

TRILL Working Group  
Internet Draft  
Intended status: Standards Track

Yizhou Li  
Weiguo Hao  
Huawei Technologies  
Jon Hudson  
Brocade  
Naveen Nimmu  
Broadcom  
Anoop Ghanwani  
DELL

Expires: January 2014

May 2, 2013

Aware Spanning Tree Topology Change on R Bridges  
draft-yizhou-trill-tc-awareness-02.txt

#### 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), 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/ietf/lid-abstracts.txt>

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>

This Internet-Draft will expire on November 2, 2013.

#### Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

Internet-Draft

STP Topology Change Awareness

May 2013

This document is subject to [BCP 78](http://trustee.ietf.org/license-info) 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.

## Abstract

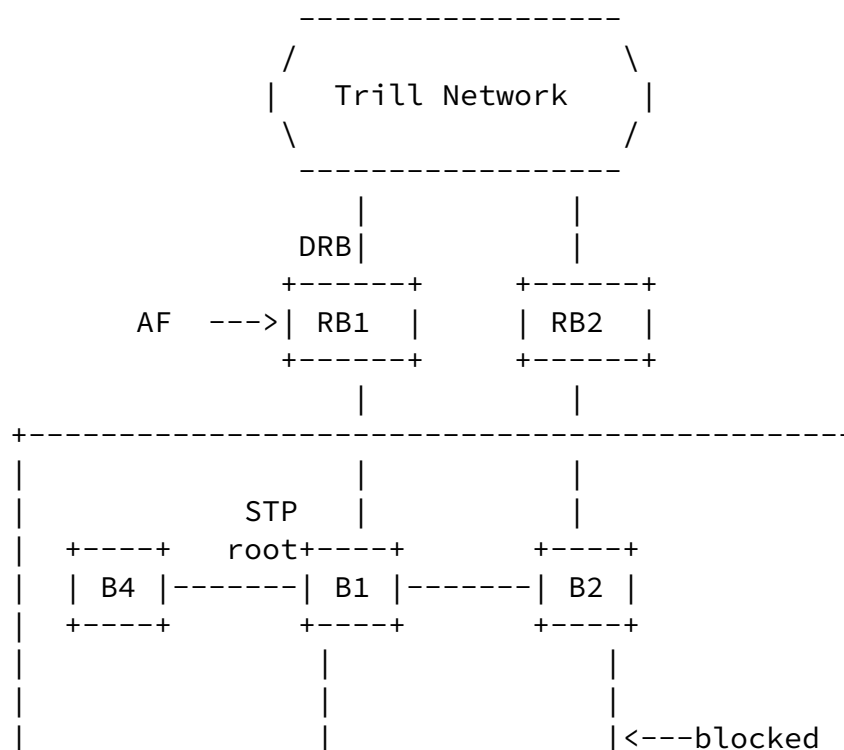
When a local LAN running spanning tree protocol connecting to TRILL campus via more than one RBridge, there are several ways to perform loop avoidance. One of them illustrated by [RFC6325](#) [RFC6325] A.3 was to make relevant ports on edge RBridges involving in spanning tree calculation. When edge RBridges are emulated as a single highest priority root, the local bridged LAN will be naturally partitioned after running spanning tree protocol. This approach achieves better link utilization and intra-VLAN load balancing in some scenarios. This document describes how the edge RBridges react to topology change occurring in bridged LAN in order to make the abovementioned spanning tree approach function correct.

