

TRILL WG
Internet-Draft
Intended status: Standards Track
Expires: February 3, 2018

Radia. Perlman
EMC Corporation
Fangwei. Hu
ZTE Corporation
Donald. Eastlake 3rd
Huawei technology
Kesava. Krupakaran
Dell
Ting. Liao
August 2, 2017

TRILL Smart Endnodes
draft-ietf-trill-smart-endnodes-06.txt

Abstract

This draft addresses the problem of the size and freshness of the endnode learning table in edge Rbridges, by allowing endnodes to volunteer for endnode learning and encapsulation/decapsulation. Such an endnode is known as a "Smart Endnode". Only the attached edge Rbridge can distinguish a "Smart Endnode" from a "normal endnode". The smart endnode uses the nickname of the attached edge Rbridge, so this solution does not consume extra nicknames. The solution also enables Fine Grained Label aware endnodes.

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 February 3, 2018.

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.

Table of Contents

1. Introduction	2
2. Solution Overview	3
3. Terminology	4
4. Smart-Hello Mechanism between Smart Endnode and RBridge . . .	5
4.1. Smart-Hello Encapsulation	5
4.2. Edge RBridge's Smart-Hello	6
4.3. Smart Endnode's Smart-Hello	7
5. Data Packet Processing	8
5.1. Data Packet Processing for Smart Endnode	8
5.2. Data Packet Processing for Edge RBridge	9
6. Multi-homing Scenario	10
7. Security Considerations	11
8. IANA Considerations	11
9. Acknowledgements	12
10. References	12
10.1. Informative References	12
10.2. Normative References	12
Authors' Addresses	14

1. Introduction

The IETF TRILL (Transparent Interconnection of Lots of Links) protocol [RFC6325] [RFC7780] provides optimal pair-wise data frame forwarding without configuration, safe forwarding even during periods of temporary loops, and support for multipathing of both unicast and multicast traffic. TRILL accomplishes this by using IS-IS [IS-IS] [RFC7176] link state routing and encapsulating traffic using a header that includes a hop count. Devices that implement TRILL are called "RBridges" (Routing Bridges) or "TRILL Switches".

An RBridge that attaches to endnodes is called an "edge RBridge" or "edge TRILL Switch", whereas one that exclusively forwards encapsulated frames is known as a "transit RBridge" or "transit TRILL Switch". An edge RBridge traditionally is the one that encapsulates a native Ethernet frame with a TRILL header, or that receives a TRILL-encapsulated packet and decapsulates the TRILL header. To

encapsulate efficiently, the edge RBridge must keep an "endnode table" consisting of (MAC, Data Label, TRILL egress switch nickname) sets, for those remote MAC addresses in Data Labels currently communicating with endnodes to which the edge RBridge is attached.

These table entries might be configured, received from ESADI [RFC7357], looked up in a directory [RFC7067], or learned from decapsulating received traffic. If the edge RBridge has attached endnodes communicating with many remote endnodes, this table could become very large. Also, if one of the MAC addresses and Data Labels in the table has moved to a different remote TRILL switch, it might be difficult for the edge RBridge to notice this quickly, and because the edge RBridge is encapsulating to the incorrect egress RBridge, the traffic will get lost.

2. Solution Overview

The Smart Endnode solution proposed in this document addresses the problem of the size and freshness of the endnode learning table in edge RBridges. An endnode E, attached to an edge RBridge R, tells R that E would like to be a "Smart Endnode", which means that E will encapsulate and decapsulate the TRILL frame, using R's nickname. Because E uses R's nickname, this solution does not consume extra nicknames.

Take the below figure as the example Smart Endnode scenario: RB1, RB2 and RB3 are the RBridges in the TRILL domain, and smart SE1 and SE2 are the smart endnodes which can encapsulate and decapsulate the TRILL packets. RB1 is the edge RB and it is been attached by SE1 and SE2. RB1 assigns its nickname to SE1 and SE2.

Each Smart Endnode, SE1 and SE2, uses RB1's nickname when encapsulating, and maintains an endnode table of (MAC, label, TRILL egress switch nickname) for remote endnodes that it (SE1 or SE2) is corresponding with. RB1 does not decapsulate packets destined for SE1 or SE2, and does not learn (MAC, label, TRILL egress switch nickname) for endnodes corresponding with SE1 or SE2, but RB1 does decapsulate, and does learn (MAC, label, TRILL egress switch nickname) for any endnodes attached to RB1 that have not declared themselves to be Smart Endnodes.

Just as an RBridge learns and times out (MAC, label, TRILL egress switch nickname), Smart Endnodes SE1 and SE2 also learn and time out endnode entries. However, SE1 and SE2 might also determine, through ICMP messages or other techniques that an endnode entry is not successfully reaching the destination endnode, and can be deleted, even if the entry has not timed out.

If SE1 wishes to correspond with destination MAC D, and no endnode entry exists, SE1 will encapsulate the packet as an unknown destination, or consulting a directory [RFC7067] (just as an RBridge would do if there was no endnode entry).

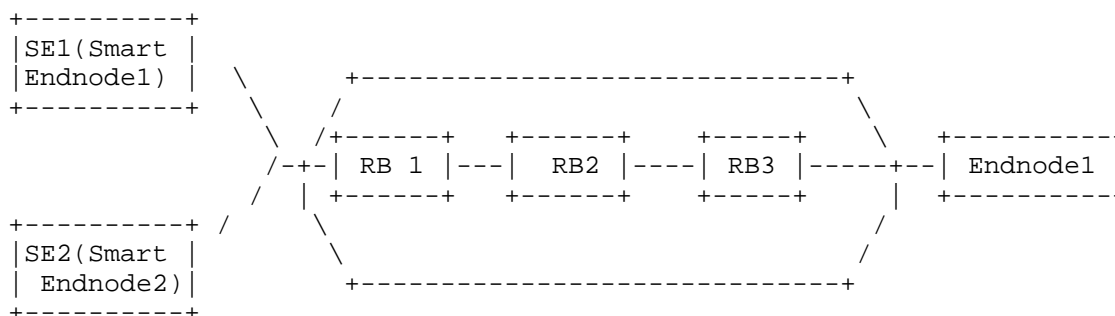


Figure 1 Smart Endnode Scenario

The mechanism in this draft is that the Smart Endnode SE1 issues a Smart-Hello, indicating SE1's desire to act as a Smart Endnode, together with the set of MAC addresses and Data Labels that SE1 owns. The Smart-Hello is used to announce the Smart Endnode capability and parameters (such as MAC address, Data Label etc.). The Smart-Hello is a type of TRILL ES-IS, which is specified in section 5 of [RFC8171]. The detailed content for a smart endnode's Smart-Hello is defined in section 4.

