

TCP Maintenance and Minor  
Extensions Working Group  
Internet-Draft  
Obsoletes: [3782](#) (if approved)  
Intended status: Standards Track  
Expires: April 22, 2012

T. Henderson  
Boeing  
S. Floyd  
ICSI  
A. Gurtov  
HIIT  
Y. Nishida  
WIDE Project  
October 22, 2011

**The NewReno Modification to TCP's Fast Recovery Algorithm**  
**draft-ietf-tcpm-rfc3782-bis-03.txt**

Abstract

[RFC 5681](#) documents the following four intertwined TCP congestion control algorithms: slow start, congestion avoidance, fast retransmit, and fast recovery. [RFC 5681](#) explicitly allows certain modifications of these algorithms, including modifications that use the TCP Selective Acknowledgement (SACK) option ([RFC 2883](#)), and modifications that respond to "partial acknowledgments" (ACKs which cover new data, but not all the data outstanding when loss was detected) in the absence of SACK. This document describes a specific algorithm for responding to partial acknowledgments, referred to as NewReno. This response to partial acknowledgments was first proposed by Janey Hoe. This document obsoletes [RFC 3782](#).

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

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."

This Internet-Draft will expire on April 22, 2012.

Copyright Notice

Copyright (c) 2011 IETF Trust and the persons identified as

the document authors. All rights reserved.

This document is subject to [BCP 78](http://trustee.ietf.org/license-info) and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

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.



## 1. Introduction

For the typical implementation of the TCP Fast Recovery algorithm described in [\[RFC5681\]](#) (first implemented in the 1990 BSD Reno release, and referred to as the Reno algorithm in [\[FF96\]](#)), the TCP data sender only retransmits a packet after a retransmit timeout has occurred, or after three duplicate acknowledgments have arrived triggering the Fast Retransmit algorithm. A single retransmit timeout might result in the retransmission of several data packets, but each invocation of the Fast Retransmit algorithm in [RFC 5681](#) leads to the retransmission of only a single data packet.

Two problems arise with Reno TCP when multiple packet losses occur in a single window. First, Reno will often take a timeout, as has been documented in [\[Hoe95\]](#). Second, even if a retransmission timeout is avoided, multiple fast retransmits and window reductions can occur, as documented in [\[F94\]](#). When multiple packet losses occur, if the SACK option [\[RFC2883\]](#) is available, the TCP sender has the information to make intelligent decisions about which packets to retransmit and which packets not to retransmit during Fast Recovery. This document applies to TCP connections that are unable to use the TCP Selective Acknowledgement (SACK) option, either because the option is not locally supported or because the TCP peer did not indicate a willingness to use SACK.

In the absence of SACK, there is little information available to the TCP sender in making retransmission decisions during Fast Recovery. From the three duplicate acknowledgments, the sender infers a packet loss, and retransmits the indicated packet. After this, the data sender could receive additional duplicate acknowledgments, as the data receiver acknowledges additional data packets that were already in flight when the sender entered Fast Retransmit.

In the case of multiple packets dropped from a single window of data, the first new information available to the sender comes when the sender receives an acknowledgment for the retransmitted packet (that is, the packet retransmitted when Fast Retransmit was first entered). If there is a single packet drop and no reordering, then the acknowledgment for this packet will acknowledge all of the packets transmitted before Fast Retransmit was entered. However, if there are multiple packet drops, then the acknowledgment for the retransmitted packet will acknowledge some but not all of the packets transmitted before the Fast Retransmit. We call this acknowledgment a partial acknowledgment.

Along with several other suggestions, [\[Hoe95\]](#) suggested that during Fast Recovery the TCP data sender responds to a partial



acknowledgment by inferring that the next in-sequence packet has been lost, and retransmitting that packet. This document describes a modification to the Fast Recovery algorithm in [RFC 5681](#) that incorporates a response to partial acknowledgments received during Fast Recovery. We call this modified Fast Recovery algorithm NewReno, because it is a slight but significant variation of the basic Reno algorithm in [RFC 5681](#). This document does not discuss the other suggestions in [\[Hoe95\]](#) and [\[Hoe96\]](#), such as a change to the ssthresh parameter during Slow-Start, or the proposal to send a new packet for every two duplicate acknowledgments during Fast Recovery. The version of NewReno in this document also draws on other discussions of NewReno in the literature [\[LM97, Hen98\]](#).

