packet drop rate for TCP ACK packets should be small (because small ACK packets can often find space on the congested queue when large data packets can't find space). In this case, ACK Congestion Control should not present any problems; the TCP ACK packets aren't contributing significantly to congestion, and aren't experiencing significant packet drop rates.

(3) The RED queue is in packet mode, and is measured in packets. This is similar to case (1) above. Because the queue is measured in packets, small TCP ACK packets contribute to the limitation in queue space, but not to the limitation in link bandwidth. Because the queue is in packet mode, large data packets and small TCP ACK packets should see similar packet drop rates.

(4) The RED queue is in packet mode, but is measured in bytes. Because the queue is measured in bytes, small TCP ACK packets don't contribute significantly to either the limitation in queue space or to the limitation in link bandwidth. Because the queue is in packet mode, large data packets and small TCP ACK packets should see similar packet drop rates. If it existed, this case would be problematic, because the TCP ACK packets would not be contributing significantly to the congestion, but they would see a similar packet drop rate as the large data packets that are contributing to congestion.

(5) The RED queue is in byte mode, and is measured in bytes. This is similar to case (2) above. Because the queue is measured in bytes, small TCP ACK packets don't contribute significantly to either the limitation in queue space or to the limitation in link bandwidth. At the same time, because the queue is in byte mode, small TCP ACK packets see much smaller packet drop rates than those of large data packets.

(6) The RED queue is in byte mode, but is measured in packets. Because the queue is measured in packets, small TCP ACK packets contribute to the limitation in queue space, but not to the limitation in link bandwidth. Because the queue is in byte mode, small TCP ACK packets see much smaller packet drop rates than those of large data packets. If this case existed, TCP ACK packets would contribute somewhat to congestion, but would see a much smaller packet drop rate than that of large data packets.

Next, we consider scenarios where the limitation on the congested link is in CPU cycles at the router in packets per second, not in bandwidth in bytes per second.

(7) The CPU load imposed by TCP ACK packets is similar to the load imposed by other packets (e.g., TCP data packets). ACK Congestion Control would be useful in this scenario, particularly if TCP ACK packets saw the same packet drop rates as TCP data packets.

(8) The CPU load imposed by TCP ACK packets is much less than the load imposed by other packets (e.g., TCP data packets). If TCP ACK packets saw a smaller packet drop rate than TCP data packets, then the TCP ACK packet drop rate would roughly match the TCP ACK packets' contribution to congestion, and this would be good. If TCP ACK

packets saw the same packet drop rate as TCP data packets, this case would be problematic, because the TCP ACK packets would not be contributing significantly to the congestion, but they would see a similar packet drop rate as the large data packets that are contributing to congestion.

5.8. Possible Complication: TCP Implementations that Skip ACK Packets

It has been reported in IETF meetings that current TCP implementations do not always acknowledge at least every other data packet, as required by the TCP specifications. In particular, it has been reported that if a TCP receiver receives many data packets in a burst, before it is able to send an acknowledgement, then it might send a single acknowledgement for the burst of packets. We note that such a behavior would cause complications for a TCP connection that used ACK congestion control, as the sender would not be able to determine when an ACK packet had been dropped in the network, and when the packet had been skipped by the receiver because it was processing a burst of data packet arrivals.

One possibility for addressing this problem would be for TCP receivers using ACK congestion control to be required to send an acknowledgement for each R packets, for ACK Ratio R. In this case, if the receiver received a large burst of data packets back-to-back, the receiver would be required to send a responding burst of ACK packets, one for each set of R data packets.

A second possibility for addressing this problem would be to define a TCP option or flag that the TCP receiver could use when sending an ACK packet to inform the sender that the TCP receiver 'skipped' some ACK packets, so that the sender should not infer ACK loss if some previous ACK packets seem to be missing.

Future work will explore the costs and benefits of these two approaches.

5.9. Possible Complication: Router or Middlebox-based ACK Mechanisms

One possible complication would be the interaction of ACK Congestion Control with router-based or middlebox-based ACK mechanisms such as ACK filtering along the reverse path [BPK97] [WWCM99] [BA03] [KLS07]. We are not aware of the deployment of ACK filtering in the Internet, but any testing of ACK congestion control would have to look for interactions with any middlebox-based mechanisms regarding ACK packets. In particular, we would consider interactions of ACK congestion control with the possible deployment of ACK filtering on satellite links, cable modems, or the like.

