Network Working Group                                          X. Xu
Internet-Draft                                    Huawei Technologies
Intended status: Informational                          C. Jacquenet
Expires: June 11, 2016                                         Orange
                                                           R. Raszuk
                                                            T. Boyes
                                                        Bloomberg LP
                                                              B. Fee
                                                    Extreme Networks
                                                   December 9, 2015


      **Virtual Subnet: A BGP/MPLS IP VPN-based Subnet Extension Solution**
                   **draft-ietf-bess-virtual-subnet-07**

Abstract

   This document describes a BGP/MPLS IP VPN-based subnet extension
   solution referred to as Virtual Subnet, which can be used for
   building Layer 3 network virtualization overlays within and/or
   between data centers.

Status of This Memo

Copyright Notice

Table of Contents

## 1.  Introduction

For business continuity purposes, Virtual Machine (VM) migration across data centers is commonly used in situations such as data center maintenance, migration, consolidation, expansion, or disaster avoidance.  The IETF community has recognized that IP renumbering of servers (i.e., VMs) after the migration is usually complex and costly.  To allow the migration of a VM from one data center to another without IP renumbering, the subnet on which the VM resides needs to be extended across these data centers.

To achieve subnet extension across multiple cloud data centers in a scalable way, the following requirements and challenges must be considered:

a.  VPN Instance Space Scalability: In a modern cloud data center
    environment, thousands or even tens of thousands of tenants could
    be hosted over a shared network infrastructure.  For security and
    performance isolation purposes, these tenants need to be isolated
    from one another.

b.  Forwarding Table Scalability: With the development of server
    virtualization technologies, it's not uncommon for a single cloud
    data center to contain millions of VMs.  This number already
    implies a big challenge to the forwarding table scalability of
    data center switches.  Provided multiple data centers of such
    scale were interconnected at Layer 2, this challenge would become
    even worse.

c.  ARP/ND Cache Table Scalability: [RFC6820] notes that the Address
    Resolution Protocol (ARP)/Neighbor Discovery (ND) cache tables
    maintained by default gateways within cloud data centers can
    raise scalability issues.  Therefore, mastering the size of the
    ARP/ND cache tables is critical as the number of data centers to
    be connected increases.

d.  ARP/ND and Unknown Unicast Flooding: It's well-known that the
    flooding of ARP/ND broadcast/multicast messages as well as
    unknown unicast traffic within large Layer 2 networks is likely
    to affect network and host performance.  When multiple data
    centers that each hosts millions of VMs are interconnected at
    Layer 2, the impact of such flooding would become even worse.  As
    such, it becomes increasingly important to avoid the flooding of
    ARP/ND broadcast/multicast as well as unknown unicast traffic
    across data centers.

e.  Path Optimization: A subnet usually indicates a location in the
    network.  However, when a subnet has been extended across
    multiple geographically-dispersed data center locations, the
    location semantics of such subnet is not retained any longer.  As
    a result, traffic exchanged between a specific user and a server
    that would be located in different data centers, may first be
    forwarded through a third data center.  This suboptimal routing
    would obviously result in an unnecessary consumption of the
    bandwidth resources between data centers.  Furthermore, in the
    case where traditional VPLS technology [RFC4761] [RFC4762] is
    used for data center interconnect, return traffic from a server
    may be forwarded to a default gateway located in a different data
    center due to the configuration of a virtual router redundancy
    group.  This suboptimal routing would also unnecessarily consume
    the bandwidth resources between data centers.

This document describes a BGP/MPLS IP VPN-based subnet extension
solution referred to as Virtual Subnet, which can be used for data
center interconnection while addressing all of the aforementioned
requirements and challenges.  Here the BGP/MPLS IP VPN means both
BGP/MPLS IPv4 VPN [RFC4364] and BGP/MPLS IPv6 VPN [RFC4659].  In
addition, since Virtual Subnet is mainly built on proven technologies
such as BGP/MPLS IP VPN and ARP/ND proxy [RFC0925][RFC1027][RFC4389],
those service providers that provide Infrastructure as a Service
(IaaS) cloud services can rely upon their existing BGP/MPLS IP VPN
infrastructure and take advantage of their BGP/MPLS VPN operational
experience to interconnect data centers.

Although Virtual Subnet is described in this document as an approach
for data center interconnection, it can be used within data centers
as well.

Note that the approach described in this document is not intended to
achieve an exact emulation of Layer 2 connectivity and therefore it
can only support a restricted Layer 2 connectivity service model with
limitations that are discussed in Section 4.  As for the discussion
about where this service model can apply, it's outside the scope of
this document.

## 2.  Terminology

This memo makes use of the terms defined in [RFC4364].

