

BESS Workgroup
INTERNET-DRAFT
Intended Status: Standards Track
Updates: 7385

A. Sajassi, Ed.
S. Salam
Cisco
J. Drake
Juniper
J. Uttaro
ATT
S. Boutros
VMware
J. Rabadan
Nokia

Expires: April 28, 2018

October 28, 2017

E-TREE Support in EVPN & PBB-EVPN
draft-ietf-bess-evpn-etree-14

Abstract

The Metro Ethernet Forum (MEF) has defined a rooted-multipoint Ethernet service known as Ethernet Tree (E-Tree). A solution framework for supporting this service in MPLS networks is described in RFC7387 ("A Framework for Ethernet-Tree (E-Tree) Service over a Multiprotocol Label Switching (MPLS) Network"). This document discusses how those functional requirements can be met with a solution based on RFC7432, BGP MPLS Based Ethernet VPN (EVPN), with some extensions and how such a solution can offer a more efficient implementation of these functions than that of RFC7796, E-Tree Support in Virtual Private LAN Service (VPLS). This document makes use of the most significant bit of the "Tunnel Type" field (in PMSI Tunnel Attribute) governed by the IANA registry created by RFC7385, and hence updates RFC7385 accordingly.

Status of this Memo

This Internet-Draft is submitted to IETF 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/lid-abstracts.html>

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

Copyright and License 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	4
1.1	Specification of Requirements	4
1.2	Terminology	5
2	E-Tree Scenarios	5
2.1	Scenario 1: Leaf or Root Site(s) per PE	6
2.2	Scenario 2: Leaf or Root Site(s) per AC	6
2.3	Scenario 3: Leaf or Root Site(s) per MAC Address	8
3	Operation for EVPN	9
3.1	Known Unicast Traffic	9
3.2	Broadcast, Unknown, and Multicast (BUM) Traffic	10
3.2.1	BUM Traffic Originated from a Single-homed Site on a Leaf AC	11
3.2.2	BUM Traffic Originated from a Single-homed Site on a Root AC	11
3.2.3	BUM Traffic Originated from a Multi-homed Site on a Leaf AC	11
3.2.4	BUM Traffic Originated from a Multi-homed Site on a Root AC	11
3.3	E-Tree Traffic Flows for EVPN	12

3.3.1 E-Tree with MAC Learning	12
3.3.2 E-Tree without MAC Learning	13
4 Operation for PBB-EVPN	13
4.1 Known Unicast Traffic	14
4.2 Broadcast, Unkonwn, and Multicast (BUM) Traffic	14
4.3 E-Tree without MAC Learning	15
5 BGP Encoding	15
5.1 E-Tree Extended Community	15
5.2 PMSI Tunnel Attribute	17
6 Acknowledgement	18
7 Security Considerations	18
8 IANA Considerations	18
8.1 Considerations for PMSI Tunnel Types	19
9 References	19
9.1 Normative References	19
9.2 Informative References	20
Appendix-A	20
Authors' Addresses	21

1 Introduction

The Metro Ethernet Forum (MEF) has defined a rooted-multipoint Ethernet service known as Ethernet Tree (E-Tree) [MEF6.1]. In an E-Tree service, a customer site that is typically represented by an Attachment Circuits (AC) (e.g., a 802.1Q VLAN tag but may also be represented by a MAC address) is labeled as either a Root or a Leaf site. Root sites can communicate with all other customer sites (both Root and Leaf sites). However, Leaf sites can communicate with Root sites but not with other Leaf sites. In this document unless explicitly mentioned otherwise, a site is always represented by an AC.

[RFC7387] describes a solution framework for supporting E-Tree service in MPLS networks. The document identifies the functional components of an overall solution to emulate E-Tree services in MPLS networks in addition to multipoint-to-multipoint Ethernet LAN (E-LAN) services specified in [RFC7432] and [RFC7623].

[RFC7432] defines EVPN, a solution for multipoint L2VPN services with advanced multi-homing capabilities, using BGP for distributing customer/client MAC address reach-ability information over the MPLS/IP network. [RFC7623] combines the functionality of EVPN with [802.1ah] Provider Backbone Bridging (PBB) for MAC address scalability.

This document discusses how the functional requirements for E-Tree service can be met with a solution based on (PBB-)EVPN (i.e., [RFC7432] and [RFC7623]) with some extensions to their procedures and BGP attributes. Such (PBB-)EVPN based solution can offer a more efficient implementation of these functions than that of RFC7796, E-Tree Support in Virtual Private LAN Service (VPLS). This efficiency is achieved by performing filtering of unicast traffic at the ingress PE nodes as opposed to egress filtering where the traffic is sent through the network and gets filtered and discarded at the egress PE nodes. The details of this ingress filtering is described in section 3.1. Since this document specifies a solution based on [RFC7432], it requires the readers to have the knowledge of [RFC7432] as prerequisite. This document makes use of the most significant bit of the "Tunnel Type" field (in PMSI Tunnel Attribute) governed by the IANA registry created by RFC7385, and hence updates RFC7385 accordingly. Section 2 discusses E-Tree scenarios. Section 3 and 4 describe E-Tree solutions for EVPN and PBB-EVPN respectively, and section 5 covers BGP encoding for E-Tree solutions.