We do not claim that the NewReno version of Fast Recovery described here is an optimal modification of Fast Recovery for responding to partial acknowledgments, for TCP connections that are unable to use SACK. Based on our experiences with the NewReno modification in the NS simulator [\[NS\]](#) and with numerous implementations of NewReno, we believe that this modification improves the performance of the Fast Retransmit and Fast Recovery algorithms in a wide variety of scenarios. Previous versions of this RFC [\[RFC2582, RFC3782\]](#) provide simulation-based evidence of the possible performance gains.

## 2. Terminology and Definitions

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC 2119](#) [\[RFC2119\]](#).

This document assumes that the reader is familiar with the terms SENDER MAXIMUM SEGMENT SIZE (SMSS), CONGESTION WINDOW (cwnd), and FLIGHT SIZE (FlightSize) defined in [\[RFC5681\]](#). FLIGHT SIZE is defined as in [\[RFC5681\]](#) as follows:

### FLIGHT SIZE:

The amount of data that has been sent but not yet cumulatively acknowledged.

This document defines an additional sender-side state variable called RECOVER:

### RECOVER:

When in Fast Recovery, this variable records the send sequence number that must be acknowledged before the Fast Recovery procedure is declared to be over.



### **3. The Fast Retransmit and Fast Recovery Algorithms in NewReno**

#### **3.1. Protocol Overview**

The basic idea of these extensions to the Fast Retransmit and Fast Recovery algorithms described in [Section 3.2 of \[RFC5681\]](#) is as follows. The TCP sender can infer, from the arrival of duplicate acknowledgments, whether multiple losses in the same window of data have most likely occurred, and avoid taking a retransmit timeout or making multiple congestion window reductions due to such an event.

The NewReno modification applies to the Fast Recovery procedure that begins when three duplicate ACKs are received and ends when either a retransmission timeout occurs or an ACK arrives that acknowledges all of the data up to and including the data that was outstanding when the Fast Recovery procedure began.

#### **3.2. Specification**

The procedures specified in [Section 3.2 of \[RFC5681\]](#) are followed with the following modifications.

- 1) Initialization of TCP protocol control block:  
When the TCP protocol control block is initialized, Recover is set to the initial send sequence number.
- 2) Three duplicate ACKs:  
When the third duplicate ACK is received, the TCP sender first checks the value of Recover to see if the Cumulative Acknowledgment field covers more than Recover. If so, the value of Recover is incremented to the value of the highest sequence number transmitted by the TCP so far. The TCP then enters Fast Retransmit (step 2 of [Section 3.2 of \[RFC5681\]](#)). If not, the TCP does not enter fast retransmit and does not reset ssthresh.
- 3) Response to newly acknowledged data:  
Step 6 of [\[RFC5681\]](#) specifies the response to the next ACK that acknowledges previously unacknowledged data. When an ACK arrives that acknowledges new data, this ACK could be the acknowledgment elicited by the retransmission from step 2, or elicited by a later retransmission. There are two cases.

Full acknowledgments:

If this ACK acknowledges all of the data up to and including Recover, then the ACK acknowledges all the intermediate segments sent between the original transmission of the lost segment and the receipt of the third duplicate ACK. Set cwnd to





either (1)  $\min(\text{ssthresh}, \max(\text{FlightSize}, \text{SMSS}) + \text{SMSS})$  or (2)  $\text{ssthresh}$ , where  $\text{ssthresh}$  is the value set when Fast Retransmit was entered, and where  $\text{FlightSize}$  in (1) is the amount of data presently outstanding. This is termed "deflating" the window. If the second option is selected, the implementation is encouraged to take measures to avoid a possible burst of data, in case the amount of data outstanding in the network is much less than the new congestion window allows. A simple mechanism is to limit the number of data packets that can be sent in response to a single acknowledgment. Exit the Fast Recovery procedure.