## Table of Contents

<a href="#">1.</a>	<a href="#">Introduction .....</a>	<a href="#">3</a>
<a href="#">1.1.</a>	<a href="#">Motivations .....</a>	<a href="#">5</a>
<a href="#">2.</a>	<a href="#">Conventions used in this document.....</a>	<a href="#">6</a>
<a href="#">3.</a>	<a href="#">BPDU RBridge Channel and TLVs.....</a>	<a href="#">6</a>
<a href="#">3.1.</a>	<a href="#">BPDU TLV .....</a>	<a href="#">7</a>
<a href="#">3.2.</a>	<a href="#">Authentication Information TLV.....</a>	<a href="#">7</a>
<a href="#">4.</a>	<a href="#">Operations .....</a>	<a href="#">8</a>
<a href="#">4.1.</a>	<a href="#">Sending BPDU using RBridge channel .....</a>	<a href="#">9</a>
<a href="#">4.2.</a>	<a href="#">Receiving BPDU in RBridge channel.....</a>	<a href="#">10</a>
<a href="#">4.3.</a>	<a href="#">Informing the remote site.....</a>	<a href="#">11</a>
<a href="#">5.</a>	<a href="#">Security Considerations.....</a>	<a href="#">13</a>
<a href="#">6.</a>	<a href="#">IANA Considerations .....</a>	<a href="#">13</a>
<a href="#">7.</a>	<a href="#">References .....</a>	<a href="#">13</a>
<a href="#">7.1.</a>	<a href="#">Normative References.....</a>	<a href="#">13</a>
<a href="#">7.2.</a>	<a href="#">Informative References.....</a>	<a href="#">14</a>
<a href="#">8.</a>	<a href="#">Acknowledgments .....</a>	<a href="#">14</a>

## 1. Introduction

The TRILL protocol [RFC6325] provides the appointed forwarder mechanism [RFC6439] for loop avoidance where, for part of the loop, the frame would be in TRILL encapsulated format, for example in the scenario shown by Figure 1. Only one of the RBridges is responsible for encapsulating/decapsulating a given VLAN's data frames on a link. Bridges in the local bridged LAN runs normal spanning tree protocol for local loop avoidance. RBridges keeps track of the root bridge by listening to BPDUs received on the local port. This information is reported per VLAN by the RBridge in its LSP and is used to detect a root bridge change. Root bridge changes trigger the reset of the inhibition timer of the appointed forwarder. When an RBridge ceases to be appointed forwarder for a VLAN on a port, it sends topology change BPDUs to purge the MAC table on local bridged LAN switches. An RBridge conformable to [RFC6325] never encapsulates or forwards any BPDUs it receives.



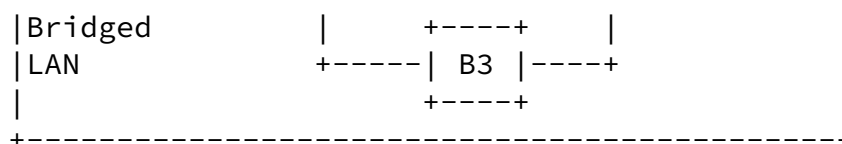
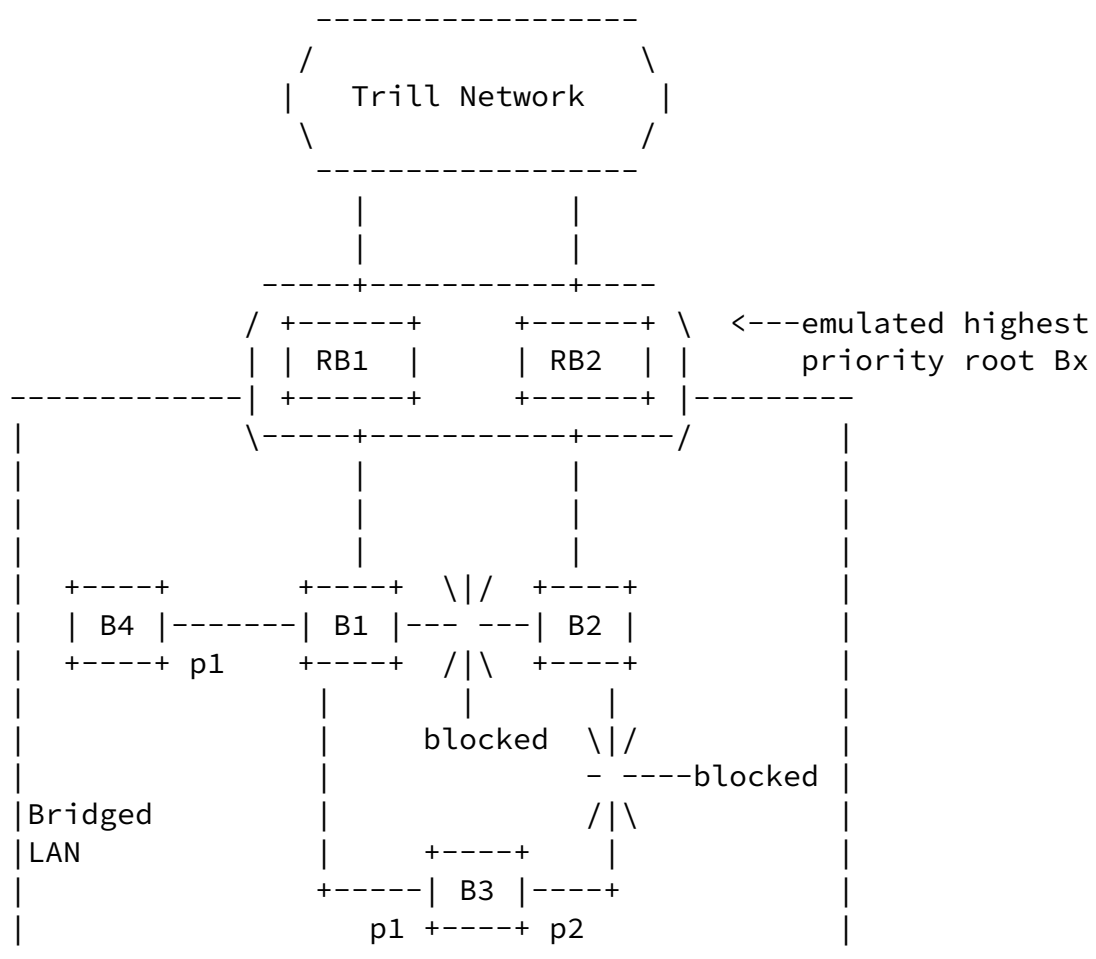


