Internet Draft Document                            Marc Lasserre
draft-lasserre-vkompella-ppvpn-vpls-01.txt    Riverstone Networks
                                                   Vach Kompella
                                                     Nick Tingle
                                                 Sunil Khandekar
                                                 Timetra Networks


Pascal Menezes                                     Loa Andersson
Terabeam Networks                                          Utfors

Andrew Smith                                          Pierre Lin
Consultant                                    Yipes Communication

Juha Heinanen                                        Giles Heron
Song Networks                               PacketExchange Ltd.

Ron Haberman                                       Tom S.C. Soon
Masergy, Inc.                               SBC Communications

Nick Slabakov                                       Luca Martini
Rob Nath                                                 Level 3
Riverstone Networks                               Communications

Expires: September 2002                              March 2002

### Virtual Private LAN Services over MPLS
draft-lasserre-vkompella-ppvpn-vpls-01.txt


1.  Status of this Memo

This document is an Internet-Draft and is in full conformance
with all provisions of Section 10 of RFC2026.

Internet-Drafts are working documents of the Internet Engineering
Task Force (IETF), its areas, and its working groups.  Note that
other groups may also distribute working documents as Internet-
Drafts.

Internet-Drafts are draft documents valid for a maximum of six
months and may be updated, replaced, or obsoleted by other documents
at any time.  It is inappropriate to use Internet-Drafts as
reference material or to cite them other than as "work in progress."

The list of current Internet-Drafts can be accessed at
     http://www.ietf.org/ietf/1id-abstracts.txt
The list of Internet-Draft Shadow Directories can be accessed at
     http://www.ietf.org/shadow.html.

2.  Abstract

This document describes a  virtual private LAN service (VPLS)
solution over MPLS, also known as Transparent LAN Services (TLS).
VPLS simulates an Ethernet virtual 802.1D bridge [802.1D-ORIG]
[802.1D-REV] for a given set of users.  It delivers a layer 2
broadcast domain that is fully capable of learning and forwarding on
Ethernet MAC addresses that is closed to a given set of users.  Many
VLS services can be supported from a single PE node.

3.  Conventions

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

Placement of this Memo in Sub-IP Area

RELATED DOCUMENTS

http:// search.ietf.org/internet-drafts/draft-martini-l2circuit-
trans-mpls-06.txt

http://search.ietf.org/internet-drafts/draft-martini-l2circuit-
encap-mpls-02.txt

http://search.ietf.org/internet-drafts/draft-augustyn-ppvpn-vpls-
reqmts-00.txt

WHERE DOES THIS FIT IN THE PICTURE OF THE SUB-IP WORK

PPVPN

WHY IS IT TARGETTED AT THIS WG

The charter of the PPVPN WG includes L2 VPN services and this draft
specifies a model for Ethernet L2 VPN services over MPLS.

JUSTIFICATION

Existing Internet drafts specify how to provide point-to-point
Ethernet L2 VPN services over MPLS. This draft defines how
multipoint Ethernet services can be provided.

Table of Contents

4.  Overview

Ethernet has become a predominant technology initially for Local
Area Networks (LANs) and now as an access technology, specifically
in metropolitan networks. Ethernet ports or IEEE VLANs are dedicated
to customers on Provider Edge (PE) routers acting as LERs. Customer
traffic gets mapped to a specific MPLS L2 VPN by configuring L2 FECs
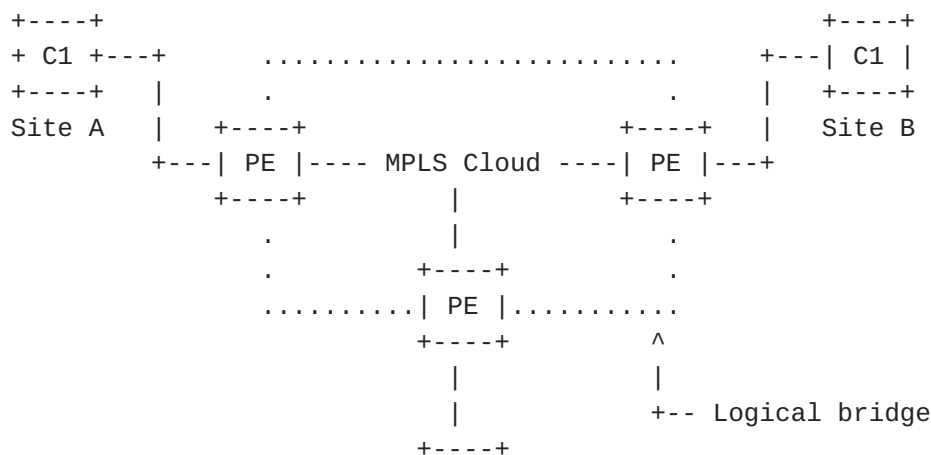based upon the input port and/or VLAN.

Broadcast and multicast services are available over traditional
LANs. MPLS does not support such services currently. Sites that
belong to the same broadcast domain and that are connected via an
MPLS network expect broadcast, multicast and unicast traffic to be
forwarded to the proper location(s). This requires MAC address
learning/aging on a per LSP basis, packet replication across LSPs
for multicast/broadcast traffic and for flooding of unknown unicast
destination traffic.

[MARTINI-ENCAP] defines how to carry L2 PDUs over point-to-point
MPLS LSPs, called VC LSPs. Such VC LSPs can be carried across MPLS
or GRE tunnels. This document describes extensions to [MARTINI-
ENCAP] for transporting Ethernet/802.3 and VLAN [802.1Q] traffic
across multiple sites that belong to the same L2 broadcast domain.
Note that the same model can be applied to other 802.1 technologies.
It describes a simple and scalable way to offer Virtual LAN
services, including the appropriate flooding of Broadcast, Multicast
and unknown unicast destination traffic over MPLS, without the need
for address resolution servers or other external servers, as
discussed in [VPLS-REQ].