Partial acknowledgments:

If this ACK does *not* acknowledge all of the data up to and including `Recover`, then this is a partial ACK. In this case, retransmit the first unacknowledged segment. Deflate the congestion window by the amount of new data acknowledged by the cumulative acknowledgment field. If the partial ACK acknowledges at least one `SMSS` of new data, then add back `SMSS` bytes to the congestion window. This artificially inflates the congestion window in order to reflect the additional segment that has left the network. Send a new segment if permitted by the new value of `cwnd`. This "partial window deflation" attempts to ensure that, when Fast Recovery eventually ends, approximately  $\text{ssthresh}$  amount of data will be outstanding in the network. Do not exit the Fast Recovery procedure (i.e., if any duplicate ACKs subsequently arrive, execute Step 4 of [Section 3.2 of \[RFC5681\]](#)).

For the first partial ACK that arrives during Fast Recovery, also reset the retransmit timer. Timer management is discussed in more detail in [Section 4](#).

4) Retransmit timeouts:

After a retransmit timeout, record the highest sequence number transmitted in the variable `Recover` and exit the Fast Recovery procedure if applicable.

Step 2 above specifies a check that the Cumulative Acknowledgment field covers more than `Recover`. Because the acknowledgment field contains the sequence number that the sender next expects to receive, the acknowledgment "`ack_number`" covers more than `Recover` when:

`ack_number - 1 > Recover;`

i.e., at least one byte more of data is acknowledged beyond the highest byte that was outstanding when Fast Retransmit was last entered.



Note that in Step 3 above, the congestion window is deflated after a partial acknowledgment is received. The congestion window was likely to have been inflated considerably when the partial acknowledgment was received. In addition, depending on the original pattern of packet losses, the partial acknowledgment might acknowledge nearly a window of data. In this case, if the congestion window was not deflated, the data sender might be able to send nearly a window of data back-to-back.

This document does not specify the sender's response to duplicate ACKs when the Fast Retransmit/Fast Recovery algorithm is not invoked. This is addressed in other documents, such as those describing the Limited Transmit procedure [[RFC3042](#)]. This document also does not address issues of adjusting the duplicate acknowledgment threshold, but assumes the threshold specified in the IETF standards; the current standard is [[RFC5681](#)], which specifies a threshold of three duplicate acknowledgments.

As a final note, we would observe that in the absence of the SACK option, the data sender is working from limited information. When the issue of recovery from multiple dropped packets from a single window of data is of particular importance, the best alternative would be to use the SACK option.

#### **4. Handling Duplicate Acknowledgments After A Timeout**

After each retransmit timeout, the highest sequence number transmitted so far is recorded in the variable "recover". If, after a retransmit timeout, the TCP data sender retransmits three consecutive packets that have already been received by the data receiver, then the TCP data sender will receive three duplicate acknowledgments that do not cover more than "recover". In this case, the duplicate acknowledgments are not an indication of a new instance of congestion. They are simply an indication that the sender has unnecessarily retransmitted at least three packets.

However, when a retransmitted packet is itself dropped, the sender can also receive three duplicate acknowledgments that do not cover more than "recover". In this case, the sender would have been better off if it had initiated Fast Retransmit. For a TCP that implements the algorithm specified in [Section 3](#) of this document, the sender does not infer a packet drop from duplicate acknowledgments in this scenario. As always, the retransmit timer is the backup mechanism for inferring packet loss in this case.

There are several heuristics, based on timestamps or on the amount of advancement of the cumulative acknowledgment field, that allow the sender to distinguish, in some cases, between three duplicate



acknowledgments following a retransmitted packet that was dropped, and three duplicate acknowledgments from the unnecessary retransmission of three packets [[Gur03](#), [GF04](#)]. The TCP sender MAY use such a heuristic to decide to invoke a Fast Retransmit in some cases, even when the three duplicate acknowledgments do not cover more than "recover".

