

Network Working Group
Internet-Draft
Intended status: Informational
Expires: December 1, 2017

I. Johansson
Ericsson AB
May 30, 2017

ECN support in QUIC
draft-johansson-quic-ecn-03

Abstract

This memo outlines the ECN (Explicit Congestion Notification) support in QUIC. The draft specifies the ECN negotiation and the ECN echo and in addition, different aspects of fallback in case of ECN failure as well as OS specific issues with ECN and monitoring for ECN capability. The intention is that most of the material ends up updating other new or existing QUIC protocol specifications, thus it may be possible that this draft does not warrant a working group status.

Requirements Language

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].

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 December 1, 2017.

Copyright Notice

Copyright (c) 2017 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.

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.

Table of Contents

1. Introduction	2
2. Elements of ECN support	3
2.1. ECN negotiation	3
2.1.1. Challenge/Response	4
2.1.2. Determine degree of ECN support	5
2.2. ECN bits in the IP header, semantics	6
2.3. ECN echo	6
2.4. Fallback in case of ECN fault	7
2.5. OS socket specifics, access to the ECN bits	7
2.6. Monitoring	8
3. IANA Considerations	8
4. Open questions	8
5. Security Considerations	9
6. Acknowledgements	9
7. References	9
7.1. Normative References	9
7.2. Informative References	9
Author's Address	11

1. Introduction

ECN support in transport protocols is a fundamental feature that should be included in the QUIC specification as a mandatory element. ECN has the key benefit that it allows for non-destructive congestion

notification by network node, i.e packets are marked instead discarded. This is particularly beneficial for realtime applications with requirements on latency, ECN also has the benefit that it provides with a congestion signal that is unambiguous. The benefits with ECN is described in more detail in [I-D.ietf-aqm-ecn-benefits]. The ECN support should be implemented to support both present and future ECN, the latter is outlined in [I-D.ietf-tsvwg-ecn-experimentation], of particular interest is the ability to discriminate between classic ECN and L4S ECN by means of differentiation between the use of the ECT(0) and ECT(1) code points. This draft does however not delve into the details of the congestion control implementation.

2. Elements of ECN support

This draft covers the following aspects of ECN support:

- o ECN negotiation
- o ECN echo
- o ECN bits in the IP header, semantics
- o Fallback in case of ECN fault
- o OS socket specifics, access to the ECN bits
- o Monitoring

2.1. ECN negotiation

ECN support in QUIC needs to be negotiated. The reasons is that network elements may not support ECN and may either clear the ECN bits or simply discard packets that have the ECN bits set. In addition, a QUIC implementation may not have access to the ECN bits in the IP header due to OS dependent restrictions, investigations (Piers O'Hanlon) have indicated that this is in certain cases an asymmetric property, for instance while it is possible to set the ECN bits it is not possible to read them.

It is also required that the ECN negotiation does not interfere with the connection setup, in other words a failed ECN negotiation should not cause an extra roundtrip for the connection setup.

The suggested method in this draft is to send an ECN negotiation frame when connection setup is completed. Both peers MUST transmit the ECN negotiation frame. The ECN negotiation frame is shown below.

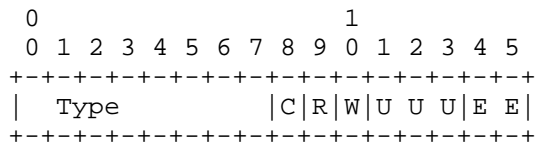


Figure 1: ECN negotiation frame

The 2nd byte contains the flags:

- o C: Challenge bit, indicates that the transmitted ECN negotiation frame is a challenge, if bit is not set then it is a response.
- o R: Possible to read ECN bits in IP header
- o W: Possible to write ECN bits in IP header
- o EE : Echo of ECN bits
- o U: Unused

The ECN negotiation has two steps.

- o Challenge/response
- o Determine degree of ECN support

2.1.1.1. Challenge/Response

A peer transmits the ECN negotiation frame with the R,W and EE bits in the 2nd byte set to '0' and the C bit set to '1'. This frame is echoed back with the flags set according to the degree of ECN support and with the ECN bits in the IP header of the received ECN negotiation frame copied to the EE field, the C bit is '0'. As both peers MUST transmit an ECN negotiation frame there will be a total of 4 ECN negotiation frames transmitted, two challenges and two responses.

