Network Working Group INTERNET-DRAFT Expires: January 2003 Reiner Ludwig Ericsson Research Andrei Gurtov Sonera Corporation July, 2002

# The Eifel Response Algorithm for TCP <<u>draft-ietf-tsvwg-tcp-eifel-response-00.txt</u>>

Status of this memo

This document is an Internet-Draft and is in full conformance with all provisions of <u>Section 10 of RFC2026</u>.

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 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 <a href="http://www.ietf.org/shadow.html">http://www.ietf.org/shadow.html</a>

# Abstract

The Eifel response algorithm uses the Eifel detection algorithm to detect a posteriori whether the TCP sender has entered loss recovery unnecessarily. In response to a spurious timeout it avoids the go-back-N retransmits that would otherwise be sent, and reinitializes the RTT estimators to avoid further spurious timeouts. Likewise, it adapts the duplicate acknowledgement threshold in response to a spurious fast retransmit. In both cases, the Eifel response algorithm restores the congestion control state in way that avoids packet bursts. Ludwig & Gurtov

[Page 1]

# Terminology

The keywords MUST, MUST NOT, REQUIRED, SHALL, SHALL NOT, SHOULD, SHOULD NOT, RECOMMENDED, MAY, and OPTIONAL, when they appear in this document, are to be interpreted as described in [RFC2119].

We refer to the first-time transmission of an octet as the 'original transmit'. A subsequent transmission of the same octet is referred to as a 'retransmit'. In most cases this terminology can likewise be applied to "data segments" as opposed to "octets". However, when repacketization occurs, a segment can contain both first-time transmissions and retransmissions of octets. In that case this terminology is only consistent when applied to "octets". For the Eifel detection and response algorithms this makes no difference as they also operate correctly when repacketization occurs.

We use the term 'acceptable ACK' as defined in [RFC793]. That is an ACK that acknowledges previously unacknowledged data. We use the term 'duplicate ACK', and the variable 'dupacks' as defined in [WS95]. The variable 'dupacks' is a counter of duplicate ACKs that have already been received by the TCP sender before the fast retransmit is sent. We use the variable 'DupThresh' to refer to the so-called duplicate acknowledgement threshold, i.e., the number of duplicate ACKs that need to arrive at the TCP sender to trigger a fast retransmit. Currently, DupThresh is specified as a fixed value of three [RFC2581].

Furthermore, we use the TCP sender state variables 'SND.UNA' and 'SND.NXT' as defined in [RFC793]. SND.UNA holds the segment sequence number of the oldest outstanding segment. SND.NXT holds the segment sequence number of the next segment the TCP sender will (re-)transmit. In addition, we define as 'SND.MAX' the segment sequence number of the next original transmit to be sent. The definition of SND.MAX is equivalent to the definition of snd\_max in [WS95].

We use the TCP sender state variables 'cwnd' (congestion window), and 'ssthresh' (slow start threshold), and the terms 'SMSS', and 'FlightSize' as defined in [RFC2581]. FlightSize is the amount of outstanding data in the network, or alternatively, the difference between SND.MAX and SND.UNA at a given point in time. We use the TCP sender state variables 'SRTT' and 'RTTVAR', and the term 'RTO' as defined in [RFC2988]. In addition, we assume that the TCP sender maintains in the variable 'RTT-SAMPLE' the value of the latest round-trip time (RTT) measurement.

# **<u>1</u>**. Introduction

The Eifel response algorithm uses the Eifel detection algorithm to detect a posteriori whether the TCP sender has entered loss recovery

Ludwig & Gurtov

[Page 2]

unnecessarily. In response to a spurious timeout it avoids the go-back-N retransmits that would otherwise be sent, and reinitializes the RTT estimators to avoid further spurious timeouts. Likewise, it adapts the duplicate acknowledgement threshold in response to a spurious fast retransmit. In both cases, the Eifel response algorithm restores the congestion control state in way that avoids packet bursts.

The Eifel response algorithm relies on the Eifel detection algorithm defined in [ $\underline{LM02}$ ]. That document discusses the relevant background and motivation that also applies to this document. Hence, the reader is expected to be familiar with [ $\underline{LM02}$ ], and should view this document as a companion document.

# 2. The Eifel Response Algorithm

The complete algorithm is specified in <u>section 2.1</u>. In sections 2.2 to 2.4, we motivate the different steps of the algorithm.

#### **<u>2.1</u>**. The Algorithm

Given that a TCP sender has enabled the Eifel detection algorithm  $[\underline{LM02}]$  for a connection, a TCP sender MAY use the Eifel response algorithm as defined in this subsection. Note that this implies that the TCP Timestamps option  $[\underline{RFC1323}]$  is used for that connection. Since the Eifel response algorithm is dependent on the Eifel detection algorithm, we describe it as an extension of the latter.

