

Network working group
Internet Draft
Category: Informational

X. Xu
Huawei

R. Raszuk

S. Hares

Y. Fan
China Telecom

C. Jacquenet
Orange

T. Boyes
Bloomberg LP

B Fee
Extreme Networks

Expires: July 2014

January 18, 2014

Virtual Subnet: A L3VPN-based Subnet Extension Solution

[draft-xu-l3vpn-virtual-subnet-03](#)

Abstract

This document describes a Layer3 Virtual Private Network (L3VPN)-based subnet extension solution referred to as Virtual Subnet, which can be used as a kind of Layer3 network virtualization overlay approach for data center interconnect.

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

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

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any

time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/1id-abstracts.txt>.

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

This Internet-Draft will expire on July 18, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Conventions used in this document

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](#) [[RFC2119](#)].

Table of Contents

1. Introduction	4
2. Terminology	6
3. Solution Description	6
3.1. Unicast	6
3.1.1. Intra-subnet Unicast	6
3.1.2. Inter-subnet Unicast	7
3.2. Multicast	9
3.3. CE Host Discovery	10
3.4. ARP/ND Proxy	10
3.5. CE Host Mobility	10
3.6. Forwarding Table Scalability on Data Center Switches	11
3.7. ARP/ND Cache Table Scalability on Default Gateways	11
3.8. ARP/ND and Unknown Uncast Flood Avoidance	11

3.9. Path Optimization	11
4. Limitations	12
4.1. Non-support of Non-IP Traffic	12
4.2. Non-support of IP Broadcast and Link-local Multicast ..	12
4.3. TTL and Traceroute	13
5. Security Considerations	13
6. IANA Considerations	13
7. Acknowledgements	13
8. References	13
8.1. Normative References	13
8.2. Informative References	14
Authors' Addresses	14

1. Introduction

For business continuity purposes, Virtual Machine (VM) migration across data centers is commonly used in those situations such as data center maintenance, data center migration, data center consolidation, data center expansion, and data center disaster avoidance. It's generally admitted that IP renumbering of servers (i.e., VMs) after the migration is usually complex and costly at the risk of extending the business downtime during the process of migration. 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.

In Infrastructure-as-a-Service (IaaS) cloud data center environments, to achieve subnet extension across multiple data centers in a scalable way, the following requirements SHOULD be considered for any data center interconnect solution:

1) 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. Hence, the data center interconnect solution SHOULD be capable of providing a large enough Virtual Private Network (VPN) instance space for tenant isolation.

2) Forwarding Table Scalability

With the development of server virtualization technologies, a single cloud data center containing millions of VMs is not uncommon. This number already implies a big challenge for data center switches, especially for core/aggregation switches, from the perspective of forwarding table scalability. Provided that multiple data centers of such scale were interconnected at layer2, this challenge would be even worse. Hence an ideal data center interconnect solution SHOULD prevent the forwarding table size of data center switches from growing by folds as the number of data centers to be interconnected increases.

3) ARP/ND Cache Table Scalability on Default Gateways

[RFC6820] notes that the Address Resolution Protocol (ARP)/Neighbor Discovery (ND) cache tables maintained by data center default gateways in cloud data centers can raise both

scalability and security issues. Therefore, an ideal data center interconnect solution SHOULD prevent the ARP/ND cache table size from growing by multiples as the number of data centers to be connected increases.

4) ARP/ND and Unknown Unicast Flood Suppression or Avoidance

It's well-known that the flooding of Address Resolution Protocol (ARP)/Neighbor Discovery (ND) broadcast/multicast and unknown unicast traffic within a large Layer2 network are likely to affect performances of networks and hosts. As multiple data centers each containing millions of VMs are interconnected together across the Wide Area Network (WAN) at layer2, the impact of flooding as mentioned above will become even worse. As such, it becomes increasingly desirable for data center operators to suppress or even avoid the flooding of ARP/ND broadcast/multicast and unknown unicast traffic across data centers.

5) 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, the traffic from a cloud user (i.e., a VPN user) which is destined for a given server located at one data center location of such extended subnet may arrive at another data center location firstly according to the subnet route, and then be forwarded to the location where the service is actually located. This suboptimal routing would obviously result in the unnecessary consumption of the bandwidth resources which are intended for data center interconnection. Furthermore, in the case where the traditional VPLS technology [RFC4761, [RFC4762](#)] is used for data center interconnect and default gateways of different data center locations are configured within the same virtual router redundancy group, the returning traffic from that server to the cloud user may be forwarded at layer2 to a default gateway located at one of the remote data center premises, rather than the one placed at the local data center location. This suboptimal routing would also unnecessarily consume the bandwidth resources which are intended for data center interconnect.

