

L2VPN Workgroup
INTERNET-DRAFT
Intended Status: Standards Track

Ali Sajassi
Samer Salam
Sami Boutros
Cisco

Wim Henderickx
Jorge Rabadan
Alcatel-Lucent

Jim Uttaro
AT&T

John Drake
Wen Lin
Juniper

Aldrin Isaac
Bloomberg

Expires: December 18, 2015

June 18, 2015

**E-TREE Support in EVPN & PBB-EVPN
draft-ietf-bess-evpn-etree-01**

Abstract

The Metro Ethernet Forum (MEF) has defined a rooted-multipoint Ethernet service known as Ethernet Tree (E-Tree). [[ETREE-FMWK](#)] proposes a solution framework for supporting this service in MPLS networks. This document discusses how those functional requirements can be easily met with (PBB-)EVPN and how (PBB-)EVPN offers a more efficient implementation of these functions.

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) 2013 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	Terminology	4
2	E-Tree Scenarios and EVPN / PBB-EVPN Support	4
2.1	Scenario 1: Leaf OR Root site(s) per PE	4
2.2	Scenario 2: Leaf AND Root site(s) per PE	5
2.3	Scenario 3: Leaf AND Root site(s) per Ethernet Segment	5
3	Operation for EVPN	6
3.1	Known Unicast Traffic	7
3.2	BUM Traffic	7
3.2.1	BUM Traffic supported by P2MP Tunnels	7
3.2.1	BUM Traffic supported by Ingress Replication	9
3.3	E-TREE Traffic Flows for EVPN	10
3.3.1	E-Tree with MAC Learning	11
3.3.2	E-Tree without MAC Learning	11
4	Operation for PBB-EVPN	12
4.1	Known Unicast Traffic	12
4.2	BUM Traffic	13
5	BGP Encoding	13
5.1	E-TREE Extended Community	13
6	Acknowledgement	14
7	Security Considerations	14
8	IANA Considerations	14
9	References	14
9.1	Normative References	14
9.2	Informative References	14
	Authors' Addresses	15

1 Introduction

The Metro Ethernet Forum (MEF) has defined a rooted-multipoint Ethernet service known as Ethernet Tree (E-Tree). In an E-Tree service, endpoints are labeled as either Root or Leaf sites. Root sites can communicate with all other sites. Leaf sites can communicate with Root sites but not with other Leaf sites.

[ETREE-FMWK] proposes the solution framework for supporting E-Tree service in MPLS networks. The document identifies the functional components of the overall solution to emulate E-Tree services in addition to Ethernet LAN (E-LAN) services on an existing MPLS network.

[EVPN] is 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. [PBB-EVPN] combines the functionality of EVPN with [802.1ah] Provider Backbone Bridging for MAC address scalability.

This document discusses how the functional requirements for E-Tree service can be easily met with (PBB-)EVPN and how (PBB-)EVPN offers a more efficient implementation of these functions.

1.1 Terminology

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](#)].

2 E-Tree Scenarios and EVPN / PBB-EVPN Support

In this section, we will categorize support for E-Tree into three different scenarios, depending on the nature of the site association (Root/Leaf) per PE or per Ethernet Segment:

- Leaf OR Root site(s) per PE
- Leaf AND Root site(s) per PE
- Leaf AND Root site(s) per Ethernet Segment

2.1 Scenario 1: Leaf OR Root site(s) per PE

In this scenario, a PE may have Root sites OR Leaf sites for a given VPN instance, but not both concurrently. The PE may have both Root

and Leaf sites albeit for different VPNs. Every Ethernet Segment connected to the PE is uniquely identified as either a Root or a Leaf site.