The following discussion applies to devices that serve as Label Edge
Routers (LERs) on an MPLS network that is VPLS capable. It will not
discuss the behavior of transit Label Switch Routers (LSRs) that are
considered a part of MPLS network. The MPLS network provides a
number of Label Switch Paths (LSPs) that form the basis for
connections between LERs attached to the same MPLS network. The
resulting set of interconnected LERs forms a private MPLS VPN where
each LSP is uniquely identified at each MPLS interface by a label.

5.  Bridging Model for MPLS

An MPLS interface acting as a bridge must be able to flood, forward,
and filter bridged frames.

```
+----+                                          +----+
+ C1 +---+      ...........................      +---| C1 |
+----+   |       .                        .      |   +----+
Site A   |   +----+                    +----+     |    Site B
         +---| PE |---- MPLS Cloud ----| PE |---+
             +----+         |          +----+
               .            |            .
               .         +----+          .
         ..........| PE |...........
                      +----+         ^
                        |            |
                        |            +-- Logical bridge
                      +----+
```

```
                         | C1 |
                         +----+
                         Site C
```

The set of PE devices interconnected via transport tunnels appears
as a single 802.1D bridge/switch to customer C1. Each PE device will
learn remote MAC addresses on LSPs (and keeps learning directly
attached MAC addresses on customer facing ports).  We note here that
while this document shows specific examples using MPLS transport
tunnels, other tunnels that can be used by pseudo-wires, e.g., GRE,
L2TP, IPSEC, etc., can also be used, as long as the sender PE can be
identified, since this is used in the learning algorithm.

The scope of the VPLS lies within the PEs in the service provider
network, highlighting the fact that apart from customer service
delineation, the form of access to a customer site is not relevant
to the VPLS [VPLS-REQ].

The PE device is typically an edge router capable of running a
signaling protocol and/or routing protocols to exchange VC label
information.  In addition, it is capable of setting up transport
tunnels to other PEs to deliver VC LSP traffic.


5.1.  Flooding and Forwarding

Flooding within the service provider network is performed by sending
unknown unicast and multicast frames to all relevant PE nodes
participating in the VPLS. In the MPLS environment this means
sending the PDU through each relevant VC LSP.

Note that multicast frames do not necessarily have to be sent to all
VPN members. For simplicity, the default approach of broadcasting
multicast frames can be used. Extensions explaining how to interact
with 802.1 GMRP protocol, IGMP snooping and static MAC multicast
filters will be discussed in a future revision.

To forward a frame, a bridge must be able to associate a destination
MAC address with a VC LSP. It is unreasonable and perhaps impossible
to require bridges to statically configure an association of every
possible destination MAC address with a VC LSP. Therefore, VPLS
bridges must provide enough information to allow an MPLS interface
to dynamically learn about foreign destinations beyond the set of
LSRs. To accomplish dynamic learning, a bridged PDU MUST conform to
the encapsulation described within [MARTINI-ENCAP].


5.2.  Address Learning

Unlike BGP VPNs [BGP-VPN], reachability information does not need to
be advertised and distributed via a control plane.  Reachability is
obtained by standard learning bridge functions in the data plane.

Since VC LSPs are uni-directional, two LSPs of opposite directions
are required to form a logical bi-directional link. When a new MAC

address is learned on an inbound LSP, it needs to be associated with
the outbound LSP that is part of the same pair. The state of this
logical link can be considered as up as soon as both incoming and
outgoing LSPs are established. Similarly, it can be considered as
down as soon as one of these two LSPs is torn down.
Standard learning, filtering and forwarding actions, as defined in
[[802.1D-ORIG](802.1D-ORIG)], [[802.1D-REV](802.1D-REV)] and [[802.1Q](802.1Q)], are required when a
logical link state changes.

5.3.  LSP Topology

PE routers typically run an IGP between them, and are assumed to
have the capability to establish MPLS tunnels.  Tunnel LSPs are set
up between PEs to aggregate traffic.  VC LSPs are signaled to
demultiplex the L2 encapsulated packets that traverse the tunnel
LSPs.

In this Ethernet L2VPN, it becomes the responsibility of the service
provider to create the loop free topology, since the PEs have to
examine the Layer 2 fields of the packets, unlike Frame Relay or
ATM, where the termination point becomes the CE node.  Therefore,
for the sake of simplicity, we assume that the topology of a VPLS is
a full mesh of tunnel and VC LSPs.

5.4.  Loop free L2 VPN

For simplicity, a full mesh of LSPs is established between PEs.

Each PE MUST create a rooted tree to every other PE router that
serve the same L2 VPN. Each PE MUST support a "split-horizon" scheme
in order to prevent loops, that is, a PE MUST NOT forward traffic
from one VC LSP to another in the same VPN (since each PE has direct
connectivity to all other PEs in the same VPN).

Note that customers are allowed to run STP such as when a customer
has a back door link used for backup. In such a case STP BPDUs are
simply tunneled through the MPLS cloud.

5.5.  LDP Based Signaling

In order to establish a full mesh of VC LSPs, all PEs in a VPLS must
have a full mesh of LDP sessions.

Once an LDP session has been formed between two PEs, all VC LSPs are
signaled over this session.