5.10. Possible Complication: Data-Limited Senders

The mechanism for adjusting the ACK Ratio is designed with the goal of having the TCP receiver send at least two ACKs in response to each window of at least four full-sized data packets. However, with ACK Congestion Control in combination with a data-limited sender, it is possible for the sender to send at least four full-sized data packets in a round-trip time, with the receiver sending less than two ACKs in response.

As an example, consider a connection where the the sender's congestion window W is greater than four and the ACK Ratio R is at its maximum value of $W/2$. If the sender becomes data-limited and sends less than W data packets in a round-trip time, then the receiver can send less than two ACK packets in response. This behavior makes the connection more sensitive to the loss of an occasional ACK packet.

Of course, there is still the safety mechanism of the receiver sending an ACK packet when the delayed ACK timer expires. However, more work would be useful to explore the conflicting goals of a congestion-controlled ACK flow and a timely ACK response to the sender for the specific case of a connection with a data-limited sender and a congested ACK path.

5.11. Other Issues

Are there any problems caused by the combination of two-way traffic and reordering? Or other issues that have not yet been addressed?

6. Evaluating ACK Congestion Control

Evaluating ACK Congestion Control will have two components: (1) evaluating the effects of ACK Congestion Control on an individual TCP connection; and (2) evaluating the effects of ACK Congestion Control on aggregate traffic (including the effects of ACK Congestion Control on the aggregate congestion of the path).

The first part, evaluating ACK Congestion Control on the performance of an individual TCP connection, will have to examine those scenarios where ACK Congestion Control might help the performance of a TCP connection, and those scenarios where the use of ACK Congestion Control might cause problems.

The second part, evaluating the effects of ACK Congestion Control on aggregate traffic, should consider scenarios where the use of ACK Congestion Control helps all of the connections sharing a path by

reducing the aggregate congestion on the path. This part should also see if there are scenarios where ACK Congestion Control causes problems by increasing the burstiness of aggregate traffic, or by otherwise changing traffic dynamics.

6.1. Contention in Wireless Links or in non-switched Ethernet

One possible benefit of ACK congestion control is that it could reduce contention in wireless links, shared Ethernet, or other environments with contention between forward-path and reverse-path traffic [AJ03] [KIA07]. At the same time, contention on the shared medium won't necessarily result in dropped ACK packets, and therefore wouldn't necessarily be detected by ACK congestion control.

6.2. Keep-alive and Other Special ACK Packets

Some TCP hosts send keep-alive packets when no data or ACK packets have been received over a long period of time [KEEP-ALIVE]. This keep-alive mechanism is not addressed in TCP specifications. However, such keep-alive packets, if used, should not interact with ACK congestion control one way or another. For ACK congestion control, the ACK Ratio is set small enough to allow the receiver to generally send at least two ACKs for a window of data. In addition, the receiver uses a delayed ACK timer with the ACK Ratio, always sending an acknowledgement if the delayed ACK timer expires. Thus, ACK congestion control will never cause the receiver to delay indefinitely in sending an acknowledgement for a received data packet.

Some TCP implementations send pure ACK packets as window probes, to solicit an ACK packet from the other end with current window information. Such ACK packets will generally be orthogonal to the ACK Congestion Control specified in this draft.

TCP receivers also can send pure ACK packets as window update packets announcing a new value for the receive window, even when the acknowledgement number and SACK options in the ACK packet are not new. The receiver may send window update packets even if the ACK Congestion Control mechanism would say that it is not time yet to send a pure ACK. The sender will not necessarily be able to detect the loss of a window update ACK packet.

7. Measurements of ACK Traffic and Congestion

There are a number of studies about the traffic composition on various links in the Internet, reporting the fraction of bandwidth used by TCP data and by TCP ACK traffic [Studies].

Are there any studies that show the relative packet drop rates for TCP data and ACK traffic, for particular links or for particular TCP connections?

Are there any studies of congested links that show the fraction of traffic on the congested link, or in the congested queue, that consist of TCP ACK packets?

8. Acknowledgement Congestion Control in DCCP's CCID 2

In the transport protocol DCCP [RFC4340], the congestion control mechanism for the CCID 2 profile is based on that of TCP. This section briefly discusses some of the issues that have been addressed in the acknowledgement congestion control already standardized in CCID 2.

Rate-based pacing: For CCID 2, RFC 4341 says that "senders MAY use a form of rate-based pacing when sending multiple data packets liberated by a single ACK packet, rather than sending all liberated data packets in a single burst." However, rate-based pacing is not required in CCID 2.