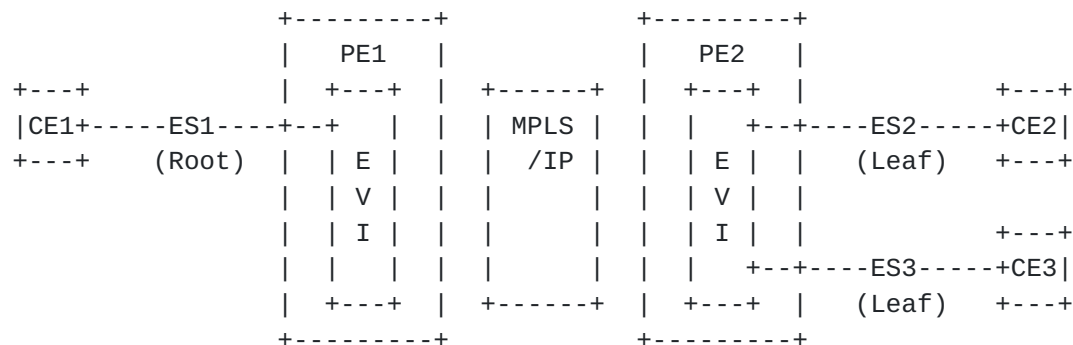


Figure 1: Scenario 1

2.2 Scenario 2: Leaf AND Root site(s) per PE

In this scenario, a PE may have a set of one or more Root sites AND a set of one or more Leaf sites for a given VPN instance. Every Ethernet Segment connected to the PE is uniquely identified as either a Root or a Leaf site.

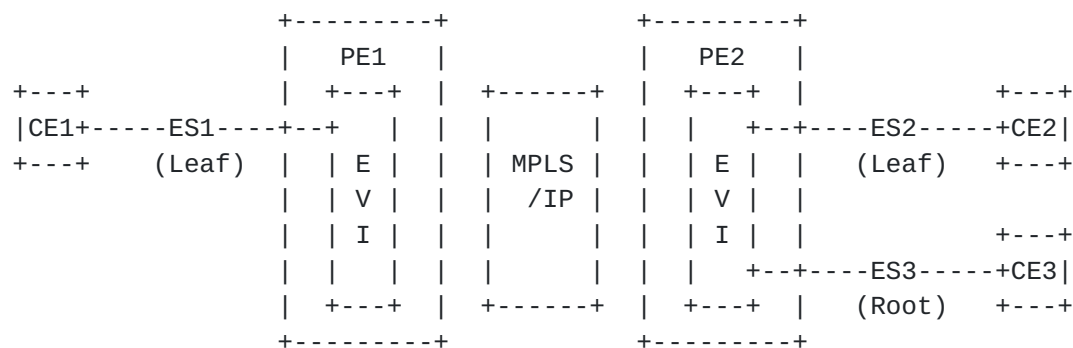


Figure 2: Scenario 2

2.3 Scenario 3: Leaf AND Root site(s) per Ethernet Segment

In this scenario, a PE may have a set of one or more Root sites AND a set of one or more Leaf sites for a given VPN instance. An Ethernet Segment connected to the PE may be identified as both a Root and a Leaf site concurrently.

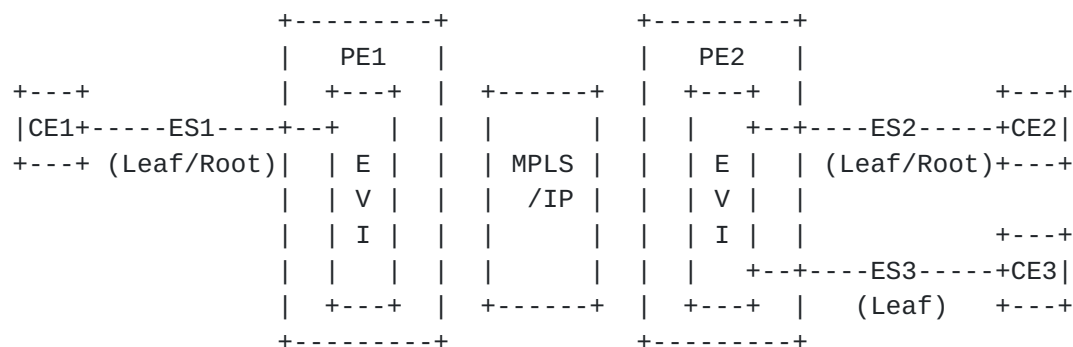


Figure 3: Scenario 3

3 Operation for EVPN

[EVPN] defines the notion of an Ethernet Segment which can be readily used to identify a Root and/or Leaf site in E-TREE services. In other words, [EVPN] has inherent capability to support E-TREE services without defining any new BGP routes. It only requires a minor modification to the existing procedures and a new BGP Extended Community for leaf indication as shown later in this document.

In addition to the procedures below (which is a MUST requirement), an EVPN PE implementation MAY provide topology constraint among the PEs belonging to the same EVI associated with an E-TREE service. The purpose of this topology constraint is to avoid having PEs with only Leaf sites (e.g., scenario 1 in [section 2.1](#)) importing and processing BGP MAC routes from each other, thereby unnecessarily exhausting their RIB tables. However, when a Root site is added to a Leaf PE (e.g., scenario 2 and 3 in [sections 2.2](#) and [2.3](#)), then that PE needs to process MAC routes from all other Leaf PEs and add them to its forwarding table. 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 and the other is associated with the Leaf sites. 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. If for a given EVI, the PEs can eventually have both Leaf and Root sites attached, even though they may start as Root-only or Leaf-only PEs, then it is recommended to use a single RT per EVI and avoid additional configuration and operational overhead. If the number of EVIs is very large (e.g., more than 32K or 64K), then RT type 0 as defined in [\[RFC4360\]](#) SHOULD be used; otherwise, RT type 2 is sufficient.