In [[MARTINI-SIG](MARTINI-SIG)], the L2 VPN information is carried in a Label

Mapping message sent in downstream unsolicited mode, which contains

the following VC FEC TLV.  VC, C, VC Info Length, Group ID,
Interface parameters are as defined in [[MARTINI-SIG](MARTINI-SIG)].

```
 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
 +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
 |    VC tlv     |C|          VC Type          |VC info Length |
 +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
 |                         Group ID                           |
 +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
 |                          VC ID                             |
 +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
 |                    Interface parameters                    |
 |                             "                              |
 |                             "                              |
 +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

This document defines a new VC type value in addition to the
following values already defined in [[MARTINI-SIG](MARTINI-SIG)]:

VC Type   Description

0x0001    Frame Relay DLCI
0x0002    ATM AAL5 VCC transport
0x0003    ATM transparent cell transport
0x0004    Ethernet VLAN
0x0005    Ethernet
0x0006    HDLC
0x0007    PPP
0x8008    CEM [8]
0x0009    ATM VCC cell transport
0x000A    ATM VPC cell transport
0x000B    Ethernet VPLS

VC types 0x0004 and 0x0005 identify VC LSPs that carry VLAN tagged
and untagged Ethernet traffic respectively, for point-to-point
connectivity.

We define a new VC type, Ethernet VPLS, with codepoint 0x000B to
identify VC LSPs that carry Ethernet traffic for multipoint
connectivity.  The Ethernet VC Type is described below.

For VC types 0x0001 to 0x000A, The VC ID identifies a particular VC.
For the VPLS VC type, the VC ID is a VPN identifier globally unique
within a service provider domain.

Note that the VCID as specified in [MARTINI_SIG] is a service

identifier, identifying a service emulating a point-to-point virtual

circuit.  In a VPLS, the VCID is a single service identifier,
identifying an emulated LAN segment.


5.6.  Ethernet VPLS VC Type
5.6.1.  VPLS Encapsulation actions

In a VPLS, a customer Ethernet packet without preamble is
encapsulated with a header as defined in [MARTINI-ENCAP].  A
customer Ethernet packet is defined as follows:

   - If the packet, as it arrives at the PE, has an encapsulation
     that is used by the local PE as a service delimiter, then that
     encapsulation is stripped before the packet is sent into the
     VPLS.  As the packet exits the VPLS, the packet may have a
     service-delimiting encapsulation inserted.

   - If the packet, as it arrives at the PE, has an encapsulation
     that is not service delimiting, then it is a customer packet
     whose encapsulation should not be modified by the VPLS.  This
     covers, for example, a packet that carries customer specific
     VLAN-Ids that the service provider neither knows about nor
     wants to modify.

By following the above rules, the Ethernet packet that traverses a
VPLS is always a customer Ethernet packet.  Note that the two
actions, at ingress and egress, of dealing with service delimiters
are local actions that neither PE has to signal to the other.  They
allow, for example, a mix-and-match of VLAN tagged and untagged
services at either end, and do not carry across a VPLS a VLAN tag
that may have only local significance.  The service delimiter may be
a VC label also, whereby an Ethernet VC given by [MARTINI-ENCAP] can
serve as the access side connection into a PE.  An RFC1483 PVC
encapsulation could be another service delimiter.  By limiting the
scope of locally significant encapsulations to the edge,
hierarchical VPLS models can be developed that provide the
capability to network-engineer VPLS deployments, as described below.


5.6.2.  VPLS Learning actions

Learning is done based on the customer Ethernet packet, as defined
above.  The Forwarding Information Base (FIB) keeps track of the
mapping of customer Ethernet packet addressing and the appropriate
VC label to use.  We define two modes of learning: qualified and
unqualified learning.

In qualified learning, the learning decisions at the PE are based on
the customer Ethernet packet's MAC address and VLAN tag, if one

exists.  If no VLAN tag exists, the default VLAN is assumed.

Effectively, within one VPLS, there are multiple logical FIBs, one
for each customer VLAN tag identified in a customer packet.

In unqualified learning, learning is based on a customer Ethernet
packet's MAC address only.  In other words, at any PE, there is only
one FIB per VPLS, which maps the MAC address in a customer Ethernet
packet to a VC label.

5.6.3.  VPLS Forwarding actions

The forwarding decisions taken at a PE couple with the learning
mode.  When using unqualified learning, unknown destination packets
are flooded to the entire VPLS.  When using qualified learning, the
scope of the flooding domain may be reduced (to the scope of the
customer VLAN).  How this may be achieved is outside the scope of
this draft.

It is important to ensure that the above learning and forwarding
modes are used consistently across the VPLS.  For example, when the
intention is to use qualified learning, duplicate MAC addresses with