## 3.  Solution Description

### 3.1.  Unicast

### 3.1.1.  Intra-subnet Unicast

```
                        +-------------------+
    +------------------+   |                   |    +------------------+
    |VPN_A:192.0.2.1/24|   |                   |    |VPN_A:192.0.2.1/24|
    |            \    |   |                   |    |  | /              |
    |    +------+   \ ++---+-+           +-+---++/   +------+          |
    |    |Host A+-----+ PE-1 |           | PE-2 +----+Host B|          |
    |    +------+\    ++-+-+-+           +-+-+-++   /+------+          |
    |     192.0.2.2/24 | | |             | | |   192.0.2.3/24          |
    |                | | |             | | |                          |
    |     DC West    | | |  IP/MPLS Backbone  | | |   DC East          |
    +------------------+ | |             | | +------------------+
                      | +-------------------+ |
                      |                     |
VRF_A :               V          VRF_A : V
+-----------+---------+--------+       +-----------+---------+--------+
|   Prefix  | Nexthop |Protocol|       |   Prefix  | Nexthop |Protocol|
+-----------+---------+--------+       +-----------+---------+--------+
|192.0.2.1/32|127.0.0.1| Direct |      |192.0.2.1/32|127.0.0.1| Direct |
+-----------+---------+--------+       +-----------+---------+--------+
|192.0.2.2/32|192.0.2.2| Direct |      |192.0.2.2/32|  PE-1   |  IBGP  |
+-----------+---------+--------+       +-----------+---------+--------+
|192.0.2.3/32|  PE-2   |  IBGP  |      |192.0.2.3/32|192.0.2.3| Direct |
+-----------+---------+--------+       +-----------+---------+--------+
|192.0.2.0/24|192.0.2.1| Direct |      |192.0.2.0/24|192.0.2.1| Direct |
+-----------+---------+--------+       +-----------+---------+--------+
```

Figure 1: Intra-subnet Unicast Example

   As shown in Figure 1, two hosts (i.e., Hosts A and B) belonging to
   the same subnet (i.e., 192.0.2.0/24) are located in different data
   centers (i.e., DC West and DC East) respectively.  PE routers (i.e.,
   PE-1 and PE-2) that are used for interconnecting these two data
   centers create host routes for their own local hosts respectively and
   then advertise these routes by means of the BGP/MPLS IP VPN
   signaling.  Meanwhile, an ARP proxy is enabled on Virtual Routing and
   Forwarding (VRF) attachment circuits of these PE routers.

   Let's now assume that host A sends an ARP request for host B before
   communicating with host B.  Upon receiving the ARP request, PE-1
   acting as an ARP proxy returns its own MAC address as a response.
   Host A then sends IP packets for host B to PE-1.  PE-1 tunnels such
   packets towards PE-2 which in turn forwards them to host B.  Thus,
   hosts A and B can communicate with each other as if they were located
   within the same subnet.

[3.1.2](#).  **Inter-subnet Unicast**

```
                          +--------------------+
    +------------------+   |                    |   +------------------+
    |VPN_A:192.0.2.1/24|   |                    |   |VPN_A:192.0.2.1/24|
    |            \   |  |                    |   |   | /             |
    |  +------+    \ ++---+-+            +-+---++/    +------+   |
    |  |Host A+-------+ PE-1 |          | PE-2 +-+----+Host B|   |
    |  +------+\      ++-+-+-+          +-+-+-++ |   /+------+   |
    |   192.0.2.2/24  | | |            | | |  | 192.0.2.3/24  |
    |   GW=192.0.2.4  | | |            | | |  | GW=192.0.2.4  |
    |                 | | |            | | |  |    +------+   |
    |                 | | |            | | |  +----+  GW  +-- |
    |                 | | |            | | |       /+------+   |
    |                 | | |            | | |      192.0.2.4/24  |
    |                 | | |            | | |                   |
    |     DC West     | | |  IP/MPLS Backbone  | | |     DC East     |
    +------------------+ | |            | | +------------------+
                      | +--------------------+ |
                      |                        |
VRF_A :               V          VRF_A : V
+------------+---------+--------+      +------------+---------+--------+
|  Prefix    | Nexthop |Protocol|      |  Prefix    | Nexthop |Protocol|
+------------+---------+--------+      +------------+---------+--------+
|192.0.2.1/32|127.0.0.1| Direct |      |192.0.2.1/32|127.0.0.1| Direct |
+------------+---------+--------+      +------------+---------+--------+
|192.0.2.2/32|192.0.2.2| Direct |      |192.0.2.2/32|  PE-1   |  IBGP  |
+------------+---------+--------+      +------------+---------+--------+
|192.0.2.3/32|   PE-2  |  IBGP  |      |192.0.2.3/32|192.0.2.3| Direct |
+------------+---------+--------+      +------------+---------+--------+
|192.0.2.4/32|   PE-2  |  IBGP  |      |192.0.2.4/32|192.0.2.4| Direct |
+------------+---------+--------+      +------------+---------+--------+
|192.0.2.0/24|192.0.2.1| Direct |      |192.0.2.0/24|192.0.2.1| Direct |
+------------+---------+--------+      +------------+---------+--------+
| 0.0.0.0/0  |   PE-2  |  IBGP  |      | 0.0.0.0/0  |192.0.2.4| Static |
+------------+---------+--------+      +------------+---------+--------+
```
                   Figure 2: Inter-subnet Unicast Example (1)

