Network Working Group                     R. Aggarwal (Editor)
Internet Draft                                Juniper Networks
Category: Standards Track
Expiration Date: December 2010                       A. Isaac
                                                    Bloomberg

                                                    J. Uttaro
                                                         AT&T

                                                   R. Shekhar
                                             Juniper Networks

                                                     F. Balus
                                               Alcatel-Lucent

                                                W. Henderickx
                                               Alcatel-Lucent

                                                June 2, 2010

                     **BGP MPLS Based MAC VPN**


                   draft-raggarwa-mac-vpn-01.txt

Status of this Memo

Abstract

   This document describes procedures for BGP MPLS based MAC VPNs (MAC-
   VPN).

Table of Contents

**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 [RFC2119].

**2**. **Contributors**

   In addition to the authors listed above, the following individuals
   also contributed to this document.

   Quaizar Vohra
   Kireeti Kompella
   Apurva Mehta
   Juniper Networks

**3**. **Introduction**

   This document describes procedures for BGP MPLS based MAC VPNs (MAC-
   VPN).

   There is a desire by Service Providers (SP) and data center providers
   to provide MPLS based bridged / LAN services or/and infrastructure
   such that they meet the requirements listed below. An example of such
   a service is a VPLS service offered by a SP. Another example is a
   MPLS based infrastructure in a data center. Here are the
   requirements:

   - Minimal or no configuration required. MPLS implementations have
   reduced the amount of configuration over the years. There is a need
   for greater auto-configuration.

   - Support of multiple active points of attachment for CEs which may
   be hosts, switches or routers. Current MPLS technologies such as
   VPLS, currently do not support this. This allows load-balancing among
   multiple active paths. Regular Ethernet switching technologies, based
   on MAC learning  do not allow the same MAC to be learned from two
   different PEs and be active at the same time in the same switching
   instance

   - Ability to span a VLAN across multiple racks in different
   geographic locations, which may not be in the same data center.

   - Minimize or eliminate flooding of unknown unicast traffic.

   - Allow hosts and Virtual Machines (VMs) in a data center to relocate
   without requiring renumbering. For instnace VMs may be moved for load
   or failure reasons.

   - Ability to scale up to hundreds of thousands of hosts or more
   across multiple data centers, where connectivity is required between
   hosts in different data centers.

   - Support for virtualization. This includes the ability to separate
   hosts and VMs working together from other such groups, and the
   ability to have overlapping IP and MAC addresses/

   - Fast convergence

   This document proposes a MPLS based technology, referred to as MPLS-
   based MAC VPN (MAC-VPN) for meeting the requirements described in
   this section. MAC-VPN requires extensions to existing IP/MPLS
   protocols as described in section 5. In addition to these extensions
   MAC-VPN uses several building blocks from existing MPLS technologies.


**4. Terminology**

   MES: MPLS Edge Switch
   CE: Host or router or switch
   MVI: MAC VPN Instance
   ESI: Ethernet segment identifier

**5**. **BGP MPLS Based MAC-VPN**

   This section describes the framework of MAC-VPN to meet the
   requirements described in section 3.

   An MAC-VPN comprises CEs that are connected to PEs or MPLS Edge
   Switches (MES) that comprise the edge of the MPLS infrastructure. A
   CE may be a host, a router or a switch. The MPLS Edge Switches
   provide layer 2 virtual bridge connectivity between the CEs. There
   may be multiple MAC VPNs in the provider's network. This document
   uses the terms MAC-VPN and MAC VPN inter-changeably. A MAC VPN
   routing and forwarding instance on a MES is referred to as a MAC VPN
   Instance (MVI).

   The MESes are connected by a MPLS LSP infrastructure which provides
   the benefits of MPLS such as fast-reroute, resiliency etc.

   In a MAC VPN, learning between MESes occurs not in the data plane (as
   happens with traditional bridging) but in the control plane. Control
   plane learning offers much greater control over the learning process,
   such as restricting who learns what, and the ability to apply
   policies.  Furthermore, the control plane chosen for this is BGP
   (very similar to IP VPNs (RFC 4364)), providing much greater scale,
   and the ability to "virtualize" or isolate groups of interacting
   agents (hosts, servers, Virtual Machines) from each other. In MAC
   VPNs MESes advertise the MAC addresses learned from the CEs that are
   connected to them, along with a MPLS label, to other MESes in the
   control plane. Control plane learning enables load balancing and
   allows CEs to connect to multiple active points of attachment. It
   also improves convergence times in the event of certain network
   failures.

   However, learning between MESes and CEs is done by the method best
   suited to the CE: data plane learning, IEEE 802.1x, LLDP, 802.1aq or
   other protocols.

   It is a local decision as to whether the Layer 2 forwarding table on
   a MES contains all the MAC destinations known to the control plane or
   implements a cache based scheme. For instance the forwarding table
   may be populated only with the MAC destinations of the active flows
   transiting a specific MES.

   The policy attributes of a MAC VPN are very similar to an IP VPN. A
   MAC-VPN instance requires a Route-Distinguisher (RD) and a MAC-VPN
   requires one or more Route-Targets (RTs). A CE attaches to a MAC-VPN
   on a MES in a particular MVI on a VLAN or simply an ethernet
   interface. When the point of attachment is a VLAN there may be one or
   more VLANs in a particular MAC-VPN. Some deployment scenarios

guarantee uniqueness of VLANs across MAC-VPNs: all points of
attachment of a given MAC VPN use the same VLAN, and no other MAC VPN
uses this VLAN. This document refers to this case as a "Default
Single VLAN MAC-VPN" and describes simplified procedures to optimize
for it.


## 6. Ethernet Segment Identifier

If a CE is multi-homed to two or more MESes, the set of attachment
circuits constitutes an "Ethernet segment". An Ethernet segment may
appear to the CE as a Link Aggregation Group (LAG).  Ethernet
segments have an identifier, called the "Ethernet Segment Identifier"
(ESI).  A single-homed CE is considered to be attached to a Ethernet
segment with ESI 0.  Otherwise, an Ethernet segment MUST have a
unique non-zero ESI.  The ESI can be assigned using various
mechanisms:

1. The ESI may be configured. For instance when MAC VPNs are used to
provide a VPLS service the ESI is fairly analogous to the VE ID used
for the procedures in [BGP-VPLS] or the Multi-homing site ID in [BGP-
VPLS-MH].

2. If LACP is used, between the MESes and CEs that are hosts, then
the ESI is determined by LACP. This is the 48 bit virtual MAC address
of the host for the LACP link bundle. As far as the host is concerned
it would treat the multiple MESes that it is homed to as the same
switch.  This allows the host to aggregate links to different MESes
in the same bundle.

3. If LLDP is used, between the MESes and CEs that are hosts, then
the ESI is determined by LLDP. The ESI will be specified in a
following version.