Figure 1 TRILL and bridged LAN topology

[RFC6325] A.2 & A.3 presented the problems using the conventional approach shown in Figure 1. Native frames enter and leave a link via the link's appointed forwarder for the VLAN of the frame can cause congestion or suboptimal routing. Four methods was illustrated in [RFC6325] to solve the problem,

1. Use RBridge instead of conventional bridge
2. Re-arrange network topology
3. Carefully select the different appointed forwarders for VLANs if end stations on local bridged LAN can be separated into multiple VLANs
4. Configure the RBridges to be like one STP tree root in local bridged LAN. The RBridge ports that are connected to the bridged LAN send spanning tree configuration BPDUs. Then the bridged LAN is forced into partitions. Figure 2 shows its network topology.

Method 1 and 2 highly depends on the network topology and equipment types and therefore have very limited applicability. Method 3 and 4 have broader applicability. Method 4 is more applicable than method 3 if all end stations in bridged LAN are on the same VLAN or intra VLAN load balancing is required to avoid per VLAN congestion and suboptimal routing. The traffic discontinuity was caused by inhibition timer setting in case of root change in method 3. Proper timeout value has to be carefully chosen for tradeoff between unnecessary traffic continuity and potential loop. Method 4 eliminates the requirement of setting inhibition timer in case of root change. Therefore method 4 is considered as a very common practice in real deployment.



## 1.1. Motivations

Bridged LANs may have topology changes at any time. When RB1 & RB2 serve as one single STP tree root as shown in Figure 2, it is required that RB1 and RB2 have to tunnel some BPDUs to help the bridged LAN convergence in certain circumstances. Figure 2 will be used to illustrate such motivation for rest of this subsection.

RB1 & RB2 use the same bridge ID to emit spanning tree BPDUs as the highest priority root Bx. All bridges in LAN see RB1 and RB2 as a single tree root. Therefore B1-B2 and B2-B3 links are blocked for loop avoidance by the spanning tree protocol. RB1 and RB2 will not receive TRILL-Hello from each other. Bridged LAN is logically partitioned into two parts. RB1 is DRB and AF for all VLANs in left partition and RB2 is DRB and AF in right partition.