different VLAN tags should not trigger re-learn events, which will
lead to incorrect forwarding decisions.  We propose that signaling
an optional parameter in the VC FEC will provide an adequate guard
against such misconfigurations.  By default, the behavior is
unqualified learning.

In order to signal the learning mode, we introduce a new interface
parameter [MARTINI-SIG].

Optional Interface Parameter
     0x06     VPLS Learning Mode
              Length: 1 byte.
              Value: 0 - unqualified learning
                     1 - qualified learning

6.  MAC Address Withdrawal

It MAY be desirable to remove MAC addresses that have been
dynamically learned for faster convergence.

We introduce an optional MAC TLV that is used to specify a list of
MAC addresses that can be removed using the Address Withdraw
Message.

The Address Withdraw message with MAC TLVs MAY be supported in order
to uninstall learned MAC addresses that have moved or gone away more
quickly.  Once a MAC address is unlearned, re-learning occurs
through flooding.

6.1.  MAC TLV

MAC addresses to be unlearned can be signaled using an LDP Address
Withdraw Message.  We define a new TLV, the MAC TLV.  Its format is
described below.  The encoding of a MAC TLV address is a 2-byte
802.1Q tag, followed by the 6-byte MAC address encoding specified by
IEEE 802 documents [g-ORIG] [802.1D-REV].  The 802.1Q tag and the
MAC address MUST appear in pairs.  If no tag is required, the value
of the tag field MUST be zero.

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|U|F|      Type             |              Length               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|         Reserved          |        802.1Q Tag #1              |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                      MAC address #1                           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          ...             |        802.1Q Tag #n               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                      MAC address #n                           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

U bit
     Unknown bit.  This bit MUST be set to 0.  If the MAC address
format is not understood, then the TLV is not understood, and MUST
be ignored.

F bit
     Forward bit.  This bit MUST be set to 0.  Since the LDP
mechanism used here is Targeted, the TLV MUST NOT be forwarded.

Type
     Type field.  This field MUST be set to 0x0404 (subject to IANA
approval).  This identifies the TLV type as MAC TLV.

Length
     Length field.  This field specifies the total length of the
TLV, including the Type and Length fields.

Reserved
     Reserved bits.  They MUST NOT be interpreted at the receiver,
and MUST be set to zero by the sender.

802.1Q Tag
     The 802.1Q Tag.  The value MUST be zero if the Ethernet VLAN
encapsulation is used.  If the Ethernet encapsulation is used, and
the Ethernet address is associated with a VLAN, it MUST be set to

the VLAN tag.  If the Ethernet encapsulation is used, and the MAC

address is not associated with a VLAN, it MUST be set to zero.
Since an 802.1Q tag is 12-bits, the high 4 bits of the field MUST be
set to zero.

MAC Address
     The MAC address being removed.

The LDP Address Withdraw Message contains a FEC TLV (to identify the
VPLS in consideration), a MAC Address TLV and optional parameters.
No optional parameters have been defined for the MAC Address
Withdraw signaling.

6.2.  Address Withdraw Message Containing MAC TLV

When MAC addresses are being removed explicitly, e.g., an adjacent
CE router has been disconnected, an Address Withdraw Message can be
sent with the list of MAC addresses to be withdrawn.

The processing for MAC TLVs received in an Address Withdraw Message
is:
  For each (VLAN tag, MAC address) pair in the TLV:
  - Remove the association between the (VLAN tag, MAC address) pair
     and VC label.  It does not matter whether the MAC address was
     installed as a static or dynamic address.

The scope of a MAC TLV is the VPLS specified in the FEC TLV in the
Address Withdraw Message.

The number of MAC addresses can be deduced from the length field in
the TLV.  The address list MAY be empty.  This tells the receiving
LSR to delete any MAC addresses learned from the sending LSR for the
VPLS specified by the FEC TLV.

7.  Operation of a VPLS

We show here an example of how a VPLS works.  The following
discussion uses the figure below, where a VPLS has been set up
between PE1, PE2 and PE3.

Initially, the VPLS is set up so that PE1, PE2 and PE3 have a full-
mesh of tunnels between them for carrying tunneled traffic.  The
VPLS service is assigned a VCID (a 32-bit quantity that is unique
across the provider network across all VPLSs). (Allocation of
domain-wide unique VCIDs is outside the scope of this draft.)

For the above example, say PE1 signals VC Label 102 to PE2 and 103
to PE3, and PE2 signals VC Label 201 to PE1 and 203 to PE3.

Assume a packet from A1 is bound for A2.  When it leaves CE1, say it
has a source MAC address of M1 and a destination MAC of M2.  If PE1
does not know where M2 is, it will multicast the packet to PE2 and
PE3.  When PE2 receives the packet, it will have an inner label of
201.  PE2 can conclude that the source MAC address M1 is behind PE1,
since it distributed the label 201 to PE1.  It can therefore
associate MAC address M1 with VC Label 102.