If RB1 supports having a Smart Endnode neighbor it also sends Smart-Hellos. The smart endnode learns from RB1's Smart-Hellos what RB1's nickname is and which trees RB1 can use when RB1 ingresses multi-destination frames. Although Smart Endnode SE1 transmits Smart-Hellos, it does not transmit or receive LSPs or E-L1FS FS-LSPs [RFC7780].

Since a Smart Endnode can encapsulate TRILL Data packets, it can cause the Inner.Label to be a Fine Grained Label [RFC7172], thus this method supports FGL aware endnodes. When and how a smart endnode decides to use the FGL instead of VLANs to encapsulate the TRILL Data packet is out of scope in this document.

3. Terminology

Edge RBridge: An RBridge providing endnode service on at least one of its ports. It is also called an edge TRILL Switch.

Data Label: VLAN or FGL.

DRB: Designated RBridge [RFC6325].

ESADI: End Station Address Distribution Information [RFC7357].

FGL: Fine Grained Label [RFC7172].

IS-IS: Intermediate System to Intermediate System [IS-IS].

RBridge: Routing Bridge, an alternative name for a TRILL switch.

Smart Endnode: An endnode that has the capability specified in this document including learning and maintaining (MAC, Data Label, Nickname) entries and encapsulating/decapsulating TRILL frame.

Transit RBridge: An RBridge exclusively forwards encapsulated frames. It is also named as transit RBridge.

TRILL: Transparent Interconnection of Lots of Links [RFC6325][RFC7780].

TRILL ES-IS: TRILL End System to Intermediate System, is a variation of TRILL IS-IS designed to operate on a TRILL link among and between one or more TRILL switches and end stations on that link[RFC8171].

TRILL Switch: a device that implements the TRILL protocol; an alternative term for an RBridge.

4. Smart-Hello Mechanism between Smart Endnode and RBridge

The subsections below describe Smart-Hello messages.

4.1. Smart-Hello Encapsulation

Although a Smart Endnode is not an RBridge, does not send LSPs or maintain a copy of the link state database, and does not perform routing calculations, it is required to have a "Hello" mechanism (1) to announce to edge RBridges that it is a Smart Endnode and (2) to tell them what MAC addresses it is handling in what Data Labels. Similarly, an edge RBridge that supports Smart Endnodes needs a message (1) to announce that support, (2) to inform Smart Endnodes what nickname to use for ingress and what nickname(s) can be used as egress nickname in a multi-destination TRILL Data packet, and (3) the list of smart end nodes it knows about on that link.

The messages sent by Smart Endnodes and by edge RBridges that support Smart Endnodes are called "Smart-Hellos". The Smart-Hello is a type of TRILL ES-IS, which is specified in [RFC8171].

The Smart-Hello Payload, both for Smart-Hellos sent by Smart Endnodes and for Smart-Hellos sent by Edge RBridges, consists of TRILL IS-IS

TLVs as described in the following two sub-sections. The non-extended format is used so TLVs, sub-TLVs, and APPsub-TLVs have an 8-bit size and type field. Both types of Smart-Hellos MUST include a Smart-Parameters APPsub-TLV as follows inside a TRILL GENINFO TLV:

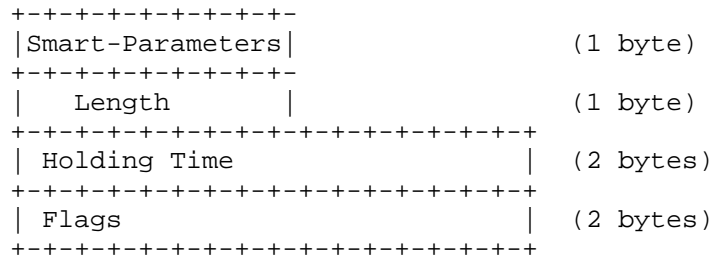


Figure 2 Smart Parameters APPsub-TLV

Type: APPsub-TLV type Smart-Parameters, value is TBD1.

Length: 4.

Holding Time: A time in seconds as an unsigned integer. It has the same meaning as the Holding Time field in IS-IS Hellos [IS-IS]. A Smart Endnode and an Edge RBridge supporting Smart Endnodes MUST send a Smart-Hello at least three times during their Holding Time. If no Smart-Hellos is received from a Smart Endnode or Edge RBridge within the most recent Holding Time it sent, it is assumed that it is no longer available.

Flags: At this time all of the Flags are reserved and MUST be send as zero and ignored on receipt.

If more than one Smart Parameters APPsub-TLV appears in a Smart-Hello, the first one is used and any following ones are ignored. If no Smart Parameters APPsub-TLV appears in a Smart-Hello, that Smart-Hello is ignored.

4.2. Edge RBridge’s Smart-Hello

The edge RBridge’s Smart-Hello contains the following information in addition to the Smart-Parameters APPsub-TLV:

- o RBridge’s nickname. The nickname sub-TLV, specified in section 2.3.2 in [RFC7176], is reused here carried inside a TLV 242 (IS-IS router capability) in a Smart-Hello frame. If more than one nickname appears in the Smart-Hello, the first one is used and the following ones are ignored.

- o Trees that RBl can use when ingressing multi-destination frames. The Tree Identifiers Sub-TLV, specified in section 2.3.4 in [RFC7176], is reused here.
- o Smart Endnode neighbor list. The TRILL Neighbor TLV, specified in section 2.5 in [RFC7176], is reused for this purpose.
- o An Autentication TLV MAY also be included.