1.1 Specification of Requirements

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 [KEYWORDS].

1.2 Terminology

Broadcast Domain: In a bridged network, the broadcast domain corresponds to a Virtual LAN (VLAN), where a VLAN is typically represented by a single VLAN ID (VID) but can be represented by several VIDs where Shared VLAN Learning (SVL) is used per [802.1Q].

Bridge Table: An instantiation of a broadcast domain on a MAC-VRF.

CE: Customer Edge device, e.g., a host, router, or switch.

EVI: An EVPN instance spanning the Provider Edge (PE) devices participating in that EVPN.

MAC-VRF: A Virtual Routing and Forwarding table for Media Access Control (MAC) addresses on a PE.

Ethernet Segment (ES): When a customer site (device or network) is connected to one or more PEs via a set of Ethernet links, then that set of links is referred to as an 'Ethernet segment'.

Ethernet Segment Identifier (ESI): A unique non-zero identifier that identifies an Ethernet segment is called an 'Ethernet Segment Identifier'.

Ethernet Tag: An Ethernet tag identifies a particular broadcast domain, e.g., a VLAN. An EVPN instance consists of one or more broadcast domains.

P2MP: Point to Multipoint.

PE: Provider Edge device.

2 E-Tree Scenarios

This document categorizes E-Tree scenarios into the following three scenarios, depending on the nature of the Root/Leaf site association:

- Either Leaf or Root site(s) per PE
- Either Leaf or Root site(s) per Attachment Circuit (AC)
- Either Leaf or Root site(s) per MAC address

2.1 Scenario 1: Leaf or Root Site(s) per PE

In this scenario, a PE may receive traffic from either Root ACs or Leaf ACs for a given MAC-VRF/bridge table, but not both. In other words, a given EVPN Instance (EVI) on a Provider Edge (PE) device is either associated with Root(s) or Leaf(s). The PE may have both Root and Leaf ACs albeit for different EVIs.

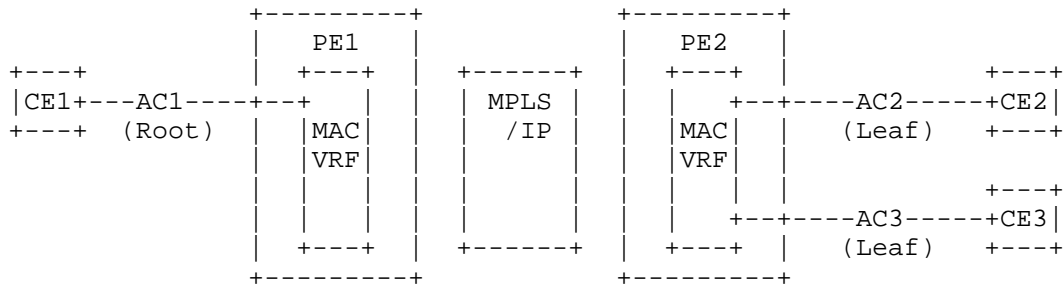


Figure 1: Scenario 1

In this scenario, tailored BGP Route Target (RT) import/export policies among the PEs belonging to the same EVI can be used to prevent the communications among Leaf PEs. To prevent the communications among Leaf ACs connected to the same PE and belonging to the same EVI, split-horizon filtering is used to block traffic from one Leaf AC to another Leaf AC on a MAC-VRF for a given E-Tree EVI. The purpose of this topology constraint is to avoid having PEs with only Leaf sites importing and processing BGP MAC routes from each other. To support such topology constrain in EVPN, two BGP Route-Targets (RTs) are used for every EVPN Instance (EVI): one RT is associated with the Root sites (Root ACs) and the other is associated with the Leaf sites (Leaf ACs). On a per EVI basis, every PE exports the single RT associated with its type of site(s). Furthermore, a PE with Root site(s) imports both Root and Leaf RTs, whereas a PE with Leaf site(s) only imports the Root RT.

For this scenario, if it is desired to use only a single RT per EVI (just like E-LAN services in [RFC7432]), then the approach B in scenario 2 (described below) needs to be used.

2.2 Scenario 2: Leaf or Root Site(s) per AC

In this scenario, a PE can receive traffic from both Root ACs and Leaf ACs for a given EVI. In other words, a given EVI on a PE can be associated with both Root(s) and Leaf(s).