```
                                                     -----
                                                    /  A1 \
        ----                                 ----CE1     |
       /    \         --------       ------- /     |     |
      | A2 CE2-      /        \      /      PE1     \     /
       \    / \     /          \---/          \      -----
        ----      ---PE2                       |
                   | Service Provider Network  |
                    \           /    \         /
             -----  PE3        /      \       /
             |Agg|_/ --------       -------
              -|    |
       ----  / -----   ----
      /    \/    \   /     \            CE = Customer Edge Router
     | A3 CE3    --C4 A4 |             PE = Provider Edge Router
      \    /        \    /             Agg = Layer 2 Aggregation
       ----          ----
```


7.1.  MAC Address Aging

PEs that learn remote MAC addresses need to have an aging mechanism
to remove unused entries associated with a VC Label.  This is
important both for conservation of memory as well as for
administrative purposes.  For example, if a customer site A is shut
down, eventually, the other PEs should unlearn A's MAC address.

As packets arrive, MAC addresses are remembered.  The aging timer
for MAC address M SHOULD be reset when a packet is received with
source MAC address M.


8.  A Hierarchical VPLS Model

The solution described above requires a full mesh of tunnel LSPs
between all the PE routers that participate in the VPLS service.
For each VPLS service, n*(n-1) VCs must be setup between the PE
routers.  While this creates signaling overhead, the real detriment
to large scale deployment is the packet replication requirements for
each provisioned VCs on a PE router.  Hierarchical connectivity,
described in this document reduces signaling and replication

overhead to allow large scale deployment.


Lasserre, Kompella et al.                          [Page 12]

In many cases, service providers place smaller edge devices in
multi-tenant buildings and aggregate them into a PE device in a
large Central Office (CO) facility. In some instances, standard IEEE
802.1q (Dot 1Q) tagging techniques may be used to facilitate mapping
CE interfaces to PE VPLS access points.  When this is done, a
hierarchical architecture is created outside the context of VPLS; no
service level signaling is present between the PE router and the MTU
bridge.

It is often beneficial to extend the VPLS service tunneling
techniques into the MTU domain.  This can be accomplished by
treating the MTU device as a PE device and provisioning VCs between
it and every other edge, as an basic VPLS.  An alternative is to
utilize [MARTINI-ENCAP] VCs between the MTU and selected VPLS
enabled PE routers.  This section focuses on this alternative
approach.  The [VPLS] mesh core tier VCs (Hub) are augmented with
access tier VCs (Spoke) to form a two tier hierarchical VPLS (H-
VPLS).

Spoke VCs may be expanded to include any L2 tunneling mechanism,
expanding the scope of the first tier to include non-bridging VPLS
PE routers. The non-bridging PE router would extend a Spoke VC from
a Layer-2 switch that connects to it, through the service core
network, to a bridging VPLS PE router supporting Hub VCs.  We also
describe how VPLS-challenged nodes and low-end CEs without MPLS
capabilities may participate in a hierarchical VPLS.


8.1.  Hierarchical connectivity

This section describes the hub and spoke connectivity model and
describes the requirements of the bridging capable and non-bridging
MTU devices for supporting the spoke connections.

For rest of this discussion we will refer to a bridging capable MTU
device as MTU-s and a non-bridging capable PE device as PE-r.  A
routing and bridging capable device will be referred to as PE-rs.


8.1.1.  Spoke connectivity for bridging-capable devices

As shown in the figure below, consider the case where an MTU-s
device has a single connection to the PE-rs device placed in the CO.
The PE-rs devices are connected in a basic VPLS full mesh.   To
participate in the VPLS service, MTU-s device creates a single
point-to-point tunnel LSP to the PE-rs device in the CO.  We will
call this the spoke connection.  For each VPLS service, a single
spoke VC is setup between the MTU-s and the PE-rs based on [MARTINI-
SIG] and [MARTINI-ENCAP].  Unlike traditional [MARTINI-ENCAP] VCs

that terminate on a physical (or a VLAN-tagged logical) port at each
end, the spoke VC terminates on a virtual bridge instance on the

MTU-s and the PE-rs devices.  The MTU-s device and the PE-rs device
treat each spoke connection like an access port of the VPLS service.
On access ports, the combination of the physical port and the VLAN
tag is used to associate the traffic to a VPLS instance while the VC
label is used to associate the traffic from the virtual spoke port
with a VPLS instance, followed by a standard L2 lookup to identify
which customer port the frame needs to be sent to.

The signaling and association of the spoke connection to the VPLS
service may be done by introducing extensions to the LDP signaling
as specified in [SHAH-PECE].