If the combined Eifel detection and response algorithm is used, the following steps MUST be taken by the TCP sender, but only upon initiation of loss recovery, i.e., when either the timeout-based retransmit or the fast retransmit is sent. Note: The algorithm MUST NOT be reinitiated after loss recovery has already started. In particular, it may not be reinitiated upon subsequent timeouts for the same segment, and not upon retransmitting segments other than the oldest outstanding segment.

Note that steps (1)-(5) are an one-to-one copy of the Eifel detection algorithm specified in [ $\underline{LM02}$ ], step (0) has been added, and step (RESP) from [ $\underline{LM02}$ ] has been replaced by steps (RESP)-(ReCC) given below.

- (0) Before the variables cwnd and ssthresh get updated when loss recovery is initiated, set a "cwnd\_prev" variable to the current value of FlightSize, and set a "ssthresh\_prev" variable to the value of ssthresh.
- (1) Set a "RetransmitTS" variable to the value of the

Timestamp Value field of the Timestamps option included in the retransmit sent when loss recovery is initiated. A

Ludwig & Gurtov

[Page 3]

TCP sender must ensure that RetransmitTS does not get overwritten as loss recovery progresses, e.g., in case of a second timeout and subsequent second retransmit of the same octet.

- (2) Set a "SpuriousRecovery" variable to FALSE.
- (3) Wait for the arrival of an acceptable ACK. If an acceptable ACK has arrived, then proceed to step (4).
- (4) If the value of the Timestamp Echo Reply field of the acceptable ACK's Timestamps option is smaller than the value of the variable RetransmitTS, then proceed to step (5),

else proceed to step (DONE).

(5) If the loss recovery has been initiated with a timeoutbased retransmit, then set SpuriousRecovery <- SPUR\_TO,</pre>

```
else set
SpuriousRecovery <- dupacks+1</pre>
```

(RESP) If SpuriousRecovery equals SPUR\_TO, then proceed to step (ST0.1),

else (spurious fast retransmit) proceed to step (SFR).

(ST0.1) Resume transmission off the top:

Set SND.NXT <- SND.MAX

(ST0.2) Reinitialize the RTT estimators:

Set

SRTT <- RTT-SAMPLE RTTVAR <- RTT-SAMPLE/2, recalcualte the RTO, and restart the retransmission timer.

Proceed to step (ReCC).

(SFR) Adapt the duplicate acknowledgement threshold:

Set

DupThresh <- max (DupThresh, SpuriousRecovery)</pre>

Proceed to step (ReCC).

Ludwig & Gurtov

[Page 4]

(ReCC) Revert the congestion control state:

If the acceptable ACK has the ECN-Echo flag [<u>RFC3168</u>] set OR the TCP sender has already taken more than three timeouts for the oldest outstanding segment, then proceed to step (DONE),

else set cwnd <- FlightSize + SMSS ssthresh <- max (cwnd\_prev, ssthresh\_prev)

Note: At this point in the algorithm, the value of FlightSize might be different from the value of FlightSize in step (0).

Proceed to step (DONE).

(DONE) No further processing.

# 2.2 Responding to Spurious Timeouts

# 2.2.1 Suppressing the Spurious go-back-N Retransmits (step ST0.1)

Without the Eifel detection algorithm, the TCP sender suffers from the retransmission ambiguity problem [KP87]. This means that when the first acceptable ACK arrives after a spurious timeout, the TCP sender must believe that that ACK was sent in response to the retransmit when in fact it was sent in response to the original transmit. Furthermore, the TCP sender must also believe that all other segments outstanding at that point have been lost. Note: the mentioned ACK cannot carry any SACK option [RFC2018].

Consequently, once the TCP sender's state has been updated after the first acceptable ACK has arrived, SND.NXT equals SND.UNA. This is what causes the often unnecessary go-back-N retransmits. Any newly arriving acceptable ACK that was sent in response to an original transmit will now clock out the segment pointed at by SND.UNA; whether it was lost or not. In fact, during this phase the TCP sender breaks 'packet conservation' [Jac88]. This is because the unnecessary go-back-N retransmits are sent during slow start, i.e., for each original packet leaving the network, two useless retransmits are sent into the network (see [LK00] for more detail).

The Eifel detection algorithm reliably eliminates the retransmission ambiguity problem. Once it detected that a timeout was spurious, it is therefore safe to let the TCP sender resume the transmission with new data. Thus, the Eifel response algorithm changes the TCP sender's state by setting SND.NXT to SND.MAX in that case.