This document describes a L3VPN-based subnet extension solution referred to as Virtual Subnet (VS), which can meet all of the requirements of cloud data center interconnect as described above. Since VS mainly reuses existing technologies including BGP/MPLS IP VPN [[RFC4364](#)] and ARP/ND proxy [[RFC925](#)][RFC1027][[RFC4389](#)], it allows

those service providers offering IaaS public cloud services to interconnect their geographically dispersed data centers in a much scalable way, and more importantly, data center interconnection design can rely upon their existing MPLS/BGP IP VPN infrastructures and their experiences in the delivery and the operation of MPLS/BGP IP VPN services.

Although Virtual Subnet is described as a data center interconnection solution in this document, there is no reason to assume that this technology couldn't be used within data centers.

Note that the approach described in this document is not intended to achieve an exact emulation of L2 connectivity and therefore can only support a restricted L2 connectivity service model with limitations declared in [Section 4](#). As for the discussion about in which environment this service model should be suitable, it's outside the scope of this document.

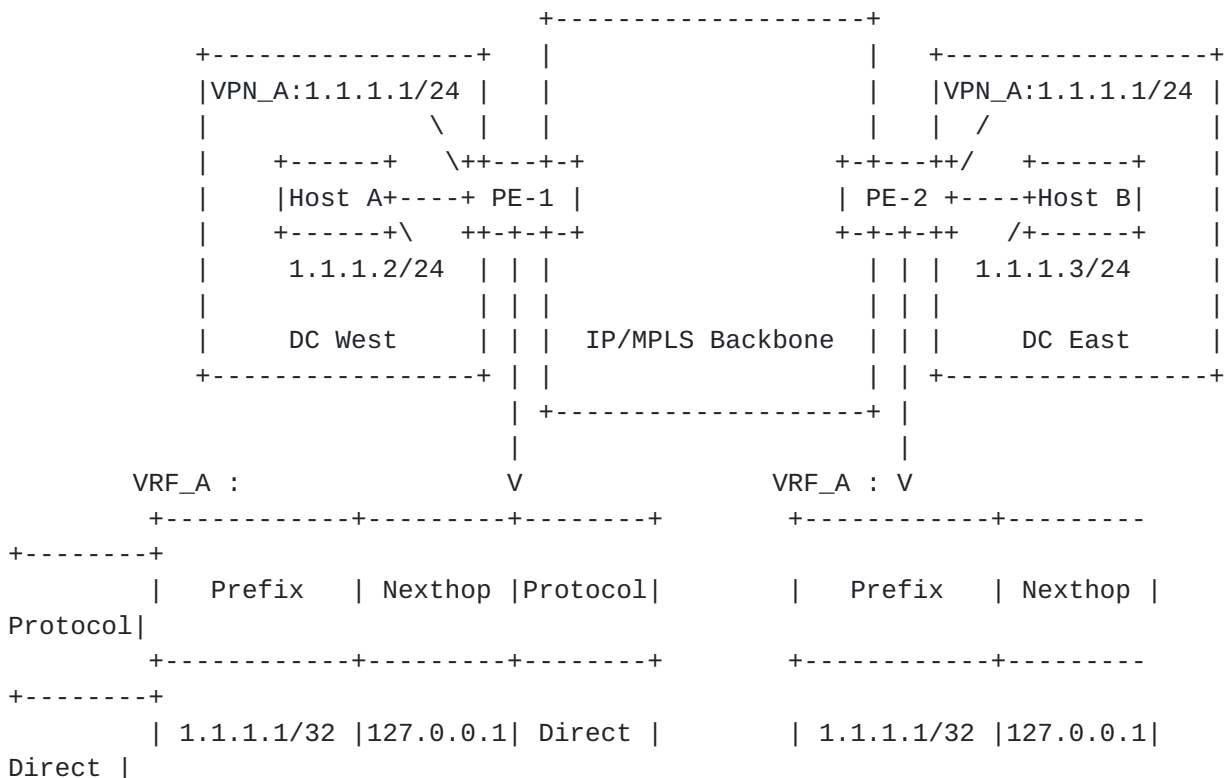
2. Terminology

This memo makes use of the terms defined in [\[RFC4364\]](#).