```
                                                    PE2-rs
                                                    ------
                                                   /      \
                                                  |   --   |
                                                  |  / \   |
   CE-1                                           |  \B /  |
     \                                             \  --  /
      \                                            /------
       \    MTU-s                   PE1-rs        /   |
        \ ------                    ------       /    |
        /      \                   /      \     /     |
       | \ --   |      VC-1       |   --   |---/      |
       |  / \--|- - - - - - - - - - - |--/  \  |      |
       |  \B / |                  |  \B /  |         |
        \ /--  /                   \  --  / ---\     |
         /-----                     ------      \    |
        /                                        \   |
      ----                                        \ ------
     |Agg |                                       /      \
      ----                                       |   --   |
     /    \                                      | / \   |
   CE-2   CE-3                                   |  \B /  |
                                                  \ --   /
  MTU-s = Bridging capable MTU                     ------
  PE-rs = VPLS capable PE                          PE3-rs

  --
 /  \
 \B / = Virtual VPLS(Bridge)Instance
  --
 Agg = Layer-2 Aggregation
```

8.1.1.1.  MTU-s Operation

MTU-s device is defined as a device that supports layer-2 switching
functionality and does all the normal bridging functions of learning

and replication on all its ports, including the virtual spoke port.
Packets to unknown destination are replicated to all ports in the
service including the virtual spoke port.  Once the MAC address is

learned, traffic between CE1 and CE2 will be switched locally by the
MTU-s device saving the link capacity of the connection to the PE-
rs.  Similarly traffic between CE1 or CE2 and any remote destination
is switched directly on to the spoke connection and sent to the PE-
rs over the point-to-point VC LSP.

Since the MTU-s is bridging capable, only a single VC is required
per VPLS instance for any number of access connections in the same
VPLS service.  This further reduces the signaling overhead between
the MTU-s and PE-rs.


8.1.1.2.  PE-rs Operation

The PE-rs device is a device that supports all the bridging
functions for VPLS service and supports the routing and MPLS
encapsulation, i.e. it supports all the functions described in
[VPLS].   The operation on the PE-rs node is identical to that
described in [VPLS] with one addition.  A point-to-point VC
associated with the VPLS is regarded as a virtual port (see
discussion in Section 5.6.1 on service delimiting).  The operation
on the virtual spoke port is identical to the operation on an access
port as described in the earlier section.  As shown in the figure
above, each PE-rs device switches traffic between aggregated
[MARTINI-ENCAP] VCs that look like virtual ports and the network
side VPLS VCs.


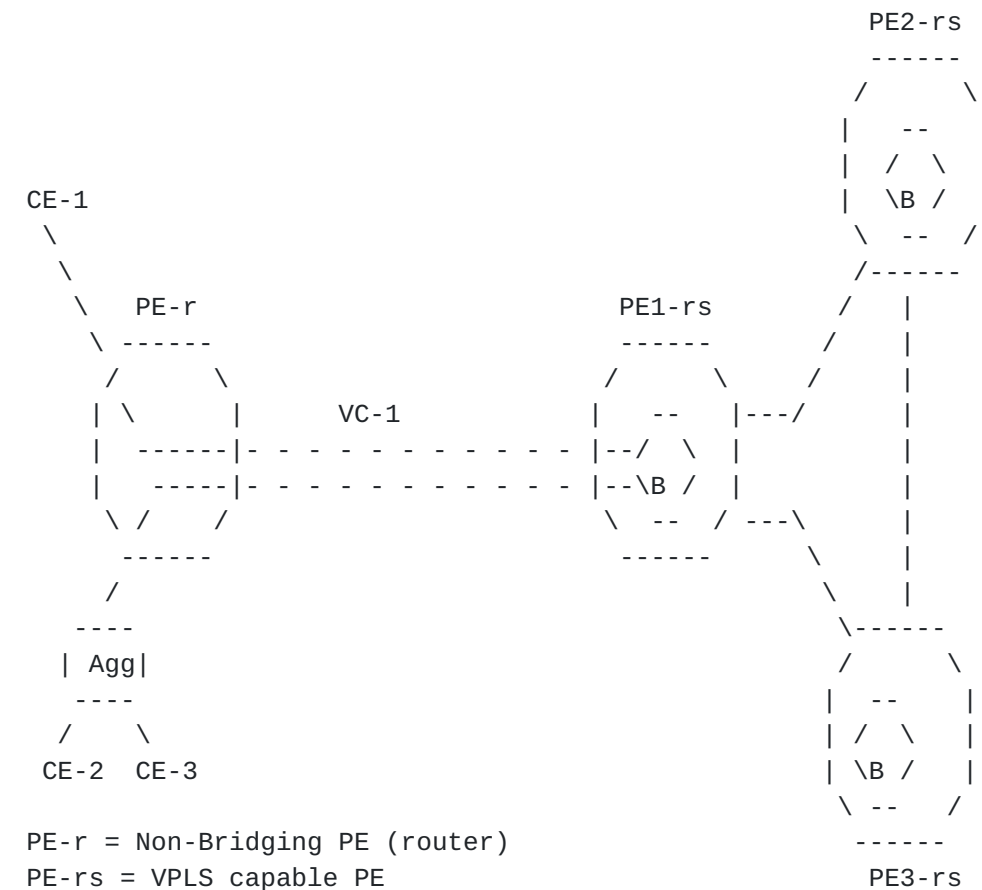8.1.2.  Advantages of spoke connectivity

Spoke connectivity offers several scaling and operational advantages
for creating large scale VPLS implementations, while retaining the
ability to offer all the functionality of the VPLS service.

- Eliminates the need for a full mesh of tunnels and full mesh of
  VCs per service between all devices participating in the VPLS
  service.
- Minimizes signaling overhead since fewer VC-LSPs are required for
  the VPLS service.
- Segments VPLS nodal discovery.  MTU-s needs to be aware of only
  the PE-rs node although it is participating in the VPLS service
  that spans multiple devices.  On the other hand, every VPLS PE-rs
  must be aware of every other VPLS PE-rs device and all of it s
  locally connected MTU-s and PE-r.
- Addition of other sites requires configuration of the new MTU-s
  device but does not require any provisioning of the existing MTU-s
  devices on that service.
- Hierarchical connections can be used to create VPLS service that
  spans multiple service provider domains. This is explained in a

later section.

8.1.3.  Spoke connectivity for non-bridging devices

In some cases, a bridging PE-rs device may not be deployed in some
CO while a PE-r might already be deployed.  If there is a need to
provide VPLS service from the CO where the PE-rs device is not
available, the service provider may prefer to use the PE-r device in
the interim.  In this section, we explain how a PE-r device that
does not support any of the bridging functionality as described in
[VPLS] can participate in the VPLS service.

```
                                                       PE2-rs
                                                       ------
                                                      /      \
                                                     |   --   |
                                                     | /  \   |
    CE-1                                             | \B /   |
      \                                               \  --  /
       \                                              /------
        \    PE-r                      PE1-rs        /   |
         \ ------                      ------       /    |
          /      \                    /      \     /     |
         | \      |      VC-1        |   --   |---/      |
         |  ------|- - - - - - - - - |--/  \  |          |
         |   -----|- - - - - - - - - |--\B /  |          |
          \ /    /                    \  --  / ---\      |
           ------                      ------      \     |
           /                                        \    |
         ----                                        \------
        | Agg|                                       /      \
         ----                                       |   --   |
        /     \                                     | /  \   |
     CE-2   CE-3                                    | \B /   |
                                                     \  --  /
     PE-r = Non-Bridging PE (router)                  ------
     PE-rs = VPLS capable PE                          PE3-rs