4. In the case of indirectly connected hosts and a bridged LAN
between the hosts and the MESes, the ESI is determined based on the
Layer 2 bridge protocol as follows:

   If STP is used then the value of the ESI is derived by listening
to BPDUs on the ethernet segment. The MES does not run STP. However
it does learn the Switch ID, MSTP ID and Root Bridge ID by listening
to BPDUs.  The ESI is as follows:

     {Switch ID (6 bits), MSTP ID (6 bits), Root Bridge ID (48
bits)}

**7**. **BGP MAC-VPN NLRI**

   This document defines a new BGP NLRI, called the MAC-VPN NLRI.

   Following is the format of the MAC-VPN NLRI:

```
               +-----------------------------------+
               |     Route Type (1 octet)          |
               +-----------------------------------+
               |      Length (1 octet)             |
               +-----------------------------------+
               | Route Type specific (variable)    |
               +-----------------------------------+
```

   The Route Type field defines encoding of the rest of MAC-VPN NLRI
   (Route Type specific MAC-VPN NLRI).

   The Length field indicates the length in octets of the Route Type
   specific field of MAC-VPN NLRI.

   This document defines the following Route Types:

   + 1 - Ethernet Tag Auto-Discovery (A-D) route
   + 2 - MAC advertisement route
   + 3 - Inclusive Multicast Ethernet Tag Route
   + 4 - Ethernet Segment Route
   + 5 - Selective Multicast Auto-Discovery (A-D) Route
   + 6 - Leaf Auto-Discovery (A-D) Route

   The detailed encoding and procedures for these route types are
   described in subsequent sections.

   The MAC-VPN NLRI is carried in BGP [RFC4271] using BGP Multiprotocol
   Extensions [RFC4760] with an AFI of TBD and an SAFI of MAC-VPN (To be
   assigned by IANA). The NLRI field in the
   MP_REACH_NLRI/MP_UNREACH_NLRI attribute contains the MAC-VPN NLRI
   (encoded as specified above).

   In order for two BGP speakers to exchange labeled MAC-VPN NLRI, they
   must use BGP Capabilities Advertisement to ensure that they both are
   capable of properly processing such NLRI. This is done as specified
   in [RFC4760], by using capability code 1 (multiprotocol BGP) with an
   AFI of TBD and an SAFI of MAC-VPN.

**8. Auto-Discovery of Ethernet Tags on Ethernet Segments**

If a CE is multi-homed to two or more MESes on a particular ethernet
segment, each MES MUST advertise to other MSEs in the MAC VPN, the
information about the Ethernet Tags (e.g., VLANs) on that ethernet
segment.  If a CE is not multi-homed, then the MES that it is
attached to MAY advertise the information about Ethernet Tags (e.g.,
VLANs) on the ethernet segment connected to the CE.

The information about an Ethernet Tag on a particular ethernet
segment is advertised using a "Ethernet Tag Auto-Discovery route
(Ethernet Tag A-D route)". This route is advertised using the MAC-VPN
NLRI.

MAC VPNs support both the non-qualified and qualified learning model.
When non-qualified learning is used the Ethernet Tag Identifier
specified in this section and in other places in this document MUST
be set to a default value. When qualified learning is used the
Ethernet Tag Identifier, when required, MUST be set to a MAC VPN
provider assigned tag that maps locally on the advertising MES to an
ethernet broadcast domain identifier such as a VLAN ID.

The Ethernet Tag Auto-discovery information is used for Designated
Forwarder (DF) election as described in section 10. It is also used
to enable equal cost multi-path as described in section 15. Further,
it can be used to optimize withdrawl of MAC addresses as described in
section 18.

A Ethernet Tag A-D route type specific MAC-VPN NLRI consists of the
following:

```
            +---------------------------------------+
            |      RD   (8 octets)                  |
            +---------------------------------------+
            | Ethernet Segment Identifier (8 octets)|
            +---------------------------------------+
            |  Ethernet Tag ID (4 octets)           |
            +---------------------------------------+
            |  MPLS Label (3 octets)                |
            +---------------------------------------+
            |   Originating Router's IP Addr        |
            +---------------------------------------+
```

Route-Distinguisher (RD) MUST be set to the RD of the MAC-VPN
instance that is advertising the NLRI. A RD MUST be assigned for a
given MAC-VPN instance on a MES. This RD MUST be unique across all
MAC-VPN instances on a MES. This can be accomplished by using a Type
1 RD [RFC4364]. The value field comprises an IP address of the MES

(typically, the loopback address) followed by a number unique to the
MES.  This number may be generated by the MES, or, in the Default
Single VLAN MAC-VPN case, may be the 12 bit VLAN ID, with the
remaining 4 bits set to 0.

Ethernet Segment Identifier MUST be an 8 octet entity as described in
section 6.

The Ethernet Tag ID is the identifier of a Ethernet Tag on the
ethernet segment. This value may be a two octet VLAN ID or it may be
another Ethernet Tagused by the MAC VPN provider. It MAY be set to
the default Ethernet Tag on the ethernet segment.

The usage of the MPLS label is described in section 15.

The Originating Router's IP address MUST be set to an IP address of
the PE.  This address SHOULD be common for all the MVIs on the PE
(e.,g., this address may be PE's loopback address).

The Next Hop field of the MP_REACH_NLRI attribute of the route MUST
be set to the same IP address as the one carried in the Originating
Router's IP Address field.

The Ethernet Tag A-D route MUST carry one or more Route Target (RT)
attributes. RTs may be configured (as in IP VPNs), or may be derived
automatically from the Ethernet Tag ID associated with the
advertisement.

The following is the procedure for deriving the RT attribute
automatically from the Ethernet Tag ID associated with the
advertisement:

  +      The Global Administrator field of the RT MUST
         be set to the Autonomous System (AS) number that the MES
    belongs to.

  +      The Local Administrator field of the RT contains a 4
         octets long number that encodes the Ethernet Tag-ID.

The above auto-configuration of the RT implies that a different RT is
used for every Ethernet Tag in a MAC-VPN, if the MAC-VPN contains
multiple Ethernet Tags. For the "Default Single VLAN MAC-VPN" this
results in auto-deriving the RT from the Ethernet Tag for that MAC-
VPN.

**9**. Determining Reachability to Unicast MAC Addresses

   MESes forward packets that they receive based on the destination MAC
   address. This implies that MESes must be able to learn how to reach a
   given destination unicast MAC address.

   There are two components to MAC address learning, "local learning"
   and "remote learning":

**9.1**. Local Learning

   A particular MES must be able to learn the MAC addresses from the CEs
   that are connected to it. This is referred to as local learning.

   The MESes in a particular MAC-VPN MUST support local data plane
   learning using vanilla ethernet learning procedures. A MES must be
   capable of learning MAC addresses in the data plane when it receives
   packets such as the following from the CE network:

      - DHCP requests

      - gratuitous ARP request for its own MAC.

