CoRE Working Group                                          C. Bormann
Internet-Draft                                    Universitaet Bremen TZI
Intended status: Informational                              A. Betzler
Expires: January 9, 2017                                  Fundacio i2CAT
                                                             C. Gomez
                                                           I. Demirkol
                        Universitat Politecnica de Catalunya/Fundacio i2CAT
                                                          July 08, 2016

                    CoAP Simple Congestion Control/Advanced
                        draft-bormann-core-cocoa-04

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.  It is making use of input from
   simulations and experiments in real networks.

Status of This Memo

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 [RFC7252].  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
   [RFC7252] 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 [RFC7252]
   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
   [RFC7641]).  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.  Area of Applicability

   The present algorithm is intended to be generally applicable.  The
   objective is to be "better" than default CoAP congestion control in a
   number of characteristics, including achievable goodput for a given
   offered load, latency, and recovery from bursts, while providing more
   predictable stress to the network and the same level of safety from
   catastrophic congestion.  It does require three state variables per
   scope plus the state needed to do RTT measurements, so it may not be
   applicable to the most constrained devices (class 1 as per
   [RFC7228]).

   The scope of each instance of the algorithm in the current set of
   evaluations has been the five-tuple, i.e., CoAP + endpoint (transport
   address) for Initiator and Responder.  Potential applicability to
   larger scopes needs to be examined.

   Aggregate Congestion Control (Appendix A) is not yet supported by
   research as well as the other algorithms in this specification.  Its
   use is more interesting on the cloud side, where a single CoAP
   endpoint may need to talk to thousands of other endpoints and may
   need to control the burstiness of the resulting aggregate traffic.

4.  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 4.2.1.  Note that such a
   mechanism must, during idle periods, decay RTO estimates that are
   shorter or longer than the basic RTO estimate back to the basic RTO
   estimate, until fresh measurements become available again, as
   proposed in Section 4.3.

   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.  Also, for communications with networks of
   constrained devices that apply radio duty cycling, large and variable
   round-trip times are likely to be observed.  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 4.2.1.

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.

## 4.1.  Blind RTO Estimate

The initial RTO estimate for an endpoint is set to 2 seconds (the
initial RTO estimate is used as the initial value for both E_weak_
and E_strong_ below).

If only the initial RTO estimate is available, the RTO estimate for
each of up to NSTART exchanges started in parallel is set to 2 s
times the number of parallel exchanges, e.g. if two exchanges are
already running, the initial RTO estimate for an additional exchange
is 6 seconds.

## 4.2.  Measured RTO Estimate

The RTO estimator runs two copies of the algorithm defined in
[RFC6298], as modified in Section 4.2.1: One copy for exchanges that
complete on initial transmissions (the "strong estimator",
E_strong_), and one copy for exchanges that have run into
retransmissions, where only the first two retransmissions are
considered (the "weak estimator", E_weak_).  For the latter, there is
some ambiguity whether a response is based on the initial
transmission or the retransmissions.  For the purposes of the weak
estimator, the time from the initial transmission counts.  Responses
obtained after the third retransmission are not used to update an
estimator.

The overall RTO estimate is an exponentially weighted moving average
(alpha = 0.5 and 0.25, respectively) computed of the strong and the
weak estimator, which is evolved after each contribution to the weak
estimator (1) or to the strong estimator (2), from the estimator that
made the most recent contribution:

$$RTO := 0.25 * E\_weak\_ + 0.75 * RTO \quad (1)$$

$$RTO := 0.5 * E\_strong\_ + 0.5 * RTO \quad (2)$$

(Splitting this update into the two cases avoids making the
contribution of the weak estimator too big in naturally lossy
networks.)

4.2.1.  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 3 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 3 s,
the RTO is multiplied only by 1.5 (this updated 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).

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.

4.2.2.  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.

Some evaluation has been done on earlier versions of this
specification [Betzler2013].  A more recent (and more comprehensive)
reference is [Betzler2015].  Additional investigation is required.

4.3.  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 16 times its current value, the RTO
   estimate is doubled.  If an estimator has a value that is higher than
   3 s, and it is left without further update for 4 times its current
   value, the RTO estimate is set to be

      RTO := 1 s + (0.5 * RTO)

   (Note that, instead of running a timer, it is possible to implement
   these RTO aging calculations cumulatively at the time the estimator
   is used next.)

5.  Advanced CoAP Congestion Control: Non-Confirmables

   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., [RFC7641] 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 4.

   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.

5.1.  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.

   [RFC7641], Section 4.5.1, specifies that the rate of NONs SHOULD NOT
   exceed 1/RTT on average, if the server can maintain an RTT estimate
   for a client.  CoCoA limits the packet rate of NONs in this situation
   to 1/RTO.  Assuming that the RTO estimation in CoCoA works as
   expected, RTO[k] should be slightly greater than the RTT[k], thus
   CoCoA would be more conservative.  The expectation therefore is that
   complying with the NON rate set by CoCoA leads to complying with
   [RFC7641].