An ECN negotiation frame should be transmitted in a unique packet, this to avoid that possible loss of ECN negotiation packets cause loss of other frames than the ECN negotiation frame.

The IP header for the ECN negotiation frame should set the ECN bits to CE '11'. When the corresponding response is received then an EE pattern of '11' indicates that ECN is likely supported in the network. This does not give a full guarantee that ECN is supported in the network. Monitoring of the ECN field in the ACK-frame serves to give further indication of ECN support once ECN is turned on.

An ECN negotiation is declared successful when an ECN negotiation response is received that indicates ECN support. A peer is not allowed to set ECT on outgoing data packets until a successful ECN negotiation is done. In other words it is only the ECN negotiation frame that is allowed to set the ECN bits in the IP header until ECN negotiation is concluded and successful.

A lack of an ECN negotiation response may indicate that the ECN challenge frame or the ECN response frame was lost or that a node in the network deliberately discards ECN-CE marked packets. The peer should transmit an additional ECN challenge within an RTO interval in case a negotiation response is not received, a maximum of retransmissions are attempted.

A failed challenge/response phase indicates that ECN should not be used in the connection. [NOTE, a special case is where one peer does not receive an ECN negotiation response but still receives ECT and CE marked packets from the other peer. It is T.B.D how this should be handled]

2.1.2. Determine degree of ECN support

If the ECN challenge/response is successful, the degree of ECN capability depends on how the R, W and EE bits are set.

- o R='1' and EE= '11': It is possible to set the ECN bits in outgoing packets.
- o R='0' or EE <> '11': ECN support is not certain as it is either not possible for remote peer to read the ECN bits or that the ECN bits are altered.
- o W='1' : It is meaningful to send ECN feedback
- o W='0' : It is not meaningful to send ECN feedback as the remote peer cannot set (write) the ECN bits in the IP header.

The mode mechanism in [RFC6679] can serve as input to a solution for the support of ECN in the case that OS ECN support is asymmetric. It is however unclear how a QUIC implementation can determine asymmetric ECN support in the underlying OS. For instance the method to send ECN marked packets to the local host to determine OS support does not reveal if the OS ECN support is asymmetric.

2.2. ECN bits in the IP header, semantics

The ECN bits in the IP header should be set according to the recommendations in [I-D.ietf-tsvwg-ecn-experimentation]. This means that the meaning of ECT(0) and ECT(1) differ.

2.3. ECN echo

The ECN echo should go into the ACK frame [I-D.ietf-quic-transport], this is beneficial as the ECN information can then use some of the already existing data in the ACK frame for improved efficiency.

The proposed alternative use one byte to encode how many bits that encode each of the ECT/CE fields.