4.3. Smart Endnode’s Smart-Hello

A new APPsub-TLV (Smart-MAC TLV) is defined for use by Smart Endnodes as defined below. In addition, there will be a Smart-Parameters APPsub-TLV and there MAY be an Authentication TLV in a Smart Endnode Smart-Hello.

If there are several VLANs/FGL Data Labels for that Smart Endnode, the Smart-MAC APPsub-TLV is included several times in Smart Endnode’s Smart-Hello. This APPsub-TLV appears inside a TRILL GENINFO TLV.

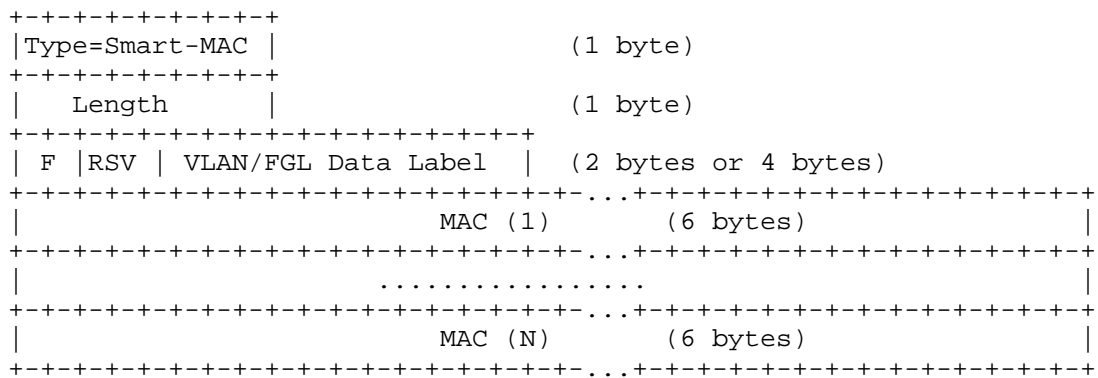


Figure 3 Smart-MAC APPsub-TLV

- o Type: TRILL APPsub-TLV Type Smart-MAC, value is TBD2.
- o Length: Total number of bytes contained in the value field.
- o F: one bit. If it sets to 1, which indicates that the endnode supports FGL data label, otherwise, the VLAN/FGL Data Labels [RFC7172] and that this Smart-MAC APPsub-TLV has an FGL in the following VLAN/FGL field. Otherwise, the VLAN/FGL Data Label field is a VLAN ID.

- o RSV: 2 bits or 6 bits, is reserved for the future use. If VLAN/FGL Data Label indicates the VLAN ID(F flag sets to 0), the RESV field is 2 bits long. Otherwise it is 6 bits.
- o VLAN/FGL Data Label: This carries a 12-bits VLAN identifier or 24-bits FGL Data Label that is valid for all subsequent MAC addresses in this APPsub-TLV, or the value zero if no VLAN/FGL data label is specified.
- o MAC(i): This is a 48-bit MAC address reachable in the Data Label given from the Smart Endnode that is announcing this APPsub-TLV.

5. Data Packet Processing

The subsections below specify Smart Endnode data packet processing. All TRILL Data packets sent to or from Smart Endnodes are sent in the Designated VLAN [RFC6325] of the local link but do not necessarily have to be VLAN tagged.

5.1. Data Packet Processing for Smart Endnode

A Smart Endnode does not issue or receive LSPs or E-L1FS FS-LSPs or calculate topology. It does the following:

- o Smart Endnode maintains an endnode table of (the MAC address of remote endnode, Data Label, the nickname of the edge RBridge's attached) entries of end nodes with which the Smart Endnode is communicating. Entries in this table are populated the same way that an edge RBridge populates the entries in its table:
 - * learning from (source MAC address ingress nickname) on packets it decapsulates.
 - * by querying a directory [RFC7067].
 - * by having some entries configured.
- o When Smart Endnode SE1 wishes to send unicast frame to remote node D, if (MAC address of remote endnode D, Data Label, nickname) entry is in SE1's endnode table, SE1 encapsulates the ingress nickname as the nickname of the RBridge(RB1), egress nickname as indicated in D's table entry. If D is unknown, SE1 either queries a directory or encapsulates the packet as a multi-destination frame, using one of the trees that RB1 has specified in RB1's Smart-Hello. The mechanism for querying a directory is given in [RFC8171].

- o When SE1 wishes to send a a multi-destination (multicast, unknown unicast, or broadcast) to the TRILL campus, SE1 encapsulates the packet using one of the trees that RB1 has specified.

If the Smart Endnode SE1 sends a multi-destination TRILL Data packet, the destination MAC of the outer Ethernet is All-RBridges multicast address.

The Smart Endnode SE1 need not send Smart-Hellos as frequently as normal RBridges. These Smart-Hellos could be periodically unicast to the Appointed Forwarder RB1 through native RBridge Channel messages. In case RB1 crashes and restarts, or the DRB changes and SE1 receives the Smart-Hello without mentioning SE1, SE1 SHOULD send a Smart-Hello immediately. If RB1 is Appointed Forwarder for any of the VLANs that SE1 claims, RB1 MUST list SE1 in its Smart-Hellos as a Smart Endnode neighbor.