The following procedures are used consistently for all the scenarios highlighted in the previous section.

3.1 Known Unicast Traffic

For known unicast traffic, the PE must advertise a Leaf indication along with each MAC Advertisement route, to indicate that the associated MAC address was learnt from a Leaf Attachment Circuit (AC). The lack of a Leaf indication, indicates the MAC address is learnt from a root AC. In other words, the default mode of operation in an EVPN is that all ACs are root (can transmit and receive traffic to/from other ACs in an EVI) unless the AC is explicitly identified as a leaf.

Tagging MAC addresses with a leaf indication when they are associated with a leaf AC, 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 Leafs, then the packet is not forwarded.

The PE places all Leaf ACs of a given bridge domain in a single split-horizon group in order to prevent intra-PE forwarding among Leaf ACs. This split-horizon function applies to both known unicast and BUM traffic.

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 is advertised with each EVPN MAC/IP Advertisement route.

3.2 BUM Traffic

For BUM traffic, it is not possible to perform filtering on the ingress PE, as is the case with known unicast, because of the multi-destination nature of the traffic. 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 Root or a Leaf AC, the MPLS-encapsulated frames MUST be tagged with an indication of whether they originated from a Root or a Leaf AC. This can be achieved in EVPN through the use of the ESI MPLS label. Therefore, the ESI MPLS label not only identifies the Ethernet segment of origin for a given frame, but also it identifies its type (e.g., Leaf or Root).

3.2.1 BUM Traffic supported by P2MP Tunnels

For multi-homing use cases where BUM traffic uses P2MP LSP, the ingress PE adds an upstream-assigned ESI MPLS label to the frame per

[RFC7432] procedures and sends it to all the intended ingress PE devices. Two ESI MPLS labels are used for each multi-homed Ethernet segment that has both Root and Leaf sites: one ESI MPLS label that only identifies the Ethernet segment of origin per [RFC7432] and another one that not only identifies the Ethernet segment of origin but also its type (which is Leaf). If an Ethernet segment has only Root sites, then the former ESI MPLS label is used and if an Ethernet segment has only Leaf sites, then the latter ESI MPLS label is used.

It should be noted that the former ESI MPLS label implicitly identifies a Root Ethernet segment - i.e., an ESI MPLS label that is signaled without the new E-TREE Extended Community (defined in section [5.1]), is assumed to be of type Root. When advertising the ESI MPLS label for an Ethernet Segment that has Leaf sites, the PE MUST indicate that the corresponding ESI is of type Leaf. This is achieved by advertising the Ethernet A-D per ES route with with the ESI MPLS label Extended Community along with the new E-TREE Extended Community that has a Leaf indication flag.

The egress PE can determine whether or not to forward a particular frame to the destination Ethernet Segment depending on the following rules:

- If the ESI MPLS label indicates that the source Ethernet Segment is the same as destination Ethernet segment, then the frame is blocked according to the split-horizon rule in [RFC7432].
- If the ESI MPLS label indicates that the source Ethernet Segment is not the same as destination Ethernet segment and it doesn't have any Leaf indication, then the frame is forwarded to the destination AC according to the split-horizon rule in [RFC7432].
- If the ESI MPLS label indicates that the source Ethernet Segment is not the same as destination Ethernet segment but it has a Leaf indication, then the frame is blocked if the destination AC is of type Leaf and it is forwarded if the destination AC is of type Root.

The ingress PE imposes the right ESI MPLS label depending on whether the Ethernet frame originated from the Root or Leaf site 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 (e.g., whether the frame come from a leaf or a root AC). Other mechanisms of identification, beyond the Ethernet Tag, are outside the scope of this document. It should be noted that support for both Root and Leaf sites on a single Ethernet Segment requires that the PE performs the Ethernet Segment split-horizon check on a per Ethernet Tag basis. In the case where a multi-homed Ethernet Segment has only either Root or

Leaf sites attached, then a single ESI MPL label is allocated and advertised.

