                   CoAP Simple Congestion Control/Advanced
                        draft-bormann-core-cocoa-01

Abstract

   The CoAP protocol needs to be implemented in such a way that it does
   not cause persistent congestion on the network it uses.  The CoRE
   CoAP specification defines basic behavior that exhibits low risk of
   congestion with minimal implementation requirements.  It also leaves
   room for combining the base specification with advanced congestion
   control mechanisms with higher performance.

   This specification defines some simple advanced CoRE Congestion
   Control mechanisms, Simple CoCoA.  In the present version -01, it is
   already making use of some input from simulations, but might still
   benefit from simplifying it further.

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 August 18, 2014.

Copyright Notice

Table of Contents

1.  Introduction

   (See Abstract.)

   Extended rationale for this specification can be found in
   [I-D.bormann-core-congestion-control] and
   [I-D.eggert-core-congestion-control], as well as in the minutes of
   the IETF 84 CoRE WG meetings.

1.1.  Terminology

   This specification uses terms from [I-D.ietf-core-coap].  In
   addition, it defines the following terminology:

   Initiator:  The endpoint that sends the message that initiates an
      exchange.  E.g., the party that sends a confirmable message, or a
      non-confirmable message conveying a request.

   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 [RFC2119] when they
appear in ALL CAPS.  These words may also appear in this document in
lower case as plain English words, absent their normative meanings.

(Note that this document is itself informational, but it is
discussing normative statements.)

The term "byte", abbreviated by "B", is used in its now customary
sense as a synonym for "octet".

## 2.  Context

In the Vancouver IETF 84 CoRE meeting, a path forward was defined
that includes a very simple basic scheme (lock-step with a number of
parallel exchanges of 1) in the base specification together with
performance-enhancing advanced mechanisms.

The present specification is based on the approved text in the
[I-D.ietf-core-coap] base specification.  It is making use of the
text that permits advanced congestion control mechanisms and allows
them to change protocol parameters, including NSTART and the binary
exponential backoff mechanism.  Note that Section 4.8 of
[I-D.ietf-core-coap] limits the leeway that implementations have in
changing the CoRE protocol parameters.

The present specification also assumes that, outside of exchanges,
non-confirmable messages can only be used at a limited rate without
an advanced congestion control mechanism (this is mainly relevant for
-observe).  It is also intended to address the [RFC5405] guideline
about combining congestion control state for a destination; and to
clarify its meaning for CoAP using the definition of an endpoint.

The present specification does not address multicast or dithering
beyond basic retransmission dithering.

## 3.  Advanced CoAP Congestion Control: RTO Estimation

For an initiator that plans to make multiple requests to one
destination endpoint, it may be worthwhile to make RTT measurements
in order to obtain a better RTO estimation than that implied by the
default initial timeout of 2 to 3 s.  This is based on the usual
algorithms for RTO estimation [RFC6298], with appropriately extended
default/base values, as proposed in Section 3.2.1.  Note that such a
mechanism must, during idle periods, decay RTT estimates that are
shorter than the basic RTT estimate back to the basic RTT estimate,
until fresh measurements become available again, as proposed in
Section 3.3.

One important consideration not relevant for TCP is the fact that a
CoAP round-trip may include application processing time, which may be
hard to predict, and may differ between different resources available
at the same endpoint.  Servers will only trigger their early ACKs
(with a non-piggybacked response to be sent later) based on the
default timers, e.g. after 1 s. A client that has arrived at a RTO
estimate shorter than 1 s SHOULD therefore use a larger backoff
factor for retransmissions to avoid expending all of its
retransmissions in the default interval of 2 to 3 s.  A proposal for
a mechanism with variable backoff factors is presented in
Section 3.2.1.

It may also be worthwhile to do RTT estimates not just based on
information measured from a single destination endpoint, but also
based on entire hosts (IP addresses) and/or complete prefixes (e.g.,
maintain an RTT estimate for a whole /64).  The exact way this can be
used to reduce the amount of state in an initiator is for further
study.