If B1-B3 link fails for some reason, alternate port p2 on B3 will send topology change (TC) BPDU to B2 as RSTP specifies [[802.1D](#)]. B2-B3 link will start forwarding frames. TC BPDU is then sent from B2 to RB2. As RB2 never forwards BPDU frame to TRILL campus, left partition has no way to know the topology change. Therefore B4 will not be able to correctly purge the MACs learnt from port p1 for end stations connected to B3. MAC table entry aging is the last resort in this case. In addition, a remote end station may keep sending traffic to an end station connected to B3 via RB1-B1 which causes frame loss. Therefore some mechanism must be introduced to purge the MACs learned both in the left partition of the bridged LAN and at the remote Rbridges when topology changes.

This draft proposes to use RBridge channel [[TRILLChannel](#)] to tunnel the TC BPDU to solve the issue. It complements [Appendix A.3 of RFC6325](#) to improve correctness and efficiency.

## 2. Conventions used in this document

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](#)].

This document uses the terminologies defined in [[RFC6325](#)] along with the following:

Root Bridge Group - A group of RBridges acting as an emulated single tree root in a spanning tree instance in local bridged LAN. The group has at least two RBridges.

### 3. BPDU RBridge Channel and TLVs

A new channel protocol is defined to carry BPDU.

Channel protocol code: TBD (BPDU)

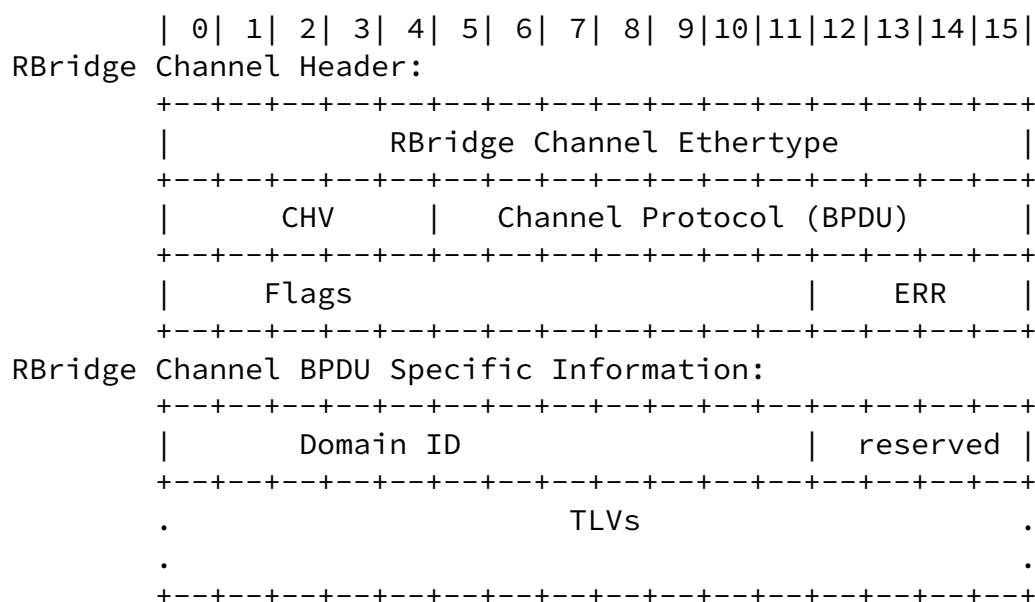


Figure 3 RBridge Channel Format for BPDU

Domain ID: Unique ID to identify a single spanning tree domain.

Reserved: 4 bits reserved field.

The fields of TRILL header and inner Ethernet header SHOULD be set as per [\[TRILLChannel\]](#) unless specified in this draft.

3.1. BPDU TLV

```
+---+---+---+---+---+
|   Type=           |           (1 byte)
+---+---+---+---+---+
|   Length           |           (1 byte)
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   BPDU             |           (variable length)
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
```

BPDU: field used to put the original BPDU frame, typically is TC (topology awareness) BPDU.