5.2. Data Packet Processing for Edge RBridge

The attached edge RBridge processes and forwards TRILL Data packets based on the endnode property rather than for encapsulation and forwarding the native frames the same way as the traditional RBridges. There are several situations for the edge RBridges as follows:

- o If receiving an encapsulated unicast TRILL Data packet from a port with a Smart Endnode, with RB1's nickname as ingress, the edge RBridge RB1 forwards the frame to the specified egress nickname, as with any encapsulated frame. However, RB1 MAY filter the encapsulation frame based on the inner source MAC and Data Label as specified for the Smart Endnode. If the MAC (or Data Label) are not among the expected entries of the Smart Endnode, the frame would be dropped by the edge RBridge.
- o If receiving a unicast TRILL Data packet with RB1's nickname as egress from the TRILL campus, and the destination MAC address in the enclosed packet is listed as "smart endnode", RB1 leaves the packet encapsulated when forwarding to the smart endnode, and both the outer and inner Ethernet destination MAC is the destination smart endnod's MAC address, and the outer Ethernet source MAC address is the RB1's port MAC address. The edge RBridge still decreases the Hop count value by 1, for there is one hop between the RB1 and Smart Endnode.
- o If receiving an multi-destination TRILL Data packet from a port with a Smart Endnode, RBridge RB1 forwards the TRILL encapsulation to the TRILL campus based on the distribution tree indicated by the egress nickname. If the egress nickname does not correspond

to a distribution tree, the packet is discarded. If there are any normal endnodes (i.e, non-Smart Endnodes) attached to the edge RBridge RB1, RB1 decapsulates the frame and sends the native frame to these ports possibly pruned based on multicast listeners, in addition to forwarding the multi-destination TRILL frame to the rest of the campus.

- o If RB1 receives a native multi-destination data frame, which is sent by a non-smart endnode, from a port, including hybrid endnodes (smart endnodes and non-smart endnodes), RB1 will encapsulate it as multi-destination TRILL Data packet , and send the encapsulated multi-destination TRILL Data Packet out that same port to the smart endnodes attached to the port, and also send the encapsulated multi-destination TRILL Data Packet to the TRILL campus through other ports .
- o If RB1 receives a multi-destination TRILL Data packet from a remote RBridge, and the exit port includes hybrid endnodes(Smart Endnodes and non-Smart Endnodes), it sends two copies of multicast frames out the port, one as native and the other as TRILL encapsulated frame. When Smart Endnode receives multi-destination TRILL Data packet, it learns the remote (MAC address, Data Label, Nickname) entry, A Smart Endnodes ignores native data frames. A normal (non-smart) endnode receives the native frame and learns the remote MAC address and ignores the TRILL data packet. This transit solution may bring some complexity for the edge RBridge and waste network bandwidth resource, so avoiding the hybrid endnodes scenario by attaching the Smart Endnodes and non-Smart Endnodes to different ports is RECOMMENDED.

6. Multi-homing Scenario

Multi-homing is a common scenario for the Smart Endnode. The Smart Endnode is on a link attached to the TRILL domain in two places: to edge RBridge RB1 and RB2. Take the figure below as example. The Smart Endnode SE1 is attached to the TRILL domain by RB1 and RB2 separately. Both RB1 and RB2 could announce their nicknames to SE1.

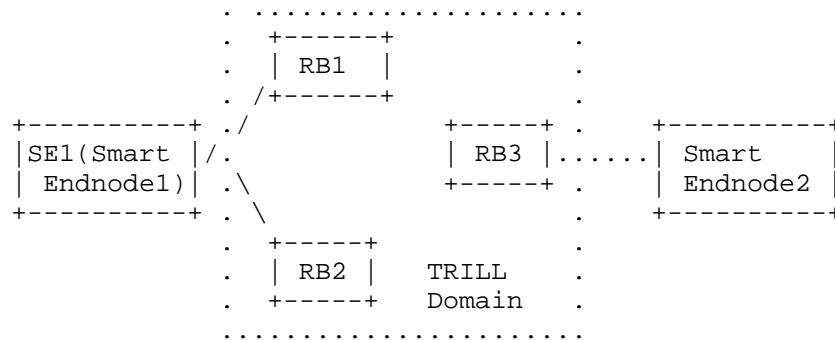


Figure 4 Multi-homing Scenario

Smart Endnode SE1 can choose either RB1 or RB2's nickname, when encapsulating and forwarding a TRILL data packet. If the active-active load balance is considered for the multi-homing scenario, the Smart Endnode SE1 could use both RB1 and RB2's nickname to encapsulate and forward TRILL Data packet. SE1 uses RB1's nickname when forwarding through RB1, and RB2's nickname when forwarding through RB2. this will cause MAC flip-flopping(see [RFC7379]) of the endnode table entry in the remote RBridges (or Smart Endnodes). The solution for the MAC flip-flopping issue is to set a multi-homing bit in the RSV field of the TRILL data packet. When remote RBridge RB3 or Smart Endnodes receives a data packet with the multi-homed bit set, the endnode entries (SE1's MAC address, label, RB1's nickname) and (SE1's MAC address, label, RB2's nickname) will coexist as endnode entries in the remote RBridge. Another solution is to use the ESADI protocol to distribute multiple attachments of a MAC address of a multi-homing group,The ESADI is deployed among the edge RBridges (See section 5.3 of [RFC7357]).

7. Security Considerations

Smart-Hellos can be secured by using Authentication TLVs based on [RFC5310].

For general TRILL Security Considerations, see [RFC6325].

For TRILL ES-IS Security Considerations, see [RFC8171].

8. IANA Considerations

IANA is requested to allocate APPsub-TLV type numbers for the Smart-MAC and Smart-Parameters APPsub-TLVs from the range below 256 and

update the "TRILL APPsub-TLV Types under IS-IS TLV 251 Application Identifier 1" registry as follows.

Protocol	Description	Reference
TBD1	Smart-Parameters	[this document]
TBD2	Smart-MAC	[this document]

Table 1

9. Acknowledgements

The contributions of the following persons are gratefully acknowledged: Mingui Zhang, Weiguo Hao, Linda Dunbar, and Andrew Qu.

10. References

10.1. Informative References

- [RFC7067] Dunbar, L., Eastlake 3rd, D., Perlman, R., and I. Gashinsky, "Directory Assistance Problem and High-Level Design Proposal", RFC 7067, DOI 10.17487/RFC7067, November 2013, <<http://www.rfc-editor.org/info/rfc7067>>.
- [RFC7379] Li, Y., Hao, W., Perlman, R., Hudson, J., and H. Zhai, "Problem Statement and Goals for Active-Active Connection at the Transparent Interconnection of Lots of Links (TRILL) Edge", RFC 7379, DOI 10.17487/RFC7379, October 2014, <<http://www.rfc-editor.org/info/rfc7379>>.