Increasing the congestion window: For CCID 2, RFC 4341 says that "when $cwnd < ssthresh$, meaning that the sender is in slow-start, the congestion window is increased by one packet for every two newly acknowledged data packets with ACK Vector State 0 (not ECN-marked), up to a maximum of $ACK\ Ratio/2$ packets per acknowledgement. This is a modified form of Appropriate Byte Counting [RFC3465] that is consistent with TCP's current standard (which does not include byte counting), but allows CCID 2 to increase as aggressively as TCP when CCID 2's ACK Ratio is greater than the default value of two. When $cwnd \geq ssthresh$, the congestion window is increased by one packet for every window of data acknowledged without lost or marked packets."

9. Security Considerations

What are the sender's incentives to cheat on ACK congestion control? What are the receiver's incentives to cheat? What are the avenues open for cheating?

As long as ACK congestion control is optional, neither host can be forced to use ACK congestion control if it doesn't want to. So ACK congestion control will only be used if the sender or receiver have some chance of receiving some benefit.

As long as ACK congestion control is optional for TCP, there is

little incentive for the TCP end nodes to cheat on non-ECN-based ACK congestion control. There is nothing now that requires TCP hosts to use congestion control in response to dropped ACK packets.

What avenues for cheating are opened by the use of ECN-Capable ACK packets? If the end nodes can use ECN to have ACK packets marked rather than dropped, and if the end nodes can then avoid the use of ACK congestion control that goes along with the use of ECN on ACK packets, then the end nodes could have an incentive to cheat. Senders could cheat by not instructing the receiver to use a higher ACK Ratio; the receiver would have a hard time detecting this cheating. Receivers could cheat by not using the ACK Ratio they were instructed to use, but senders could easily detect this cheating. However, receivers could also cheat by not using ACK congestion control and still sending ACK packets as ECN-capable, so ACK congestion control is not a necessary component for receivers to cheat about sending ECN-capable ACK packets. One question would be whether there is any way for receivers to cheat about sending ECN-Capable ACK packets and not using appropriate ACK congestion control without this cheating being easily detected by the sender.

What about the ability of routers or middleboxes to detect TCP receivers that cheat by inappropriately sending ACK packets as ECN-capable? The router will only know if the receiver is authorized to send ACK packets as ECN-Capable if the router can see traffic on both the forward and reverse paths, and monitored both the SYN and SYN/ACK packets (and was able to read the TCP options in the packet headers). If ACK congestion control has been negotiated, the router will only know if ACK congestion control is being used correctly by the receiver if it can monitor the ACK Ratio options sent from the sender to the receiver. If ACK congestion control is being used, the router will not necessarily be able to tell if ACK congestion control is being used correctly by the sender, because drops of ACK packets might be occurring after the ACK packets have left the router. However, if the router sees the ACK Ratio options sent from the sender, the router will be able to tell if the sender is correctly accounting for those ACK packets that are dropped or ECN-marked on the path from the receiver to the router.

10. IANA Considerations

No IANA action is needed at this time. If this document was advanced as Experimental or Proposed Standard, then IANA would allocate the option numbers for the two TCP options, the ACK Congestion Control Permitted Option, and the ACK Ratio Option. In this case, the following two lines would be added to the TCP Option Numbers registry (currently located at <http://www.iana.org/assignments/tcp->

parameters):

Kind	Length	Meaning	Reference
TBD1	2	ACK Congestion Control Permitted	[RFC{this}]
TBD2	3	ACK Ratio	[RFC{this}]

In the absence of TCP option numbers allocated by IANA, experimenters may use the TCP Option Numbers set aside for Experimentation in RFC4727 [RFC4727]. As stressed in Section 1 of RFC 3692 [RFC3692], the TCP Option Numbers in the experimental range are intended for experimentation and testing and not for wide or general deployments; these option numbers could be in use by other experimentors for other purposes.

11. Conclusions

This document specifies a congestion control mechanism for acknowledgement traffic (ACKs) for TCP, and discusses the possible complications. We are deferring a recommendation on the use of this mechanism for TCP until it has been evaluated more fully.

12. Acknowledgements

Many thanks for feedback from Mark Allman, Armando Caro, Alfred Hoenes, Ilpoo Jarvinen, Sara Landstrom, Anantha Ramaiah, and Michael Welzl, and for contributed text from Michael Welzl.