3.2. Authentication Information TLV

```
+---+---+---+---+---+
|   Type=           |           (1 byte)
+---+---+---+---+---+
|   Length           |           (1 byte)
+---+---+---+---+---+
|   Auth Mode        |           (1 byte)
+---+---+---+---+---+.....+---+---+---+---+---+---+
|   Auth Value       |           (variable length)
+---+---+---+---+---+.....+---+---+---+---+---+---+
```

Auth Mode:

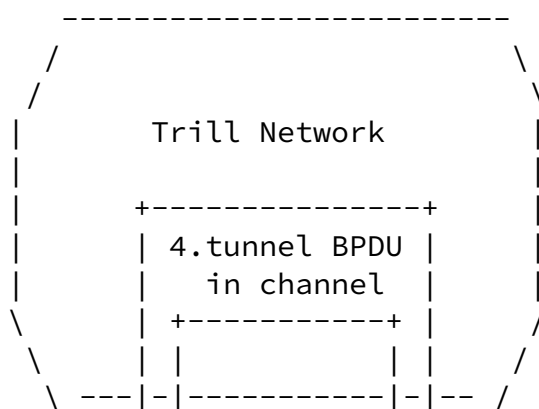
Authentication Mode Code	Value	Reference
Reserved	0	[ISO-10589]
Cleartext Password	1	[ISO-10589]
ISO 10589 Reserved	2	[ISO-10589]
HMAC-MD5 Authentication	54	<a href="#">RFC 5304</a>
Routing Domain private authentication method	255	[ISO-10589]

Auth Value: Value used to authenticate the RBridge Channel BPDUs  
Specific Information using the algorithm specified by Auth Mode

This TLV is used for channel message authentication. When an RBridge receives a channel message with Authentication Information, it MUST discard it if the Authentication Value is incorrect.

#### 4. Operations

Figure 4 shows TC BPDUs tunneled from RB2 to RB1 using RBridge Channel.