      - ARP request for a peer.

        Alternatively if a CE is a host then MESes MAY learn the MAC
        addresses of the host in the control plane.

        In the case where a CE is a host or a switched network connected
        on ESI X to hosts, the MAC address that is reachable via a given
        MES may move such that it becomes reachable via the same MES on
        another MES on ESI Y.  This is referred to as a "MAC Move"
        Procedures to support this are described in section 16.

**9.2**. Remote learning

   A particular MES must be able to determine how to send traffic to MAC
   addresses that belong to or are behind CEs connected to other MESes
   i.e. to remote CEs or hosts behind remote CEs. We call such MAC
   addresses as "remote" MAC addresses.

   This document requires a MES to learn remote MAC addresses in the
   control plane. In order to achieve this each MES advertises the MAC
   addresses it learns from its locally attached CEs in the control
   plane, to all the other MESes in the MAC-VPN, using BGP.

**9.2.1**. **BGP MAC-VPN MAC Address Advertisement**

BGP is extended to advertise these MAC addresses using the MAC
advertisement route type in the MAC-VPN-NLRI.

A MAC advertisement route type specific MAC-VPN NLRI consists of the
following:

```
            +----------------------------------------+
            |      RD    (8 octets)                   |
            +----------------------------------------+
            | Ethernet Segment Identifier (8 octets)|
            +----------------------------------------+
            |  Ethernet Tag ID (4 octets)            |
            +----------------------------------------+
            |  MAC Address (6 octets)                |
            +----------------------------------------+
            |  MPLS Label (3 octets)                 |
            +----------------------------------------+
            |  Originating Router's IP Addr          |
            +----------------------------------------+
```

The RD MUST be the RD of the MAC-VPN instance that is advertising the
NLRI. The procedures for setting the RD for a given MAC VPN are
described in section 8.

The Ethernet Segment Identifier is set to the eight octet ESI
identifier described in section 6.

If qualified learning is used and the MAC address that is learned
from the CE is associated with an Ethernet Tag, the Ethernet Tag ID
MUST be the Ethernet Tag Identifier, assigned by the MAC VPN provider
and mapped to the CE's ethernet tag. If non-qualified learning is
used the Ethernet Tag identifier SHOULD be set to the default
Ethernet Tag on the ethernet segment.

The encoding of a MAC address is the 6-octet MAC address specified by
IEEE 802 documents [802.1D-ORIG] [802.1D-REV].

The MPLS label MUST be the downstream assigned MAC-VPN MPLS label
that is used by the MES to forward MPLS encapsulated ethernet packets
received from remote MESes, where the destination MAC address in the
ethernet packet is the MAC address advertised in the above NLRI. The
forwarding procedures are specified in section 13. A MES may
advertise the same MAC-VPN label for all MAC addresses in a given
MAC-VPN instance. This label assignment methodology is referred to as
a per MVI label assigment. Or a MES may advertise a unique MAC-VPN
label per <ESI, Ethernet Tag> combination.  This label methodology is

referred to as a per <ESI, Ethernet Tag> label assignment. Or a MES
may advertise a unique MAC-VPN label per MAC address.  All of these
methodologies have their tradeoffs.

Per MVI label assignment requires the least number of MAC-VPN labels,
but requires a MAC lookup in addition to a MPLS lookup on an egress
MES for forwarding. On the other hand a unique label per <ESI,
Ethernet Tag> or a unique label per MAC allows an egress MES to
forward a packet that it receives from another MES, to the connected
CE, after looking up only the MPLS labels and not having to do a MAC
lookup.

The Originating Router's IP address MUST be set to an IP address of
the PE.  This address SHOULD be common for all the MVIs on the PE
(e.,g., this address may be PE's loopback address).

The Next Hop field of the MP_REACH_NLRI attribute of the route MUST
be set to the same IP address as the one carried in the Originating
Router's IP Address field.

The BGP advertisement that advertises the MAC advertisement route
MUST also carry one or more Route Target (RT) attributes. The
assignemnt of RTs described in section 8 MUST be followed.

It is to be noted that this document does not require MESes to create
forwarding state for remote MACs when they are learned in the control
plane. When this forwarding state is actually created is a local
implementation matter.


**10. Designated Forwarder Election**

Consider a CE that is a host or a router that is multi-homed directly
to more than one MES in a MAC-VPN on a given ethernet segment. One or
more Ethernet Tags may be configured on the ethernet segment. In this
scenario only one of the MESes, referred to as the Designated
Forwarder (DF), is responsible for certain actions:

    -       Sending multicast and broadcast traffic, on a given Ethernet
       Tag
            on a particular ethernet segment, to the CE. Note that
            this behavior, which allows selecting a DF at the
            granularity of <ESI, Ethernet Tag> for multicast and
       broadcast
            traffic is the default behavior in this specification.
            Optional mechanisms, which will be specified in the
            future, will allow selecting a DF at the granularity of
            <ESI, Ethernet Tag, S, G>.

-          Flooding unknown unicast traffic (i.e. traffic for
           which a MES does not know the destination MAC address),
           on a given Ethernet Tag on a particular ethernet segment to
      the CE,
           if the environment requires flooding of unknown unicast
           traffic.


   Note that a CE always sends packets using a single link. For instance
   if the CE is a host then, as mentioned earlier, the host treats the
   multiple links that it uses to reach the MESes as a Link Aggregation
   Group (LAG).

   If a bridge network is multi-homed to more than one MES in a MAC-VPN
   via switches, then the support of active-active points of attachments
   as described in this specification requires the bridge network to be
   connected to two or more MESes using a LAG. In this case the reasons
   for doing DF election are the same as those described above when a CE
   is a host or a router.

   If a bridge network does not connect to the MESes using LAG, then
   only one of the links between a CE that is a switch and the MESes
   must be the active link. Procedures for supporting active-active
   points of attachments, when a bridge network does not connect to the
   MESes using LAG, are for further study.

   The granularity of the DF election MUST be at least the ethernet
   segment via which the CE is multi-homed to the MESes. If the DF
   election is done at the ethernet segment granularity then a single
   MES MUST be elected as the DF on the ethernet segment.

   If there are one or more Ethernet Tags (e.g., VLANs) on the ethernet
   segment then the granularity of the DF election SHOULD be the
   combination of the ethernet segment and Ethernet Tag on that ethernet
   segment. In this case the same MES MUST be elected as the DF for a
   particular Ethernet Tag on that ethernet segment.

   The MESes perform a designated forwarder (DF) election, for an
   ethernet segment, or ethernet segment, Ethernet Tag combination using
   the Ethernet Tag A-D BGP route described in section 8.