For single-homing use cases where BUM traffic uses P2MP LSP, the ingress PE adds a special ESI MPLS label to the frame if the frame is originated from a Leaf site. This special ESI MPLS label used for single-homing scenarios is not on a per ES basis but rather on a per PE basis - i.e., a single ESI MPLS label is used for all single-homed segments on that PE. If the frame is originated from a Root site, then the ingress PE does not add any ESI MPLS label per [\[RFC7432\]](#) procedures. The egress PE, when receiving this special ESI MPLS label, it blocks the frame if the destination AC is of type Leaf and it forwards the frame if the destination AC is of type Root.

When a PE wants to advertise this special ESI label to other PE devices, it advertises it using ESI MPLS label Extended Community with the Ethernet A-D per ES route. The ESI for the Ethernet A-D per ES route, can be of type 3, 4, or 5.

[3.2.1](#) BUM Traffic supported by Ingress Replication

The procedures for supporting BUM traffic using ingress replication, are similar to the ones in the previous section. The main differences are that the ESI label is downstream assigned and not all egress PE devices need to receive the ESI label just like ingress replication procedures defined in [\[RFC7432\]](#).

For frames received from a multi-homed Ethernet segment, the ingress PE may or may not add an ESI MPLS label based on the following criteria:

- If the frame is forwarded to a PE that participates in the same multi-homed Ethernet Segment and the frame is received on a Root AC, then the ingress PE adds a per-ES downstream-assigned ESI MPLS label to the frame per [\[RFC7432\]](#).
- If the frame is forwarded to a PE that participates in the same multi-homed Ethernet Segment and the frame is received on a Leaf AC, then the ingress PE adds the a per-ES downstream-assigned ESI MPLS label indicating Leaf to the frame.
- If the frame is forwarded to a PE that does not participate in the same multi-homed Ethernet Segment and the frame is received on a Leaf AC, then the ingress PE adds a per-PE downstream-assigned special ESI MPLS label indicating Leaf to the frame. This special ESI MPLS label is per PE.

- If the frame is forwarded to a PE that does not participate in the same multi-homed Ethernet Segment and the frame is received on a Root AC, then the ingress PE does not add any ESI MPLS label to the frame per [\[RFC7432\]](#).

For frames received from a single-homed Ethernet segment, the ingress PE may or may not add an ESI MPLS label based on the following criteria:

- If the frame is received on a Root AC, then the ingress PE does not add any ESI MPLS label to the frame.
- If the frame is received on a Leaf AC, then the ingress PE adds a special downstream-assigned ESI MPLS label indicating Leaf to the frame.

Just as described in the previous section, the Leaf indication is signaled using the new E-TREE extended community defined in section [5.1] along with the ESI MPLS label extended community with the Ethernet A-D per ES route.

The egress PE can determine whether or not to forward a particular frame to the destination Ethernet Segment depending on the following rules:

- If the ESI MPLS label indicates that the source Ethernet Segment is the same as destination Ethernet segment, then the frame is blocked according to the split-horizon rule in [\[RFC7432\]](#).
- If the ESI MPLS label indicates that the source Ethernet Segment is not the same as destination Ethernet segment and it doesn't have any Leaf indication, then the frame is forwarded to the destination AC according to the split-horizon rule in [\[RFC7432\]](#).
- If the ESI MPLS label indicates that the source Ethernet Segment is not the same as destination Ethernet segment but it has a Leaf indication, then the frame is blocked if the destination AC is of type Leaf and it is forwarded if the destination AC is of type Root.
- If the ESI label is a special ESI MPLS label, then the frame is blocked if the destination AC is of type Leaf and it is forwarded if the destination AC is of type Root.

[3.3](#) E-TREE Traffic Flows for EVPN

Per [\[ETREE-FMWK\]](#), a generic E-Tree service supports all of the following traffic flows:

- Ethernet Unicast from Root to Roots & Leaf
- Ethernet Unicast from Leaf to Root
- Ethernet Broadcast/Multicast from Root to Roots & Leafs
- Ethernet Broadcast/Multicast 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 unicast flows need not be supported, the L2VPN PE can avoid performing any MAC learning function.

In the subsections that follow, we 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 from Root to Leaf or from Leaf to Root sites. In this case, the PE with Root sites performs MAC learning in the data-path over the Ethernet Segments, and advertises reachability in EVPN MAC Advertisement routes. These routes will be imported by PEs that have Leaf sites as well as by PEs that have Root sites, in a given EVI. Similarly, the PEs with Leaf sites perform MAC learning in the data-path over their Ethernet Segments, and advertise reachability in EVPN MAC Advertisement routes which are imported only by PEs with at least one Root site in the EVI. A PE with only Leaf sites will not import these routes. PEs with Root and/or Leaf sites may use the Ethernet A-D routes for aliasing (in the case of multi-homed segments) and for mass MAC withdrawal.