## 3.1.  Blind RTO Estimate

The initial RTO estimate for an endpoint is set to 2 seconds.

If only the initial RTO estimate is available, the RTO estimate for
each of up to NSTART exchanges started in parallel is set to 1 s
times 2 to the power of the number of parallel exchanges, e.g. if two
exchanges are already running, the initial RTO estimate for an
additional exchange is 8 seconds.

The binary exponential backoff is truncated at 32 seconds.  Similar
to the way retransmissions are handled in the base specification,
they are dithered between 1 x RTO and ACK_RANDOM_FACTOR x RTO.

## 3.2.  Measured RTO Estimate

The RTO estimator runs two copies of the algorithm defined in
[RFC6298], as modified in Section 3.2.1: One copy for exchanges that
complete on initial transmissions (the "strong estimator"), and one
copy for exchanges that have run into retransmissions (the "weak
estimator").  For the latter, there is some ambiguity whether a
response is based on the initial transmission or any retransmission.
For the purposes of the weak estimator, the time from the initial
transmission counts.

The overall RTO estimate is a exponentially weighted moving average
(alpha = 0.5) computed of the strong and the weak estimator, which is
evolved after each contribution to the weak estimator or to the
strong estimator, from the estimator that made the most recent
contribution:

    RTO_overall_ := 0.5 * RTO_recent_ + 0.5 * RTO_overall_

(Note that the contributionof the weak estimator may be too big in
naturally lossy networks.  TBD.)

3.2.1.  Modifications to the algorithm of RFC 6298

   This subsection presents three modifications that must be applied to
   the algorithm of [RFC6298] as per this document.  The first two
   recommend new parameter settings.  The third one is the variable
   backoff factor mechanism.

   The initial value for each of the two RTO estimators is 2 s.

   For the weak estimator, the factor K (the RTT variance multiplier) is
   set to 1 instead of 4.  This is necessary to avoid a strong increase
   of the RTO in the case that the RTTVAR value is very large, which may
   be the case if a weak RTT measurement is obtained after one or more
   retransmissions.

   If an RTO estimation is lower than 1 s or higher than 8 s, instead of
   applying a binary backoff factor in both cases, a variable backoff
   factor is used.  For RTO estimations below 1 s, the RTO for a
   retransmission is multiplied by 3, while for estimations above 8 s,
   the RTO is multiplied only by 1.3 (this initial choice of numbers to
   be verified by more simulations).  This helps to avoid that exchanges
   with small initial RTOs use up all retransmissions in a short
   interval of time and exchanges with large initial RTOs may not be
   able to carry out all retransmissions within MAX_TRANSMIT_WAIT
   (93 s).

3.2.2.  Discussion

   In contrast to [RFC6298], this algorithm attempts to make use of
   ambiguous information from retransmissions.  This is motivated by the
   high non-congestion loss rates expected in constrained node networks,
   and the need to update the RTO estimators even in the presence of
   loss.  Additional investigation is required to determine whether this
   is indeed justified.

3.3.  Lifetime, Aging

The state of the RTO estimators for an endpoint SHOULD be kept as long as possible.  If other state is kept for the endpoint (such as a DTLS connection), it is very strongly RECOMMENDED to keep the RTO state alive at least as long as this other state.  It MUST be kept for at least 255 s.

If an estimator has a value that is lower than 1 s, and it is left without further update for a time that is more than 16 times its current value, its value is doubled.

(It is allowed to implement this cumulatively at the time it is used next, possibly approximating multiple doublings by replacing the value with 1/8th of the time that has elapsed since the last update. Alternatively, simple estimators can be simply updated to 1 s after being without update for a time that is more than 16 times its value, or, even simpler, be clamped at 1 s or above.)

4.  Advanced CoAP Congestion Control: Non-Confirmables

(TO DO: Align this with final consensus on -observe!)

A CoAP endpoint MUST NOT send non-confirmables to another CoAP endpoint at a rate higher than defined by this document.  Independent of any congestion control mechanisms, a CoAP endpoint can always send non-confirmables if their rate does not exceed 1 B/s.