   The DF election for a particular ESI or a particular <ESI, Ethernet
   Tag> combination proceeds as follows. First a MES constructs a
   candidate list of MESes. This comprises all the Ethernet Tag A-D
   routes with that particular ESI or <ESI, Ethernet Tag> tuple that a
   MES imports in a MAC-VPN instance, including the Ethernet Tag A-D
   route generated by the MES itself, if any.  The DF MES is chosen from
   this candidate list. Note that DF election is carried out by all the

MESes that import the DF route.

The default procedure for choosing the DF is the MES with the highest
IP address, of all the MESes in the candidate list. This procedure
MUST be implemented. It ensures that except during routing transients
each MES chooses the same DF MES for a given ESI and Ethernet Tag
combination.

Other alternative procedures for performing DF election are possible
and will be described in the future.

## 11. Handling of Broadcast, Multicast and Unknown Unicast Traffic

Procedures are required for a given MES to send broadcast or
multicast traffic, received from a CE encapsulated in a given
Ethernet Tag in a MAC VPN, to all the other MESes that span that
Ethernet Tag in the MAC VPN. In certain scenarios, described in
section 12, a given MES may also need to flood unknown unicast
traffic to other MESes.

The MESes in a particular MAC-VPN may use ingress replication or P2MP
LSPs to send unknown unicast, broadcast or multicast traffic to other
MESes.

Each MES MUST advertise an "Inclusive Multicast Ethernet Tag Route"
to enable the above. This section provides the encoding and the
overview of the Inclusive Multicast Ethernet Tag route. Subsequent
sections describe in further detail its usage.

An Inclusive Multicast Ethernet Tag route type specific MAC-VPN NLRI
consists of the following:

```
            +---------------------------------------+
            |      RD   (8 octets)                   |
            +---------------------------------------+
            | Ethernet Segment Identifier (8 octets)|
            +---------------------------------------+
            |  Ethernet Tag ID (4 octets)           |
            +---------------------------------------+
            |    Originating Router's IP Addr       |
            +---------------------------------------+
```

The RD MUST be the RD of the MAC-VPN instance that is advertising the
NLRI. The procedures for setting the RD for a given MAC VPN are

described in section 8.

The Ethernet Segment Identifier MAY be set to the eight octet ESI
identifier described in section 6. Or it MAY be set to 0. It MUST be
set to 0 if the Ethernet Tag is set to 0.

The Ethernet Tag ID is the identifier of the Ethernet Tag. It MAY be
set to 0 in which case an egress MES MUST perform a MAC lookup to
forward the packet.

The Originating Router's IP address MUST be set to an IP address of
the PE.  This address SHOULD be common for all the MVIs on the PE
(e.,g., this address may be PE's loopback address).

The Next Hop field of the MP_REACH_NLRI attribute of the route MUST
be set to the same IP address as the one carried in the Originating
Router's IP Address field.

The BGP advertisement that advertises the Inclusive Multicast
Ethernet Tag route MUST also carry one or more Route Target (RT)
attributes. The assignemnt of RTs described in section 8 MUST be
followed.


## 11.1. P-Tunnel Identification

In order to identify the P-Tunnel used for sending broadcast, unknown
unicast or multicast traffic, the Inclusive Multicast Ethernet Tag
route MUST carry a "PMSI Tunnel Attribute" specified in [BGP MVPN].

Depending on the technology used for the P-tunnel for the MAC VPN on
the PE, the PMSI Tunnel attribute of the Inclusive Multicast Ethernet
Tag route is constructed as follows.

  + If the PE that originates the advertisement uses a P-Multicast
    tree for the P-tunnel for the MAC VPN, the PMSI Tunnel attribute
    MUST contain the identity of the tree (note that the PE could
    create the identity of the tree prior to the actual instantiation
    of the tree).

  + A PE that uses a P-Multicast tree for the P-tunnel MAY aggregate
    two or more Ethernet Tags in the same or different MAC VPNs
    present on the PE onto the same tree. In this case in addition to
    carrying the identity of the tree, the PMSI Tunnel attribute MUST
    carry an MPLS upstream assigned label which the PE has bound
    uniquely to the <ESI, Ethernet Tag> for MAC VPN associated with
    this update (as determined by its RTs).

If the PE has already advertised Inclusive Multicast Ethernet Tag
routes for two or more Ethernet Tags that it now desires to
aggregate, then the PE MUST re-advertise those routes. The re-
advertised routes MUST be the same as the original ones, except
for the PMSI Tunnel attribute and the label carried in that
attribute.

+ If the PE that originates the advertisement uses ingress
  replication for the P-tunnel for the MAC VPN, the route MUST
  include the PMSI Tunnel attribute with the Tunnel Type set to
  Ingress Replication and Tunnel Identifier set to a routable
  address of the PE. The PMSI Tunnel attribute MUST carry a
  downstream assigned MPLS label. This label is used to demultiplex
  the broadcast, multicast or unknown unicast MAC VPN traffic
  received over a unicast tunnel by the PE.

+ The Leaf Information Required flag of the PMSI Tunnel attribute
  MUST be set to zero, and MUST be ignored on receipt.

## [11.2](11.2). Ethernet Segment Identifier and Ethernet Tag

As described above the encoding rules allow setting the Ethernet
Segment Identifier and Ethernet Tag to either valid values or to 0.
If the Ethernet Tag is set to a valid value, then an egress MES can
forward the packet to the set of egress ESIs in the Ethernet Tag, in
the MAC VPN, by performing a MPLS lookup alone. Further if the ESI is
also set to non zero then the egress MES does not need to replicate
the packet as it is destined for a given ethernet segment. If both
Ethernet Tag and ESI are set to 0 then an egress MES MUST perform a
MAC lookup in the MVI determined by the MPLS label, after the MPLS
lookup, to forward the packet.

If a MES advertises multiple Inclusive Ethernet Tag routes for a
given MAC VPN then the PMSI Tunnel Attributes for these routes MUST
be distinct.

## [12](12). Processing of Unknown Unicast Packets

The procedures in this document do not require MESes to flood unknown
unicast traffic to other MESes. If MESes learn CE MAC addresses via a
control plane, the MESes can then distribute MAC addresses via BGP,
and all unicast MAC addresses will be learnt prior to traffic to
those destinations.

However, if a destination MAC address of a received packet is not
known by the MES, the MES may have to flood the packet. Flooding must

take into account "split horizon forwarding" as follows. The
principles behind the following procedures are borrowed from the
split horizon forwarding rules in VPLS solutions [RFC 4761, RFC
4762].  When a MES capable of flooding (say MESx) receives a
broadcast Ethernet frame, or one with an unknown destination MAC
address, it must flood the frame.  If the frame arrived from an
attached CE, MESx must send a copy of the frame to every other
attached CE, as well as to all other MESs participating in the MAC
VPN. If, on the other hand, the frame arrived from another MES (say
MESy), MESx must send a copy of the packet only to attached CEs. MESx
MUST NOT send the frame to other MESs, since MESy would have already
done so. Split horizon forwarding rules apply to broadcast and
multicast packets, as well as packets to an unknown MAC address.

