TRILL Working Group INTERNET-DRAFT Intended status: Proposed Standard Donald Eastlake Huawei Bob Briscoe Simula Research Lab March 21, 2016

Expires: September 20, 2016

Abstract

Explicit congestion notification (ECN) allows a forwarding element to notify downstream devices, including the destination, of the onset of congestion without having to drop packets. This document extends this capability to TRILL switches, including integration with IP ECN.

Status of This Memo

This Internet-Draft is submitted to IETF in full conformance with the provisions of $\underline{BCP 78}$ and $\underline{BCP 79}$.

Distribution of this document is unlimited. Comments should be sent to the TRILL working group mailing list <trill@ietf.org>.

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 to cite them other than as "work in progress."

The list of current Internet-Drafts can be accessed at http://www.ietf.org/lid-abstracts.html. The list of Internet-Draft Shadow Directories can be accessed at http://www.ietf.org/shadow.html.

[Page 1]

INTERNET-DRAFT

Table of Contents

$\underline{1}$. Introduction
<u>2</u> . The ECN Specific Extended Header Flags5
3. ECN Support
4. IANA Considerations
5.Security Considerations106.Acknowledgements10
Normative References
Authors' Addresses <u>12</u>

[Page 2]

1. Introduction

Explicit congestion notification (ECN [RFC3168]) allows a forwarding element, such as a router, to notify downstream devices, including the destination, of the onset of congestion without having to drop packets. Instead, the forwarding element can explicitly mark a proportion of packets in a two-bit ECN field. For example, a two-bit field in IP headers is available for ECN marking.

The transit of user data through a TRILL campus is similar to transport through a tunnel with the ingress and egress RBridges equivalent to the ends of the tunnel. Thus, existing ECN tunneling recommendatons, particularly [<u>RFC6040</u>], apply.

		• • •		• • • •								
		•										
	+ -		+									
++		Ingre	ess									
Source	+->	RBrid	dge						+		+	
++		RB1	1						Foi	rward	ing	
	+-		-++	+-			+		E.	Lement	t	
V	1				Trans	it	1			Y	ĺ	
+	-++		+	->	RBrid	ges	1		+		-+-+	
Forwar	ding			Í	RBn	-	Ì		/	^	1	
Eleme	nt			+-		-+	+ +		+		V	
X	i					1	1	Egres	s	+-·		+
+	+					+	>	RBrid	lge +	+ De	estina	tion
							i	RBS)	+		+
			TRIL	L			+		· +			
			camp	us								

In the figure above, if ECN is implemented and assuming IP traffic, RB1 is effectivley a tunnel entrance and RB9 a tunnel exit. Traffic from Source to RB1 might or might not get marked as having experienced congestion in forwarding elements, such as X, before being encapsulated at ingress RB1. Any such ECN marking is encapsulated with a TRILL Header and provision is made in the TRILL Header extension Flags Word for ECN marking by the RBridges through which this traffic passes.

Any ECN marking in the traffic at the ingress is copied out to the TRILL Header Flags Word. At RB9, the TRILL egress, any ECN markings in the TRILL Header Flags Word and in the encapsulated traffic are combined so that subsequent forwarding elements, such as Y and the Destination, can see if congestion was experienced at any previous point in the path from Source if the forwarding elements are ECN capable and the Source marked packets as ECT (ECN Capabile Transport).

D. Eastlake & B.Briscoe

[Page 3]

<u>1.1</u> Conventions used in this document

The terminology and acronyms defined in $[\underline{RFC6325}]$ are used herein with the same meaning.

In this documents, "IP" refers to both IPv4 and IPv6.

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 [<u>RFC2119</u>].

Acronyms:

CE - Congestion Experienced

- ECN Explicit Congestion Notification
- ECT ECN Capabile Transport

[Page 4]

2. The ECN Specific Extended Header Flags

RBridges MAY implement ECN (Explicit Congestion Notification) [<u>RFC3168</u>] through a two-bit field in the TRILL Header extension Flags Word [<u>RFC7780</u>]. If implemented, it SHOULD be enabled by default but can be disable on a per RBridge basis by configuration.

This field is show below as "ECN" and consists of bits 12 and 13 which are in the range reserved for non-critical hop-by-hop bits. See [RFC7780] and [RFC7179] for the meaining of the other bits.