3. Solution Description

3.1. Unicast

3.1.1. Intra-subnet Unicast



	1.1.1.2/32	1.1.1.2	Direct		1.1.1.2/32	PE-1	
IBGP							
	1.1.1.3/32	PE-2	IBGP		1.1.1.3/32	1.1.1.3	
Direct							

Figure 1: Intra-subnet Unicast Example

Now assume 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.

```

+-----+ | +-----+
|VPN_A:1.1.1.1/24 | | |VPN_A:1.1.1.1/24 | | | |
| \ | | / |
| +-----+ \ +-----+ | +-----+ / +-----+ |
| |Host A+-----+ PE-1 | | PE-2 +-----+Host B| |
| +-----+\ +-----+ | +-----+ / +-----+ |
| 1.1.1.2/24 | | | 1.1.1.3/24 | |
| GW=1.1.1.4 | | | GW=1.1.1.4 | |
| | | | +-----+ |
| | | | +-----+ GW +---+ |
| | | | /+-----+ |
| | | | 1.1.1.4/24 |
| | | |
| DC West | | | IP/MPLS Backbone | | | DC East |
+-----+ | +-----+
| | |
VRF_A : V VRF_A : V
+-----+ +-----+
+-----+ | Prefix | Nexthop |Protocol| | Prefix | Nexthop |
Protocol| | | | | |
+-----+ +-----+
+-----+ | 1.1.1.1/32 |127.0.0.1| Direct | | 1.1.1.1/32 |127.0.0.1|
Direct | |

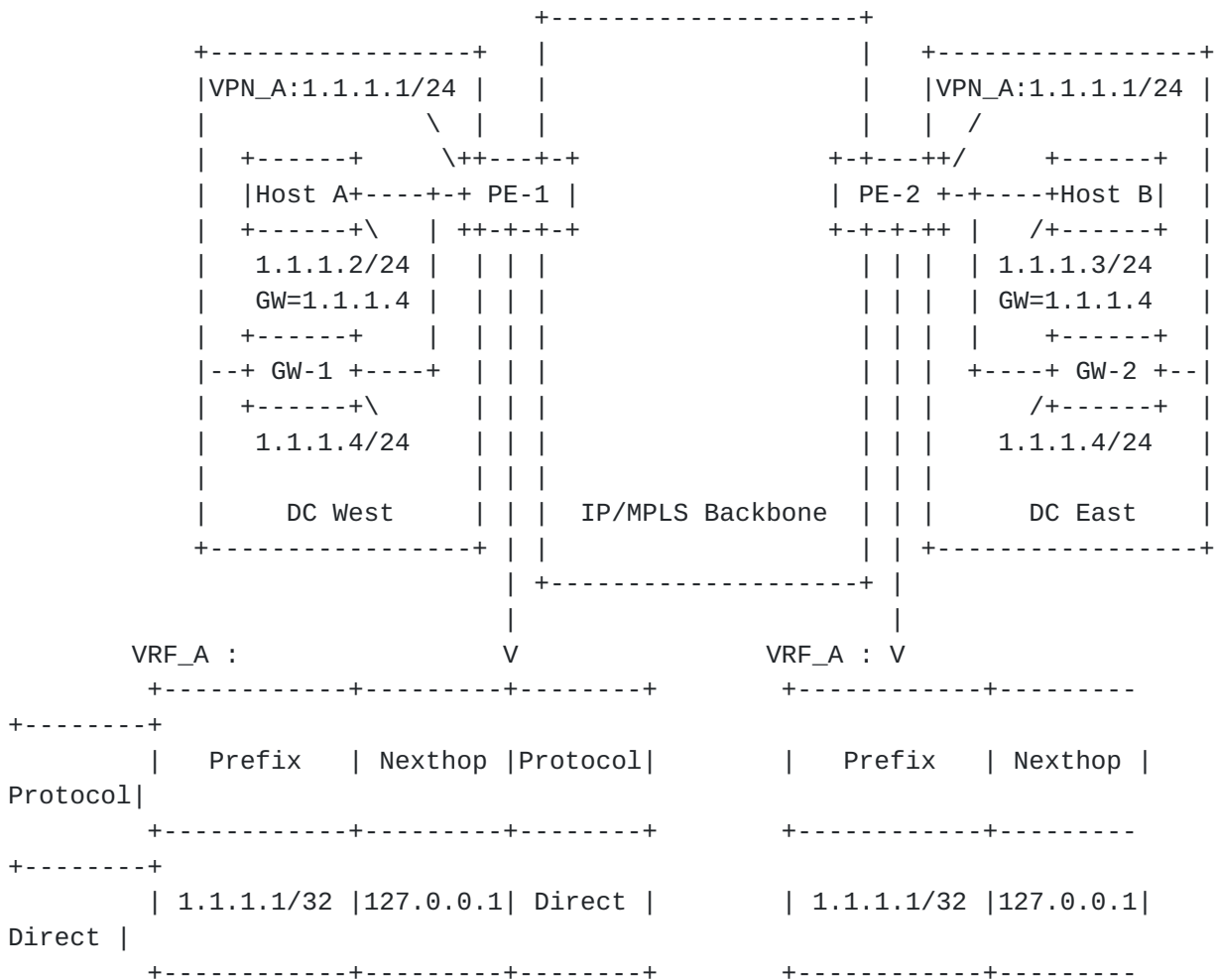
```