Whether or not to flood packets to unknown destination MAC addresses
should be an administrative choice, depending on how learning happens
between CEs and MESes.

The MESes in a particular MAC VPN may use ingress replication using
RSVP-TE P2P LSPs or LDP MP2P LSPs for sending broadcast, multicast
and unknown unicast traffic to other MESes. Or they may use RSVP-TE
or LDP P2MP LSPs for sending such traffic to other MESes.

## 12.1. Ingress Replication

If ingress replication is in use, the P-Tunnel attribute, carried in
the Inclusive Multicast Ethernet Tag routes (section 11) for the MAC
VPN, specifies the downstream label that the other MESes can use to
send unknown unicast, multicast or broadcast traffic for the MAC VPN
to this particular MES.

The MES that receives a packet with this particular MPLS label MUST
treat the packet as a broadcast, multicast or unknown unicast packet.
Further if the MAC address is a unicast MAC address, the MES MUST
treat the packet as an unknown unicast packet.

## 12.2. P2MP MPLS LSPs

The procedures for using P2MP LSPs are very similar to VPLS
procedures [VPLS-MCAST]. The P-Tunnel attribute used by a MES for
sending unknown unicast, broadcast or multicast traffic for a
particular ethernet segment, is advertised in the Inclusive Ethernet
Tag Multicast route as described in section 11.

The P-Tunnel attribute specifies the P2MP LSP identifier. This is the
equivalent of an Inclusive tree in [VPLS-MCAST]. Note that multiple

   Ethernet Tags, which may be in different MAC-VPNs, may use the same
   P2MP LSP, using upstream labels [VPLS-MCAST]. When P2MP LSPs are used
   for flooding unknown unicast traffic, packet re-ordering is possible.

   The MES that receives a packet on the P2MP LSP specified in the PMSI
   Tunnel Attribute MUST treat the packet as a broadcast, multicast or
   unknown unicast packet. Further if the MAC address is a unicast MAC
   address, the MES MUST treat the packet as an unknown unicast packet.


**13. Forwarding Unicast Packets**

**13.1. Forwarding packets received from a CE**

   When a MES receives a packet from a CE, on a given Ethernet Tag, it
   must first look up the source MAC address of the packet. In certain
   environments the source MAC address may be used to authenticate the
   CE and determine that traffic from the host can be allowed into the
   network.

   If the MES decides to forward the packet the destination MAC address
   of the packet must be looked up. If the MES has received MAC address
   advertisements for this destination MAC address from one or more
   other MESes or learned it from locally connected CEs, it is
   considered as a known MAC address. Else the MAC address is considered
   as an unknown MAC address.

   For known MAC addresses the MES forwards this packet to one of the
   remote MESes. The packet is encapsulated in the MAC-VPN MPLS label
   advertised by the remote MES, for that MAC address, and in the MPLS
   LSP label stack to reach the remote MES.

   If the MAC address is unknown then, if the administrative policy on
   the MES requires flooding of unknown unicast traffic:
      - The MES MUST flood the packet to other MESes. If the ESI over
   which the MES receives the packet is multi-homed, then the MES MUST
   first encapsulate the packet in the ESI MPLS label as described in
   section 14.  If ingress replication is used the packet MUST be
   replicated one or more times to each remote MES with the bottom label
   of the stack being a MPLS label determined as follows. This is the
   MPLS label advertised by the remote MES in a PMSI Tunnel Attribute in
   the Inclusive Multicast Ethernet Tag route for an <ESI, Ethernet Tag>
   combination. The Ethernet Tag in the route must be the same as the
   Ethernet Tag advertised by the ingress MES in its Ethernet Tag A-D
   route associated with the interface on which the ingress MES receives
   the packet. If P2MP LSPs are being used the packet MUST be sent on
   the P2MP LSP that the MES is the root of for the Ethernet Tag in the
   MAC-VPN. If the same P2MP LSP is used for all Ethernet Tags then all

the MESes in the MAC VPN MUST be the leaves of the P2MP LSP. If a
distinct P2MP LSP is used for a given Ethernet Tag in the MAC VPN
then only the MESes in the Ethernet Tag MUST be the leaves of the
P2MP LSP. The packet MUST be encapsulated in the P2MP LSP label
stack.

If the MAC address is unknown then, if the admnistrative policy on
the MES does not allow flooding of unknown unicast traffic:
    - The MES MUST drop the packet.


**13.2. Forwarding packets received from a remote MES**

**13.2.1. Unknown Unicast Forwarding**

When a MES receives a MPLS packet from a remote MES then, after
processing the MPLS label stack, if the top MPLS label ends up being
a P2MP LSP label associated with a MAC-VPN or the downstream label
advertised in the P-Tunnel attribute and after performing the split
horizon procedures described in section 14:

    - If the MES is the designated forwarder of unknown unicast,
broadcast or multicast traffic, on a particular set of ESIs for the
Ethernet Tag, the default behavior is for the MES to flood the packet
on the ESIs. In other words the default behavior is for the MES to
assume that the destination MAC address is unknown unicast, broadcast
or multicast and it is not required to do a destination MAC address
lookup, as long as the granularity of the MPLS label included the
Ethernet Tag. As an option the MES may do a destination MAC lookup to
flood the packet to only a subset of the CE interfaces in the
Ethernet Tag. For instance the MES may decide to not flood an unknown
unicast packet on certain ethernet segments even if it is the DF on
the ethernet segment, based on administrative policy.

    - If the MES is not the designated forwarder on any of the ESIs
for the Ethernet Tag, the default behavior is for it to drop the
packet.


**13.2.2. Known Unicast Forwarding**

If the top MPLS label ends up being a MAC-VPN label that was
advertised in the unicast MAC advertisements, then the MES either
forwards the packet based on CE next-hop forwarding information
associated with the label or does a destination MAC address lookup to
forward the packet to a CE.

14. **Split Horizon**

   Consider a CE that is multi-homed to two or more MESes on an ethernet
   segment ES1. If the CE sends a multicast, broadcast or unknown
   unicast packet to a particular MES, say MES1, then MES1 will forward
   that packet to all or subset of the other MESes in the MAC VPN. In
   this case the MESes, other than MES1, that the CE is multi-homed to
   MUST drop the packet and not forward back to the CE. This is referred
   to as "split horizon" in this document.

   In order to accomplish this each MES distributes to other MESes that
   are connected to the ethernet segment an "Ethernet Segment Route".

