

Network Working Group
Internet-Draft
Intended status: Experimental
Expires: May 3, 2017

N. Khademi
M. Welzl
University of Oslo
G. Armitage
Swinburne University of
Technology
G. Fairhurst
University of Aberdeen
October 30, 2016

TCP Alternative Backoff with ECN (ABE)
draft-khademi-tcpm-alternativebackoff-ecn-01

Abstract

This memo updates the TCP sender-side reaction to a congestion notification received via Explicit Congestion Notification (ECN). The updated method reduces FlightSize in Congestion Avoidance by a smaller amount than the TCP reaction to loss. The intention is to achieve good throughput when the queue at the bottleneck is smaller than the bandwidth-delay-product of the connection. This is more likely when an Active Queue Management (AQM) mechanism has used ECN to CE-mark a packet, than when a packet was lost. Future versions of this document will also describe a corresponding method for SCTP.

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 May 3, 2017.

Copyright Notice

Copyright (c) 2016 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 (<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.

Table of Contents

1. Definitions	3
2. Introduction	3
3. Discussion	4
3.1. Why Use ECN to Vary the Degree of Backoff?	4
3.2. Focus on ECN as Defined in RFC3168	5
3.3. Discussion: Choice of ABE Multiplier	5
4. Specification	6
5. Status of the Update	6
6. Acknowledgements	6
7. IANA Considerations	7
8. Implementation Status	7
9. Security Considerations	7
10. Revision Information	7
11. References	8
11.1. Normative References	8
11.2. Informative References	8
Authors' Addresses	10

1. 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].

2. Introduction

Complementing [I-D.AQM-ECN-benefits], [I-D.ECN-exp] enables wider ECN deployment by updating rules in [RFC3168] that prohibited certain experiments. Specifically, [I-D.ECN-exp] allows for experiments to specify a congestion control response to a CE-marked packet that differs from the response to a dropped packet. This memo defines such a different congestion control response, called "ABE" (Alternative Backoff with ECN). ABE is thus an Experiment in accordance with [I-D.ECN-exp].

[RFC5681] stipulates that TCP congestion control sets "sssthresh" to $\max(\text{FlightSize} / 2, 2 * \text{SMSS})$ in response to packet loss. This corresponds to a backoff multiplier of 0.5 (halving cwnd and sssthresh after packet loss). Consequently, a standard TCP flow using this reaction needs significant network queue space: it can only fully utilise a bottleneck when the length of the link queue (or the AQM dropping threshold) is at least the bandwidth-delay product (BDP) of the flow.

A backoff multiplier of 0.5 is not the only available strategy. As defined in [I-D.CUBIC], CUBIC multiplies the current cwnd by 0.7 in response to loss (the Linux implementation of CUBIC has used a multiplier of 0.7 since kernel version 2.6.25 released in 2008). Consequently, CUBIC utilises paths well even when the bottleneck queue is shorter than the bandwidth-delay product of the flow. However, in the case of a DropTail (FIFO) queue without AQM, such less-aggressive backoff increases the risk of creating a standing queue [CODEL2012].

The standard TCP backoff behaviour defined in [RFC5681] entails reduced link utilisation in situations with short queues and low statistical multiplexing. This memo proposes a concrete sender-side-only congestion control response that remedies this problem.

Devices implementing AQM are likely to be the dominant (and possibly only) source of ECN CE-marking for packets from ECN-capable senders. AQM mechanisms typically strive to maintain a small average queue length, regardless of the bandwidth-delay product of flows passing through them. Receipt of an ECN CE-mark might therefore reasonably be taken to indicate that a small bottleneck queue exists in the

path, and hence the TCP flow would benefit from using a less aggressive backoff multiplier.

Much of the background to this proposal can be found in [ABE2015]. Using a mix of experiments, theory and simulations with standard NewReno and CUBIC, [ABE2015] recommends enabling ECN and letting individual TCP senders use a larger multiplicative decrease factor as a reaction to the receiver reporting ECN CE-marks from AQM-enabled bottlenecks. Such a change is noted to result in "...significant performance gains in lightly-multiplexed scenarios, without losing the delay-reduction benefits of deploying CoDel or PIE" [I-D.CoDel] [I-D.PIE]. This is achieved when reacting to ECN-Echo in Congestion Avoidance by multiplying `cwnd` and `ssthresh` with a value in the range [0.7..0.85].