   As shown in Figure 2, only one data center (i.e., DC East) is
   deployed with a default gateway (i.e., GW).  PE-2 that is connected
   to GW would either be configured with or learn from GW a default
   route with the next-hop being pointed at GW.  Meanwhile, this route
   is distributed to other PE routers (i.e., PE-1) as per normal
   [RFC4364] operation.  Assume host A sends an ARP request for its
   default gateway (i.e., 192.0.2.4) prior to communicating with a
   destination host outside of its subnet.  Upon receiving this ARP
   request, PE-1 acting as an ARP proxy returns its own MAC address as a
   response.  Host A then sends a packet for Host B to PE-1.  PE-1

tunnels such packet towards PE-2 according to the default route
learnt from PE-2, which in turn forwards that packet to GW.

```
                            +-------------------+
    +-----------------+     |                   |    +-----------------+
    |VPN_A:192.0.2.1/24|    |                   |    |VPN_A:192.0.2.1/24|
    |            \    |     |                   |    |   | /           |
    |  +------+    \ ++---+-+                 +-+---++/   +------+     |
    |  |Host A+----+--+ PE-1 |                | PE-2 +-+----+Host B|   |
    |  +------+\   |  ++-+-+-+                 +-+-+-++ |   /+------+   |
    |  192.0.2.2/24 | | | |                    | | | | 192.0.2.3/24   |
    |  GW=192.0.2.4 | | | |                    | | | | GW=192.0.2.4   |
    |  +------+    |  | | |                    | | | |   +------+     |
    |--+ GW-1 +----+  | | |                    | | |  +----+ GW-2 +-- |
    |  +------+\      | | |                    | | |    /+------+     |
    |  192.0.2.4/24   | | |                    | | |    192.0.2.4/24  |
    |                 | | |                    | | |                  |
    |     DC West     | | |  IP/MPLS Backbone  | | |     DC East      |
    +-----------------+ | |                    | | +-----------------+
                      | +--------------------+ |
                      |                        |
VRF_A :               V            VRF_A : V
+-----------+--------+--------+         +-----------+--------+--------+
|  Prefix   | Nexthop |Protocol|        |  Prefix   | Nexthop |Protocol|
+-----------+--------+--------+         +-----------+--------+--------+
|192.0.2.1/32|127.0.0.1| Direct |       |192.0.2.1/32|127.0.0.1| Direct |
+-----------+--------+--------+         +-----------+--------+--------+
|192.0.2.2/32|192.0.2.2| Direct |       |192.0.2.2/32| PE-1   | IBGP   |
+-----------+--------+--------+         +-----------+--------+--------+
|192.0.2.3/32|  PE-2  | IBGP   |        |192.0.2.3/32|192.0.2.3| Direct |
+-----------+--------+--------+         +-----------+--------+--------+
|192.0.2.4/32|192.0.2.4| Direct |       |192.0.2.4/32|192.0.2.4| Direct |
+-----------+--------+--------+         +-----------+--------+--------+
|192.0.2.0/24|192.0.2.1| Direct |       |192.0.2.0/24|192.0.2.1| Direct |
+-----------+--------+--------+         +-----------+--------+--------+
| 0.0.0.0/0 |192.0.2.4| Static |        | 0.0.0.0/0 |192.0.2.4| Static |
+-----------+--------+--------+         +-----------+--------+--------+
```