0	1		2	3
0123	4 5 6 7 8 9 0	1 2 3 4 5 6	78901234	5678901
+ - + - + - + - +	+ - + - + - + - + - + - + - +	-+-+-+-+-+	-+-+-+-+-+-+-	+ - + - + - + - + - + - + - + - +
Crit.	CHbH NCH	bH CRSV	NCRSV CITE	E NCITE
C C C	C N			
R R R	R C	ECN Ext		Ext
H I R	C C	Hop		Clr
b t s	A A	Cnt		
H E V	F F			
+ - + - + - + - +	+ - + - + - + - + - + - + - +	-+-+-+-+-+	-+-+-+-+-+-+-	+ - + - + - + - + - + - + - +

The following table is modified from [<u>RFC3168</u>] and shows the meaning of bit values in TRILL Header extended flags 12 and 13. These are also the meanings of bits 6 and 7 of the DS field in the IPv4 and IPv6 heders as defined in [<u>RFC3168</u>]:

Binary	Meaning
00	Not-ECT (Not ECN-Capable Transport)
01	ECT(1) (ECN-Capable Transport(1))
10	ECT(0) (ECN-Capable Transport(0))
11	CE (Congestion Experienced)

Table 1. ECN Field Bit Combinations

[Page 5]

3. ECN Support

An RBridge that has ECN support as specified herein advertises this through bit TBD in the Extended RBridge Capabilities APPsub-TLV [<u>RFC7782</u>] (see <u>Section 4.2</u>). On encapsulation, transit, and decapsulation it behaves as described in the subsections below, which correspond to the recommended provisions of [<u>RFC6040</u>].

3.1 Ingress ECN Support

Behavior at the ingress depends on whether the egress RBridge supports ECN. If it does, then the behavior is as follows (called "normal mode" in [<u>RFC6040</u>]):

- o When encapsulating an IP frame that is ECN enabled (non-zero ECN field), the ingress RBridge MUST create a flags word as part of the TRILL Header, setting the F flag, and copy the two ECN bits from the IP header into flag word bits 12 and 13.
- o When encapsulating a frame for a non-IP protocol, where that protocol has a means of indicating ECN that is understood by the ingress RBridge, it MAY add a flags word to the TRILL Header with the ECN bits set from the encapsulated native frame.

If the egress RBridge does not support ECN, the behavior is as follows (called "compatibility mode" in [<u>RFC6040</u>]):

- o A TRILL Header Flags Word need not be created unless there is some reason other than ECN to do so.
- o If a Flags Word is created, the ECN bits are set to zero (the Non-ECT value).

3.2 Transit ECN Support

When forwarding a TRILL Data packet encountering congestion at an RBridge, if the TRILL Header flags word is present, bits 12 and 13 are updated in the usual ECN manner [<u>RFC3168</u>]. An RBridge detects congestion either by monitoring its own queue depths or from participation in a link-specific protocol.

If, for reasons other than ECN, conditions at a transit RBridge require the insertion of a TRLL Header Flags Word into a TRILL Data packet, this implies that the egress RBridge is not ECN capable -- if it was, the Flags Word would have been included in the TRILL Data packet at the ingress. Thus, when a transit RBridge creates such a

D. Eastlake & B.Briscoe

[Page 6]

Flags Word, it sets bits 12 and 13 to zero.

3.3 Egress ECN Support

Egress RBridge support of ECN is determined by looking at the Extended Capabilities APPsub-TLV that RBridge advertises. If bit TBD is zero, or the APPsub-TLV is absent, that RBridge does not support ECN. If the APPsub-TLV is present and bit TBD is one, then it does support ECN. If there are inconsistent APPsub-TLVs, the egress RBridge is assumed to support ECN if any of those APPsub-TLVs indicate that it does.

If the egress RBridge does not support ECN, it will ignore bits 12 and 13 of any Flags Word that is present, because it does not contain any special ECN logic.

If the egress RBridge supports ECN, it does the following:

- o When decapsulating an IP frame, the RBridge MUST set the outgoing native IP frame ECN field to the code point at the intersection of the values for that field in the encapsulated IP frame (row) and the TRILL Header flags word ECN field (column) in Table 2 below or drop the frame in the case where the TRILL header indicates congestion experienced but the encapsulated native IP frame indicates a not ECN-capable transport. (Such frame dropping is necessary because IP transport that is not ECN-capable requires dropped frames to sense congestion.)
- o When decapsulating a non-IP protocol frame with a means of indicating ECN that is understood by the RBridge, it MAY set the ECN information in the decapsulated native frame by combining that information in the TRILL Header flags word and the encapsulated non-IP native frame as specified in Table 2.

Table 2 below (adapted from [RFC6040]) shows how, at the egress, to combine the ECN information in the extended TRILL Header ECN field with the ECN information in an encapsulated frame to produce the ECN information to be carried in the resulting native frame.

[Page 7]

+		. + .							+
	Inner Native	 +.	Arrivin	g TRILL Hea	der	Flag Word I		N Field	 +
	Header		Not-ECT	ECT(0)		ECT(1)	 _	CE	 _
+.		• + •			+		- + -		T
Ι	Not-ECT		Not-ECT	Not-ECT(*)	Not-ECT(*)		<drop>(*)</drop>	
	ECT(0)		ECT(0)	ECT(0)		ECT(1)	Ι	CE	I
	ECT(1)		ECT(1)	ECT(1)(*)	ECT(1)		CE	I
Ι	CE		CE	CE	- 1	CE(*)	Ι	CE	I
+ -		. + .		+	+		-+-		+
 +.	CÈ	 .+.	CÉ	CE	 +	CE(*)	 - + -	CE	 +

Table 2: Egress ECN Behavior

An asterisk in the above table indicates a probably erroneous condition that SHOULD be logged.