   An Ethernet Segment route type specific MAC-VPN NLRI consists of the
   following:

```
                +----------------------------------------+
                |       RD   (8 octets)                  |
                +----------------------------------------+
                | Ethernet Segment Identifier (8 octets)|
                +----------------------------------------+
                |   MPLS Label (3 octets)                |
                +----------------------------------------+
                |    Originating Router's IP Addr        |
                +----------------------------------------+
```

   The RD MUST be the RD of the MAC-VPN instance that is advertising the
   NLRI. The procedures for setting the RD for a given MAC VPN are
   described in section 8.

   The Ethernet Segment Identifier MUST be set to the eight octet ESI
   identifier described in section 6.

   The MPLS label is referred to as an "ESI label". This label MUST be a
   downstream assigned MPLS label if the advertising MES is using
   ingress replication for sending multicast, broadcast or unknown
   unicast traffic, to other MESes. If the advertising MES is using P2MP
   MPLS LSPs for the same, then this label MUST be an upstream assigned
   MPLS label. The usage of this label is described below.

   The Originating Router's IP address MUST be set to an IP address of
   the PE.  This address SHOULD be common for all the MVIs on the PE
   (e.,g., this address may be PE's loopback address).

   The Next Hop field of the MP_REACH_NLRI attribute of the route MUST
   be set to the same IP address as the one carried in the Originating
   Router's IP Address field.

   The BGP advertisement that advertises the MAC advertisement route
   MUST also carry one Route Target (RT) attribute. The construction of
   this RT will be specified in the next version.

   This route will be enhanced to carry LAG specific information such as
   LACP parameters in the future.


**14.1. ESI MPLS Label: Ingress Replication**

   An MES that is using ingress replication for sending broadcast,
   multicast or unknown unicast traffic, distributes to other MESes,
   that belong to the ethernet segment, a downstream assigned "ESI MPLS
   label" in the Ethernet Segment route. This label MUST be programmed
   in the platform label space by the advertising MES. Further the
   forwarding entry for this label must result in NOT forwarding packets
   received with this label onto the ethernet segment that the label was
   distributed for.

   Consider MES1 and MES2 that are multi-homed to CE1 on ES1. Further
   consider that MES1 is using P2P or MP2P LSPs to send packets to MES2.
   Consider that MES1 receives a a multicast, broadcast or unknown
   unicast packet from CE1 on VLAN1 on ESI1.

   First consider the case where MES2 distributes an unique Inclusive
   Multicast Ethernet Tag route for VLAN1, for each ethernet segment on
   MES2. In this case MES1 MUST NOT replicate the packet to MES2 for
   <ESI1, VLAN1>.

   Next consider the case where MES2 distributes a single Inclusive
   Multicast Ethernet Tag route for VLAN1 for all ethernet segments on
   MES2. In this case when MES1 sends a multicast, broadcast or unknown
   unicast packet, that it receives from CE1, it MUST first push onto
   the MPLS label stack the ESI label that MES2 has distributed for
   ESI1. It MUST then push on the MPLS label distributed by MES2 in the
   Inclusive Ethernet Tag Multicast route for Ethernet Tag1. The
   resulting packet is further encapsulated in the P2P or MP2P LSP label
   stack required to transmit the packet to MES2.  When MES2 receives
   this packet it determines the set of ESIs to replicate the packet to
   from the top MPLS label, after any P2P or MP2P LSP labels have been
   removed. If the next label is the ESI label assigned by MES2 then
   MES2 MUST NOT forward the packet onto ESI1.

**14.2. ESI MPLS Label: P2MP MPLS LSPs**

   An MES that is using P2MP LSPs for sending broadcast, multicast or
   unknown unicast traffic, distributes to other MESes, that belong to
   the ethernet segment, an upstream assigned "ESI MPLS label" in the
   Ethernet Segment route. This label is upstream assigned by the MES
   that advertises the route. This label MUST be programmed by the other
   MESes, that are connected to the ESI advertised in the route, in the
   context label space for the advertising MES. Further the forwarding
   entry for this label must result in NOT forwarding packets received
   with this label onto the ethernet segment that the label was
   distributed for.

   Consider MES1 and MES2 that are multi-homed to CE1 on ES1. Further
   assume that MES1 is using P2MP MPLS LSPs to send broadcast, multicast
   or uknown unicast packets. When MES1 sends a multicast, broadcast or
   unknown unicast packet, that it receives from CE1, it MUST first push
   onto the MPLS label stack the ESI label that it has assigned for the
   ESI that the packet was received on. The resulting packet is further
   encapsulated in the P2MP MPLS label stack necessary to transmit the
   packet to the other MESes. Penultimate hop popping MUST be disabled
   on the P2MP LSPs used in the MPLS transport infrastructure for MAC
   VPN. When MES2 receives this packet it decapsulates the top MPLS
   label and forwards the packet using the context label space
   determined by the top label. If the next label is the ESI label
   assigned by MES1 then MES2 MUST NOT forward the packet onto ESI1.


**15. Load Balancing of Unicast Packets**

   This section specifies how load balancing is achieved to/from a CE
   that has more than one interface that is directly connected to one or
   more MESes. The CE may be a host or a router or it may be a switched
   network that is connected via LAG to the MESes.


**15.1. Load balancing of traffic from a MES to remote CEs**

   Whenever a remote MES imports a MAC advertisement for a given <ESI,
   Ethernet Tag> in a MAC VPN instance, it MUST consider the MAC as
   reachahable via all the MESes from which it has imported Ethernet Tag
   A-D routes for that <ESI, Ethernet Tag>. Further the remote MES MUST
   use these MAC advertisement and Ethernet Tag A-D routes to constuct
   the set of next-hops that it can use to send the packet to the
   destination MAC. Each next-hop comprises a MPLS label, that is to be
   used by the egress MES to forward the packet. This label is
   determined as follows. If the next-hop is constructed as a result of
   a MAC route which has a valid MPLS label, then this label MUST be

used. However if the MAC route doesn't have a valid MPLS label or if
the next-hop is constructed as a result of a Ethernet Tag A-D route
then the MPLS label from the Ethernet Tag A-D route MUST be used.

Consider a CE, CE1, that is dual homed to two MESes, MES1 and MES2 on
a LAG interface, ES1, and is sending packets with MAC address MAC1 on
VLAN1. Based on MAC-VPN extensions described in sections [8](8) and [9](9), a
remote MES say MES3 is able to learn that a MAC1 is reachable via
MES1 and MES2. Both MES1 and MES2 may advertise MAC1 in BGP if they
receive packets with MAC1 from CE1. If this is not the case and if
MAC1 is advertised only by MES1, MES3 still considers MAC1 as
reachable via both MES1 and MES2 as both MES1 and MES2 advertise a
Ethernet Tag A-D route for <ESI1, VLAN1>.

The MPLS label stack to send the packets to MES1 is the MPLS LSP
stack to get to MES1 and the MAC-VPN label advertised by MES1 for
CE1's MAC.