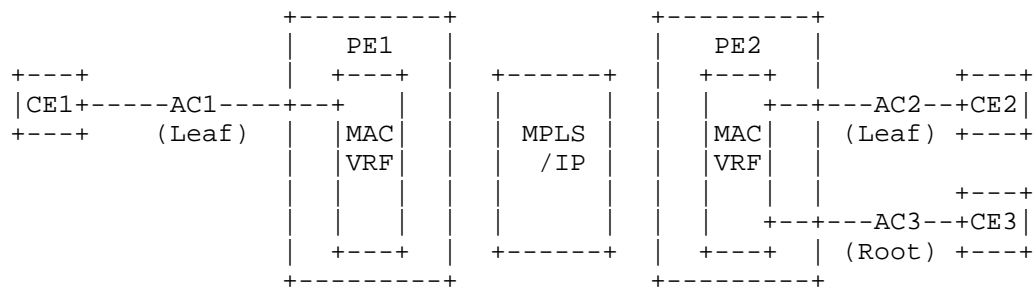


Figure 2: Scenario 2

In this scenario, just like the previous scenario (in section 2.1), two Route Targets (one for Root and another for Leaf) can be used. However, the difference is that on a PE with both Root and Leaf ACs, all remote MAC routes are imported and thus there needs to be a way to differentiate remote MAC routes associated with Leaf ACs versus the ones associated with Root ACs in order to apply the proper ingress filtering.

In order to recognize the association of a destination MAC address to a Leaf or Root AC and thus support ingress filtering on the ingress PE with both Leaf and Root ACs, MAC addresses need to be colored with Root or Leaf indication before advertisements to other PEs. There are two approaches for such coloring:

A) To always use two RTs (one to designate Leaf RT and another for Root RT)

B) To allow for a single RT be used per EVI just like [RFC7432] and thus color MAC addresses via a "color" flag in a new extended community as detailed in section 5.1.

Approach (A) would require the same data plane enhancements as approach (B) if MAC-VRF and bridge tables used per VLAN, are to remain consistent with [RFC7432] (section 6). In order to avoid data-plane enhancements for approach (A), multiple bridge tables per VLAN may be considered; however, this has major drawbacks as described in appendix-A and thus is not recommended.

Given that both approaches (A) and (B) would require the same data-plane enhancements, approach (B) is chosen here in order to allow for RT usage consistent with baseline EVPN [RFC7432] and for better generality. It should be noted that if one wants to use RT constraints in order to avoid MAC advertisements associated with a Leaf AC to PEs with only Leaf ACs, then two RTs (one for Root and another for Leaf) can still be used with approach (B); however, in

such applications Leaf/Root RTs will be used to constrain MAC advertisements and they are not used to color the MAC routes for ingress filtering - i.e., in approach (B), the coloring is always done via the new extended community.

If, for a given EVI, a significant number of PEs have both Leaf and Root sites attached (even though they may start as Root-only or Leaf-only PEs), then a single RT per EVI should be used. The reason for such recommendation is to alleviate the configuration overhead associated with using two RTs per EVI at the expense of having some unwanted MAC addresses on the Leaf-only PEs.

2.3 Scenario 3: Leaf or Root Site(s) per MAC Address

In this scenario, a customer Root or Leaf site is represented by a MAC address and a PE may receive traffic from both Root AND Leaf sites on a single Attachment Circuit (AC) of an EVI. This scenario is not covered in either [RFC7387] or [MEF6.1]; however, it is covered in this document for the sake of completeness. In this scenario, since an AC carries traffic from both Root and Leaf sites, the granularity at which Root or Leaf sites are identified is on a per MAC address. This scenario is considered in this document for EVPN service with only known unicast traffic because the Designated Forwarding (DF) filtering per [RFC7432] would not be compatible with the required egress filtering - i.e., Broadcast, Unknown, and Multicast (BUM) traffic is not supported in this scenario and it is dropped by the ingress PE.

For this scenario, the approach B in scenario 2 (described above) is used in order to allow for single RT usage by service providers.

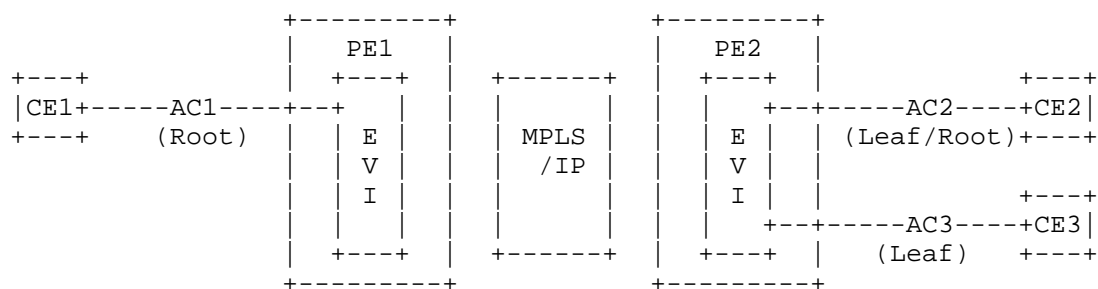


Figure 3: Scenario 3

In conclusion, the approach B in scenario 2 is the recommended approach across all the above three scenarios and the corresponding solution is detailed in the following sections.

3 Operation for EVPN

[RFC7432] defines the notion of Ethernet Segment Identifier (ESI) MPLS label used for split-horizon filtering of BUM traffic at the egress PE. Such egress filtering capabilities can be leveraged in provision of E-Tree services as it will be seen shortly for BUM traffic. For known unicast traffic, additional extensions to [RFC7432] is needed (i.e., a new BGP Extended Community for Leaf indication described in section 5.1) in order to enable ingress filtering as described in detail in the following sections.