10.2. Normative References

- [IS-IS] ISO/IEC 10589:2002, Second Edition,, "Intermediate System to Intermediate System Intra-Domain Routing Exchange Protocol for use in Conjunction with the Protocol for Providing the Connectionless-mode Network Service (ISO 8473)", 2002.
- [RFC5310] Bhatia, M., Manral, V., Li, T., Atkinson, R., White, R., and M. Fanto, "IS-IS Generic Cryptographic Authentication", RFC 5310, DOI 10.17487/RFC5310, February 2009, <<http://www.rfc-editor.org/info/rfc5310>>.

- [RFC6325] Perlman, R., Eastlake 3rd, D., Dutt, D., Gai, S., and A. Ghanwani, "Routing Bridges (RBridges): Base Protocol Specification", RFC 6325, DOI 10.17487/RFC6325, July 2011, <<http://www.rfc-editor.org/info/rfc6325>>.
- [RFC7172] Eastlake 3rd, D., Zhang, M., Agarwal, P., Perlman, R., and D. Dutt, "Transparent Interconnection of Lots of Links (TRILL): Fine-Grained Labeling", RFC 7172, DOI 10.17487/RFC7172, May 2014, <<http://www.rfc-editor.org/info/rfc7172>>.
- [RFC7176] Eastlake 3rd, D., Senevirathne, T., Ghanwani, A., Dutt, D., and A. Banerjee, "Transparent Interconnection of Lots of Links (TRILL) Use of IS-IS", RFC 7176, DOI 10.17487/RFC7176, May 2014, <<http://www.rfc-editor.org/info/rfc7176>>.
- [RFC7178] Eastlake 3rd, D., Manral, V., Li, Y., Aldrin, S., and D. Ward, "Transparent Interconnection of Lots of Links (TRILL): RBridge Channel Support", RFC 7178, DOI 10.17487/RFC7178, May 2014, <<http://www.rfc-editor.org/info/rfc7178>>.
- [RFC7357] Zhai, H., Hu, F., Perlman, R., Eastlake 3rd, D., and O. Stokes, "Transparent Interconnection of Lots of Links (TRILL): End Station Address Distribution Information (ESADI) Protocol", RFC 7357, DOI 10.17487/RFC7357, September 2014, <<http://www.rfc-editor.org/info/rfc7357>>.
- [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", RFC 7780, DOI 10.17487/RFC7780, February 2016, <<http://www.rfc-editor.org/info/rfc7780>>.
- [RFC7783] Senevirathne, T., Pathangi, J., and J. Hudson, "Coordinated Multicast Trees (CMT) for Transparent Interconnection of Lots of Links (TRILL)", RFC 7783, DOI 10.17487/RFC7783, February 2016, <<http://www.rfc-editor.org/info/rfc7783>>.
- [RFC8171] Eastlake 3rd, D., Dunbar, L., Perlman, R., and Y. Li, "Transparent Interconnection of Lots of Links (TRILL): Edge Directory Assistance Mechanisms", RFC 8171, DOI 10.17487/RFC8171, June 2017, <<http://www.rfc-editor.org/info/rfc8171>>.

Authors' Addresses

Radia Perlman
EMC Corporation
2010 156th Ave NE, suite #200
Bellevue, WA 98007
USA

Phone: +1-206-291-367
Email: radiaperlman@gmail.com

Fangwei Hu
ZTE Corporation
No.889 Bibo Rd
Shanghai 201203
China

Phone: +86 21 68896273
Email: hu.fangwei@zte.com.cn

Donald Eastlake, 3rd
Huawei technology
155 Beaver Street
Milford, MA 01757
USA

Phone: +1-508-634-2066
Email: d3e3e3@gmail.com

Kesava Vijaya Krupakaran
Dell
Olympia Technology Park
Guindy Chennai 600 032
India

Phone: +91 44 4220 8496
Email: Kesava_Vijaya_Krupak@Dell.com

Ting Liao
Nanjing, Jiangsu 210012
China

Email: liaoting82@163.com

TRILL WG
INTERNET-DRAFT
Intended Status:Informational
Expires: October 20, 2017

R. Parameswaran,
Brocade Communications, Inc.
April 22, 2017

TRILL: Parent node Shifts in Tree Construction, Mitigation.
<draft-rp-trill-parent-selection-03.txt>

Abstract

This draft documents a known problem in the TRILL tree construction mechanism and offers an approach requiring no change to the TRILL protocol in order to solve the problem.

Status of This Memo

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

Distribution of this document is unlimited. Comments should be sent to the authors or 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>.

Terminology and Acronyms.

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].

Table of Contents

1. Introduction.....	1
2. Tree construction in TRILL.....	2
3. Issues with the TRILL tree construction algorithm.....	2
4. Solution using the Affinity sub-TLV.....	4
5. Network wide selection of computation algorithm.....	7
6. Relationship to draft-ietf-trill-resilient-trees.....	7
7. Security Considerations.....	9
8. IANA Considerations.....	9
9. Informative References.....	9

1. Introduction.

TRILL is a data center technology that uses link-state routing mechanisms in a layer 2 setting, and serves as a replacement for spanning-tree. TRILL uses trees rooted at pre-determined nodes as a way to distribute multi-destination traffic. Multi-destination traffic includes traffic such as layer-2 broadcast frames, unknown unicast flood frames, and layer 2 traffic with multicast MAC addresses (collectively referred to as BUM traffic). Multi-destination traffic is typically hashed onto one of the available trees and sent over the tree, potentially reaching all nodes in the network (hosts behind which may own/need the packet in question).

2. Tree construction in TRILL.

Tree construction in TRILL is defined by [RFC6325], with additional corrections defined in [RFC7780].