For example, when three duplicate acknowledgments are caused by the unnecessary retransmission of three packets, this is likely to be accompanied by the cumulative acknowledgment field advancing by at least four segments. Similarly, a heuristic based on timestamps uses the fact that when there is a hole in the sequence space, the timestamp echoed in the duplicate acknowledgment is the timestamp of the most recent data packet that advanced the cumulative acknowledgment field [[RFC1323](#)]. If timestamps are used, and the sender stores the timestamp of the last acknowledged segment, then the timestamp echoed by duplicate acknowledgments can be used to distinguish between a retransmitted packet that was dropped and three duplicate acknowledgments from the unnecessary retransmission of three packets.

#### **[4.1.](#) ACK Heuristic**

If the ACK-based heuristic is used, then following the advancement of the cumulative acknowledgment field, the sender stores the value of the previous cumulative acknowledgment as `prev_highest_ack`, and stores the latest cumulative ACK as `highest_ack`. In addition, the following step is performed if Step 1 in [Section 3](#) fails, before proceeding to Step 1B.

- 1\*) If the Cumulative Acknowledgment field didn't cover more than "recover", check to see if the congestion window is greater than SMSS bytes and the difference between `highest_ack` and `prev_highest_ack` is at most  $4 \times \text{SMSS}$  bytes. If true, duplicate ACKs indicate a lost segment (proceed to Step 1A in [Section 3](#)). Otherwise, duplicate ACKs likely result from unnecessary retransmissions (proceed to Step 1B in [Section 3](#)).

The congestion window check serves to protect against fast retransmit immediately after a retransmit timeout.

If several ACKs are lost, the sender can see a jump in the cumulative ACK of more than three segments, and the heuristic can fail. [[RFC5681](#)] recommends that a receiver should send duplicate ACKs for every out-of-order data packet, such as a data packet received during Fast Recovery. The ACK heuristic is more likely to fail if the receiver does not follow this advice, because then a smaller number of ACK losses are needed to produce a



sufficient jump in the cumulative ACK.

#### **4.2. Timestamp Heuristic**

If this heuristic is used, the sender stores the timestamp of the last acknowledged segment. In addition, the second paragraph of step 1 in [Section 3](#) is replaced as follows:

1\*\*) If the Cumulative Acknowledgment field didn't cover more than "recover", check to see if the echoed timestamp in the last non-duplicate acknowledgment equals the stored timestamp. If true, duplicate ACKs indicate a lost segment (proceed to Step 1A in [Section 3](#)). Otherwise, duplicate ACKs likely result from unnecessary retransmissions (proceed to Step 1B in [Section 3](#)).

The timestamp heuristic works correctly, both when the receiver echoes timestamps as specified by [\[RFC1323\]](#), and by its revision attempts. However, if the receiver arbitrarily echoes timestamps, the heuristic can fail. The heuristic can also fail if a timeout was spurious and returning ACKs are not from retransmitted segments. This can be prevented by detection algorithms such as [\[RFC3522\]](#).

#### **5. Implementation Issues for the Data Receiver**

[\[RFC5681\]](#) specifies that "Out-of-order data segments SHOULD be acknowledged immediately, in order to accelerate loss recovery." Neal Cardwell has noted that some data receivers do not send an immediate acknowledgment when they send a partial acknowledgment, but instead wait first for their delayed acknowledgment timer to expire [\[C98\]](#). As [\[C98\]](#) notes, this severely limits the potential benefit of NewReno by delaying the receipt of the partial acknowledgment at the data sender. Echoing [\[RFC5681\]](#), our recommendation is that the data receiver send an immediate acknowledgment for an out-of-order segment, even when that out-of-order segment fills a hole in the buffer.