      --
     /  \
     \B / = Virtual VPLS(Bridge)Instance
      --
     Agg = Layer-2 Aggregation
```

As shown in this figure, the PE-r device creates a point-to-point
tunnel LSP to a PE-rs device.  Then for every access port that needs
to participate in a VPLS service, the PE-r device creates a point-
to-point [MARTINI-ENCAP] VC that terminates on the physical port at
the PE-r and terminates on the virtual bridge instance of the VPLS

service at the PE-rs.

8.1.3.1.  PE-r Operation

The PE-r device is defined as a device that supports routing but
does not support any bridging functions.  However, it is capable of
setting up [Martini-Encap] VCs between itself and the PE-rs.  For
every port that is supported in the VPLS service, a [MARTINI-ENCAP]
VC is setup from the PE-r to the PE-rs.  Once the VCs are setup,
there is no learning or replication function required on part of the
PE-r.  All traffic received on any of the access ports is
transmitted on the VC.  Similarly all traffic received on a VC is
transmitted to the access port where the VC terminates.  Thus
traffic from CE1 destined for CE2 is switched at PE-rs and not at
PE-r.

This approach adds more overhead than the bridging capable (MTU-s)
spoke approach since a VC is required for every access port that
participates in the service versus a single VC required per service
(regardless of access ports) when a MTU-s type device is used.
However, this approach offers the advantage of offering a VPLS
service in conjunction with a routed internet service without
requiring the addition of new MTU device.


8.1.3.2.  PE-rs Operation

The operation of PE-rs is independent of the type of device at the
other end of the spoke connection.  Whether there is a bridging
capable device (MTU-s) at the other end of the spoke connection or
there is a non-bridging device (PE-r) at the other end of the spoke
connection, the operation of PE-rs is exactly the same.  Thus, the
spoke connection from the PE-r is treated as a virtual port and the
PE-rs device switches traffic between the virtual port, access ports
and the network side VPLS VCs once it has learned the MAC addresses.


8.2.  Redundant Spoke Connections

An obvious weakness of the hub and spoke approach described thus far
is that the MTU device has a single connection to the PE-rs device.
In case of failure of the connection or the PE-rs device, the MTU
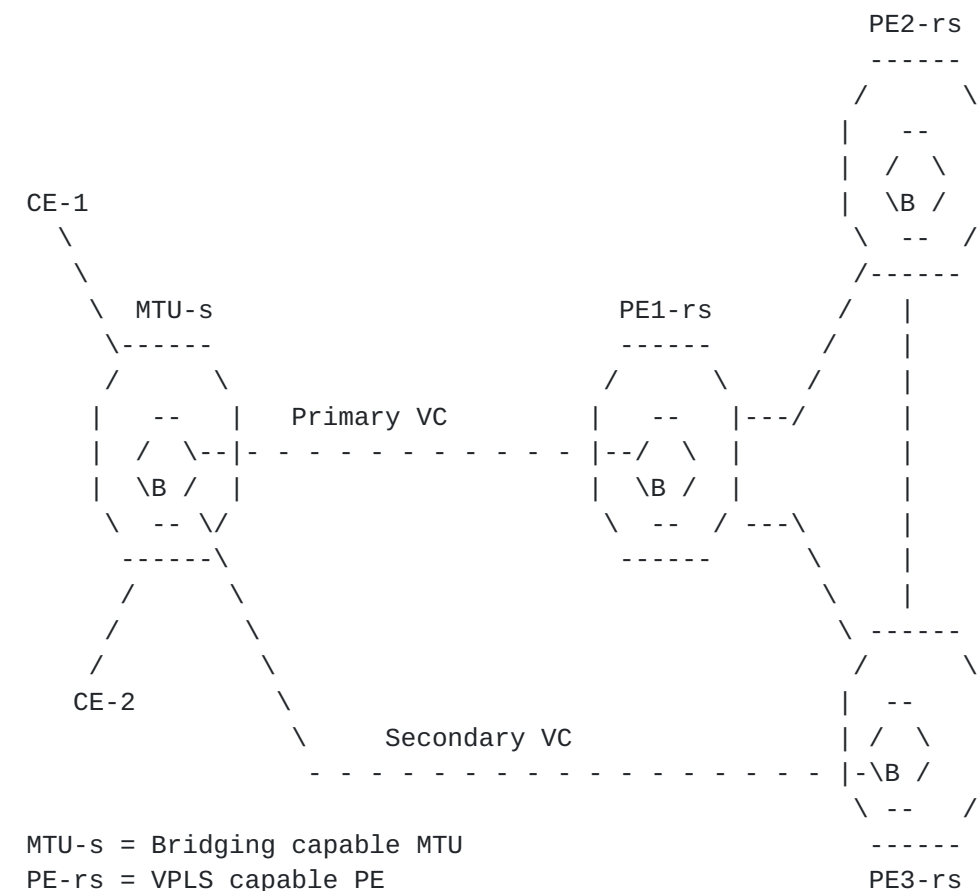device suffers total loss of connectivity.

In this section we describe how the redundant connections can be
provided to avoid total loss of connectivity from the MTU device.
The mechanism described is identical for both, MTU-s and PE-r type
of devices

8.2.1.  Dual-homed MTU device

To protect from connection failure of the VC or the failure of the
PE-rs device, the MTU-s device or the PE-r is dual-homed into two

PE-rs devices, as shown in figure-3.  The PE-rs devices must be part
of the same VPLS service instance.

An MTU-s device will setup two [MARTINI-ENCAP] VCs (one each to PE-
rs1 and PE-rs2) for each VPLS instances. One of the two VC is
designated as primary and is the one that is actively used under
normal conditions, while the second VC is designated as secondary
and is held in a standby state.  The MTU device negotiates the VC-
labels for both the primary and secondary VC, but does not use the
secondary VC unless the primary VC fails.  Since only one link is
active at a given time, a loop does not exist and hence 802.1D
spanning tree is not required.

```
                                                    PE2-rs
                                                    ------
                                                   /      \
                                                  |   --   |
                                                  |  / \   |
    CE-1                                          |  \B /  |
      \                                           \   --  /
       \                                           /------
        \   MTU-s                    PE1-rs       /    |
         \------                     ------      /     |
         /      \                   /      \    /      |
        |   --   |    Primary VC   |   --   |---/      |
        |  / \--|- - - - - - - - - - - |--/ \  |          |
        |  \B /  |                  |  \B /  |          |
         \  -- \/                    \   --  / ---\     |
          ------\                     ------      \     |
          /      \                                 \    |
         /        \                                 \ ------
        /          \                                /      \
     CE-2           \                              |   --   |
                     \        Secondary VC         | / \   |
                - - - - - - - - - - - - - - - - - - |-\B /   |
                                                   \  --   /
 MTU-s = Bridging capable MTU                       ------
 PE-rs = VPLS capable PE                            PE3-rs