A. Appendix: Related Work

There exist several papers dealing with controlling congestion in the reverse path of a TCP connection, especially in the context of networks with bandwidth asymmetry. Some of these proposals require explicit support from routers or middleboxes, whereas others are "pure" end-to-end schemes.

RFC 3449 [RFC3449] discusses TCP performance problems that arise in TCP connections over asymmetric paths. Section 3 of RFC 3449 describes in detail how congestion on the ACK path can affect overall TCP performance. RFC 3449 also outlines a number of proposed mitigations, including ACK Congestion Control. The experimental ACK Congestion Control mechanism discussed in the RFC relies on ECN, with the TCP sender detecting congestion on the ACK path from ECN-marked packets. The RFC also discusses two receiver-based mechanisms, Window Prediction Mechanism (WPM) [CLP98] and Acknowledgement based on Cwnd Estimation (ACE) [MJW00], for using a dynamic ACK Ratio. RFC 3449 also considers link and network layer techniques that

address congestion on the upstream path. These include header compression, and bandwidth management techniques for the upstream link including ACK filtering and ACK reconstruction.

RFC 3135 [RFC3135] on "Performance Enhancing Proxies Intended to Mitigate Link-Related Degradations" surveys a range of Performance Enhancing Proxies used to improve TCP behavior, including proxies for ACK filtering and reconstruction. RFC 2760 [RFC2760] on "Ongoing TCP Research Related to Satellites" discusses both ACK Congestion Control and ACK Filtering and Reconstruction, with detailed descriptions of the mechanisms proposed by Balakrishnan et al. in [BPK97].

Landstrom et al. in [LL07] explore a mechanism where the receiver sends only four acknowledgements per window of data, along with the sender using a form of Appropriate Byte Counting. In addition, the receiver reverts to a lower acknowledgement frequency after a timeout, to avoid unnecessary Retransmit Timeouts. One conclusion of the paper is that pacing at the sender introduces an additional delay, and might not be necessary. A key result of the paper is that with the use of some form of byte counting at the sender, it is possible for TCP to use a lower acknowledgement frequency than that of delayed acknowledgements.

A.1. ECN-only Mechanisms

Balakrishnan et al. in [BPK97] describe the use of ECN to detect congestion in the return path, in order to reduce the sending rate of ACKs. The use of a RED queue in the reverse path allows for marking of ACK packets. The sender echoes back ECN congestion marks to the receiver. The receiver keeps an ACK ratio d (called the "delayed-ACK factor"), specifying the number of data segments that have to be received before the receiver sends a new ACK. The ACK ratio d is managed using multiplicative-increase, additive-decrease; upon reception of a congestion mark, the receiver doubles the value of d (hence dividing the ACK sending rate by two). The ACK ratio decreases linearly for each RTT in which no ECN-marked ACKs are received. Multiple congestion marks received in an RTT are treated as a single congestion event, i.e., d can be doubled at most once per RTT. The TCP timestamp option is used to keep track of the RTT values.

A.2. Receiver-only Mechanisms

In [MJW00], Tam Ming-Chit et al. propose a receiver-based method for calculating an "appropriate" number of ACKs per congestion window ($cwnd$) of data, in order to alleviate congestion on the reverse path. The sender's $cwnd$ is estimated at the receiver by counting the number of received packets per RTT (which also has to be estimated by the

receiver). From this estimate, a simple algorithm is used to compute the number of ACKs to be sent per cwnd. The algorithm enforces a lower bound on the number of ACKs per cwnd, aiming at minimizing the probability of timeout at the sender due to ACK loss. Similarly, the ACK ratio is upper-bounded so as to avoid excessive ACK delay.

Blandford et al. [BGG+07] propose an end-to-end, receiver-oriented scheme called "smartacking". The algorithm is based upon the receiver monitoring the inter-segment arrival time for data packets and adapting the ACK sending rate in response. When the bottleneck link is underutilized, ACKs are sent frequently (up to one ACK per received segment) to promote fast growth of the congestion window. On the other hand, when the bottleneck is close to full utilization, the algorithm tries to reduce control traffic overhead and slow congestion window growth by generating ACKs at the minimum rate needed to keep the data pipe full.

Reducing the number of ACKs (or, equivalently, increasing the amount of bytes acknowledged by each ACK) can increase the burstiness of the TCP sender. Hence, any mechanism as those cited above should be coupled with a burst mitigation technique, Rate-Based Pacing, that paces the sending of data segments [AB05] [ASA00] [BPK97].