#### **6. Implementation Issues for the Data Sender**

In [Section 3](#), Step 5 above, it is noted that implementations should take measures to avoid a possible burst of data when leaving Fast Recovery, in case the amount of new data that the sender is eligible to send due to the new value of the congestion window is large. This can arise during NewReno when ACKs are lost or treated as pure window updates, thereby causing the sender to underestimate the number of new segments that can be sent during the recovery procedure. Specifically, bursts can occur when the FlightSize is much less than the new congestion window when exiting from Fast Recovery. One





simple mechanism to avoid a burst of data when leaving Fast Recovery is to limit the number of data packets that can be sent in response to a single acknowledgment. (This is known as "maxburst\_" in the ns simulator.) Other possible mechanisms for avoiding bursts include rate-based pacing, or setting the slow-start threshold to the resultant congestion window and then resetting the congestion window to FlightSize. A recommendation on the general mechanism to avoid excessively bursty sending patterns is outside the scope of this document.

An implementation may want to use a separate flag to record whether or not it is presently in the Fast Recovery procedure. The use of the value of the duplicate acknowledgment counter for this purpose is not reliable because it can be reset upon window updates and out-of-order acknowledgments.

When updating the Cumulative Acknowledgment field outside of Fast Recovery, the "recover" state variable may also need to be updated in order to continue to permit possible entry into Fast Recovery ([Section 3](#), step 1). This issue arises when an update of the Cumulative Acknowledgment field results in a sequence wraparound that affects the ordering between the Cumulative Acknowledgment field and the "recover" state variable. Entry into Fast Recovery is only possible when the Cumulative Acknowledgment field covers more than the "recover" state variable.

It is important for the sender to respond correctly to duplicate ACKs received when the sender is no longer in Fast Recovery (e.g., because of a Retransmit Timeout). The Limited Transmit procedure [[RFC3042](#)] describes possible responses to the first and second duplicate acknowledgments. When three or more duplicate acknowledgments are received, the Cumulative Acknowledgment field doesn't cover more than "recover", and a new Fast Recovery is not invoked, it is important that the sender not execute the Fast Recovery steps (3) and (4) in [Section 3](#). Otherwise, the sender could end up in a chain of spurious timeouts. We mention this only because several NewReno implementations had this bug, including the implementation in the NS simulator.

It has been observed that some TCP implementations enter a slow start or congestion avoidance window updating algorithm immediately after the cwnd is set by the equation found in ([Section 3](#), step 5), even without a new external event generating the cwnd change. Note that after cwnd is set based on the procedure for exiting Fast Recovery ([Section 3](#), step 5), cwnd SHOULD NOT be updated until a further event occurs (e.g., arrival of an ack, or timeout) after this adjustment.



## **7. Security Considerations**

[RFC5681] discusses general security considerations concerning TCP congestion control. This document describes a specific algorithm that conforms with the congestion control requirements of [\[RFC5681\]](#), and so those considerations apply to this algorithm, too. There are no known additional security concerns for this specific algorithm.

## **8. IANA Considerations**

This document has no actions for IANA.

## **9. Conclusions**

This document specifies the NewReno Fast Retransmit and Fast Recovery algorithms for TCP. This NewReno modification to TCP can even be important for TCP implementations that support the SACK option, because the SACK option can only be used for TCP connections when both TCP end-nodes support the SACK option. NewReno performs better than Reno ([RFC5681](#)) in a number of scenarios discussed in previous versions of this RFC ([\[RFC2582\]](#), [\[RFC3782\]](#)).

A number of options to the basic algorithm presented in [Section 3](#) are also referenced in [Appendix A](#) to this document. These include the handling of the retransmission timer, the response to partial acknowledgments, and whether or not the sender must maintain a state variable called Recover. Our belief is that the differences between these variants of NewReno are small compared to the differences between Reno and NewReno. That is, the important thing is to implement NewReno instead of Reno, for a TCP connection without SACK; it is less important exactly which of the variants of NewReno is implemented.