6.  IANA Considerations

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

7.  Security Considerations

   (TBD.  The security considerations of, e.g., [RFC5681], [RFC2914],
   and [RFC5405] apply.  Some issues are already discussed in the
   security considerations of [RFC7252].)

8.  References

8.1.  Normative References

   [RFC2119]  Bradner, S., "Key words for use in RFCs to Indicate
              Requirement Levels", BCP 14, RFC 2119,
              DOI 10.17487/RFC2119, March 1997,
              <http://www.rfc-editor.org/info/rfc2119>.

   [RFC2914]  Floyd, S., "Congestion Control Principles", BCP 41,
              RFC 2914, DOI 10.17487/RFC2914, September 2000,
              <http://www.rfc-editor.org/info/rfc2914>.

   [RFC5405]  Eggert, L. and G. Fairhurst, "Unicast UDP Usage Guidelines
              for Application Designers", BCP 145, RFC 5405,
              DOI 10.17487/RFC5405, November 2008,
              <http://www.rfc-editor.org/info/rfc5405>.

   [RFC6298]  Paxson, V., Allman, M., Chu, J., and M. Sargent,
              "Computing TCP's Retransmission Timer", RFC 6298,
              DOI 10.17487/RFC6298, June 2011,
              <http://www.rfc-editor.org/info/rfc6298>.

   [RFC7252]  Shelby, Z., Hartke, K., and C. Bormann, "The Constrained
              Application Protocol (CoAP)", RFC 7252,
              DOI 10.17487/RFC7252, June 2014,
              <http://www.rfc-editor.org/info/rfc7252>.

8.2.  Informative References

   [Betzler2013]
              Betzler, A., Gomez, C., Demirkol, I., and J. Paradells,
              "Congestion control in reliable CoAP communication",
              ACM MSWIM'13 p. 365-372, DOI 10.1145/2507924.2507954,
              2013.

   [Betzler2015]
              Betzler, A., Gomez, C., Demirkol, I., and J. Paradells,
              "CoCoA+: an Advanced Congestion Control Mechanism for
              CoAP", Ad Hoc Networks Vol. 33 pp. 126-139,
              DOI 10.1016/j.adhoc.2015.04.007, October 2015.

   [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.

   [RFC5348]  Floyd, S., Handley, M., Padhye, J., and J. Widmer, "TCP
              Friendly Rate Control (TFRC): Protocol Specification",
              RFC 5348, DOI 10.17487/RFC5348, September 2008,
              <http://www.rfc-editor.org/info/rfc5348>.

   [RFC5681]  Allman, M., Paxson, V., and E. Blanton, "TCP Congestion
              Control", RFC 5681, DOI 10.17487/RFC5681, September 2009,
              <http://www.rfc-editor.org/info/rfc5681>.

   [RFC7228]  Bormann, C., Ersue, M., and A. Keranen, "Terminology for
              Constrained-Node Networks", RFC 7228,
              DOI 10.17487/RFC7228, May 2014,
              <http://www.rfc-editor.org/info/rfc7228>.

   [RFC7641]  Hartke, K., "Observing Resources in the Constrained
              Application Protocol (CoAP)", RFC 7641,
              DOI 10.17487/RFC7641, September 2015,
              <http://www.rfc-editor.org/info/rfc7641>.

Appendix A.  Advanced CoAP Congestion Control: Aggregate Congestion
            Control

   (The mechanism defined in this appendix has received less research
   than the ones in the main body of this specification.)

A.1.  Proposed Algorithm

   To avoid possible congestion when sending many packets to different
   destination endpoints in parallel, the overall number of outstanding
   interactions towards different destination endpoints should be
   limited.  An upper limit PLIMIT determines the maximum number of
   outstanding interactions towards different destination endpoints that
   are allowed in parallel.  When a request is to be sent to a
   destination endpoint, PLIMIT is determined according to Equation (3)
   in the case that no RTO information is already available for the
   destination endpoint, or using Equation (4) in case that valid RTO
   information is available for the destination endpoint.  Both formulas
   use LAMBDA, as defined in Equation (5).

      PLIMIT = LAMBDA        (3)

      PLIMIT = max(LAMBDA, LAMBDA*ACK_TIMEOUT/mean(RTO))       (4)

      LAMBDA = max(4, KNOWN_DEST_ENDPOINTS/4)        (5)

   mean(RTO) is the average value of all valid RTO estimations
   maintained by the device.  LAMBDA is the maximum of a constant value
   (4 by default) and the rounded up value of KNOWN_DEST_ENDPOINTS/4 ,
   where KNOWN_DEST_ENDPOINTS is the overall number of "known"
   destination endpoints (i.e. destination endpoints for which an RTO
   estimate is maintained).

   A new interaction may only be processed if the current overall number
   of outstanding interactions is lower than the PLIMIT calculated when
   the request is initiated.