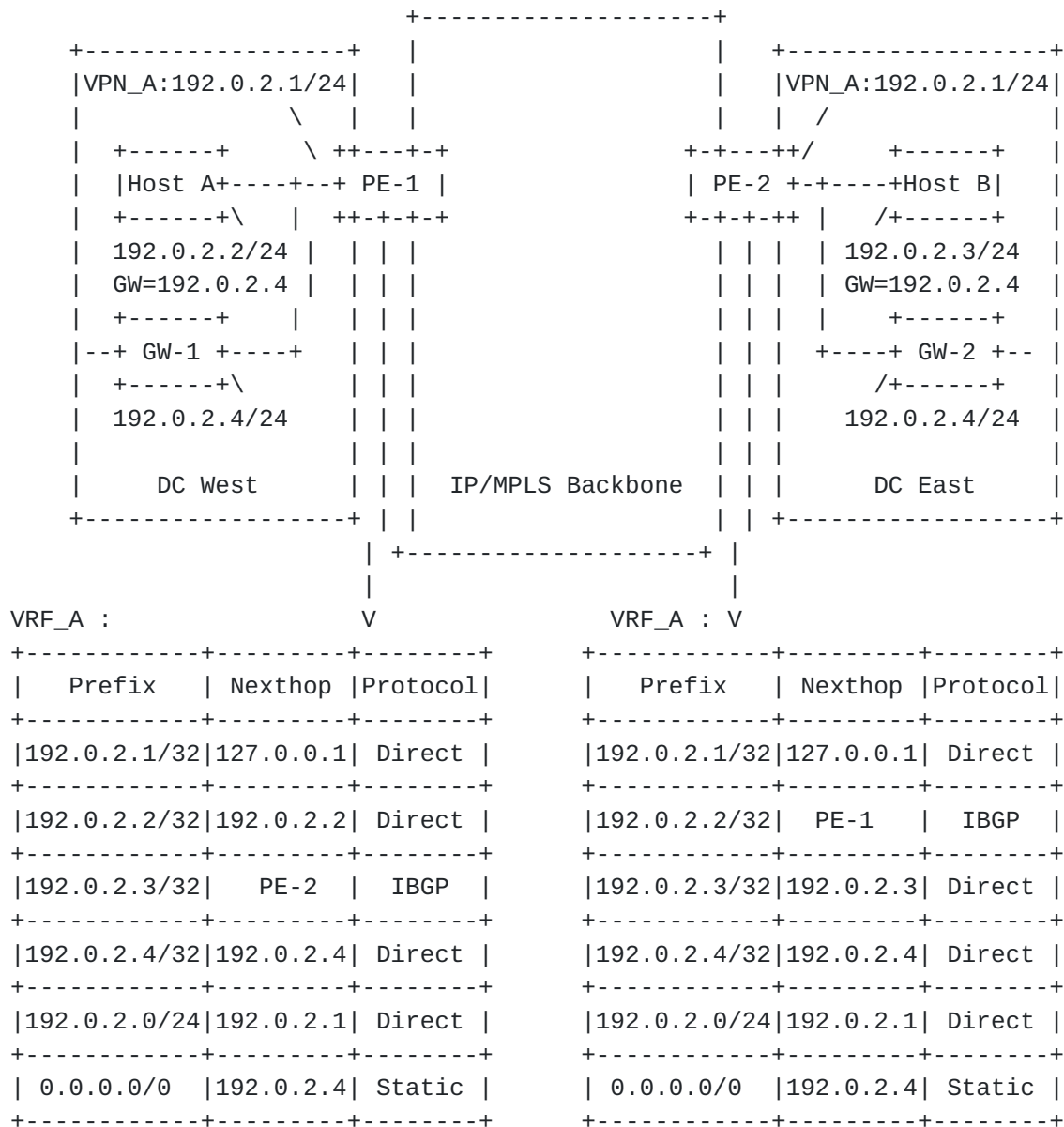                Figure 3: Inter-subnet Unicast Example (2)