## **10. Acknowledgments**

Many thanks to Anil Agarwal, Mark Allman, Armando Caro, Jeffrey Hsu, Vern Paxson, Kacheong Poon, Keyur Shah, and Bernie Volz for detailed feedback on this document or on its precursor, [RFC 2582](#). Jeffrey Hsu provided clarifications on the handling of the recover variable that were applied to [RFC 3782](#) as errata, and now are in [Section 8](#) of this document. Yoshifumi Nishida contributed a modification to the fast recovery algorithm to account for the case in which flightsize is 0 when the TCP sender leaves fast recovery, and the TCP receiver uses delayed acknowledgments. Alexander Zimmermann provided several suggestions to improve the clarity of the document.

## **11. References**



### **11.1. Normative References**

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC5681] Allman, M., Paxson, V. and E. Blanton, "TCP Congestion Control", [RFC 5681](#), September 2009.
- [RFC6298] Paxson, V., Allman, M., Chu, J., and Sargent, M., "Computing TCP's Retransmission Timer", [RFC 6298](#), June 2011.

### **11.2. Informative References**

- [C98] Cardwell, N., "delayed ACKs for retransmitted packets: ouch!". November 1998, Email to the tcpimpl mailing list, Message-ID "Pine.LNX.4.02A.9811021421340.26785-100000@sake.cs.washington.edu", archived at "<http://tcp-impl.lerc.nasa.gov/tcp-impl>".
- [F98] Floyd, S., Revisions to [RFC 2001](#), "Presentation to the TCPIMPL Working Group", August 1998. URLs "<ftp://ftp.ee.lbl.gov/talks/sf-tcpimpl-aug98.ps>" and "<ftp://ftp.ee.lbl.gov/talks/sf-tcpimpl-aug98.pdf>".
- [F03] Floyd, S., "Moving NewReno from Experimental to Proposed Standard? Presentation to the TSVWG Working Group", March 2003. URLs "<http://www.icir.org/floyd/talks/newreno-Mar03.ps>" and "<http://www.icir.org/floyd/talks/newreno-Mar03.pdf>".
- [FF96] Fall, K. and S. Floyd, "Simulation-based Comparisons of Tahoe, Reno and SACK TCP", Computer Communication Review, July 1996. URL "<ftp://ftp.ee.lbl.gov/papers/sacks.ps.Z>".
- [F94] Floyd, S., "TCP and Successive Fast Retransmits", Technical report, October 1994. URL "<ftp://ftp.ee.lbl.gov/papers/fastretrans.ps>".
- [GF04] Gurtov, A. and S. Floyd, "Resolving Acknowledgment Ambiguity in non-SACK TCP", Next Generation Teletraffic and Wired/Wireless Advanced Networking (NEW2AN'04), February 2004. URL "<http://www.cs.helsinki.fi/u/gurtov/papers/heuristics.html>".
- [Gur03] Gurtov, A., "[Tsvwg] resolving the problem of unnecessary fast retransmits in go-back-N", email to the tsvwg mailing list, message ID <3F25B467.9020609@cs.helsinki.fi>, July 28, 2003. URL "<http://www1.ietf.org/mail-archive/>

working-groups/ tsvwg/current/msg04334.html".

Henderson, et al.

Expires April 22, 2012

[Page 12]