3.1 Known Unicast Traffic

Since in EVPN, MAC learning is performed in the control plane via advertisement of BGP routes, the filtering needed by E-Tree service for known unicast traffic can be performed at the ingress PE, thus providing very efficient filtering and avoiding sending known unicast traffic over the MPLS/IP core to be filtered at the egress PE as done in traditional E-Tree solutions - i.e., E-Tree for VPLS [RFC7796].

To provide such ingress filtering for known unicast traffic, a PE MUST indicate to other PEs what kind of sites (Root or Leaf) its MAC addresses are associated with. This is done by advertising a Leaf indication flag (via an Extended Community) along with each of its MAC/IP Advertisement routes learned from a Leaf site. The lack of such flag indicates that the MAC address is associated with a Root site. This scheme applies to all scenarios described in section 2.

Tagging MAC addresses with a Leaf indication enables remote PEs to perform ingress filtering for known unicast traffic - i.e., on the ingress PE, the MAC destination address lookup yields, in addition to the forwarding adjacency, a flag which indicates whether the target MAC is associated with a Leaf site or not. The ingress PE cross-checks this flag with the status of the originating AC, and if both are leaves, then the packet is not forwarded.

In situation where MAC moves are allowed among Leaf and Root sites (e.g., non-static MAC), PEs can receive multiple MAC/IP advertisements routes for the same MAC address with different Leaf/Root indications (and possibly different ESIs for multi-homing scenarios). In such situations, MAC mobility procedures (section 15 of [RFC7432]) take precedence to first identify the location of the MAC before associating that MAC with a Root or a Leaf site.

To support the above ingress filtering functionality, a new E-Tree Extended Community with a Leaf indication flag is introduced [section 5.1]. This new Extended Community MUST be advertised with MAC/IP Advertisement routes learned from a Leaf site. Besides MAC/IP Advertisement route, no other EVPN routes are required to carry this new extended community.

3.2 Broadcast, Unknown, and Multicast (BUM) Traffic

This specification does not provide support for filtering BUM (Broadcast, Unknown, and Multicast) traffic on the ingress PE; due to the multi-destination nature of BUM traffic, it is not possible to perform filtering of the same on the ingress PE. As such, the solution relies on egress filtering. In order to apply the proper egress filtering, which varies based on whether a packet is sent from a Leaf AC or a Root AC, the MPLS-encapsulated frames MUST be tagged with an indication when they originated from a Leaf AC - i.e., to be tagged with a Leaf label as specified in section 5.1. This Leaf label allows for destination PE (e.g., egress PE) to perform the necessary egress filtering function in data-plane similar to ESI label in [RFC7432]. The allocation of the Leaf label is on a per PE basis (e.g., independent of ESI and EVI) as described in the following sections.

The Leaf label can be upstream assigned for P2MP LSP or downstream assigned for ingress replication tunnels. The main difference between downstream and upstream assigned Leaf label is that in case of downstream assigned not all egress PE devices need to receive the label in MPLS encapsulated BUM packets just like ESI label for ingress replication procedures defined in [RFC7432].

On the ingress PE, the PE needs to place all its Leaf ACs for a given bridge domain in a single split-horizon group in order to prevent intra-PE forwarding among its Leaf ACs. This intra-PE split-horizon filtering applies to BUM traffic as well as known-unicast traffic.

There are four scenarios to consider as follows. In all these scenarios, the ingress PE imposes the right MPLS label associated with the originated Ethernet Segment (ES) depending on whether the Ethernet frame originated from a Root or a Leaf site on that Ethernet Segment (ESI label or Leaf label). The mechanism by which the PE identifies whether a given frame originated from a Root or a Leaf site on the segment is based on the AC identifier for that segment (e.g., Ethernet Tag of the frame for 802.1Q frames). Other mechanisms for identifying Root or Leaf sites such as the use of source MAC address of the receiving frame are optional. The scenarios below are described in context of Root/Leaf AC; however, they can be extended to Root/Leaf MAC address if needed.

3.2.1 BUM Traffic Originated from a Single-homed Site on a Leaf AC

In this scenario, the ingress PE adds a Leaf label advertised using the E-Tree Extended Community (Section 5.1) indicating a Leaf site. This Leaf label, used for single-homing scenarios, is not on a per ES basis but rather on a per PE basis - i.e., a single Leaf MPLS label is used for all single-homed ES's on that PE. This Leaf label is advertised to other PE devices, using the E-Tree Extended Community (section 5.1) along with an Ethernet Auto-discovery per ES (EAD-ES) route with ESI of zero and a set of Route Targets (RTs) corresponding to all EVIs on the PE where each EVI has at least one Leaf site. Multiple EAD-ES routes will need to be advertised if the number of Route Targets (RTs) that need to be carried exceed the limit on a single route per [RFC7432]. The ESI for the EAD-ES route is set to zero to indicate single-homed sites.

When a PE receives this special Leaf label in the data path, it blocks the packet if the destination AC is of type Leaf; otherwise, it forwards the packet.

3.2.2 BUM Traffic Originated from a Single-homed Site on a Root AC

In this scenario, the ingress PE does not add any ESI label or Leaf label and it operates per [RFC7432] procedures.