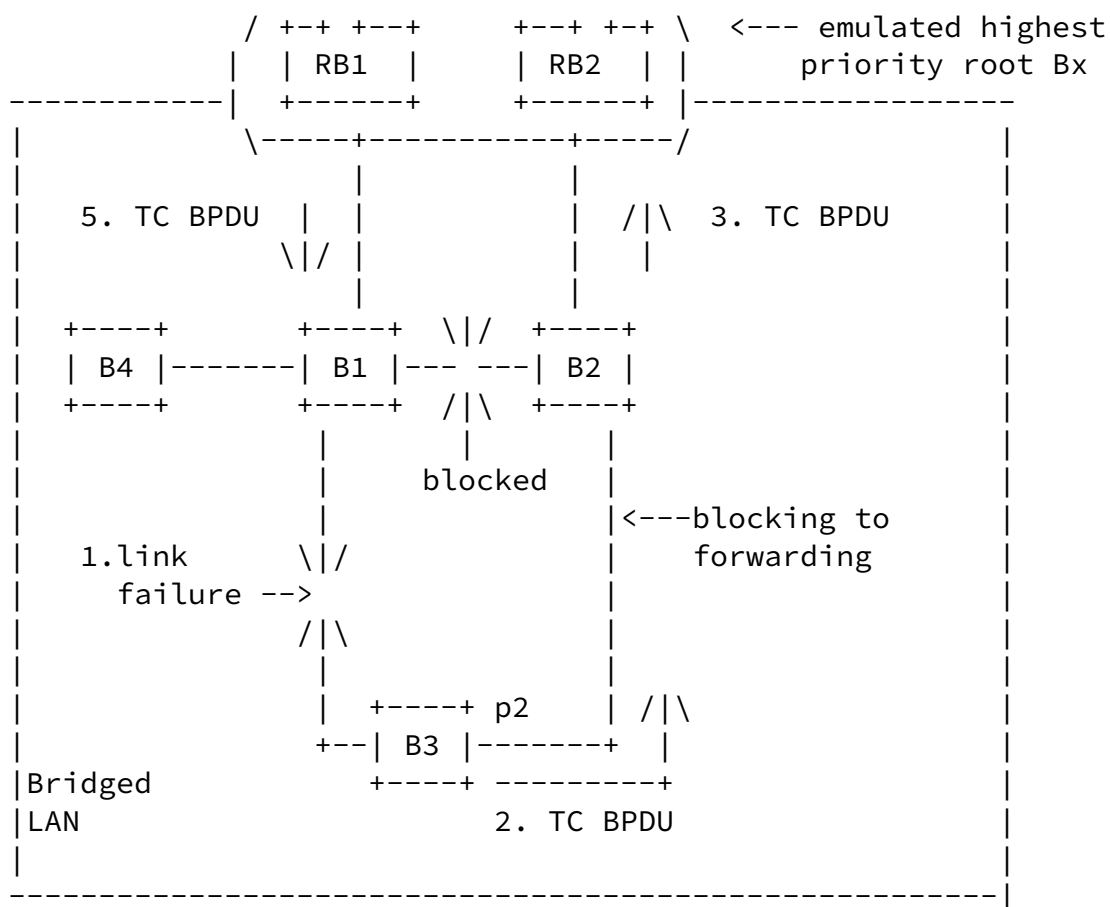


Figure 4 Tunneled TC BPDUs

#### 4.1. Sending BPDUs using RBridge channel

In figure 4, when B1-B3 link fails, alternate port p2 on B3 will start to send TC BPDUs and go to forwarding state. RB2 receives TC BPDUs from B2 sequentially. RB2 encapsulates the TC BPDUs in RBridge channel and sends it to RB1.

Interested VLANs and Spanning Tree Roots Sub-TLV [RFC6326] carries spanning tree root bridge IDs seen for all ports for which the RBridge is the appointed forwarder for a VLAN. As RB1 and RB2 use the same bridge ID and that bridge ID is the spanning tree root, RB1 and RB2 are considered as in a root bridge group. Static configuration of root bridge group is also allowed.

When RBridge receives TC BPDU from an access port, it tunnels the frame to all the other RBridges in the same root bridge group using RBridge channel protocol specified in [section 3](#). Normally the number of RBridges in a root bridge group is limited, say 2 or 3; such tunneling is performed using TRILL unicast encapsulation. N members in a root bridge group results in N-1 sequential unicast BPDU tunneled. In figure 4, RB2 knows RB1 is in the same root bridge group from LSP exchange; hence RB2 uses RB1's nickname as egress nickname and encapsulates the TC BPDU in RBridge channel. M bit in TRILL header SHOULD be 0.

If TRILL Campus was partitioned temporarily in some unusual cases, RBridges in the same root bridge group may not reach each other. For instance, if RB2 was not able to reach RB1 through TRILL campus at some transition period due to network fault, RB1 would not receive the tunneled TC BPDU from RB2. Then the approach illustrated in this document will take effect again only after RB1 and RB2 connectivity via TRILL recovers from the network fault.

It is possible to statically configure a root bridge group, especial when network is relatively small and stable. Therefore when an RBridge tunnels the TC BPDU to other members in the same root bridge group, it has to make sure the destination is reachable.

If edges RBridges configured in the same root bridge group connect to separate TRILL campus intentionally, it is not recommended to use spanning tree partition mechanism and such root bridge group provisioning is normally considered as mis-configuration.

#### [4.2](#). Receiving BPDU in RBridge channel

When an RBridge receives a TC BPDU from RBridge channel, it determines the frame was sent from an RB in the same root bridge group. Then RBridge decapsulates the frame and sends the original TC BPDU to its local bridged LAN. TC BPDU will be flooded throughout the access ports configured as the same domain ID specified by BPDU RBridge channel in the left partition to purge MAC table of bridges.

#### [4.3](#). Informing the remote site

When local topology changes, the correspondence of end station and its attaching RBridge cached by remote RB may become invalid. The RBridges who is the appointed forwarder for the specified VLAN in remote sites should be informed to update the stale correspondence table entry.

When traffic is bi-directional, the remote RBridge will receive the data frames from the newly attached RBridge of the local end station. The remote RBridge will update its MAC-Nickname correspondence table naturally through data frame learning.

When traffic is uni-directional from the remote to local site or traffic from local to remote has to be triggered by traffic from remote to local, remote RBridge will not receive the data frame from local RBridge to refresh its table. Then traffic discontinuity may last for some time until the table entry is aged out at the remote RBridge.