The tree construction mechanism used in TRILL codifies certain tree construction steps which make the resultant trees very brittle. Specifically, the parent selection mechanism in TRILL causes problems in case of node failures. TRILL uses the following rule - when constructing an SPF tree, if there are multiple possible parents for a given node (i.e. if multiple upstream nodes can potentially pull in a given node during SPF, all at the same cumulative cost, then the parent selection is imposed in the following manner):

[RFC6325]:

"When building the tree number j , remember all possible equal cost parents for node N . After calculating the entire 'tree' (actually, directed graph), for each node N , if N has ' p ' parents, then order the parents in ascending order according to the 7-octet IS-IS ID considered as an unsigned integer, and number them starting at zero. For tree j , choose N 's parent as choice $j \bmod p$."

There is an additional correction posted to this in [RFC7780]:

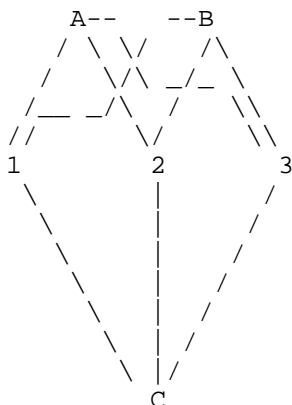
[RFC7780], Section 3.4:

"Section 4.5.1 of [RFC6325] specifies that, when building distribution tree number j , node (RBridge) N that has multiple possible parents in the tree is attached to possible parent number $j \bmod p$. Trees are numbered starting with 1, but possible parents are numbered starting with 0. As a result, if there are two trees and two possible parents, then in tree 1 parent 1 will be selected, and in tree 2 parent 0 will be selected.

This is changed so that the selected parent MUST be $(j-1) \bmod p$. As a result, in the case above, tree 1 will select parent 0, and tree 2 will select parent 1. This change is not backward compatible with [RFC6325]. If all RBridges in a campus do not determine distribution trees in the same way, then for most topologies, the RPFC will drop many multi-destination packets before they have been properly delivered."

3. Issues with the TRILL tree construction algorithm.

With this tree construction mechanism in mind, let's look at the Spine-Leaf topology presented below and consider the calculation of Tree number 2 in TRILL. Assume all the links in the tree are at the same cost.



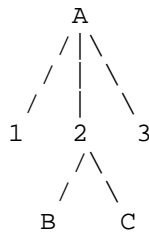
Assume that in the above topology, when ordered by 7-octet ISIS-id,

1 < 2 < 3 holds and that the root for Tree number 2 is A. Given the ordered set {1, 2, 3} , these nodes have the following indices (with a starting index of 0):

Node	Index
1	0
2	1
3	2

Given the SPF constraint and that the tree root is A, the parent for nodes 1,2, and 3 will be A. However, when the SPF algorithm tries to pull B or C into the tree, we have a choice of parents, namely 1, 2, or 3.

Given that this is tree 2, the parent will be the one with index $(2-1) \bmod 3$ (which is equal to 1). Hence the parent for node B will be node 2.

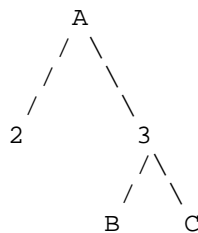


However, due to TRILL's parent selection algorithm, the sub-tree rooted at Node 2 will be impacted even if Node 1 or Node 3 go down.

Take the case where Node 1 goes down. Tree 2 must now be re-computed (this is normal) - but now, when the SPF computation is underway, when the SPF process tries to pull in B, the list of potential parents for B now are {2 and 3}. So, after ordering these by ISIS-Id as {2, 3} (where 2 is considered to be at index of 0 and 3 is considered to be at index 1), for tree 1, we apply TRILL's formula of:

$$\begin{aligned}
 \text{Parent's index} &= (\text{TreeNumber}-1) \bmod \text{Number_of_parents.} \\
 &= (2-1) \bmod 2 \\
 &= 1 \bmod 2 \\
 &= 1 \text{ (which is the index of Node 3)}
 \end{aligned}$$

The re-calculated tree now looks as shown below. The shift in parent nodes (for B) may cause disruption to live traffic in the network, and is unnecessary in absolute terms because the existing parent for node B, node 2, was not perturbed in any way.



Aside from the disruption posed by the change in the tree links, depending upon how the concerned rbridges stripe vlans/FGLs across trees and how they may prune these, additional disruption is possible if the forwarding state on the new parent rbridge is not primed to match the new tree structure. This churn could simply be avoided with a better approach.

The parent shift issue noted above can be solved by using

the Affinity sub-TLV.

While the technique identified in this draft has an immediate benefit when applied to spine/leaf networks popular in data-center designs, nothing in the approach outlined below assumes a spine-leaf network. The technique presented below will work on any connected graph. Furthermore, no directional symmetry in link-cost is assumed.

4. Solution using the Affinity sub-TLV.

At a high level, this problem can be solved by having the affected parent send out an Affinity sub-TLV identifying the children for which it wants to preserve the parent-child relationship, subject to network events which may change the structure of the tree. The affected parent node would send out an Affinity sub-TLV with multiple Affinity records, one per child node, listing the concerned tree number.

It would be sufficient to have a local configuration option (e.g. a CLI) at one of the nodes which is deemed to be the parent of choice (referred to as designated parent below). The following steps provide a way to implement this proposal:

- a. The operator locally configures the designated parent to indicate its stickiness in tree construction for a specific tree number and tree root via the Affinity sub-TLV. This can be done before tree construction if the operator consults the 7 octet ISIS-ID relative ordering of the concerned nodes and decides up-front which of the potential parent nodes should become the parent node for a given set of children on that tree number under the TRILL tree construction mechanism. The operator **MUST** configure the designated parent stickiness on only one node amongst a set of sibling (potential parent) nodes relative to the tree root for that tree number. It is suggested that the parent stickiness be configured on the node that would have been selected as the parent under default Trill parent selection rules. Parent stickiness **MUST NOT** be configured on the root of the tree, or if configured previously on a non-root node with the root for that tree shifting to that node subsequently, such configuration **MUST** be ignored on the root node.
- b. On any subsequent SPF calculation after the operator configures the designated parent as indicated above, when the designated parent node finds that it could be a potential parent for one or more child nodes during tree construction, it declares itself to be the parent for the concerned child nodes, over-riding the default TRILL parent selection rules. The configured node advertises its parent preference via the Affinity sub-TLV when it completes a tree calculation, and finds itself the parent of one or more child nodes per the SPF tree calculation. The Affinity sub-TLV **MUST** reflect the appropriate tree number and the child nodes for which the concerned node is a parent node. The Affinity sub-TLV **SHOULD** be published when the tree computation is deemed to have converged (more on this under d. below).
- c. Likewise, when any change event happens in the network, one which forces a tree re-calculation for the concerned tree, the designated parent node should run through the normal TRILL tree calculation agnostic of the fact that it has published an Affinity sub-TLV as well as agnostic of the default TRILL tree selection rules i.e the node asserts its right to be a parent without directly referencing either the default Trill parent selection rules or its own published Affinity sub-TLV in establishing parent relationships.
- d. During the SPF tree calculation, the designated parent node should react in the following manner:

- i. If the node is a potential parent for some of the children identified in an existing Affinity sub-TLV, if any, after convergence of the tree computation, the node MUST send out an (updated) Affinity sub-TLV identifying the correct sub-set of children for which the node aspires to establish/continue the parent relationship. This case would also apply if there are new child nodes for which the node is now a parent (however, see the conflicted Affinity sub-TLV rules in vii and j. below).

For its own tree computation, the designated parent node MUST use itself as parent in order to pull the set of children identified during the SPF run into the tree, barring a conflicting affinity sub-TLV seen from another node (see vii. below for handling this case).

- ii. If the tree structure changes such that the designated node is no longer a potential parent for any of the child nodes in the advertised Affinity sub-TLV, then it SHOULD retract the Affinity sub-TLV, upon convergence of the tree computation. In this case, the default TRILL tie-break rule would need to be used during SPF construction for the nodes that were children of this designated node previously. One specific case may be worth high-lighting - if a parent-child relationship inverts i.e. if the designated parent becomes a child of its former child node due to a change in the tree structure, it MUST exclude that child from its Affinity sub-TLV. In such case, if the designated parent node cannot maintain a parent relationship with any of its prior child nodes, then it MUST retract any previously published affinity sub-TLV.
- iii. Nodes SHOULD use a convergence timer to track completion of the tree computation. If there are any additional tree computations while the convergence timer is running, the timer SHOULD be re-started/extended in order to absorb the interim network events. It is possible that the intended action at the expiration of the timer may change meanwhile. The timer needs to be large enough to absorb multiple network events that may happen due to a change in the physical state of the network, and yet short enough to avoid delaying the update of the Affinity sub-TLV.
- iv. At the expiration of the convergence timer, the existing state of the tree MUST be compared with the existing Affinity sub-TLV and the intended change in the status of the Affinity sub-TLV is carried out e.g. a fresh publication, or an update to the list of children, or a retraction.
- v. Alternately, the above steps (re-examination of the Affinity sub-TLV and update) MAY be tied to/triggered from the download of the tree routes to the L2 RIB, since that typically happens upon a successful computation of the complete tree. An additional stabilization timer could be used to counteract back-to-back L2 RIB downloads due to repeated computations of the tree due to a burst of network events.
- vi. Note that this approach may cause an additional tree computation at remote nodes once the updated Affinity sub-TLV (or lack of it) is received/perceived, beyond the network events which led up to the change in the tree. In the case where an operator introduced a designated parent configuration on an existing tree, then remote nodes would need to receive the Affinity sub-TLV indicating the designated parent's Affinity for its children before the remote nodes shift away from the default TRILL parent selection rules. However, in most cases, in steady state, this mechanism should cause very little tree churn unless

a designated parent configuration was introduced, removed, or a link between the designated parent and its children changed state. In cases where the network change event originated on the designated parent node, it may be possible to optimize on the churn by packing both the data bearing the network change event and the Affinity sub-TLV into the same link-state update packet.

- vii. In situations where the designated parent node would normally originate an affinity sub-TLV to indicate affinity to a specific set of child nodes, it MUST NOT originate an Affinity sub-TLV if it sees an Affinity sub-TLV from some other node for the same tree number and for all of the same child-nodes, such that the other node's Affinity sub-TLV would win using the conflict tie-break rules in section 5.3 of [RFC7783]. Any existing Affinity sub-TLV already published by this node in such a situation MUST be retracted. If only some of the child nodes overlap between the two conflicting Affinity sub-TLVs, then this designated parent node MAY continue to publish its affinity sub-TLV listing its child nodes that are not in conflict with the other Affinity sub-TLV. Other guide-lines listed in [RFC7783] MUST be adhered to as well - the originator of the Affinity sub-TLV must name only directly adjacent nodes as children, and must not name the tree root as a child.
- e. Situations where the node advertising the Affinity sub-TLV dies or restarts SHOULD be handled using the normal handling for such scenarios relating to the parent Router Capability TLV, and as specified in [RFC4971].
- f. Situations where a parent-child link directly connected to the designated parent node constantly flaps, MUST be handled by having the designated parent node retract the Affinity sub-TLV, if it affects the parent-child relationships in consideration. The long-term state of the Affinity sub-TLV can be monitored by the designated parent node to see if it is being published and retracted repeatedly in multiple iterations or if a specific set of children are being constantly added and removed. The designated parent may resume publication of the Affinity sub-TLV once it perceives the network to be stable again in the future.
- g. If the designated parent node is forced to retract its Affinity sub-TLV due to a change in the tree structure, it can then repeat these steps in a subsequent tree construction, if the same node becomes a parent again, so long as it perceives its parent-child links to be stable (free of link/node flaps).
- h. In terms of nodes that do not support this draft, they are expected to seamlessly inter-operate with this draft, so long as they understand and honor the Affinity sub-TLV. The draft assumes that most TRILL implementations now support the Affinity sub-TLV. In any case, the guide-lines specified in section 4.1 of [RFC7783] MUST be used i.e. if all nodes in the network do not support the Affinity sub-TLV then the network must default to the Trill parent selection rules.
- i. Remote nodes MUST default to the Trill parent selection rules if they do not see an Affinity sub-TLV sent by any node in the network.
- j. At remote nodes, conflicting Affinity sub-TLVs from different originators for the same tree number and child node MUST be handled as specified in section 5.3 of [RFC7783], namely by selecting the Affinity sub-TLV originated by the node with the highest priority to be a tree root, with System-ID as tie-breaker.