3.2.3 BUM Traffic Originated from a Multi-homed Site on a Leaf AC

In this scenario, it is assumed that while different ACs (VLANs) on the same ES could have different Root/Leaf designation (some being Roots and some being Leafs), the same VLAN does have the same Root/Leaf designation on all PEs on the same ES. Furthermore, it is assumed that there is no forwarding among subnets - ie, the service is EVPN L2 and not EVPN IRB [EVPN-IRB]. IRB use cases described in [EVPN-IRB] are outside the scope of this document.

In this scenario, if a multicast or broadcast packet is originated from a Leaf AC, then it only needs to carry Leaf label described in section 3.2.1. This label is sufficient in providing the necessary egress filtering of BUM traffic from getting sent to Leaf ACs including the Leaf AC on the same Ethernet Segment.

3.2.4 BUM Traffic Originated from a Multi-homed Site on a Root AC

In this scenario, both the ingress and egress PE devices follows the procedure defined in [RFC7432] for adding and/or processing an ESI MPLS label - i.e., existing procedures for BUM traffic in [RFC7432] are sufficient and there is no need to add a Leaf label.

3.3 E-Tree Traffic Flows for EVPN

Per [RFC7387], a generic E-Tree service supports all of the following traffic flows:

- Known unicast traffic from Root to Roots & Leaf
- Known unicast traffic from Leaf to Root
- BUM traffic from Root to Roots & Leafs
- BUM traffic from Leaf to Roots

A particular E-Tree service may need to support all of the above types of flows or only a select subset, depending on the target application. In the case where only multicast and broadcast flows need to be supported, the L2VPN PEs can avoid performing any MAC learning function.

The following subsections will describe the operation of EVPN to support E-Tree service with and without MAC learning.

3.3.1 E-Tree with MAC Learning

The PEs implementing an E-Tree service must perform MAC learning when unicast traffic flows must be supported among Root and Leaf sites. In this case, the PE(s) with Root sites performs MAC learning in the data-path over the Ethernet Segments, and advertises reachability in EVPN MAC/IP Advertisement Routes. These routes will be imported by all PEs for that EVI (i.e., PEs that have Leaf sites as well as PEs that have Root sites). Similarly, the PEs with Leaf sites perform MAC learning in the data-path over their Ethernet Segments, and advertise reachability in EVPN MAC/IP Advertisement Routes. For scenarios where two different RTs are used per EVI (one to designate Root site and another to designate Leaf site), the MAC/IP Advertisement routes are imported only by PEs with at least one Root site in the EVI - i.e., a PE with only Leaf sites will not import these routes. PEs with Root and/or Leaf sites may use the Ethernet Auto-discovery per EVI (EAD-EVI) routes for aliasing (in the case of multi-homed segments) and EAD-ES routes for mass MAC withdrawal per [RFC7432].

To support multicast/broadcast from Root to Leaf sites, either a P2MP tree rooted at the PE(s) with the Root site(s) (e.g., Root PEs) or ingress replication can be used (section 16 of [RFC7432]). The multicast tunnels are set up through the exchange of the EVPN Inclusive Multicast route, as defined in [RFC7432].

To support multicast/broadcast from Leaf to Root sites, either ingress replication tunnels from each Leaf PE or a P2MP tree rooted at each Leaf PE can be used. The following two paragraphs describes

when each of these tunneling schemes can be used and how to signal them.

When there are only a few Root PEs with small amount of multicast/broadcast traffic from Leaf PEs toward Root PEs, then ingress replication tunnels from Leaf PEs toward Root PEs should be sufficient. Therefore, if a Root PE needs to support a P2MP tunnel in transmit direction from itself to Leaf PEs and at the same time it wants to support ingress-replication tunnels in receive direction, the Root PE can signal it efficiently by using a new composite tunnel type defined in section 5.2. This new composite tunnel type is advertised by the Root PE to simultaneously indicate a P2MP tunnel in transmit direction and an ingress-replication tunnel in the receive direction for the BUM traffic.

If the number of Root PEs is large, P2MP tunnels (e.g., mLDP or RSVP-TE) originated at the Leaf PEs may be used and thus there will be no need to use the modified PMSI tunnel attribute and the composite tunnel type values defined in section 5.2.

3.3.2 E-Tree without MAC Learning

The PEs implementing an E-Tree service need not perform MAC learning when the traffic flows between Root and Leaf sites are mainly multicast or broadcast. In this case, the PEs do not exchange EVPN MAC/IP Advertisement Routes. Instead, the Inclusive Multicast Ethernet Tag route is used to support BUM traffic. In such scenarios, the small amount of unicast traffic (if any) is sent as part of BUM traffic.

The fields of this route are populated per the procedures defined in [RFC7432], and the multicast tunnel setup criteria are as described in the previous section.

Just as in the previous section, if the number of Root PEs are only a few and thus ingress replication is desired from Leaf PEs to these Root PEs, then the modified PMSI attribute and the composite tunnel type values defined in section 5.2 should be used.