IBGP	1.1.1.2/32	1.1.1.2	Direct		1.1.1.2/32	PE-1	
Direct	1.1.1.3/32	PE-2	IBGP		1.1.1.3/32	1.1.1.3	
Direct	1.1.1.4/32	PE-2	IBGP		1.1.1.4/32	1.1.1.4	

Internet-Draft				Virtual Subnet				January 2014					
		1.1.1.0/24		1.1.1.1		Direct				1.1.1.0/24		1.1.1.1	
Direct													
		+-----+-----+-----+											
		+-----+											
		0.0.0.0/0		PE-2		IBGP				0.0.0.0/0		1.1.1.4	
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 which is connected to GW would either be configured with or learn from GW a default route with next-hop being pointed to 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., 1.1.1.4) 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.



IBGP	1.1.1.2/32 1.1.1.2 Direct	1.1.1.2/32 PE-1
	+-----+	+-----+
Direct	1.1.1.3/32 PE-2 IBGP	1.1.1.3/32 1.1.1.3
	+-----+	+-----+
Direct	1.1.1.4/32 1.1.1.4 Direct	1.1.1.4/32 1.1.1.4
	+-----+	+-----+
Direct	1.1.1.0/24 1.1.1.1 Direct	1.1.1.0/24 1.1.1.1
	+-----+	+-----+
Static	0.0.0.0/0 1.1.1.4 Static	0.0.0.0/0 1.1.1.4
	+-----+	+-----+

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, CE 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.