3. Discussion

3.1. Why Use ECN to Vary the Degree of Backoff?

The classic rule-of-thumb dictates that a transport provides a BDP of bottleneck buffering if a TCP connection wishes to optimise path utilisation. A single TCP connection running through such a bottleneck will have opened `cwnd` up to $2 \times \text{BDP}$ by the time packet loss occurs. [RFC5681]'s halving of `cwnd` and `ssthresh` pushes the TCP connection back to allowing only a BDP of packets in flight -- just sufficient to maintain 100% utilisation of the network path.

AQM schemes like CoDel [I-D.CoDel] and PIE [I-D.PIE] use congestion notifications to constrain the queuing delays experienced by packets, rather than in response to impending or actual bottleneck buffer exhaustion. With current default delay targets, CoDel and PIE both effectively emulate a shallow buffered bottleneck (section II, [ABE2015]) while allowing short traffic bursts into the queue. This interacts acceptably for TCP connections over low BDP paths, or highly multiplexed scenarios (many concurrent TCP connections). However, it interacts badly with lightly-multiplexed cases (few concurrent connections) over a high BDP path. Conventional TCP backoff in such cases leads to gaps in packet transmission and under-utilisation of the path.

The idea to react differently to loss upon detecting an ECN CE-mark pre-dates [ABE2015]. [ICC2002] also proposed using ECN CE-marks to modify TCP congestion control behaviour, using a larger multiplicative decrease factor in conjunction with a smaller additive increase factor to work with RED-based bottlenecks that were not necessarily configured to emulate a shallow queue.

3.2. Focus on ECN as Defined in RFC3168

Some mechanisms rely on ECN semantics that differ from the definitions in [RFC3168] -- for example, Congestion Exposure (ConEx) [RFC7713] and DCTCP [I-D.ietf-tcpm-dctcp] need more accurate ECN information than the feedback mechanism in [RFC3168] offers (defined in [I-D.ietf-tcpm-accurate-ecn]). Such mechanisms allow a sending rate adjustment more frequent than each RTT. These mechanisms are out of the scope of the current document.

3.3. Discussion: Choice of ABE Multiplier

Alternative Backoff with ECN (ABE) decouples a TCP sender's reaction to loss and ECN CE-marks in Congestion Avoidance. The description respectively uses β_{loss} and β_{ecn} to refer to the multiplicative decrease factors applied in response to packet loss, and also in response to a receiver indicating that an ECN CE-mark was received on an ECN-enabled TCP connection (based on the terms used in [ABE2015]). For non-ECN-enabled TCP connections, no ECN CE-marks are received and only β_{loss} applies.

In other words, in response to detected loss:

$$\text{FlightSize}_{(n+1)} = \text{FlightSize}_n * \beta_{\text{loss}}$$

and in response to an indication of a received ECN CE-mark:

$$\text{FlightSize}_{(n+1)} = \text{FlightSize}_n * \beta_{\text{ecn}}$$

where, as in [RFC5681], FlightSize is the amount of outstanding data in the network, upper-bounded by the sender's congestion window (cwnd) and the receiver's advertised window (rwnd). The higher the values of β_{loss} and β_{ecn} , the less aggressive the response of any individual backoff event.

The appropriate choice for β_{loss} and β_{ecn} values is a balancing act between path utilisation and draining the bottleneck queue. More aggressive backoff (smaller $\beta_{\text{*}}$) risks underutilising the path, while less aggressive backoff (larger $\beta_{\text{*}}$) can result in slower draining of the bottleneck queue.

The Internet has already been running with at least two different β_{loss} values for several years: the value in [RFC5681] is 0.5, and Linux CUBIC uses 0.7. ABE proposes no change to β_{loss} used by any current TCP implementations.