A.2.  Example 1

   In the following we give an example, with LAMBDA = 4 (our proposed
   default LAMBDA):

   Assume that a sender has so far obtained RTO estimations for two
   destination endpoints A (RTO = 0.5 s) and B (RTO = 1.5 s), and
   currently pcount (a variable which accounts for the number of
   outstanding interactions towards endpoints) is equal to 0.  Now three
   transactions are initiated consecutively in the following order: one
   for A, one for B and one for a new destination C.

When an interaction with node A is initiated, LAMBDA is calculated

    LAMBDA = max(4, 3/4) = 4.

Then PLIMIT is calculated:

    PLIMIT = max(4, (4*2 s)/mean(0.5 s, 1.5 s)) = max (4, 8 s/1 s) =
    max (4, 8) = 8

This means that with the current RTO information the sender has
obtained about the destination endpoints, up to 8 outstanding
interactions to different destination endpoints would be allowed.  By
initiating an interaction with A, pcount is increased to 1, which is
still below PLIMIT.  Thus, the interaction may be processed.  The
same applies to B: pcount increases to 2 after obtaining the same
PLIMIT value of 8.

Destination C is unknown to CoCoA, therefore the updated PLIMIT
before processing the interaction with node C is 4.  The CoAP request
may be processed (pcount = 3).  If two more interactions with
different unknown destination endpoints would have been initiated,
only the first one would have met the requirements to process it
(PLIMIT = 4, pcount = 4).  The second interaction would have
increased pcount to 5, which is not permitted, since PLIMIT is 4.  It
may occur that pcount exceeds PLIMIT in particular cases, in this
case, the interaction is not permitted as well.  If the number of
destinations exchanges are initiated with would increase further,
eventually LAMBDA could grow beyond 4, allowing for more interactions
to be sent in parallel.

A.3.  Example 2

Let us now assume that a sender has so far obtained RTO estimations
for 101 destination endpoints, their average RTO is 1 s, and
currently pcount is equal to 0.  When a new transaction is initiated
with a destination endpoint for which an RTO estimate is available,
LAMBDA is calculated

    LAMBDA = max(4, 101/4) = 26

Based on this, PLIMIT is calculated as follows:

    PLIMIT= max(26, (26*2 s)/1 s) = max (26, 52) = 52

This means that with the current RTO information that the sender has
obtained about the destination endpoints, up to 52 outstanding
interactions to known destination endpoints would be allowed.

However, if the new exchange is to be initiated with an "unknown"
destination endpoint (i.e. an endpoint for which an RTO estimate is
not available), then PLIMIT would be 26.

A.4.  Discussion

The idea of the proposal is to allow more parallel transactions to
different destination endpoints if we have low RTO estimations for
them (which can be interpreted as good connections and low degree of
congestion).  If the RTO estimations are large or interactions with
unknown destinations are initiated, the mechanism behaves more
conservatively by reducing the maximum number of parallel
interactions towards different destinations, but allowing at least
LAMBDA outstanding interactions.  The second term of the max()
statement used to calculate LAMBDA avoids behaving too restrictively
when exchanges with many different destination endpoints are
initiated.  If no RTO information is available for a destination
endpoint, PLIMIT is simply set to be LAMBDA.

If at any moment pcount would exceed PLIMIT, CoAP does not
immediately perform the transaction.  Further, it is important that
in parallel, NSTART for each destination endpoint applies (which, for
now, we assume to be 1).  The default value used for LAMBDA (equal to
4 as per this document) determines how aggressive/conservative CoCoA
behaves by default for a limited set of destination endpoints and it
should be chosen carefully.  The term KNOWN_DEST_ENDPOINTS/4 loosens
the hard limit of exchanges when large numbers of destination
endpoints are addressed.

It will be necessary to see whether this approach is effective in the
sense that it avoids congestion in use cases where transactions to a
multitude of different destination endpoints are initiated.  An
important aspect of such evaluations would be whether LAMBDA is too
conservative when dealing with few destination endpoints and whether
it allows for a dynamic adjustment of parallel exchanges with large
numbers of destination endpoints.  On the other hand, a more safe
approach would use max(RTO) instead of mean(RTO).  Other concerns
include the fact that the congestion degree of the paths to "known"
destination endpoints influence whether a new interaction is
permitted to some new endpoint which may be in very different
conditions in terms of congestion.  However, it is desirable to avoid
adding a lot of complexity to the current CoCoA mechanisms.

Authors' Addresses

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

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


   August Betzler
   Fundacio i2CAT
   Mobile and Wireless Internet Group
   C/ del Gran Capita, 2
   Barcelona  08034
   Spain

   Email: august.betzler@entel.upc.edu

Carles Gomez
Universitat Politecnica de Catalunya/Fundacio i2CAT
Escola d'Enginyeria de Telecomunicacio i Aeroespacial
de Castelldefels
C/Esteve Terradas, 7
Castelldefels  08860
Spain

Phone: +34-93-413-7206
Email: carlesgo@entel.upc.edu


Ilker Demirkol
Universitat Politecnica de Catalunya/Fundacio i2CAT
Departament d'Enginyeria Telematica
C/Jordi Girona, 1-3
Barcelona  08034
Spain

Email: ilker.demirkol@entel.upc.edu