Ludwig & Gurtov

[Page 5]

# 2.2.2 Re-Initializing the RTT Estimators (step ST0.2)

Since the timeout was spurious, the TCP sender's RTT estimators are likely to be off. On the other hand, since timestamps are used, a new and valid RTT measurement (RTT-SAMPLE) can be derived from the acceptable ACK. It is therefore suggested to reinitialize the RTT estimators from RTT-SAMPLE.

To have the new RTO become effective, the retransmission timer needs to be restarted. This is consistent with [<u>RFC2988</u>] which recommends restarting the retransmission timer with the arrival of an acceptable ACK.

## 2.3 Responding to Spurious Fast Retransmits (step SFR)

The assumption behind the fast retransmit algorithm [RFC2581] is that a segment was lost if as many duplicate ACKs have arrived at the TCP sender as indicated by DupThresh. Currently, DupThresh is specified as a fixed value of three [RFC2581]. That value is assumed to be sufficiently conservative so that packet reordering and/or packet duplication does not falsely trigger the fast retransmit algorithm. Clearly, this assumption does not hold for a particular TCP connection once the TCP sender detects that the last fast retransmit has been spurious. It is therefore suggested to dynamically adapt DupThresh to the reordering characteristics observed over the course of a particular connection.

At the beginning of a connection DupThresh is initialized with three. Then for each spurious fast retransmit that is detected, DupThresh is set to the maximum of the previous DupThresh, and the lowest value that would have avoided that spurious fast retransmit. Note that the Eifel detection algorithm records the latter value in SpuriousRecovery. This strategy ensures that the TCP sender is able to cope with the longest reordering length seen on a particular connection so far.

However, the strategy bears the risk that the retransmission timer expires before the TCP sender receives the duplicate ACK that would trigger a fast retransmit of the oldest outstanding segment. To alleviate that potential problem the TCP sender should implement the "fast timeout" algorithm proposed in [Lu02].

Also, we believe that this strategy should be implemented together with an advanced version of the Limited Transmit algorithm [<u>RFC3042</u>]. That is for each duplicate ACK that arrives until DupThresh is reached, the TCP sender should sent a new data segment if new data is available, and the TCP receiver's advertised window allows so. Although, the current Limited Transmit algorithm only allows this for the first two duplicate ACKs, we believe that this is safe. This is already implemented in widely deployed TCPs [SK02].

Ludwig & Gurtov

[Page 6]

Other alternatives for responding to spurious fast retransmits are discussed in [BA02a].

# **2.4** Reverting Congestion Control State (step ReCC)

When a TCP sender enters loss recovery, it also assumes that is has received a congestion indication. In response to that it reduces cwnd, and ssthresh. However, once the TCP sender detects that the loss recovery has been falsely triggered, this reduction was unnecessary. In fact, no congestion signal has been received. We therefore believe that it is safe to revert to the previous congestion control state.

Instead, of simply restoring cwnd, and ssthresh, it is suggested to set cwnd to one half the previous cwnd, and then enter the slow start phase. This is more conservative than the original proposal, but it avoids the packet burst that could otherwise be triggered after a spurious fast retransmit [<u>LK00</u>]. When the spurious loss recovery has been triggered during slow start, the previous slow start threshold is restored. Otherwise, the TCP sender slow starts to the FlightSize it had before the loss recovery was initiated (cwnd\_prev).

Clearly, when the acceptable ACK signals congestion through the ECN-Echo flag [RFC3168], the TCP sender MUST refrain from reverting congestion control state. The same is true if the TCP sender has already taken more than three timeouts for the oldest outstanding segment. Allowing three timeouts while still reverting congestion control state goes beyond [RFC2581]. That standard recommends setting cwnd to no more than the restart window before beginning transmission if the TCP sender has not sent data in an interval exceeding the current RTO. The motivation for doing so is to restart the ACK clock which is believed to have been lost. The case in step (ReCC) of the Eifel response algorithm is different. Since, an acceptable ACK has finally returned, the TCP has reason to believe that the ACK clock was merely interrupted but has now resumed "ticking" again.

### 3. Interoperability with Advanced Loss Recovery Schemes

We believe that there are no problems concerning interoperability with advanced loss recovery schemes such as NewReno [RFC2582], or SACK-based schemes [2018], [BA02b]. This is because in case loss recovery has been initiated unnecessarily, the Eifel response algorithm will have caused the TCP sender to back out of loss recovery before those schemes would have kicked in.

In fact, we recommend that the Eifel response algorithm is implemented together with one of those advanced loss recovery scheme; ideally a SACK-based alternative. In an environment where spurious timeouts and back-to-back packet losses often coincide, we have found

Ludwig & Gurtov

[Page 7]

that TCP's performance can even suffer if the Eifel response algorithm is operated without an advanced loss recovery scheme [<u>GL02</u>].