The MPLS label stack to send packets to MES2 is the MPLS LSP stack to
get to MES2 and the upstream assigned label in the Ethernet Tag A-D
route advertised by MES2 for <ES1, VLAN1>, if MES2 has not advertised
MAC1 in BGP.

We will refer to these label stacks as MPLS next-hops.

The remote MES, MES3, can now load balance the traffic it receives
from its CEs, destined for CE1, between MES1 and MES2.  MES3 may use
the IP flow information for it to hash into one of the MPLS next-hops
for load balancing for IP traffic. Or MES3 may rely on the source and
destination MAC addresses for load balancing.

Note that once MES3 decides to send a particular packet to MES1 or
MES2 it can pick from more than path to reach the particular remote
MES using regular MPLS procedures. For instance if the tunneling
technology is based on RSVP-TE LSPs, and MES3 decides to send a
particular packet to MES1 then MES3 can choose from multiple RSVP-TE
LSPs that have MES1 as their destination.

When MES1 or MES2 receive the packet destined for CE1 from MES3, if
the packet is a unicast MAC packet it is forwarded to CE1.  If it is
a multicast or broadcast MAC packet then only one of MES1 or MES2
must forward the packet to the CE. Which of MES1 or MES2 forward this
packet to the CE is determined by default based on which of the two
is the DF. An alternate procedure to load balance multicast packets
will be described in the future.

If the connectivity between the multi-homed CE and one of the MESes
that it is multi-homed to fails, the MES MUST withdraw the MAC

address from BGP.  This enables the remote MESes to remove the MPLS
next-hop to this particular MES from the set of MPLS next-hops that
can be used to forward traffic to the CE. For further details and
procedures on withdrawl of MAC VPN route types in the event of MES to
CE failures please section 18.4.

## 15.2. Load balancing of traffic between a MES and a local CE

A CE may be configured with more than one interface connected to
different MESes or the same MES for load balancing. The MES(s) and
the CE can load balance traffic onto these interfaces using one of
the following mechanisms.

### 15.2.1. Data plane learning

Consider that the MESes perform data plane learning for local MAC
addresses learned from local CEs. This enables the MES(s) to learn a
particular MAC address and associate it with one or more interfaces.
The MESes can now load balance traffic destined to that MAC address
on the multiple interfaces.

Whether the CE can load balance traffic that it generates on the
multiple interfaces is dependent on the CE implementation.

### 15.2.2. Control plane learning

The CE can be a host that advertises the same MAC address using a
control protocol on both interfaces. This enables the MES(s) to learn
the host's MAC address and associate it with one or more interfaces.
The MESes can now load balance traffic destined to the host on the
multiple interfaces. The host can also load balance the traffic it
generates onto these interfaces and the MES that receives the traffic
employs MAC-VPN forwarding procedures to forward the traffic.

## 16. MAC Moves

In the case where a CE is a host or a switched network connected to
hosts, the MAC address that is reachable via a given MES on a
particular ESI may move such that it becomes reachable via another
MES on another ESI.  This is referred to as a "MAC Move".

Remote MESes must be able to distinguish a MAC move from the case
where a MAC address on an ESI is reachable via two different MESes
and load balancing is performed as described in section 15. This

distinction can be made as follows. If a MAC is learned by a
particular MES from multiple MESes, then the MES performs load
balancing only amongst the set of MESes that advertised the MAC with
the same ESI. If this is not the case then the MES chooses only one
of the advertising MESes to reach the MAC as per BGP path selection.

There can be traffic loss during a MAC move. Consider MAC1 that is
advertised by MES1 and learned from CE1 on ESI1. If MAC1 now moves
behind MES2, on ESI2, MES2 advertises the MAC in BGP. Until a remote
MES, MES3, determines that the best path is via MES2, it will
continue to send traffic destined for MAC1 to MES1. This will not
occur deterministially until MES1 withdraws the advertisement for
MAC1.

One recommended optimization to reduce the traffic loss during MAC
moves is the following option. When an MES sees a MAC update from a
CE on an ESI, which is different from the ESI on which the MES has
currently learned the MAC, the corresponding entry in the local
bridge forwarding table SHOULD be immediately purged causing the MES
to withdraw its own MAC-VPN MAC advertisement route and replace it
with the update.

A future version of this specification will describe other optimized
procedures to minimize traffic loss during MAC moves.


## 17. Multicast

The MESes in a particular MAC-VPN may use ingress replication or P2MP
LSPs to send multicast traffic to other MESes.


## 17.1. Ingress Replication

The MESes may use ingress replication for flooding unknown unicast,
multicast or broadcast traffic as described in section 11. A given
unknown unicast or broadcast packet must be sent to all the remote
MESes. However a given multicast packet for a multicast flow may be
sent to only a subset of the MESes. Specifically a given multicast
flow may be sent to only those MESes that have receivers that are
interested in the multicast flow. Determining which of the MESes have
receivers for a given multicast flow is done using explicit tracking
described below.

**17.2. P2MP LSPs**

   A MES may use an "Inclusive" tree for sending an unknown unicast,
   broadcast or multicast packet or a "Selective" tree. This terminology
   is borrowed from [VPLS-MCAST].

   A variety of transport technologies may be used in the SP network.
   For inclusive P-Multicast trees, these transport technologies include
   point-to-multipoint LSPs created by RSVP-TE or mLDP. For selective P-
   Multicast trees, only unicast MES-MES tunnels (using MPLS or IP/GRE
   encapsulation) and P2MP LSPs are supported, and the supported P2MP
   LSP signaling protocols are RSVP-TE, and mLDP.


**17.2.1. Inclusive Trees**

    An Inclusive Tree allows the use of a single multicast distribution
   tree, referred to as an Inclusive P-Multicast tree, in the SP network
   to carry all the multicast traffic from a specified set of MAC VPN
   instances on a given MES. A particular P-Multicast tree can be set up
   to carry the traffic originated by sites belonging to a single MAC
   VPN, or to carry the traffic originated by sites belonging to
   different MAC VPNs. The ability to carry the traffic of more than one
   MAC VPN on the same tree is termed 'Aggregation'. The tree needs to
   include every MES that is a member of any of the MAC VPNs that are
   using the tree. This implies that a MES may receive multicast traffic
   for a multicast stream even if it doesn't have any receivers that are
   interested in receiving traffic for that stream.

   An Inclusive P-Multicast tree as defined in this document is a P2MP
   tree.  A P2MP tree is used to carry traffic only for MAC VPN CEs that
   are connected to the MES that is the root of the tree.

   The procedures for signaling an Inclusive Tree are the same as those
   in [VPLS-MCAST] with the VPLS-AD route replaced with the Inclusive
   Multicast Ethernet Tag route. The P-Tunnel attribute [VPLS-MCAST] for
   an Inclusive tree is advertised in the Inclusive Ethernet Tag A-D
   route as described in section 11.  Note that a MES can "aggregate"
   multiple inclusive trees for different MAC-VPNs on the same P2MP LSP
   using upstream labels. The procedures for aggregation are the same as
   those described in [VPLS-MCAST], with VPLS A-D routes replaced by
   MAC-VPN Inclusive Multicast Ethernet Tag A-D routes.