A lightweight method is to use RBridge channel to carry MAC purge information. In Figure 4, When RB2 receives TC BPDU from B2, it derives the corresponding VLAN list. For example, if MSTP is used, RB2 will get the VLAN IDs in the same MSTP instance as TC BPDU. RB2 sends out MAC purge information using RBridge channel with VLAN information and RBridges' nicknames in the same root bridge group. All remote RBridges received MAC purge should clear its MAC-to-nickname correspondence table for entries with the specified nicknames and VLAN IDs. If no VLAN list is specified, the remote RBridges should clear the correspondence in all VLANs relevant to the given nicknames. The MAC purge is recommended to send on the management VLAN in which all RBridges joins.

A new channel protocol code for MAC purge should be defined as follows.

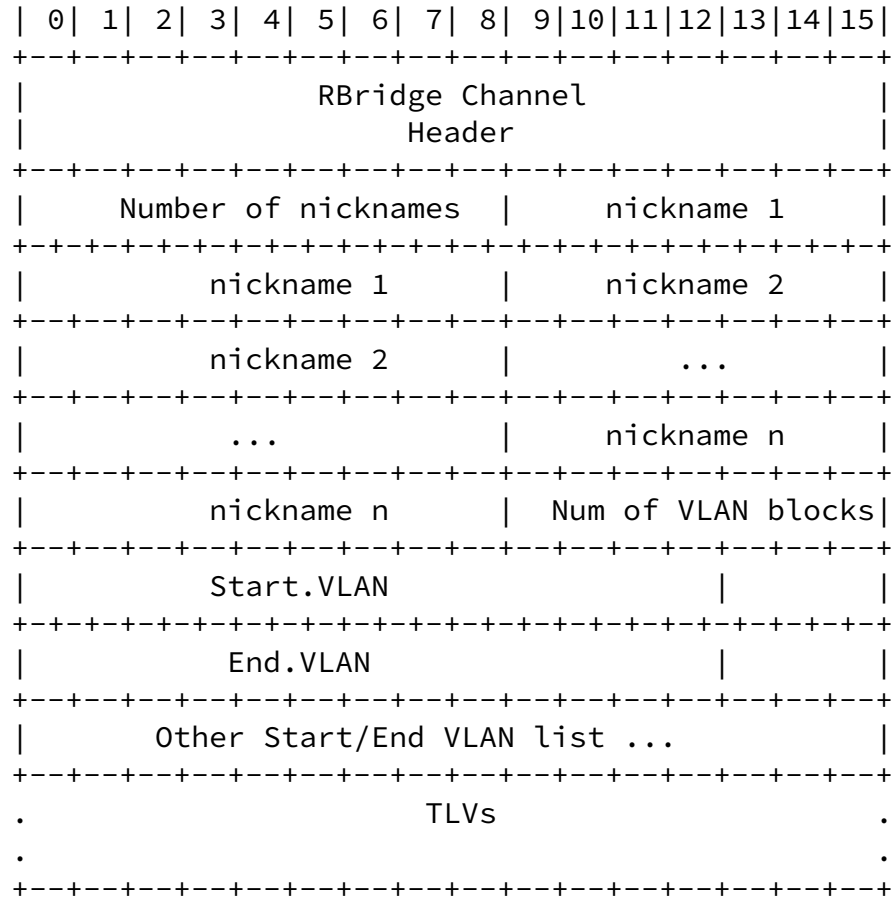


Figure 5 RBridge Channel Format for MAC Purge

Number of nicknames: number of the following nicknames, which will be used by the receivers to purge their relevant MAC-to-nickname correspondence table entries.

Num of VLAN blocks: number of the following VLAN block. A VLAN block is specified by a start and an end VLAN IDs. When start and end VLAN IDs are the same, it implies only one VLAN ID is in the block. When number of VLAN block is 0, it implies no VLAN ID is specified.

For any nickname x specified and any VLAN y specified in this TLV, the receivers should purge MAC-to-nickname correspondence table entries with (any-MAC, VLAN-y, nickname-x). When number of VLAN block is 0, the receivers should purge entries with (any-MAC, any-VLAN, nickname-x).

Authentication Information TLV specified by [section 3.2](#) can be used in MAC Purge channel message.