In that study we among other variants compared TCP-Reno with and without the Eifel response algorithm (TCP-Reno/Eifel vs. TCP-Reno), and without an advanced loss recovery scheme for both variants. The reason that TCP-Reno performed better in the mentioned scenario, is its aggressiveness after a spurious timeout. Even though it breaks 'packet conservation' (see <u>Section 2.2.1</u>) when blindly retransmitting all outstanding segments, it usually recovers the back-to-back packet losses within a single round-trip time. On the contrary, the more conservative TCP-Reno/Eifel was forced into another (backed-off) timeout in that case. In the study, we found that the best end-to-end performance was achieved when the TCP sender implemented both the Eifel response algorithm and SACK-based loss recovery. In case NewReno is chosen as the advanced loss recovery scheme, we found that it performs better if the 'bugfix' feature is disabled. That feature often lead the TCP sender to the wrong decision.

## **<u>4</u>**. Security Considerations

There is a risk that TCP receivers make a genuine retransmit appear to the TCP sender as a spurious retransmit by forging echoed timestamps. This could effectively disable congestion control at the TCP sender. A reliable method to protect against that risk is to implement the safe variant of the Eifel detection algorithm specified in [LM02].

## Acknowledgments

Many thanks to Keith Sklower, Randy Katz, Michael Meyer, Stephan Baucke, Sally Floyd, Vern Paxson, Mark Allman, and Ethan Blanton for very useful discussions that contributed to this work.

### References

- [RFC3042] M. Allman, H. Balakrishnan, S. Floyd, Enhancing TCP's Loss Recovery Using Limited Transmit, <u>RFC 3042</u>, January 2001.
- [BA02a] E. Blanton, M. Allman, On Making TCP More Robust to Packet Reordering, ACM Computer Communication Review, Vol. 32, No. 1, January 2002.

[BA02b] E. Blanton, M. Allman, A Conservative SACK-based Loss Recovery Algorithm for TCP, work in progress, July 2002.

Ludwig & Gurtov

[Page 8]

- [RFC2119] S. Bradner, Key words for use in RFCs to Indicate Requirement Levels, <u>RFC 2119</u>, March 1997.
- [RFC2582] S. Floyd, T. Henderson, The NewReno Modification to TCP's Fast Recovery Algorithm, RFC 2582, April 1999.
- [Gu01] A. Gurtov, Effect of Delays on TCP Performance, In Proceedings of IFIP Personal Wireless Conference, August 2001.
- [GL02] A. Gurtov, R. Ludwig, Evaluating the Eifel Algorithm for TCP in a GPRS Network, In Proceedings of the European Wireless Conference, February 2002.
- [RFC1323] V. Jacobson, R. Braden, D. Borman, TCP Extensions for High Performance, <u>RFC 1323</u>, May 1992.
- [RFC2018] M. Mathis, J. Mahdavi, S. Floyd, A. Romanow, TCP Selective Acknowledgement Options, <u>RFC 2018</u>, October 1996.
- [KP87] P. Karn, C. Partridge, Improving Round-Trip Time Estimates in Reliable Transport Protocols, In Proceedings of ACM SIGCOMM 87.
- [LK00] R. Ludwig, R. H. Katz, The Eifel Algorithm: Making TCP Robust Against Spurious Retransmissions, ACM Computer Communication Review, Vol. 30, No. 1, January 2000.
- [LM02] R. Ludwig, M. Meyer, The Eifel Detection Algorithm for TCP, work in progress, July 2002.
- [Lu02] R. Ludwig, Responding to Fast Timeouts in TCP, work in progress, July 2002.
- [RFC793] J. Postel, Transmission Control Protocol, <u>RFC793</u>, September 1981.
- [RFC3168] K. Ramakrishnan, S. Floyd, D. Black, The Addition of Explicit Congestion Notification (ECN) to IP, <u>RFC 3168</u>, September 2001
- [SK02] P. Sarolahti, A. Kuznetsov, Congestion Control in Linux TCP, In Proceedings of USENIX, June 2002.
- [WS95] G. R. Wright, W. R. Stevens, TCP/IP Illustrated, Volume 2

(The Implementation), Addison Wesley, January 1995.

Ludwig & Gurtov

[Page 9]

Author's Address

Reiner Ludwig Ericsson Research (EED) Ericsson Allee 1 52134 Herzogenrath, Germany Email: Reiner.Ludwig@ericsson.com

Andrei Gurtov Cellular Systems Development P.O. Box 970, FIN-00051 Sonera Helsinki, Finland Phone: +358(0)20401 +358(0)204064365 Fax: Email: andrei.gurtov@sonera.com Homepage: <u>http://www.cs.helsinki.fi/u/gurtov</u>

This Internet-Draft expires in January 2003.

Ludwig & Gurtov

[Page 10]