**17.2.2. Selective Trees**

   A Selective P-Multicast tree is used by a MES to send IP multicast
   traffic for one or IP more specific multicast streams, originated by
   CEs connected to the MES, that belong to the same or different MAC
   VPNs, to a subset of the MESs that belong to those MAC VPNs. Each of
   the MESs in the subset should be on the path to a receiver of one or
   more multicast streams that are mapped onto the tree. The ability to
   use the same tree for multicast streams that belong to different MAC
   VPNs is termed a MES the ability to create separate SP multicast
   trees for specific multicast streams, e.g. high bandwidth multicast
   streams. This allows traffic for these multicast streams to reach
   only those MES routers that have receivers in these streams. This
   avoids flooding other MES routers in the MAC VPN.

   A SP can use both Inclusive P-Multicast trees and Selective P-
   Multicast trees or either of them for a given MAC VPN on a MES, based
   on local configuration.

   The granularity of a selective tree is <RD, MES, S, G> where S is an
   IP multicast source address and G is an IP multicast group address or
   G is a multicast MAC address. Wildcard sources and wildcard groups
   are supported. Selective trees require explicit tracking as described
   below.

   A MAC-VPN MES advertises a selective tree using a MAC-VPN selective
   A-D route. The procedures are the same as those in [VPLS-MCAST] with
   S-PMSI A-D routes in [VPLS-MCAST] replaced by MAC-VPN Selective A-D
   routes. The information elements of the MAC VPN selective
    A-D route are similar to those of the VPLS S-PMSI A-D route with the
   following differences. A MAC VPN Selective A-D route includes an
   optional Ethernet Tag field. Also a MAC VPN selective A-D route may
   encode a MAC address in the Group field. The encoding details of the
   MAC VPN selective A-D route will be described in the next revision.

   Selective trees can also be aggregated on the same P2MP LSP using
   aggregation as described in [VPLS-MCAST].

**17.3. Explicit Tracking**

   [VPLS-MCAST] describes procedures for explicit tracking that rely on
   Leaf A-D routes. The same procedures are used for explicit tracking
   in this specification with VPLS Leaf A-D routes replaced with MAC-VPN
   Leaf A-D routes.  These procedures allow a root MES to request
   multicast membership information for a given (S, G), from leaf MESs.
   Leaf MESs rely on IGMP snooping or PIM snooping between the MES and
   the CE to determine the multicast membership information. Note that

   the procedures in [VPLS-MCAST] do not describe how explicit tracking
   is performed if the CEs are enabled with join suppression. The
   procedures for this case will be described in a future version.


## 18. Convergence

   This section describes failure recovery from different types of
   network failures.


### 18.1. Transit Link and Node Failures between MESes

   The use of existing MPLS Fast-Reroute mechanisms can provide failure
   recovery in the order of 50ms, in the event of transit link and node
   failures in the infrastructure that connects the MESes.


### 18.2. MES Failures

   Consider a host host1 that is dual homed to MES1 and MES2. If MES1
   fails, a remote MES, MES3, can discover this based on the failure of
   the BGP session.  This failure detection can be in the sub-second
   range if BFD is used to detect BGP session failure. MES3 can update
   its forwarding state to start sending all traffic for host1 to only
   MES2. It is to be noted that this failure recovery is potentially
   faster than what would be possible if data plane learning were to be
   used. As in that case MES3 would have to rely on re-learning of MAC
   addresses via MES2.


#### 18.2.1. Local Repair

   It is possible to perform local repair in the case of MES failures.
   Details will be specified in the future.


### 18.3. MES to CE Network Failures

   When an ethernet segment connected to a MES fails or when a Ethernet
   Tag is deconfigured on an ethernet segment, then the MES MUST
   withdraw the Ethernet Tag A-D route(s) announced for the <ESI,
   Ethernet Tags> that are impacted by the failure or de-configuration.
   In addition the MES MUST also withdraw the MAC advertisement routes
   that are impacted by the failure or de-configuration.

   The Ethernet Tag A-D routes should be used by an implementation to
   optimize the withdrawal of MAC advertisement routes. When a MES

receives a withdrawl of a particular Ethernet Tag A-D route it SHOULD
consider all the MAC advertisement routes, that are learned from the
same <ESI, Ethernet Tag> as in the Ethernet Tag A-D route, as having
been withdrawn. This optimizes the network convergence times in the
event of MES to CE failures.

## 19. Acknowledgements

We would like to thank Yakov Rekhter, Kaushik Ghosh, Nischal Sheth
and Amit Shukla for discussions that helped shape this document.  We
would also like to thank Han Nguyen for his comments and support of
this work.

## 20. References

[RFC4364] "BGP/MPLS IP VPNs", Rosen, Rekhter, et. al., February 2006

[VPLS-MCAST] "Multicast in VPLS". R. Aggarwal et.al., draft-ietf-
l2vpn-vpls-mcast-04.txt

[RFC4761] Kompella, K. and Y. Rekhter, "Virtual Private LAN Service
(VPLS) Using BGP for Auto-Discovery and Signaling", RFC 4761, January
2007.

[RFC4762] Lasserre, M. and V. Kompella, "Virtual Private LAN Service
(VPLS) Using Label Distribution Protocol (LDP) Signaling", RFC 4762,
January 2007.

[VPLS-MULTIHOMING] "BGP based Multi-homing in Virtual Private LAN
Service", K. Kompella et. al., draft-ietf-l2vpn-vpls-
multihoming-00.txt

[PIM-SNOOPING] "PIM Snooping over VPLS", V. Hemige et. al., draft-
ietf-l2vpn-vpls-pim-snooping-01

[IGMP-SNOOPING] "Considerations for Internet Group Management
Protocol (IGMP) and Multicast Listener Discovery (MLD) Snooping
Switches", M. Christensen et. al., RFC4541,

21. Author's Address

Rahul Aggarwal
Juniper Networks
1194 N. Mathilda Ave.
Sunnyvale, CA  94089 US


Email: rahul@juniper.net

Aldrin Isaac
Bloomberg
Email: aisaac71@bloomberg.net

James Uttaro
AT&T
200 S. Laurel Avenue
Middletown, NJ  07748
USA
Email: uttaro@att.com

Ravi Shekhar
Juniper Networks
1194 N. Mathilda Ave.
Sunnyvale, CA  94089 US

Wim Henderickx
Alcatel-Lucent
e-mail: wim.henderickx@alcatel-lucent.be

Florin Balus
Alcatel-Lucent
e-mail: Florin.Balus@alcatel-lucent.be