A.3. Middlebox-based Mechanisms

ACK filtering (AF) [BPK97] from Balakrishnan et al. is a router-based technique that tries to reduce the number of ACKs sent over the congested return link. With AF, an arriving ACK may replace preceding, older ACKs at the bottleneck queue. An aggressive replacement policy might guarantee that at most one ACK per connection is waiting in the queue, alleviating congestion. However, as in other proposals, care must be taken to avoid sender timeouts in case the (too few) ACKs resulting from the filtering get lost. The idea of filtering ACKs has been extended in [YMH03] to deal with SACK information.

Aweya et al. [AOM02] present a middlebox-based approach for mitigating data packet bursts and for controlling the uplink ACK congestion. The main idea is to perform pacing on ACK segments on an edge device close to the sender, so as to control the ACK arrival rate at the sender.

B. Appendix: Design Considerations

B.1. The TCP ACK Ratio Option, or an AckNow bit in data packets?

In the ACK congestion control mechanism specified in this document, the sender uses the TCP ACK Ratio Option to tell the receiver the ACK

Ratio to use. An alternate approach to the TCP ACK Ratio Option could be for the sender to use an AckNow bit in the TCP header of data packets, telling the receiver to acknowledge this data packet. In the discussion below, we call these two approaches the TCP ACK Ratio Option approach and the AckNow approach.

An advantage of an AckNow approach is that it would require less state from the receiver; the receiver would not need to maintain a variable for the current ACK Ratio, and would not need to keep track of the number of data packets unacked to date.

However, a disadvantage of the AckNow approach is that the sender does not know when packets will be reordered, delayed, or dropped on the path to the receiver. In particular, the sender does not have control over whether a data packet with the AckNow bit set is reordered, delayed, or dropped in the network. For this reason, we have chosen the approach of the sender determining the ACK Ratio that should be used, and sending the ACK Ratio to the receiver, rather than the sender telling the receiver exactly which data packets to acknowledge.

An additional disadvantage of the AckNow approach is that it would add complications and add difficulties for the default cases of the receiver using an ACK Ratio of one or two, as is done in the absence of ACK congestion control.

For these reasons, we have specified that the sender determines the ACK Ratio to use and tells the receiver, rather than the sender setting an AckNow bit in the TCP Header of selected data packets.

Normative References (if standardized)

- [RFC2119] S. Bradner, "Key Words For Use in RFCs to Indicate Requirement Levels", RFC 2119, March 1997.
- [RFC2581] Allman, M., V. Paxson, and W. Stevens, "TCP Congestion Control", RFC 2581, April 1999.
- [RFC3692] Narton, T., "Assigning Experimental and Testing Numbers Considered Useful", RFC 3692, January 2004.
- [RFC3465] Allman, M., "TCP Congestion Control with Appropriate Byte Counting (ABC)", RFC 3465, Experimental, February 2003.
- [RFC4340] Kohler, E., Handley, M., and S. Floyd, "Datagram Congestion Control Protocol (DCCP)", RFC 4340, March 2006.

[RFC4341] Floyd, S., and E. Kohler, "Profile for Datagram Congestion Control Protocol (DCCP) Congestion Control ID 2: TCP-like Congestion Control", RFC 4341, March 2006.

[RFC4727] Fenner, B., "Experimental Values in IPv4, IPv6, ICMPv4, ICMPv6, UDP, and TCP Headers", RFC 4727, November 2006.

Informative References

[RFC2760] Allman, M., Dawkins, S., Glover, D., Griner, J., Henderson, T., Heidemann, J., Kruse, H., Ostermann, S., Scott, K., Semke, J., Touch, J. and D. Tran, "Ongoing TCP Research Related to Satellites", RFC 2760, February 2000.

[RFC3135] Border, J., Kojo, M., Griner, J., Montenegro, G. and Z. Shelby, "Performance Enhancing Proxies Intended to Mitigate Link-Related Degradations", RFC 3135, June 2001.

[RFC3168] K. Ramakrishnan, S. Floyd and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, September 2001.

[RFC3449] H. Balakrishnan, V. N. Padmanabhan, G. Fairhurst, and M. Sooriyabandara, "TCP Performance Implications of Network Path Asymmetry", RFC 3449, December 2002.

[RFC4653] S. Bhandarkar, A. L. N. Reddy, M. Allman and E. Blanton, "Improving the Robustness of TCP to Non-Congestion Events", RFC 4653, August 2006.

[ASA00] A. Aggarwal, S. Savage, and T. Anderson, "Understanding the Performance of TCP Pacing", INFOCOM (3), pp. 1157-1165, 2000.