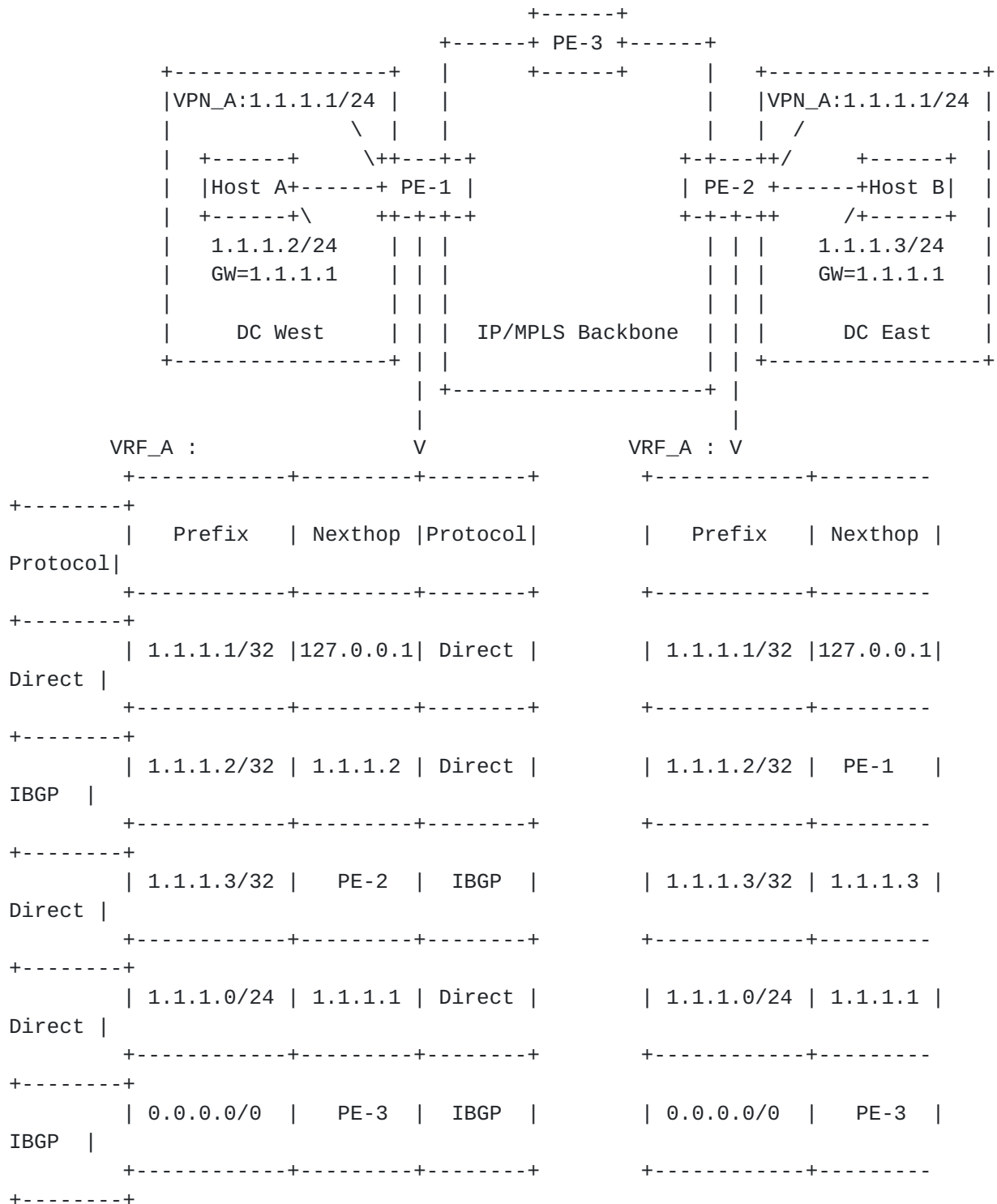


Figure 4: Inter-subnet Unicast Example (3)

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

[3.2.](#) Multicast

To support IP multicast between CE hosts of the same virtual subnet, MVPN technology [[MVPN](#)] could be directly reused. For example, PE routers attached to a given VPN join a default provider multicast distribution tree which is dedicated for that VPN. Ingress PE routers, upon receiving multicast packets from their local CE hosts, forward them towards remote PE routers through the corresponding default provider multicast distribution tree.

More details about how to support multicast and broadcast in VS will be explored in a later version of this document.

3.3. CE Host Discovery

PE routers SHOULD be able to discover their local CE 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 CE host discovery by some traditional host discovery mechanisms using ARP or ND protocols. Furthermore, Link Layer Discovery Protocol (LLDP) described in [[802.1AB](#)] or VSI Discovery and Configuration Protocol (VDP) described in [[802.1Qbg](#)], or even interaction with the data center orchestration system could also be considered as a means to dynamically discover local CE hosts.

3.4. ARP/ND Proxy

Acting as an ARP or ND proxies, a PE routers 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, 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. CE 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 CE 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 receiving a notification message of VM detachment (e.g., a VDP message about VM detachment). Meanwhile, the latter PE router could optionally broadcast a gratuitous ARP or send an unsolicited NA message on behalf of that CE host with source MAC address being one of its own. In this way, the ARP/ND entry of this CE host that moved and which has been cached on any local CE 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 determine the VM detachment event: upon learning a route

update for a local CE host from a remote PE router for the first time, the PE router could immediately check whether that local CE 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 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 the default gateway active in each data center share the same virtual MAC and virtual IP addresses.

3.6. Forwarding Table Scalability on Data Center Switches

In a VS 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

In case where data center default gateway functions are implemented on PE routers of the VS as shown in Figure 4, since the ARP/ND cache table on each PE router only needs to contain ARP/ND entries of local CE hosts, the ARP/ND cache table size will not grow as the number of data centers to be connected increases.