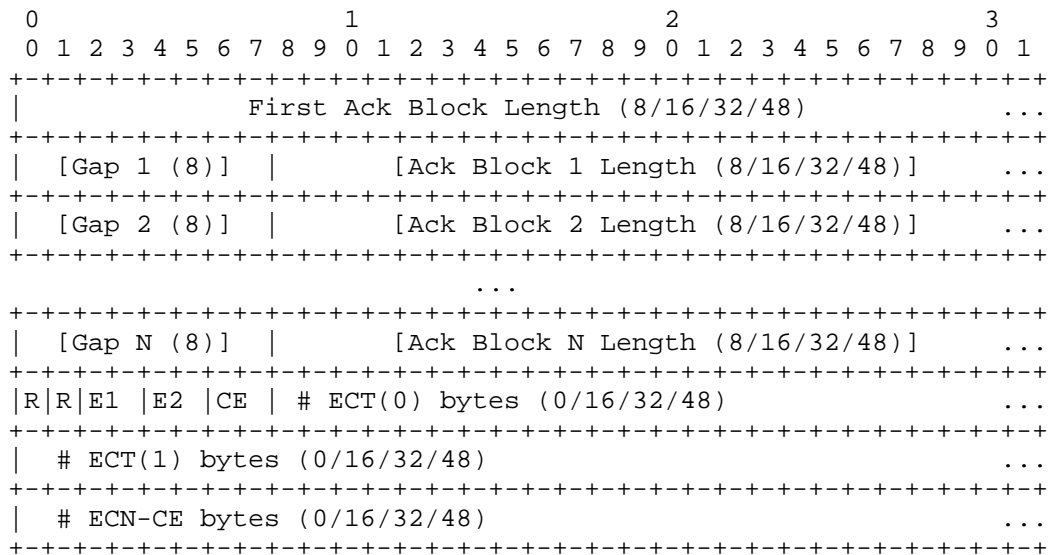


Figure 2: ECN field in ACK frame ACK block

The E1, E2 and CE fields indicate the length of each encoding for the number of ECT(0), ECT(1) and ECN-CE marked bytes. This is encoded as:

- o 00: 0 bits
- o 01: 16 bits
- o 10: 32 bits
- o 11: 48bits

R indicates reserved bits.

The proposed encoding enables flexible encoding of the ECN information, with a minimal 1 octet overhead for the cases where ECN is not supported by the connection.

2.4. Fallback in case of ECN fault

ECN can be subject to issues in network equipment, such as remarking to Not-ECN, remarking from ECT(0) to ECT(1) and vice versa or constant remarking to ECN-CE. Furthermore ECT marked packets may be discarded in the network. While these problems seem to be rare, see for instance [McQuistin-Perkins] and [APPLE-ECN], it is still necessary to safeguard against such problems.

A peer should disable ECN for its outgoing packets if ECN fault is detected, it is however still possible for the other peer to use ECN.

TODO add more information as regards to how to detect network ECN faults. [ECN-fallback](expired) gives a few examples for fault detection. Examples on how to detect ECN faults include for instance the method to set ECT and CE for outgoing packets according to a given pattern.

Fallback in case of ECN faults is not an issue only for QUIC, it is here suggested that mechanisms for this is described in a non QUIC related draft, for instance in TSVWG.

2.5. OS socket specifics, access to the ECN bits

ECN support in QUIC comes with the additional challenge that it is necessary to somehow access the ECN bits in the IP headers. In TCP this is provided without major concerns as TCP is generally implemented in OS kernel space. QUIC can however be implemented both in user space or kernel space and is layered on top of UDP, which means that access to the ECN bits is not a given, instead various tricks are needed.

The text below is copy-pasted from [OHanlon].

"To set ECN on Linux, BSD and OSX one can use IP_TOS socket option, with the setsockopt() call, to set the relevant ECN bits of the TOS byte. On Windows one can use a similar technique though firstly one has to enable TOS byte setting by enabling a particular Registry key (DisableUserTOSSetting=0 (see <https://msdn.microsoft.com/en-us/library/windows/desktop/dd874008%28v=vs.85%29.aspx> One could also probably use the libpcap write functionality."

"To obtain the ECN bits from a packet one needs a mechanism to retrieve the ECN bits from each packet. On Linux, one needs to firstly set the `IP_RECVTOS` socket option on the receiving socket, and use the `recvmsg()` call to receive a packet, and then retrieve the TOS byte from the associated `cmsg` structure returned by the `recvmsg()` call. This still works with linux-4.2.3. On OSX/BSD there are no suitable socket options to retrieve the ECN/TOS bits and one cannot use raw sockets as they do not function for UDP/TCP sockets (they do work with ICMP), so one has to use alternatives such the `bpf` interface, or a `REDIRECT` socket. Whilst on Windows it seems that the only way to retrieve the ECN bits is via a raw socket, or custom `NDIS` driver, though it's possible there's an API I'm missing."

TODO: Write a more detailed description on how to implement ECN support in QUIC for different OS stacks.

2.6. Monitoring

A QUIC implementation should monitor the ECN functionality in order to provide input to e.g. service providers to improve ECN support in the networks. Items of interest are:

- o Black holes, ECT or CE marked packets are discarded.
- o Faulty remarking, e.g. `ECT(0)` is remarked to `ECT(1)` or Not-ECT.
- o Continuous CE marking, possible indication of faulty on/off ECN marking, but can also be an effect of severe congestion.
- o Degree of L4S support. L4S should generally give low queue latency. Estimation of one way queue delay for L4S enabled QUIC connections can be used to determine if there are congested nodes along the path that are not L4S capable.