[Page 8]

<u>4</u>. IANA Considerations

This section summarizes IANA actions required.

4.1 Flags Word Bits

IANA is requested to assign bits 12 and 13 in the TRILL Header Flags Word for ECN and update the TRILL Extended Header Flags registry by replacing the line for bits 9-13 with the following"

BitsPurposeReference9-11available non-critical hop-by-hop flags12-13ECN (Explicit Congestion Notification) [this document]

4.2 Extended RBridge Capability Bit

IANA is requested to assign bit TBD in the Extended RBridge Capabilities to indicate ECN support. The Extended RBridge Capabilities registry on the TRILL Parameters page is updated by adding the folloing line and updating any "Unassigned" line that is affected.

Bit	Mnemonic	Description	Reference
TBD	ECN	ECN Support	[this document]

[Page 9]

<u>5</u>. Security Considerations

TBD

For ECN tunneling security considerations, see [RFC6040].

For general TRILL protocol security considerations, see [RFC6325].

<u>6</u>. Acknowledgements

This document was prepared with basic NROFF. All macros used were defined in the source file.

[Page 10]

Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u>, <u>RFC 2119</u>, DOI 10.17487/RFC2119, March 1997, <<u>http://www.rfc-editor.org/info/rfc2119</u>>.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", <u>RFC 3168</u>, DOI 10.17487/RFC3168, September 2001, <<u>http://www.rfc-</u> editor.org/info/rfc3168>.
- [RFC6040] Briscoe, B., "Tunnelling of Explicit Congestion Notification", <u>RFC 6040</u>, DOI 10.17487/RFC6040, November 2010, <<u>http://www.rfc-editor.org/info/rfc6040</u>>.
- [RFC6325] Perlman, R., Eastlake 3rd, D., Dutt, D., Gai, S., and A. Ghanwani, "Routing Bridges (RBridges): Base Protocol Specification", <u>RFC 6325</u>, DOI 10.17487/RFC6325, July 2011, <<u>http://www.rfc-editor.org/info/rfc6325</u>>.
- [RFC7179] Eastlake 3rd, D., Ghanwani, A., Manral, V., Li, Y., and C. Bestler, "Transparent Interconnection of Lots of Links (TRILL): Header Extension", <u>RFC 7179</u>, DOI 10.17487/RFC7179, May 2014, <<u>http://www.rfc-editor.org/info/rfc7179</u>>.
- [RFC7780] Eastlake 3rd, D., Zhang, M., Perlman, R., Banerjee, A., Ghanwani, A., and S. Gupta, "Transparent Interconnection of Lots of Links (TRILL): Clarifications, Corrections, and Updates", <u>RFC 7780</u>, DOI 10.17487/RFC7780, February 2016, <<u>http://www.rfc-editor.org/info/rfc7780</u>>.
- [RFC7782] Zhang, M., Perlman, R., Zhai, H., Durrani, M., and S. Gupta, "Transparent Interconnection of Lots of Links (TRILL) Active-Active Edge Using Multiple MAC Attachments", <u>RFC 7782</u>, DOI 10.17487/RFC7782, February 2016, <<u>http://www.rfc-</u> <u>editor.org/info/rfc7782</u>>.

Informative References

[none]

[Page 11]

Authors' Addresses

Donald E. Eastlake, 3rd Huawei Technologies 155 Beaver Street Milford, MA 01757 USA

Tel: +1-508-333-2270 Email: d3e3e3@gmail.com

Bob Briscoe (editor) Simula Research Lab

Email: ietf@bobbriscoe.net
URI: http://bobbriscoe.net/

[Page 12]

Copyright and IPR Provisions

Copyright (c) 2016 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to <u>BCP 78</u> 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. The definitive version of an IETF Document is that published by, or under the auspices of, the IETF. Versions of IETF Documents that are published by third parties, including those that are translated into other languages, should not be considered to be definitive versions of IETF Documents. The definitive version of these Legal Provisions is that published by, or under the auspices of, the IETF. Versions of these Legal Provisions that are published by third parties, including those that are translated into other languages, should not be considered to be definitive versions of these Legal Provisions. For the avoidance of doubt, each Contributor to the IETF Standards Process licenses each Contribution that he or she makes as part of the IETF Standards Process to the IETF Trust pursuant to the provisions of RFC 5378. No language to the contrary, or terms, conditions or rights that differ from or are inconsistent with the rights and licenses granted under RFC 5378, shall have any effect and shall be null and void, whether published or posted by such Contributor, or included with or in such Contribution.

[Page 13]