[AB05] M. Allman and E. Blanton, "Notes on Burst Mitigation for Transport Protocols", SIGCOMM Comput. Commun. Rev., 35(2):5360, 2005.

[AJ03] E. Altman and T. Jimenez, "Novel Delayed ACK Techniques for Improving TCP Performance in Multihop Wireless Networks", Proc. of the Personal Wireless Communications, 2003.

[AOM02] J. Aweya, M. Ouellette, and D. Y. Montuno, "A Self-regulating TCP Acknowledgement (ack) Pacing Scheme". Int. J. Netw. Manag., 12(3):145163, 2002.

[BA03] C. Barakat and E. Altman, "On ACK Filtering on a Slow Reverse Channel", International Journal of Satellite Communications and Networking, V.21 N.3, 2003.

- [BPK97] Balakrishnan, H., V. Padmanabhan, and Katz, R., "The Effects of Asymmetry on TCP Performance", Third ACM/IEEE Mobicom Conference, September 1997.
- [BGG+07] D.K. Blandford, S.A. Goldman, S. Gorinsky, Y. Zhou, and D.R. Dooley, "Smartacking: Improving TCP Performance from the Receiving End", Journal of Internet Engineering, 1(1), 2007.
- [CLP98] Calveras, A., Linares, J., and J. Paradells, "Window Prediction Mechanism for Improving TCP in Wireless Asymmetric Links". Proc. IEEE Global Communications Conference (GLOBECOM), Sydney Australia, November 1998, pp. 533-538.
- [KIA07] Keceli, F., I. Inan, and E. Ayanoglu, "TCP ACK Congestion Control and Filtering for Fairness Provision in the Uplink of IEEE 802.11 Infrastructure Basic Service Set", Proc. IEEE ICC 2007, June 2007.
- [KEEP-ALIVE] Fabio Busatto, "TCP Keepalive HOWTO", May 2007, URL "<http://tldp.org/HOWTO/TCP-Keepalive-HOWTO/index.html>".
- [KLS07] H. Kim, H. Lee, and S. Shin, "On the Cross-Layer Impact of TCP ACK Thinning on IEEE 802.11 Wireless MAC Dynamics", IEICE Transactions on Communications, 2007
- [LL07] Landstrom, S., and Larzon, L.-A., "Reducing the TCP Acknowledgement Frequency", CCR, July 2007.
- [MJW00] Ming-Chit, I.T., Jinsong, D., and W. Wang, "Improving TCP Performance Over Asymmetric Networks", ACM SIGCOMM, ACM Computer Communications Review (CCR), Vol.30, No.3, 2000.
- [Studies] Web page on "Measurement Studies of End-to-End Congestion Control in the Internet", URL "<http://www.icir.org/floyd/ccmeasure.html>".
- [WWCM99] H. Wu, J. Wu, S. Cheng, and J. Ma, "ACK Filtering on Bandwidth Asymmetry Networks", IEEE Communications, 1999.
- [YMH03] L. Yu, Y. Minhua, and Z. Huimin, "The Improvement of TCP Performance in Bandwidth Asymmetric Network", IEEE PIMRC, 1:482-486, September 2003.

Authors' Addresses

Sally Floyd
ICSI Center for Internet Research
1947 Center Street, Suite 600
Berkeley, CA 94704
USA

Email: floyd@icir.org

Andres Arcia
Networking, Security & Multimedia (RSM) Dpt.
TELECOM Bretagne
Rue de la Chataigneraie, CS 17607
35576 Cesson Sevigne Cedex
France

Email:ae <dot> arcia <at> telecom-bretagne <dot> eu

Universidad de Los Andes
Facultad de Ingenieria
Nucleo La Hechicera
Merida, Merida 5101
Venezuela

Email: amoret <at> ula <dot> ve
URI: <http://www.ula.ve>

Janardhan R. Iyengar
Math and Computer Science
Franklin & Marshall College
P. O. Box 3003
Lancaster, PA 17604-3003
USA

Email: jiyengar <at> fandm <dot> edu

David Ros
Networking, Security & Multimedia (RSM) Dpt.
TELECOM Bretagne
Rue de la Chataigneraie, CS 17607
35576 Cesson Sevigne Cedex
France

Email: David <dot> Ros <at> telecom-bretagne <dot> eu

Acknowledgement

Funding for the RFC Editor function is provided by the IETF Administrative Support Activity (IASA).