β_{ecn} depends on how the response of a TCP connection to shallow AQM marking thresholds is optimised. β_{loss} reflects the

preferred response of each TCP algorithm when faced with exhaustion of buffers (of unknown depth) signalled by packet loss. Consequently, for any given TCP algorithm the choice of β_{ecn} is likely to be algorithm-specific, rather than a constant multiple of the algorithm's existing β_{loss} .

A range of experiments (section IV, [ABE2015]) with NewReno and CUBIC over CoDel and PIE in lightly-multiplexed scenarios have explored this choice of parameter. These experiments indicate that CUBIC connections benefit from β_{ecn} of 0.85 (cf. $\beta_{\text{loss}} = 0.7$), and NewReno connections see improvements with β_{ecn} in the range 0.7 to 0.85 (cf. $\beta_{\text{loss}} = 0.5$).

4. Specification

This document RECOMMENDS that experimental deployments multiply the FlightSize by 0.8 and reduce the slow start threshold 'ssthresh' in Congestion Avoidance in response to reception of a TCP segment that sets the ECN-Echo flag.

5. Status of the Update

This update is a sender-side only change. Like other changes to congestion-control algorithms it does not require any change to the TCP receiver or to network devices (except to enable an ECN-marking algorithm [RFC3168] [RFC7567]). If the method is only deployed by some TCP senders, and not by others, the senders that use this method can gain advantage, possibly at the expense of other flows that do not use this updated method. This advantage applies only to ECN-marked packets and not to loss indications. Hence, the new method can not lead to congestion collapse.

The present specification has been assigned an Experimental status, to provide Internet deployment experience before being proposed as a Standards-Track update.

6. Acknowledgements

Authors N. Khademi, M. Welzl and G. Fairhurst were part-funded by the European Community under its Seventh Framework Programme through the Reducing Internet Transport Latency (RITE) project (ICT-317700). The views expressed are solely those of the authors.

The authors would like to thank the following people for their contributions to [ABE2015]: Chamil Kulatunga, David Ros, Stein

Gjessing, Sebastian Zander. Thanks to (in alphabetical order) Bob Briscoe, Markku Kojo, John Leslie, Dave Taht and the TCPM WG for providing valuable feedback on this document.

The authors would like to thank feedback on the congestion control behaviour specified in this update received from the IRTF Internet Congestion Control Research Group (ICCRG).

7. IANA Considerations

XX RFC ED - PLEASE REMOVE THIS SECTION XXX

This memo includes no request to IANA.

8. Implementation Status

ABE is implemented as a patch for Linux and FreeBSD. It is meant for research and available for download from <http://heim.ifi.uio.no/naeemk/research/ABE/> This code was used to produce the test results that are reported in [ABE2015].

9. Security Considerations

The described method is a sender-side only transport change, and does not change the protocol messages exchanged. The security considerations of [RFC3168] therefore still apply.

This document describes a change to TCP congestion control with ECN that will typically lead to a change in the capacity achieved when flows share a network bottleneck. Similar unfairness in the way that capacity is shared is also exhibited by other congestion control mechanisms that have been in use in the Internet for many years (e.g., CUBIC [I-D.CUBIC]). Unfairness may also be a result of other factors, including the round trip time experienced by a flow. This advantage applies only to ECN-marked packets and not to loss indications, and will therefore not lead to congestion collapse.

10. Revision Information

XX RFC ED - PLEASE REMOVE THIS SECTION XXX

-01. This I-D now refers to draft-black-tsvwg-ecn-experimentation-02, which replaces draft-khademi-tsvwg-ecn-response-00 to make a broader update to

RFC3168 for the sake of allowing experiments. As a result, some of the motivating and discussing text that was moved from draft-khademi-alternativebackoff-ecn-03 to draft-khademi-tsvwg-ecn-response-00 has now been re-inserted here.

-00. draft-khademi-tsvwg-ecn-response-00 and draft-khademi-tcpm-alternativebackoff-ecn-00 replace draft-khademi-alternativebackoff-ecn-03, following discussion in the TSVWG and TCPM working groups.