5. Network wide selection of computation algorithm.

The proposed solution above does not need any operational change to the TRILL protocol, beyond the usage of the Affinity sub-TLV (which is already in the proposed standard) for the use case identified in this draft.

6. Relationship to draft-ietf-trill-resilient-trees.

Given that both draft-ietf-trill-resilient-trees, and draft-rp-trill-parent-selection-03 drafts use the Affinity sub-TLV, it is worthwhile to examine if there is any functional overlap between the two drafts. At a high level, the two drafts have different goals and appear to solve unrelated problems.

draft-ietf-trill-resilient-trees relates to link protection, and defines the notion of a primary distribution tree and a backup distribution tree (DT), where these trees are intentionally kept link disjoint to the extent possible, and the backup tree is pre-programmed in the hardware, and activated either up front or upon failure of the primary distribution tree.

On the other hand, draft-rp-trill-parent-selection-03 protects parent-child relationships of interest on the primary DT, and has no direct notion of a backup DT.

draft-ietf-trill-resilient-trees considers the following algorithmic approaches to the building the backup distribution tree (section numbers listed below are from draft-ietf-trill-resilient-trees):

1. Operator hand-configuration for links on the backup DT/manual generation of Affinity sub-TLV - this is very tedious and unlikely to scale or be implemented in practice, and hence is disregarded in the analysis here.
2. Section 3.2.1.1a: Use of MRT algorithms (which will produce conjugate trees - link disjoint trees with roots for primary and backup trees that are coincident on the same rBridge).
3. Section 3.2.1.1b: Once the primary DT is constructed, the links used in the primary DT are additively cost re-weighted, and a second SPF is run to derive the links comprising the backup DT. Affinity sub-TLV is used to mark links on the back-up DT which are not also on the primary DT. This approach can handle conjugate trees as well as non-conjugate trees (link disjoint trees that are rooted at different rBridges).
4. Section 3.2.2: A variation on the section 3.2.1.1b approach, but without Affinity sub-TLV advertisement. Once the primary DT is constructed, costs for links on the primary DT are multiplied by a fixed multiplier to prevent them from being selected in a subsequent SPF run, unless there is no other choice, and the subsequent SPF yields links on the backup DT.

All of the approaches above yield maximally link disjoint trees, when applied as prescribed.

Approach 4 above does not seem to use Affinity sub-TLVs and instead seems to depend upon a network wide agreement on the alternative tree computation algorithm being used.

Approaches 2 and 3 use Affinity sub-TLV on the backup DT, for links that are not already on the primary DT. The primary DT does not appear to use Affinity sub-TLVs. Additionally, from an end-to-end perspective the backup DT comes into picture when the primary DT fails (this is effectively true even in the 1+1 protection mechanism

and in the local protection case), and then again, only until the primary DT is recalculated. Once the primary DT is recalculated, the backup DT is recalculated as well, and can change corresponding to the new primary DT.

draft-ietf-trill-resilient-trees cannot directly prevent/mitigate a parent node shift on the primary DT at a given parent node, and while usage of the Affinity sub-TLV on the backup DT might confer a parent affinity on some nodes on the backup DT, these are not necessarily the nodes on which the network operator may want/prefer an explicit parent affinity. Further, the backup DT is only used on a transient basis, from a forwarding perspective, until the primary DT is recomputed.

However, a parent shift can be triggered by link or node failure. In a situation where both drafts are active in the implementation, failure of a specific link may cause the backup DT to kick in, but when the primary DT is re-calculated, draft-rp-trill-parent-selection-03 can be used to preserve parent-child relationships on the primary DT, to the extent possible, during the re-calculation. So, there does not appear to be a direct functional overlap in the simultaneous usage of these drafts, and it ought to be possible to use both drafts simultaneously, so long as the primary and back-up DTs can be uniquely identified/differentiated.

7. Security Considerations.

The proposal primarily influences tree construction and tries to preserve parent-child relationships in the tree from prior computations of the same tree, without changing any of operational aspects of the protocol. Hence, no new security considerations for TRILL are raised by this proposal.

8. IANA Considerations.

No new registry entries are requested to be assigned by IANA. The Affinity Sub-TLV has been defined in [RFC7176], and this proposal does not change its semantics in any way.

9. Informative 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>>.
- [RFC6325] Perlman, R., Eastlake 3rd, D., Dutt, D., Gai, S., and A. Ghanwani, "Routing Bridges (Rbridges): Base Protocol Specification", RFC 6325, DOI 10.17487/RFC6325, July 2011, <<http://www.rfc-editor.org/info/rfc6325>>.
- [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", RFC 7780, DOI 10.17487/RFC7780, February 2016, <<http://www.rfc-editor.org/info/rfc7780>>.
- [RFC7783] Senevirathne, T., Pathangi, J., Hudson, J., "Coordinated Multicast Trees (CMT) for Transparent Interconnection of Lots of Links (TRILL)", RFC 7783, February 2016, <<http://datatracker.ietf.org/doc/rfc7783>>
- [RFC4971] Vasseur, JP., Shen, N., Aggarwal, R., "Intermediate System to Intermediate System (IS-IS) Extensions for Advertising Router Information", RFC 4971, July 2007, <<http://datatracker.ietf.org/doc/rfc4971>>

Author's Address:

R. Parameswaran,
Brocade Communications, Inc.
120 Holger Way,
San Jose, CA 95134.

Email: parameswaran.r7@gmail.com

Copyright and IPR Provisions

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. 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.