To support multicast/broadcast from Root to Leaf sites, either a P2MP tree rooted at the PE(s) with the Root site(s) or ingress replication can be used. 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, ingress replication should be sufficient for most scenarios where there is a single Root or few Roots. If the number of Roots is large, a P2MP tree rooted at the PEs with Leaf sites may be used.

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 multicast or broadcast. In this case, the PEs do not exchange EVPN MAC Advertisement routes. Instead, the Ethernet A-D routes are used to exchange the EVPN labels.

The fields of the Ethernet A-D route are populated per the procedures defined in [[RFC7432](#)], and the route import rules are as described in previous sections.

4 Operation for PBB-EVPN

In PBB-EVPN, the PE must advertise a Root/Leaf indication along with each MAC Advertisement route, to indicate whether the associated B-MAC address corresponds to a Root or a Leaf site. Similar to the EVPN case, this flag will be added to the new E-TREE extended community defined in section [5.1], and advertised with each MAC Advertisement route.

In the case where a multi-homed Ethernet Segment has both Root and Leaf sites attached, two B-MAC addresses are allocated and advertised: one B-MAC address implicitly denoting Root and the other 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.

The ingress PE uses the right B-MAC source address depending on whether the Ethernet frame originated from the Root or Leaf site 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. It should be noted that support for both Root and Leaf sites on a single Ethernet Segment requires that the PE performs the Ethernet Segment split-horizon check on a per Ethernet Tag basis.

In the case where a multi-homed Ethernet Segment has only Root or Leaf sites attached, then a single B-MAC address is allocated and advertised per segment.

Furthermore, a PE advertises two special global B-MAC addresses: one for Root and another for Leaf, and tags them as such in the MAC Advertisement routes. These B-MAC addresses are used as source addresses for traffic originating from single-homed segments.

4.1 Known Unicast Traffic

For known unicast traffic, the PEs perform ingress filtering: On the ingress PE, the C-MAC 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 cross-checks this flag with the status of the originating site, and if both are a Leaf, then the packet is not

forwarded.

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.

4.2 BUM Traffic

For BUM traffic, the PEs must perform egress filtering. When a PE receives a MAC advertisement route, it updates its Ethernet Segment egress filtering function (based on the B-MAC source address), as follows:

- If the MAC 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 the B-MAC filtering list.
- Otherwise, the B-MAC filtering list is not updated.

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 [[PBB-EVPN](#)] and does not require any additional flags in the data-plane.

5 BGP Encoding

This document defines one new BGP Extended Community for EVPN.

5.1 E-TREE Extended Community

A new EVPN BGP Extended Community called E-TREE is introduced here. This new extended community is a transitive extended community with the Type field of 0x06 (EVPN) and the Sub-Type of 0x04. This extended community is used to for leaf indication and it is advertised with an EVPN MAC/IP route or an Ethernet A-D per ES route. When advertised with an Ethernet A-D per ES route, it is sent along with ESI Label Extended Community defined in [section 7.5 of \[RFC7432\]](#).

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


```

      0              1              2              3
    0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Type=0x06      | Sub-Type=0x04 |      E-TREE Flags      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                                     E-TREE Flags          | L |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Leaf flag (L): A value of 1 indicates a leaf

6 Acknowledgement

We would like to thank Dennis Cai for his comments.

7 Security Considerations

Same security considerations as [[RFC7432](#)].

8 IANA Considerations

Allocation of Extended Community Type and Sub-Type for EVPN.

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.

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

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

9.2 Informative References

[ETREE-FMWK] Key et al., "A Framework for E-Tree Service over MPLS Network", [draft-ietf-l2vpn-etree-frwk-03](#), work in progress, September 2013.

[PBB-EVPN] Sajassi et al., "PBB-EVPN", [draft-ietf-l2vpn-pbb-evpn-05.txt](#), work in progress, October, 2013.

Authors' Addresses

Ali Sajassi
Cisco
Email: sajassi@cisco.com

Samer Salam
Cisco
Email: ssalam@cisco.com

Wim Henderickx
Alcatel-Lucent
Email: wim.henderickx@alcatel-lucent.com

Jim Uttaro
AT&T
Email: ju1738@att.com

Aldrin
Bloomberg Issac
Email: aisaac71@bloomberg.net

Sami Boutros
Cisco
Email: sboutros@cisco.com