11. References

11.1. Normative References

- [I-D.ECN-exp] Black, D., "Explicit Congestion Notification (ECN) Experimentation", Internet-draft, IETF work-in-progress draft-black-tsvwg-ecn-experimentation-02, October 2016.
- [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>>.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, DOI 10.17487/RFC3168, September 2001, <<http://www.rfc-editor.org/info/rfc3168>>.
- [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>>.
- [RFC7567] Baker, F., Ed. and G. Fairhurst, Ed., "IETF Recommendations Regarding Active Queue Management", BCP 197, RFC 7567, DOI 10.17487/RFC7567, July 2015, <<http://www.rfc-editor.org/info/rfc7567>>.

11.2. Informative References

- [ABE2015] Khademi, N., Welzl, M., Armitage, G., Kulatunga, C., Ros, D., Fairhurst, G., Gjessing, S., and S. Zander, "Alternative Backoff: Achieving Low Latency and High Throughput with ECN and AQM", CAIA Technical Report CAIA-TR-150710A, Swinburne University of Technology, July 2015, <<http://caia.swin.edu.au/reports/150710A/>>

CAIA-TR-150710A.pdf>.

[CODEL2012]

Nichols, K. and V. Jacobson, "Controlling Queue Delay", July 2012, <<http://queue.acm.org/detail.cfm?id=2209336>>.

[I-D.AQM-ECN-benefits]

Fairhurst, G. and M. Welzl, "The Benefits of using Explicit Congestion Notification (ECN)", Internet-draft, IETF work-in-progress draft-ietf-aqm-ecn-benefits-08, November 2015.

[I-D.CUBIC]

Rhee, I., Xu, L., Ha, S., Zimmermann, A., Eggert, L., and R. Scheffenegger, "CUBIC for Fast Long-Distance Networks", Internet-draft, IETF work-in-progress draft-ietf-tcpm-cubic-02, August 2016.

[I-D.CoDel]

Nichols, K., Jacobson, V., McGregor, V., and J. Iyengar, "Controlled Delay Active Queue Management", Internet-draft, IETF work-in-progress draft-ietf-aqm-codel-04, June 2016.

[I-D.PIE]

Pan, R., Natarajan, P., Baker, F., and G. White, "PIE: A Lightweight Control Scheme To Address the Bufferbloat Problem", Internet-draft, IETF work-in-progress draft-ietf-aqm-pie-10, September 2016.

[I-D.ietf-tcpm-accurate-ecn]

Briscoe, B., Kuehlewind, M., and R. Scheffenegger, "More Accurate ECN Feedback in TCP", draft-ietf-tcpm-accurate-ecn-01 (work in progress), June 2016.

[I-D.ietf-tcpm-dctcp]

Bensley, S., Eggert, L., Thaler, D., Balasubramanian, P., and G. Judd, "Datacenter TCP (DCTCP): TCP Congestion Control for Datacenters", draft-ietf-tcpm-dctcp-02 (work in progress), July 2016.

[ICC2002]

Kwon, M. and S. Fahmy, "TCP Increase/Decrease Behavior with Explicit Congestion Notification (ECN)", IEEE ICC 2002, New York, New York, USA, May 2002, <<http://dx.doi.org/10.1109/ICC.2002.997262>>.

[RFC7713]

Mathis, M. and B. Briscoe, "Congestion Exposure (ConEx) Concepts, Abstract Mechanism, and Requirements", RFC 7713,

DOI 10.17487/RFC7713, December 2015,
<<http://www.rfc-editor.org/info/rfc7713>>.

Authors' Addresses

Naeem Khademi
University of Oslo
PO Box 1080 Blindern
Oslo, N-0316
Norway

Email: naeemk@ifi.uio.no

Michael Welzl
University of Oslo
PO Box 1080 Blindern
Oslo, N-0316
Norway

Email: michawe@ifi.uio.no

Grenville Armitage
Centre for Advanced Internet Architectures
Swinburne University of Technology
PO Box 218
John Street, Hawthorn
Victoria, 3122
Australia

Email: garmitage@swin.edu.au

Godred Fairhurst
University of Aberdeen
School of Engineering, Fraser Noble Building
Aberdeen, AB24 3UE
UK

Email: gorry@erg.abdn.ac.uk