4 Operation for PBB-EVPN

In PBB-EVPN, the PE advertises a Root/Leaf indication along with each B-MAC Advertisement route to indicate whether the associated B-MAC address corresponds to a Root or a Leaf site. Just like the EVPN case, the new E-Tree Extended Community defined in section [5.1] is advertised with each EVPN MAC/IP Advertisement route.

In the case where a multi-homed Ethernet Segment has both Root and Leaf sites attached, two B-MAC addresses are advertised: one B-MAC address is per ES as specified in [RFC7623] and implicitly denoting Root, and the other B-MAC address is per PE and explicitly denoting Leaf. The former B-MAC address is not advertised with the E-Tree extended community but the latter B-MAC denoting Leaf is advertised with the new E-Tree extended community where "Leaf-indication" flag is set. In multi-homing scenarios where an Ethernet Segment has both Root and Leaf ACs, it is assumed that while different ACs (VLANs) on the same ES could have different Root/Leaf designation (some being Roots and some being Leafs), the same VLAN does have the same Root/Leaf designation on all PEs on the same ES. Furthermore, it is assumed that there is no forwarding among subnets - ie, the service is L2 and not IRB. IRB use case is outside the scope of this document.

The ingress PE uses the right B-MAC source address depending on whether the Ethernet frame originated from the Root or Leaf AC on that Ethernet Segment. The mechanism by which the PE identifies whether a given frame originated from a Root or Leaf site on the segment is based on the Ethernet Tag associated with the frame. Other mechanisms of identification, beyond the Ethernet Tag, are outside the scope of this document.

Furthermore, a PE advertises two special global B-MAC addresses: one for Root and another for Leaf, and tags the Leaf one as such in the MAC Advertisement route. These B-MAC addresses are used as source addresses for traffic originating from single-homed segments. The B-MAC address used for indicating Leaf sites can be the same for both single-homed and multi-homed segments.

4.1 Known Unicast Traffic

For known unicast traffic, the PEs perform ingress filtering: On the ingress PE, the C-MAC [RFC7623] destination address lookup yields, in addition to the target B-MAC address and forwarding adjacency, a flag which indicates whether the target B-MAC is associated with a Root or a Leaf site. The ingress PE also checks the status of the originating site, and if both are a Leaf, then the packet is not forwarded.

4.2 Broadcast, Unknown, and Multicast (BUM) Traffic

For BUM traffic, the PEs must perform egress filtering. When a PE receives an EVPN MAC/IP advertisement route (which will be used as a source B-MAC for BUM traffic), it updates its egress filtering (based on the source B-MAC address), as follows:

- If the EVPN MAC/IP Advertisement route indicates that the advertised B-MAC is a Leaf, and the local Ethernet Segment is a Leaf as well, then the source B-MAC address is added to its B-MAC list used for egress filtering - i.e., to block traffic from that B-MAC address.
- Otherwise, the B-MAC filtering list is not updated.
- If the EVPN MAC/IP Advertisement route indicates that the advertised B-MAC has changed its designation from a Leaf to a Root and the local Ethernet Segment is a Leaf, then the source B-MAC address is removed from the B-MAC list corresponding to the local Ethernet Segment used for egress filtering - i.e., to unblock traffic from that B-MAC address.

When the egress PE receives the packet, it examines the B-MAC source address to check whether it should filter or forward the frame. Note that this uses the same filtering logic as baseline [RFC7623] for an ESI and does not require any additional flags in the data-plane.

Just as in section 3.2, the PE places all Leaf Ethernet Segments of a given bridge domain in a single split-horizon group in order to prevent intra-PE forwarding among Leaf segments. This split-horizon function applies to BUM traffic as well as known-unicast traffic.

4.3 E-Tree without MAC Learning

In scenarios where the traffic of interest is only multicast and/or broadcast, the PEs implementing an E-Tree service do not need to do any MAC learning. In such scenarios the filtering must be performed on egress PEs. For PBB-EVPN, the handling of such traffic is per section 4.2 without the need for C-MAC learning (in data-plane) in I-component (C-bridge table) of PBB-EVPN PEs (at both ingress and egress PEs).

5 BGP Encoding

This document defines a new BGP Extended Community for EVPN.

5.1 E-Tree Extended Community

This Extended Community is a new transitive Extended Community [RFC4360] having a Type field value of 0x06 (EVPN) and the Sub-Type 0x05. It is used for Leaf indication of known unicast and BUM traffic. It indicates that the frame is originated from a Leaf site.

The E-Tree Extended Community is encoded as an 8-octet value as follows:

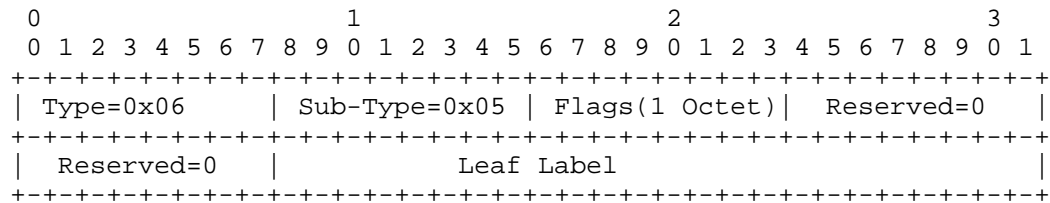
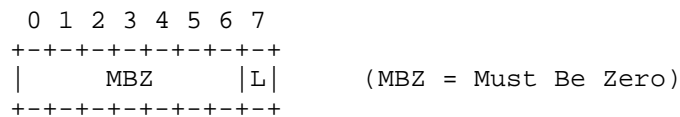


Figure 4: E-Tree Extended Community

The Flags field has the following format:



This document defines the following flags:

+ Leaf-Indication (L)

A value of one indicates a Leaf AC/Site. The rest of flag bits are reserved and should be set to zero.