   As shown in Figure 3, in the case where each data center is deployed
   with a default gateway, hosts will get ARP responses directly from
   their local default gateways, rather than from their local PE routers
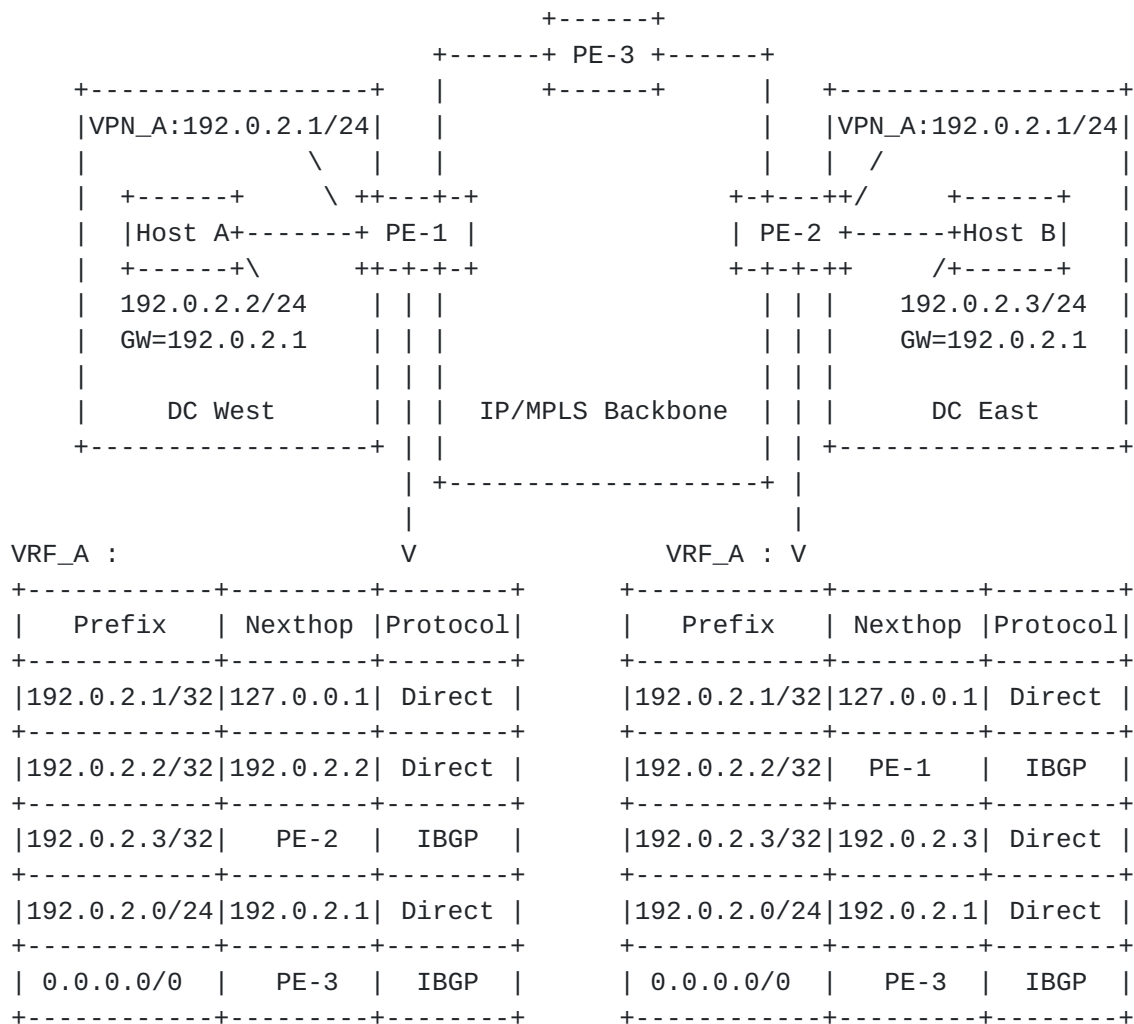   when sending ARP requests for their default gateways.

```
                                 +------+
                         +------+ PE-3 +------+
      +------------------+   |     +------+      |  +------------------+
      |VPN_A:192.0.2.1/24|   |                   |  |VPN_A:192.0.2.1/24|
      |            \   |   |                   |  |  /               |
      |   +------+    \ ++---+-+          +-+---++/     +------+    |
      |   |Host A+-------+ PE-1 |          | PE-2 +------+Host B|    |
      |   +------+\      ++-+-+-+          +-+-+-++     /+------+    |
      |   192.0.2.2/24    | | |            | | |    192.0.2.3/24   |
      |   GW=192.0.2.1    | | |            | | |    GW=192.0.2.1   |
      |                   | | |            | | |                   |
      |      DC West      | | |  IP/MPLS Backbone | | |    DC East      |
      +------------------+ | |            | | +------------------+
                          | +-------------------+ |
                          |                       |
VRF_A :                   V          VRF_A : V
+-----------+---------+--------+      +-----------+---------+--------+
|  Prefix   | Nexthop |Protocol|      |  Prefix   | Nexthop |Protocol|
+-----------+---------+--------+      +-----------+---------+--------+
|192.0.2.1/32|127.0.0.1| Direct |      |192.0.2.1/32|127.0.0.1| Direct |
+-----------+---------+--------+      +-----------+---------+--------+
|192.0.2.2/32|192.0.2.2| Direct |      |192.0.2.2/32|  PE-1   |  IBGP  |
+-----------+---------+--------+      +-----------+---------+--------+
|192.0.2.3/32|  PE-2   |  IBGP  |      |192.0.2.3/32|192.0.2.3| Direct |
+-----------+---------+--------+      +-----------+---------+--------+
|192.0.2.0/24|192.0.2.1| Direct |      |192.0.2.0/24|192.0.2.1| Direct |
+-----------+---------+--------+      +-----------+---------+--------+
| 0.0.0.0/0 |  PE-3   |  IBGP  |      | 0.0.0.0/0 |  PE-3   |  IBGP  |
+-----------+---------+--------+      +-----------+---------+--------+
```

               Figure 4: Inter-subnet Unicast Example (3)

   Alternatively, as shown in Figure 4, PE routers themselves could be
   configured as default gateways for their locally connected hosts as
   long as these PE routers have routes to reach outside networks.