Non-confirmables that form part of exchanges are governed by the rules for exchanges.

Non-confirmables outside exchanges (e.g., [I-D.ietf-core-observe] notifications sent as non-confirmables) are governed by the following rules:

1.  Of any 16 consecutive messages towards this endpoint that aren't responses or acknowledgments, at least 2 of the messages must be confirmable.

2.  The confirmable messages must be sent under an RTO estimator, as specified in Section 3.

3.  The packet rate of non-confirmable messages cannot exceed 1/RTO, where RTO is the overall RTO estimator value at the time the non-confirmable packet is sent.

4.1.  Discussion

This is relatively conservative.  More advanced versions of this
algorithm could run a TFRC-style Loss Event Rate calculator [RFC5348]
and apply the TCP equation to achieve a higher rate than 1/RTO.

5.  IANA Considerations

This document makes no requirements on IANA.  (This section to be
removed by RFC editor.)

6.  Security Considerations

(TBD.  The security considerations of, e.g., [RFC5681], [RFC2914],
and [RFC5405] apply.  Some issues are already discussed in the
security considerations of [I-D.ietf-core-coap].)

7.  Acknowledgements

The first document to examine CoAP congestion control issues in
detail was [I-D.eggert-core-congestion-control], to which this draft
owes a lot.

Michael Scharf did a review of CoAP congestion control issues that
asked a lot of good questions.  Several Transport Area
representatives made further significant inputs this discussion
during IETF84, including Lars Eggert, Michael Scharf, and David
Black.  Andrew McGregor, Eric Rescorla, Richard Kelsey, Ed Beroset,
Jari Arkko, Zach Shelby, Matthias Kovatsch and many others provided
very useful additions.  Carles Gomez, August Betzler and Ilker
Demirkol provided a first detailed review of the present document;
August provided simulation results that shaped some of the
recommendations here but also indicate a need for further effort.

8.  References

8.1.  Normative References

[RFC2119]  Bradner, S., "Key words for use in RFCs to Indicate
           Requirement Levels", BCP 14, RFC 2119, March 1997.

[RFC2914]  Floyd, S., "Congestion Control Principles", BCP 41, RFC
           2914, September 2000.

[RFC5405]  Eggert, L. and G. Fairhurst, "Unicast UDP Usage Guidelines
           for Application Designers", BCP 145, RFC 5405, November
           2008.

   [RFC6298]  Paxson, V., Allman, M., Chu, J., and M. Sargent,
              "Computing TCP's Retransmission Timer", RFC 6298, June
              2011.

8.2.  Informative References

   [I-D.bormann-core-congestion-control]
              Bormann, C. and K. Hartke, "Congestion Control Principles
              for CoAP", draft-bormann-core-congestion-control-02 (work
              in progress), July 2012.

   [I-D.eggert-core-congestion-control]
              Eggert, L., "Congestion Control for the Constrained
              Application Protocol (CoAP)", draft-eggert-core-
              congestion-control-01 (work in progress), January 2011.

   [I-D.ietf-core-coap]
              Shelby, Z., Hartke, K., and C. Bormann, "Constrained
              Application Protocol (CoAP)", draft-ietf-core-coap-18
              (work in progress), June 2013.

   [I-D.ietf-core-observe]
              Hartke, K., "Observing Resources in CoAP", draft-ietf-
              core-observe-11 (work in progress), October 2013.

   [RFC5348]  Floyd, S., Handley, M., Padhye, J., and J. Widmer, "TCP
              Friendly Rate Control (TFRC): Protocol Specification", RFC
              5348, September 2008.

   [RFC5681]  Allman, M., Paxson, V., and E. Blanton, "TCP Congestion
              Control", RFC 5681, September 2009.

Author's Address

   Carsten Bormann
   Universitaet Bremen TZI
   Postfach 330440
   Bremen  D-28359
   Germany

   Phone: +49-421-218-63921
   Email: cabo@tzi.org