3. IANA Considerations

T.B.D.

4. Open questions

A list of open questions:

- o Is it sufficient that one peer sends an ECN negotiation challenge frame?.
- o Should all packets be ECT or should there be special patterns to improve fault detection.

- o Write up a more detailed description on how to implement ECN support in QUIC for different OS stacks.
- o Is a completely new ACK frame an alternative ?
- o Should amount on ECT(0), ECT(1) and CE marked bytes account for the IP+UDP headers or is it only the QUIC header + data that counts ?
- o Outline possible connection migration actions
- o Are there any security implications with the small ECN negotiation frame ?

5. Security Considerations

T.B.D

6. Acknowledgements

The following persons have contributed with comments and suggestions for improvements: Mirja Kuehlewind, Koen De Schepper, Piers O'Hanlon, Michael Welzl, Marcelo Bagnulo Braun, Martin Duke

7. References

7.1. 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>>.

7.2. Informative References

- [APPLE-ECN] Apple Inc., "TCP ECN: Experience with Enabling ECN on the Internet", <<https://www.ietf.org/proceedings/98/slides/slides-98-maprg-tcp-ecn-experience-with-enabling-ecn-on-the-internet-padma-bhooma-00.pdf>>.
- [Bagnulo] "Adding Explicit Congestion Notification (ECN) to TCP control packets and TCP retransmissions", <<https://tools.ietf.org/id/draft-bagnulo-tcpm-generalized-ecn-00.txt>>.

- [ECN-fallback]
"A Mechanism for ECN Path Probing and Fallback",
<<https://www.ietf.org/archive/id/draft-kuehlewind-tcpm-ecn-fallback-01.txt>>.
- [I-D.ietf-aqm-ecn-benefits]
Fairhurst, G. and M. Welzl, "The Benefits of using Explicit Congestion Notification (ECN)", draft-ietf-aqm-ecn-benefits-08 (work in progress), November 2015.
- [I-D.ietf-quic-transport]
Iyengar, J. and M. Thomson, "QUIC: A UDP-Based Multiplexed and Secure Transport", draft-ietf-quic-transport-03 (work in progress), May 2017.
- [I-D.ietf-tsvwg-ecn-experimentation]
Black, D., "Explicit Congestion Notification (ECN) Experimentation", draft-ietf-tsvwg-ecn-experimentation-02 (work in progress), April 2017.
- [McQuistin-Perkins]
"Is Explicit Congestion Notification usable with UDP?",
Proceedings of the ACM Internet Measurement Conference,
Tokyo, Japan, October 2015. DOI:10.1145/2815675.2815716",
<<https://csparks.org/publications/2015/10/mcquistin2015ecn-udp.pdf>>.
- [OHanlon] "ECN support in different OS stacks",
<<https://mailarchive.ietf.org/arch/msg/rmcat/rRKF3PvmFL2zHCplbOPKimqSsbM>>.
- [RFC6679] Westerlund, M., Johansson, I., Perkins, C., O'Hanlon, P., and K. Carlberg, "Explicit Congestion Notification (ECN) for RTP over UDP", RFC 6679, DOI 10.17487/RFC6679, August 2012, <<http://www.rfc-editor.org/info/rfc6679>>.
- [RFC6789] Briscoe, B., Ed., Woundy, R., Ed., and A. Cooper, Ed., "Congestion Exposure (ConEx) Concepts and Use Cases", RFC 6789, DOI 10.17487/RFC6789, December 2012, <<http://www.rfc-editor.org/info/rfc6789>>.
- [RFC7560] Kuehlewind, M., Ed., Scheffenegger, R., and B. Briscoe, "Problem Statement and Requirements for Increased Accuracy in Explicit Congestion Notification (ECN) Feedback", RFC 7560, DOI 10.17487/RFC7560, August 2015, <<http://www.rfc-editor.org/info/rfc7560>>.

Author's Address

Ingemar Johansson
Ericsson AB
Laboratoriegränd 11
Luleå 977 53
Sweden

Phone: +46 730783289
Email: ingemar.s.johansson@ericsson.com