3.8. ARP/ND and Unknown Unicast Flood Avoidance

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

3.9. Path Optimization

Take the scenario shown in Figure 4 as an example, to optimize the forwarding path for traffic between cloud users and cloud data centers, PE routers located at cloud data centers (i.e., PE-1 and PE-2), which are also data center default gateways, propagate host routes for their local CE 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 the traffic received from its local CE hosts according to the best-match route in the corresponding VRF. As a result, traffic from data centers to cloud user sites is forwarded along the 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. That's to say, all IP traffic including both intra-subnet and inter-subnet would be processed by the Virtual Subnet process, while the non-IP traffic would be resorted to a particular Layer2 VPN approach. Such unified L2/L3 VPN approach requires ingress PE routers to classify the traffic received from CE 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 is forwarded at Layer3 in the Virtual Subnet solution. Therefore, IP broadcast and link-local multicast traffic cannot be supported by 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's to say, the IP broadcast and link-local multicast would be resorted to the L2VPN forwarding

process while the routable IP traffic would be processed by the Virtual Subnet process.

4.3. TTL and Traceroute

As illustrated before, intra-subnet traffic is forwarded at Layer3 in the Virtual Subnet context. Since it doesn't require any change to the TTL handling mechanism of the BGP/MPLS IP VPN, when doing a traceroute operation on one CE host for another CE 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 belonging to the same subnet are actually connected via an virtual subnet emulated by ARP proxy, rather than a normal LAN. In addition, for any other applications which generate intra-subnet traffic with TTL set to 1, these applications may not be workable in the Virtual Subnet context, unless special TTL processing for such case has been implemented (e.g., if the source and destination addresses of a packet whose TTL is set to 1 belong to the same extended subnet, both ingress and egress PE routers MUST NOT decrement the TTL of such packet. Furthermore, the TTL of such packet SHOULD NOT be copied into the TTL of the transport tunnel and vice versa).

5. Security Considerations

This document doesn't introduce additional security risk to BGP/MPLS IP VPN, nor does it provide any additional security feature for BGP/MPLS IP VPN.

6. IANA Considerations

There is no requirement for any IANA action.

7. Acknowledgements

Thanks to 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 and Wim Henderickx for their valuable comments and suggestions on this document.

8. References

8.1. Normative References

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

8.2. Informative References

- [RFC4364] Rosen. E and Y. Rekhter, "BGP/MPLS IP Virtual Private Networks (VPNs)", [RFC 4364](#), February 2006.
- [MVPN] Rosen. E and Aggarwal. R, "Multicast in MPLS/BGP IP VPNs", [draft-ietf-l3vpn-2547bis-mcast-10.txt](#), Work in Progress, January 2010.
- [RFC925] Postel, J., "Multi-LAN Address Resolution", [RFC-925](#), USC Information Sciences Institute, October 1984.
- [RFC1027] Smoot Carl-Mitchell, John S. Quarterman, "Using ARP to Implement Transparent Subnet Gateways", [RFC 1027](#), October 1987.
- [RFC4389] D. Thaler, M. Talwar, and C. Patel, "Neighbor Discovery Proxies (ND Proxy) ", [RFC 4389](#), April 2006.
- [RFC5798] S. Nadas., "Virtual Router Redundancy Protocol", [RFC 5798](#), March 2010.
- [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.
- [802.1AB] IEEE Standard 802.1AB-2009, "Station and Media Access Control Connectivity Discovery", September 17, 2009.
- [802.1Qbg] IEEE Draft Standard P802.1Qbg/D2.0, "Virtual Bridged Local Area Networks -Amendment XX: Edge Virtual Bridging", Work in Progress, December 1, 2011.
- [RFC6820] Narten, T., Karir, M., and I. Foo, "Problem Statement for ARMD", [RFC 6820](#), January 2013.

Authors' Addresses

Xiaohu Xu
Huawei Technologies,
Beijing, China.
Phone: +86 10 60610041
Email: xuxiaohu@huawei.com

Robert Raszuk
Email: robert@raszuk.net

Susan Hares
Email: shares@ndzh.com

Yongbing Fan
Guangzhou Institute, China Telecom
Guangzhou, China.
Phone: +86 20 38639121
Email: fanyb@gsta.com

Christian Jacquenet
Orange
Rennes France
Email: christian.jacquenet@orange.com

Truman Boyes
Bloomberg LP
Phone: +1 2126174826
Email: tboyes@bloomberg.net

Brendan Fee
Extreme Networks
9 Northeastern Blvd.
Salem, NH, 03079
Email: bfee@enterasys.com