- [Hen98] Henderson, T., Re: NewReno and the 2001 Revision. September 1998. Email to the tcpimpl mailing list, Message ID "Pine.BSI.3.95.980923224136.26134A-100000@raptor.CS.Berkeley.EDU", archived at "<http://tcp-impl.lerc.nasa.gov/tcp-impl>".
- [Hoe95] Hoe, J., "Startup Dynamics of TCP's Congestion Control and Avoidance Schemes", Master's Thesis, MIT, 1995.
- [Hoe96] Hoe, J., "Improving the Start-up Behavior of a Congestion Control Scheme for TCP", ACM SIGCOMM, August 1996. URL "<http://www.acm.org/sigcomm/sigcomm96/program.html>".
- [LM97] Lin, D. and R. Morris, "Dynamics of Random Early Detection", SIGCOMM 97, September 1997. URL "<http://www.acm.org/sigcomm/sigcomm97/program.html>".
- [NS] The Network Simulator (NS). URL "<http://www.isi.edu/nsnam/ns/>".
- [PF01] Padhye, J. and S. Floyd, "Identifying the TCP Behavior of Web Servers", June 2001, SIGCOMM 2001.
- [RFC1323] Jacobson, V., Braden, R. and D. Borman, "TCP Extensions for High Performance", [RFC 1323](#), May 1992.
- [RFC2582] Floyd, S. and T. Henderson, "The NewReno Modification to TCP's Fast Recovery Algorithm", [RFC 2582](#), April 1999.
- [RFC2883] Floyd, S., J. Mahdavi, M. Mathis, and M. Podolsky, "The Selective Acknowledgment (SACK) Option for TCP", [RFC 2883](#), July 2000.
- [RFC3042] Allman, M., Balakrishnan, H. and S. Floyd, "Enhancing TCP's Loss Recovery Using Limited Transmit", [RFC 3042](#), January 2001.
- [RFC3522] Ludwig, R. and M. Meyer, "The Eifel Detection Algorithm for TCP", [RFC 3522](#), April 2003.
- [RFC3782] Floyd, S., T. Henderson, and A. Gurtov, "The NewReno Modification to TCP's Fast Recovery Algorithm", [RFC 3782](#), April 2004.

## [Appendix A](#). Additional Information

Previous versions of this RFC ([[RFC2582](#)], [[RFC3782](#)]) contained additional informative material on the following subjects, and may be consulted by readers who may want more information about



possible variants to the algorithm and who may want references to specific [\[NS\]](#) simulations that provide NewReno test cases.

[Section 4 of \[RFC3782\]](#) discusses some alternative behaviors for resetting the retransmit timer after a partial acknowledgment.

[Section 5 of \[RFC3782\]](#) discusses some alternative behaviors for performing retransmission after a partial acknowledgment.

[Section 6 of \[RFC3782\]](#) describes more information about the motivation for the sender's state variable Recover.

[Section 9 of \[RFC3782\]](#) introduces some NS simulation test suites for NewReno. In addition, references to simulation results can be found throughout [\[RFC3782\]](#).

[Section 10 of \[RFC3782\]](#) provides a comparison of Reno and NewReno TCP.

[Section 11 of \[RFC3782\]](#) listed changes relative to [\[RFC3782\]](#).

#### **[Appendix B](#). Changes Relative to [RFC 3782](#)**

In [\[RFC3782\]](#), the cwnd after Full ACK reception will be set to (1) min (sssthresh, FlightSize + SMSS) or (2) sssthresh. However, there is a risk in the first logic which results in performance degradation. With the first logic, if FlightSize is zero, the result will be 1 SMSS. This means TCP can transmit only 1 segment at this moment, which can cause delay in ACK transmission at receiver due to delayed ACK algorithm.

The FlightSize on Full ACK reception can be zero in some situations. A typical example is where sending window size during fast recovery is small. In this case, the retransmitted packet and new data packets can be transmitted within a short interval. If all these packets successfully arrive, the receiver may generate a Full ACK that acknowledges all outstanding data. Even if window size is not small, loss of ACK packets or receive buffer shortage during fast recovery can also increase the possibility to fall into this situation.

The proposed fix in this document ensures that sender TCP transmits at least two segments on Full ACK reception.

In addition, errata for [RFC3782](#) (editorial clarification to [Section 8 of RFC2582](#), which is now [Section 6](#) of this document) has been applied.

The specification text ([Section 3.2](#) herein) was rewritten to more closely track [Section 3.2 of \[RFC5681\]](#).

Sections [4](#), [5](#), [9-11](#) of [\[RFC3782\]](#) were removed, and instead Appendix



E-Mail: [gurtov@hiit.fi](mailto:gurtov@hiit.fi)



Yoshifumi Nishida  
WIDE Project  
Endo 5322  
Fujisawa, Kanagawa 252-8520  
Japan  
  
Email: nishida@wide.ad.jp