## 5. Security Considerations

This document does not change the general RBridge security considerations of the TRILL base protocol and TRILL RBridge Channel. See [Section 6 of \[RFC6325\]](#) and section 7 of [\[TRILLChannel\]](#).

Massive TC BPDU may trigger RBridges continuously sending tunneled BPDU and MAC purges. It may cause denial-of-service in TRILL campus. Similar as the traditional bridged LAN running spanning tree, it is suggested to monitor the receiving rate of TC BPDU on bridged LAN facing port of RBridges. If the receiving rate is beyond the threshold, RBridge should only process and tunnel the TC BPDU in the configured rate.

Authentication TLV should be included if risk of injecting forged information is a concern. However keeping channel information in this document secret and private is not necessary and is not a requirement.

## 6. IANA Considerations

IANA is requested to allocate the new channel protocol codes as following.

Channel protocol code X1: BPDU

Channel protocol code X2: MAC purge

BPDU TLV type:

Authentication Information TLV type:

## 7. References

### 7.1. Normative References

- [RFC6325] Perlman, R., Eastlake 3rd, D., Dutt, D., Gai, S., and A. Ghanwani, "Routing Bridges (RBridges): Base Protocol Specification", [RFC 6325](#), July 2011.
- [6326bis] Eastlake, D. et.al., "'Transparent Interconnection of Lots of Links (TRILL) Use of IS-IS'", [draft-eastlake-isis-rfc6326bis-07.txt](#), Work in Progress, December 2011.

[RFC6439] Eastlake, D. et.al., 'RBridge: Appointed Forwarder', [RFC 6439](#), November 2011.

Li, et al.

Expires January 2, 2014

[Page 13]

---

Internet-Draft

STP Topology Change Awareness

May 2013

[TRILLChannel] - Eastlake, D., V. Manral, Y. Li, S. Aldrin, D. Ward, "RBridges: RBridge Channel Support in TRILL", [draft-ietf-trill-rbridge-channel](#), work in progress.

[RFC6327] Eastlake 3rd, D., Perlman, R., Ghanwani, A., Dutt, D., and V. Manral, "Routing Bridges (RBridges): Adjacency", [RFC 6327](#), July 2011

[802.1D] "IEEE Standard for Local and metropolitan area networks /Media Access Control (MAC) Bridges", 802.1D-2004, 9 June 2004.

## [7.2](#). Informative References

[RFC6165] Banerjee, A. and D. Ward, "Extensions to IS-IS for Layer-2 Systems", [RFC 6165](#), April 2011.

[802.1Q-2011] "IEEE Standard for Local and metropolitan area networks /Virtual Bridged Local Area Networks", 802.1Q-2011, 31 Aug 2011.

## [8](#). Acknowledgments

This document was prepared using 2-Word-v2.0.template.dot.

## Appendix: Change History

Changes from -01 to -02

1. Add domain ID field to BPDU channel message.
2. Allow TLVs for BPDU and MAC Purge channel message.
3. Change BPDU channel message. Make BPDU field from fixed field to TLV

4. Add Authentication Information TLV

6. Assorted editorial changes.

Li, et al.

Expires January 2, 2014

[Page 14]

---

Internet-Draft

STP Topology Change Awareness

May 2013

#### Authors' Addresses

Yizhou Li  
Huawei Technologies  
101 Software Avenue,  
Nanjing 210012  
China

Phone: +86-25-56625375  
Email: liyizhou@huawei.com

Weiguo Hao  
Huawei Technologies  
101 Software Avenue,  
Nanjing 210012  
China

Phone: +86-25-56623144  
Email: haoweiguo@huawei.com

Jon Hudson  
Brocade  
120 Holger Way  
San Jose, CA 95134  
USA.

Email: jon.hudson@gmail.com

Naveen Nimmu  
Broadcom  
9th Floor, Building no 9, Raheja Mind space  
Hi-Tec City, Madhapur,  
Hyderabad - 500 081, INDIA

Phone: +1-408-218-8893  
Email: naveen@broadcom.com

Anoop Ghanwani  
DELL  
350 Holger Way  
San Jose, CA 95134  
USA.

Phone: +1-408-571-3500  
Email: Anoop@alumni.duke.edu