## 3.2.  Multicast

   To support IP multicast between hosts of the same Virtual Subnet,
   MVPN technologies [RFC6513] could be used without any change.  For
   example, PE routers attached to a given VPN join a default provider
   multicast distribution tree which is dedicated to that VPN.  Ingress
   PE routers, upon receiving multicast packets from their local hosts,
   forward them towards remote PE routers through the corresponding
   default provider multicast distribution tree.  Within this context,
   the IP multicast doesn't include link-local multicast.

### 3.3.  Host Discovery

PE routers should be able to dynamically discover their local hosts
and keep the list of these hosts up-to-date in a timely manner so as
to ensure the availability and accuracy of the corresponding host
routes originated from them.  PE routers could accomplish local host
discovery by some traditional host discovery mechanisms using ARP or
ND protocols.

### 3.4.  ARP/ND Proxy

Acting as an ARP or ND proxy, a PE router should only respond to an
ARP request or Neighbor Solicitation (NS) message for a target host
when it has a best route for that target host in the associated VRF
and the outgoing interface of that best route is different from the
one over which the ARP request or NS message is received.  In the
scenario where a given VPN site (i.e., a data center) is multi-homed
to more than one PE router via an Ethernet switch or an Ethernet
network, the Virtual Router Redundancy Protocol (VRRP) [RFC5798] is
usually enabled on these PE routers.  In this case, only the PE
router being elected as the VRRP Master is allowed to perform the
ARP/ND proxy function.

### 3.5.  Host Mobility

During the VM migration process, the PE router to which the moving VM
is now attached would create a host route for that host upon
receiving a notification message of VM attachment (e.g., a gratuitous
ARP or unsolicited NA message).  The PE router to which the moving VM
was previously attached would withdraw the corresponding host route
when noticing the detachment of that VM.  Meanwhile, the latter PE
router could optionally broadcast a gratuitous ARP or send an
unsolicited NA message on behalf of that host with source MAC address
being one of its own.  In this way, the ARP/ND entry of this host
that moved and which has been cached on any local host would be
updated accordingly.  In the case where there is no explicit VM
detachment notification mechanism, the PE router could also use the
following trick to detect the VM detachment: upon learning a route
update for a local host from a remote PE router for the first time,
the PE router could immediately check whether that local host is
still attached to it by some means (e.g., ARP/ND PING and/or ICMP
PING).  It is important to ensure that the same MAC and IP are
associated to the default gateway active in each data center, as the
VM would most likely continue to send packets to the same default
gateway address after having migrated from one data center to
another.  One possible way to achieve this goal is to configure the
same VRRP group on each location so as to ensure that the default

gateway active in each data center shares the same virtual MAC and
virtual IP addresses.

## 3.6.  Forwarding Table Scalability on Data Center Switches

In a Virtual Subnet environment, the MAC learning domain associated
with a given Virtual Subnet which has been extended across multiple
data centers is partitioned into segments and each segment is
confined within a single data center.  Therefore data center switches
only need to learn local MAC addresses, rather than learning both
local and remote MAC addresses.

## 3.7.  ARP/ND Cache Table Scalability on Default Gateways

When default gateway functions are implemented on PE routers as shown
in Figure 4, the ARP/ND cache table on each PE router only needs to
contain ARP/ND entries of local hosts.  As a result, the ARP/ND cache
table size would not grow as the number of data centers to be
connected increases.

## 3.8.  ARP/ND and Unknown Unicast Flood Avoidance

In a Virtual Subnet environment, the flooding domain associated with
a given Virtual Subnet that has been extended across multiple data
centers, is partitioned into segments and each segment is confined
within a single data center.  Therefore, the performance impact on
networks and servers imposed by the flooding of ARP/ND broadcast/
multicast and unknown unicast traffic is minimized.

## 3.9.  Path Optimization

Take the scenario shown in Figure 4 as an example, to optimize the
forwarding path for the traffic between cloud users and cloud data
centers, PE routers located in cloud data centers (i.e., PE-1 and PE-
2), which are also acting as default gateways, propagate host routes
for their own local hosts respectively to remote PE routers which are
attached to cloud user sites (i.e., PE-3).  As such, traffic from
cloud user sites to a given server on the Virtual Subnet which has
been extended across data centers would be forwarded directly to the
data center location where that server resides, since traffic is now
forwarded according to the host route for that server, rather than
the subnet route.  Furthermore, for traffic coming from cloud data
centers and forwarded to cloud user sites, each PE router acting as a
default gateway would forward traffic according to the longest-match
route in the corresponding VRF.  As a result, traffic from data
centers to cloud user sites is forwarded along an optimal path as
well.