When this Extended Community (EC) is advertised along with MAC/IP Advertisement route (for known unicast traffic) per section 3.1, the Leaf-Indication flag MUST be set to one and Leaf Label SHOULD be set to zero. The receiving PE MUST ignore Leaf Label and only processes Leaf-Indication flag. A value of zero for Leaf-Indication flag is invalid when sent along with MAC/IP advertisement route and an error should be logged.

When this EC is advertised along with EAD-ES route (with ESI of zero) for BUM traffic to enable egress filtering on disposition PEs per sections 3.2.1 and 3.2.3, the Leaf Label MUST be set to a valid MPLS label (i.e., non-reserved assigned MPLS label [RFC3032]) and the Leaf-Indication flag SHOULD be set to zero. The value of the 20-bit MPLS label is encoded in the high-order 20 bits of the Leaf Label field. The receiving PE MUST ignore the Leaf-Indication flag. A non-valid MPLS label when sent along with the EAD-ES route, should be ignored and logged as an error.

The reserved bits SHOULD be set to zero by the transmitter and MUST

be ignored by the receiver.

5.2 PMSI Tunnel Attribute

[RFC6514] defines PMSI Tunnel attribute which is an optional transitive attribute with the following format:

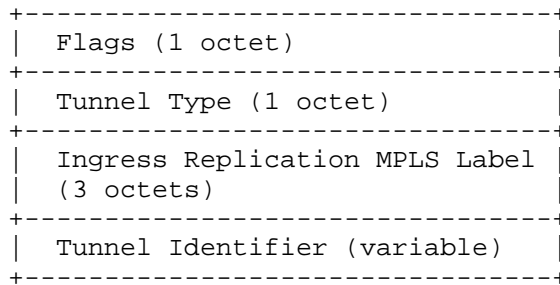


Figure 5: PMSI Tunnel Attribute

This document defines a new Composite tunnel type by introducing a new 'Composite Tunnel' bit in the Tunnel Type field and adding a MPLS label to the Tunnel Identifier field of PMSI Tunnel attribute as detailed below. All other fields remain as defined in [RFC6514]. Composite tunnel type is advertised by the Root PE to simultaneously indicate a non-(ingress replication) tunnel (e.g., P2MP tunnel) in transmit direction and an ingress-replication tunnel in the receive direction for the BUM traffic.

When receiver ingress-replication labels are needed, the high-order bit of the tunnel type field (Composite Tunnel bit) is set while the remaining low-order seven bits indicate the tunnel type as before (for the existing tunnel types). When this Composite Tunnel bit is set, the "tunnel identifier" field begins with a three-octet label, followed by the actual tunnel identifier for the transmit tunnel. PEs that don't understand the new meaning of the high-order bit treat the tunnel type as an undefined tunnel type and treat the PMSI tunnel attribute as a malformed attribute [RFC6514]. That is why the composite tunnel bit is allocated in the Tunnel Type field rather than the Flags field. For the PEs that do understand the new meaning of the high-order, if ingress replication is desired when sending BUM traffic, the PE will use the the label in the Tunnel Identifier field when sending its BUM traffic.

Using the Composite Tunnel bit for Tunnel Types 0x00 'no tunnel information present' and 0x06 'Ingress Replication' is invalid, and a

PE that receives a PMSI Tunnel attribute with such information, considers it as malformed and it SHOULD treat this Update as though all the routes contained in this Update had been withdrawn per section 5 of [RFC6514].

6 Acknowledgement

We would like to thank Eric Rosen, Jeffrey Zhang, Wen Lin, Aldrin Issac, Wim Henderickx, Dennis Cai, and Antoni Przygienda for their valuable comments and contributions. The authors would also like to thank Thomas Morin for shepherding this document and providing valuable comments.

7 Security Considerations

Since this document uses the EVPN constructs of [RFC7432] and [RFC7623], the same security considerations in these documents are also applicable here. Furthermore, this document provides an additional security check by allowing sites (or ACs) of an EVPN instance to be designated as "Root" or "Leaf" by the network operator/ service provider and thus preventing any traffic exchange among "Leaf" sites of that VPN through ingress filtering for known unicast traffic and egress filtering for BUM traffic. Since by default and for the purpose of backward compatibility, an AC that doesn't have a Leaf designation is considered as a Root AC, in order to avoid any traffic exchange among Leaf ACs, the operator SHOULD configure the AC with a proper role (Leaf or Root) before activating the AC.