 --
/  \
\B / = Virtual VPLS(Bridge)Instance
 --
```

8.2.2.   Failure detection and recovery

The MTU-s device controls the usage of the VC links to the PE-rs

nodes.  Since LDP signaling is used to negotiate the VC-labels, the
hello messages used for the LDP session are used to detect failure
of the primary VC.

Upon failure of the primary VC, MTU-s device immediately switches to
the secondary VC.  At this point the PE3-rs device that terminates
the secondary VC starts learning MAC addresses on the spoke VC.  All
other PE-rs nodes in the network think that CE-1 and CE-2 are behind
PE1-rs and may continue to send traffic to PE1-rs until they learn
that the devices are now behind PE3-rs.  The relearning process can
take a long time and may adversely affect the connectivity of higher
level protocols from CE1 and CE2.  To enable faster convergence, the
PE1-rs device where the primary VC failed sends out a flush message,
using the MAC TLV as defined in [Section 6](Section 6), to all other PE-rs
devices participating in the VPLS service.  Upon receiving the
message, all PE-rs flush the MAC addresses learned from PE1-rs.

8.3.  Multi-domain VPLS service

Hierarchy can also be used to create a large scale VPLS service
within a single domain or a service that spans multiple domains
without requiring full mesh connectivity between all VPLS capable
devices.  Two fully meshed VPLS networks are connected together
using a single LSP tunnel between the VPLS gateway devices.  A
single VC is setup per VPLS service to connect the two domains
together.  The VPLS gateway device joins two VPLS services together
to form a single multi-domain VPLS service.

9.  Acknowledgments

We wish to thank Joe Regan, Kireeti Kompella, Anoop Ghanwani, Joel
Halpern, Rick Wilder and Eric Rosen for their valuable feedback.

10.  Security Considerations

Security issues resulting from this draft will be discussed in
greater depth at a later point.  It is recommended in [[RFC3036](RFC3036)] that
LDP security (authentication) methods be applied.  This would
prevent unauthorized participation by a PE in a VPLS.  Traffic
separation for a VPLS is effected by using VC labels.  However, for
additional levels of security, the customer MAY deploy end-to-end
security, which is out of the scope of this draft.

11.  Intellectual Property Considerations

This document is being submitted for use in IETF standards
discussions.

12.  Full Copyright Statement

or assist in its implementation may be prepared, copied, published
and distributed, in whole or in part, without restriction of any

## 13.  References

[MARTINI-ENCAP] "Encapsulation Methods for Transport of Layer 2
Frames Over MPLS", draft-martini-l2circuit-encap-mpls-04.txt (Work
in progress)

[MARTINI-SIG] "Transport of Layer 2 Frames Over MPLS", draft-
martini-l2circuit-trans-mpls-08.txt (Work in progress)

[802.1D-ORIG] Original 802.1D - ISO/IEC 10038, ANSI/IEEE Std 802.1D-
1993 "MAC Bridges".

[802.1D-REV] 802.1D - "Information technology - Telecommunications
and information exchange between systems - Local and metropolitan
area networks - Common specifications - Part 3: Media Access Control
(MAC) Bridges: Revision. This is a revision of ISO/IEC 10038: 1993,
802.1j-1992 and 802.6k-1992. It incorporates P802.11c, P802.1p and
P802.12e." ISO/IEC 15802-3: 1998.

[802.1Q] 802.1Q - ANSI/IEEE Draft Standard P802.1Q/D11, "IEEE
Standards for Local and Metropolitan Area Networks: Virtual Bridged
Local Area Networks", July 1998.

[BGP-VPN] Rosen and Rekhter, "BGP/MPLS VPNs". RFC 2547, March 1999

[VPLS-REQ] "Requirements for Virtual Private LAN Services (VPLS)",
draft-augustyn-vpls-requirements-00.txt (Work in progress).

[RFC3036] "LDP Specification", L. Andersson, et al.  RFC 3036.
January 2001.

[SHAH-PECE] " Signaling between PE and L2PE/MTU for Decoupled VPLS and Hierarchical VPLS ", draft-shah-ppvpn-vpls-pe-mtu-signaling-00.txt, February, 2002. (Work in progress)

14.  Authors' Addresses

Marc Lasserre
Riverstone Networks
5200 Great America Pkwy
Santa Clara, CA 95054
Email: marc@riverstonenet.com

Vach Kompella
TiMetra Networks
274 Ferguson Dr.
Mountain View, CA 94043
Email: vkompella@timetra.com

Sunil Khandekar
TiMetra Networks
274 Ferguson Dr.
Mountain View, CA 94043
Email: sunil@timetra.com

Nick Tingle
TiMetra Networks
274 Ferguson Dr.
Mountain View, CA 94043
Email: ntingle@timetra.com

Loa Andersson
Utfors Bredband AB
Rasundavagen 12 169 29 Solna
Email: loa.andersson@utfors.se

Pascal Menezes
TeraBeam Networks
2300 Seventh Ave
Seattle, WA 98121
Email: Pascal.Menezes@Terabeam.com

Pierre Lin
Yipes Communication
114 Sansome St
San Francisco, CA 94104
Email: pierre.lin@yipes.com

Andrew Smith

Consultant
Email: ah_smith@pacbell.net

Giles Heron
PacketExchange Ltd.
The Truman Brewery
91 Brick Lane
LONDON E1 6QL
United Kingdom
Email: giles@packetexchange.net

Juha Heinanen
Song Networks, Inc.
Email: jh@lohi.eng.song.fi

Tom S. C. Soon
SBC Technology Resources Inc.
4698 Willow Road
Pleasanton, CA 94588
Email: sxsoon@tri.sbc.com

Ron Haberman
Masergy Inc.
2901 Telestar Ct.
Falls Church, VA 22042
Email: ronh@masergy.com

Luca Martini
Level 3 Communications, LLC.
1025 Eldorado Blvd.
Broomfield, CO, 80021
Email: luca@level3.net

Nick Slabakov
Riverstone Networks
5200 Great America Pkwy
Santa Clara, CA 95054
Email: nslabakov@riverstonenet.com

Rob Nath
Riverstone Networks
5200 Great America Pkwy
Santa Clara, CA 95054
Email: rnath@riverstonenet.com