4. Limitations

4.1. Non-support of Non-IP Traffic

   Although most traffic within and across data centers is IP traffic,
   there may still be a few legacy clustering applications which rely on
   non-IP communications (e.g., heartbeat messages between cluster
   nodes).  Since Virtual Subnet is strictly based on L3 forwarding,
   those non-IP communications cannot be supported in the Virtual Subnet
   solution.  In order to support those few non-IP traffic (if present)
   in the environment where the Virtual Subnet solution has been
   deployed, the approach following the idea of "route all IP traffic,
   bridge non-IP traffic" could be considered.  In other words, all IP
   traffic including both intra- and inter-subnet, would be processed
   according to the Virtual Subnet design, while non-IP traffic would be
   forwarded according to a particular Layer 2 VPN approach.  Such
   unified L2/L3 VPN approach requires ingress PE routers to classify
   packets received from hosts before distributing them to the
   corresponding L2 or L3 VPN forwarding processes.  Note that more and
   more cluster vendors are offering clustering applications based on
   Layer 3 interconnection.

4.2. Non-support of IP Broadcast and Link-local Multicast

   As illustrated before, intra-subnet traffic across PE routers is
   forwarded at Layer 3 in the Virtual Subnet solution.  Therefore, IP
   broadcast and link-local multicast traffic cannot be forwarded across
   PE routers in the Virtual Subnet solution.  In order to support the
   IP broadcast and link-local multicast traffic in the environment
   where the Virtual Subnet solution has been deployed, the unified L2/
   L3 overlay approach as described in Section 4.1 could be considered
   as well.  That is, IP broadcast and link-local multicast messages
   would be forwared at Layer 2 while routable IP traffic would be
   processed according to the Virtual Subnet design.

4.3. TTL and Traceroute

   As mentioned before, intra-subnet traffic is forwarded at Layer 3 in
   the Virtual Subnet context.  Since it doesn't require any change to
   the Time To Live (TTL) handling mechanism of the BGP/MPLS IP VPN,
   when doing a traceroute operation on one host for another host
   (assuming that these two hosts are within the same subnet but are
   attached to different sites), the traceroute output would reflect the
   fact that these two hosts within the same subnet are actually
   connected via a Virtual Subnet, rather than a Layer 2 connection
   since the PE routers to which those two hosts are connected would be
   displayed in the traceroute output.  In addition, for any other
   applications that generate intra-subnet traffic with TTL set to 1,

these applications may not work properly in the Virtual Subnet
context, unless special TTL processing and loop-prevention mechanisms
for such context have been implemented.  Details about such special
TTL processing and loop-prevention mechanisms are outside the scope
of this document.

## 5.  Acknowledgements

Thanks to Susan Hares, Yongbing Fan, Dino Farinacci, Himanshu Shah,
Nabil Bitar, Giles Heron, Ronald Bonica, Monique Morrow, Rajiv Asati,
Eric Osborne, Thomas Morin, Martin Vigoureux, Pedro Roque Marque, Joe
Touch, Wim Henderickx, Alia Atlas and Stephen Farrell for their
valuable comments and suggestions on this document.  Thanks to Loa
Andersson for his WG LC review on this document.  Thanks to Alvaro
Retana for his AD review on this document.  Thanks to Ronald Bonica
for his RtgDir review.  Thanks to Donald Eastlake for his Sec-DIR
review of this document.  Thanks to Jouni Korhonen for the OPS-Dir
review of this document.  Thanks to Roni Even for the Gen-ART review
of this document.  Thanks to Sabrina Tanamal for the IANA review of
this document.

## 6.  IANA Considerations

There is no requirement for any IANA action.

## 7.  Security Considerations

Since the BGP/MPLS IP VPN signaling is reused without any change,
those security considerations as described in [RFC4364] are
applicable to this document.  Meanwhile, since security issues
associated with the NDP are inherited due to the use of NDP proxy,
those security considerations and recommendations as described in
[RFC6583] are applicable to this document as well.

Inter data-center traffic often carries highly sensitive information
at higher layers that is not directly understood (parsed) within an
egress or ingress PE.  For example, migrating a VM will often mean
moving private keys and other sensitive configuration information.
For this reason inter data-center traffic should always be protected
for both confidentiality and integrity using a strong security
mechanism such as IPsec [RFC4301].  In future it may be feasible to
protect that traffic within the MPLS layer
[I-D.ietf-mpls-opportunistic-encrypt] though at the time of writing
the mechanism for that is not sufficiently mature to recommend.
Exactly how such security mechanisms are deployed will vary from case
to case, so securing the inter data-center traffic may or may not
involve deploying security mechanisms on the ingress/egress PEs or
further "inside" the data centers concerned.  Note though that if

security is not deployed on the egress/ingress PEs there is a
substantial risk that some sensitive traffic may be sent in clear and
therefore be vulnerable to pervasive monitoring [RFC7258] or other
attacks.