8 IANA Considerations

IANA has allocated value 5 in the "EVPN Extended Community Sub-Types" registry defined in [RFC7153] as follow:

SUB-TYPE	VALUE	NAME	Reference
	0x05	E-Tree Extended Community	This document

This document creates a one-octet registry called "E-Tree Flags". New registrations will be made through the "RFC Required" procedure defined in [RFC8126]. Initial registrations are as follows:

bit	Name	Reference
0-6	Unassigned	
7	Leaf-Indication	This document

8.1 Considerations for PMSI Tunnel Types

The "P-Multicast Service Interface Tunnel (PMSI Tunnel) Tunnel Types" registry in the "Border Gateway Protocol (BGP) Parameters" registry needs to be updated to reflect the use of the most significant bit as "Composite Tunnel" bit (section 5.2).

For this purpose, this document updates [RFC7385] by changing the previously unassigned values (i.e., 0x08 - 0xFA) as follow:

Value	Meaning	Reference
0x08-0x7A	Unassigned	
0x7B-0x7E	Experimental	this document
0x7F	Reserved	this document
0x80-0xFA	Reserved for Composite tunnel	this document
0xFB-0xFE	Experimental	[RFC7385]
0xFF	Reserved	[RFC7385]

The allocation policy for values 0x08-0x7A is per IETF Review [RFC8126]. The range for experimental has been expanded to include the previously assigned range of 0xFB-0xFE and the new range of 0x7B-0x7E. The value in these ranges are not to be assigned. The value 0x7F which is the mirror image of (0xFF) is reserved in this document.

9 References

9.1 Normative References

[KEYWORDS] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, March 1997.

[RFC8126] Cotton et al, "Guidelines for Writing an IANA Considerations Section in RFCs", June, 2017.

[RFC7387] Key et al., "A Framework for E-Tree Service over MPLS Network", October 2014.

[MEF6.1] Metro Ethernet Forum, "Ethernet Services Definitions - Phase

2", MEF 6.1, April 2008, https://mef.net/PDF_Documents/technical-specifications/MEF6-1.pdf

[RFC7432] Sajassi et al., "BGP MPLS Based Ethernet VPN", February, 2015.

[RFC7623] Sajassi et al., "Provider Backbone Bridging Combined with Ethernet VPN (PBB-EVPN)", September, 2015.

[RFC7385] Andersson et al., "IANA Registry for P-Multicast Service Interface (PMSI) Tunnel Type Code Points", October, 2014.

[RFC7153] Rosen et al., "IANA Registries for BGP Extended Communities", March, 2014.

[RFC6514] Aggarwal et al., "BGP Encodings and Procedures for Multicast in MPLS/BGP IP VPNs", February, 2012.

[RFC4360] Sangli et al., "BGP Extended Communities Attribute", February, 2006.

9.2 Informative References

[RFC4360] S. Sangli et al, "BGP Extended Communities Attribute", February, 2006.

[RFC3032] E. Rosen et al, "MPLS Label Stack Encoding", January 2001.

[RFC7796] Y. Jiang et al, "Ethernet-Tree (E-Tree) Support in Virtual Private LAN Service (VPLS)", March 2016.

[EVPN-IRB] A. Sajassi et al, "Integrated Routing and Bridging in EVPN", draft-ietf-bess-evpn-inter-subnet-forwarding-03, February 8, 2017.

[802.1ah] IEEE, "IEEE Standard for Local and metropolitan area networks - Media Access Control (MAC) Bridges and Virtual Bridged Local Area Networks", Clauses 25 and 26, IEEE Std 802.1Q, DOI 10.1109/IEEESTD.2011.6009146.

Appendix-A

When two MAC-VRFs (two bridge tables per VLANs) are used for an E-Tree service (one for Root ACs and another for Leaf ACs) on a given PE, then the following complications in data-plane path can result.

Maintaining two MAC-VRFs (two bridge tables) per VLAN (when both Leaf and Root ACs exists for that VLAN) would either require two lookups

be performed per MAC address in each direction in case of a miss, or duplicating many MAC addresses between the two bridge tables belonging to the same VLAN (same E-Tree instance). Unless two lookups are made, duplication of MAC addresses would be needed for both locally learned and remotely learned MAC addresses. Locally learned MAC addresses from Leaf ACs need to be duplicated onto Root bridge table and locally learned MAC addresses from Root ACs need to be duplicated onto Leaf bridge table. Remotely learned MAC addresses from Root ACs need to be copied onto both Root and Leaf bridge tables. Because of potential inefficiencies associated with dataplane implementation of additional MAC lookup or duplication of MAC entries, this option is not believed to be implementable without dataplane performance inefficiencies in some platforms and thus this document introduces the coloring as described in section 2.2 and detailed in section 3.1.

Authors' Addresses

Ali Sajassi
Cisco
Email: sajassi@cisco.com

Samer Salam
Cisco
Email: ssalam@cisco.com

John Drake
Juniper
Email: jdrake@juniper.net

Jim Uttaro
AT&T
Email: jul738@att.com

Sami Boutros
VMware
Email: sboutros@vmware.com

Jorge Rabadan
Nokia
Email: jorge.rabadan@nokia.com