## 8.  References

### 8.1.  Normative References

[RFC0925]  Postel, J., "Multi-LAN address resolution", RFC 925,
           DOI 10.17487/RFC0925, October 1984,
           <http://www.rfc-editor.org/info/rfc925>.

[RFC1027]  Carl-Mitchell, S. and J. Quarterman, "Using ARP to
           implement transparent subnet gateways", RFC 1027,
           DOI 10.17487/RFC1027, October 1987,
           <http://www.rfc-editor.org/info/rfc1027>.

[RFC4364]  Rosen, E. and Y. Rekhter, "BGP/MPLS IP Virtual Private
           Networks (VPNs)", RFC 4364, DOI 10.17487/RFC4364, February
           2006, <http://www.rfc-editor.org/info/rfc4364>.

[RFC4389]  Thaler, D., Talwar, M., and C. Patel, "Neighbor Discovery
           Proxies (ND Proxy)", RFC 4389, DOI 10.17487/RFC4389, April
           2006, <http://www.rfc-editor.org/info/rfc4389>.

### 8.2.  Informative References

[I-D.ietf-mpls-opportunistic-encrypt]
           Farrel, A. and S. Farrell, "Opportunistic Security in MPLS
           Networks", draft-ietf-mpls-opportunistic-encrypt-00 (work
           in progress), July 2015.

[RFC4301]  Kent, S. and K. Seo, "Security Architecture for the
           Internet Protocol", RFC 4301, DOI 10.17487/RFC4301,
           December 2005, <http://www.rfc-editor.org/info/rfc4301>.

[RFC4659]  De Clercq, J., Ooms, D., Carugi, M., and F. Le Faucheur,
           "BGP-MPLS IP Virtual Private Network (VPN) Extension for
           IPv6 VPN", RFC 4659, DOI 10.17487/RFC4659, September 2006,
           <http://www.rfc-editor.org/info/rfc4659>.

[RFC4761]  Kompella, K., Ed. and Y. Rekhter, Ed., "Virtual Private
           LAN Service (VPLS) Using BGP for Auto-Discovery and
           Signaling", RFC 4761, DOI 10.17487/RFC4761, January 2007,
           <http://www.rfc-editor.org/info/rfc4761>.

   [RFC4762]  Lasserre, M., Ed. and V. Kompella, Ed., "Virtual Private
              LAN Service (VPLS) Using Label Distribution Protocol (LDP)
              Signaling", RFC 4762, DOI 10.17487/RFC4762, January 2007,
              <http://www.rfc-editor.org/info/rfc4762>.

   [RFC5798]  Nadas, S., Ed., "Virtual Router Redundancy Protocol (VRRP)
              Version 3 for IPv4 and IPv6", RFC 5798,
              DOI 10.17487/RFC5798, March 2010,
              <http://www.rfc-editor.org/info/rfc5798>.

   [RFC6513]  Rosen, E., Ed. and R. Aggarwal, Ed., "Multicast in MPLS/
              BGP IP VPNs", RFC 6513, DOI 10.17487/RFC6513, February
              2012, <http://www.rfc-editor.org/info/rfc6513>.

   [RFC6583]  Gashinsky, I., Jaeggli, J., and W. Kumari, "Operational
              Neighbor Discovery Problems", RFC 6583,
              DOI 10.17487/RFC6583, March 2012,
              <http://www.rfc-editor.org/info/rfc6583>.

   [RFC6820]  Narten, T., Karir, M., and I. Foo, "Address Resolution
              Problems in Large Data Center Networks", RFC 6820,
              DOI 10.17487/RFC6820, January 2013,
              <http://www.rfc-editor.org/info/rfc6820>.

   [RFC7258]  Farrell, S. and H. Tschofenig, "Pervasive Monitoring Is an
              Attack", BCP 188, RFC 7258, DOI 10.17487/RFC7258, May
              2014, <http://www.rfc-editor.org/info/rfc7258>.

Authors' Addresses

   Xiaohu Xu
   Huawei Technologies
   No.156 Beiqing Rd
   Beijing  100095
   CHINA

   Email: xuxiaohu@huawei.com


   Christian Jacquenet
   Orange
   4 rue du Clos Courtel
   Cesson-Sevigne,  35512
   FRANCE

   Email: christian.jacquenet@orange.com

   Robert Raszuk
   Bloomberg LP
   731 Lexington Ave
   New York City, NY  10022
   USA

   Email: robert@raszuk.net


   Truman Boyes
   Bloomberg LP

   Email: tboyes@bloomberg.net


   Brendan Fee
   Extreme Networks

   Email: bfee@extremenetworks.com