

Network Working Group
Internet Draft
Category: Informational

L. Yong
Huawei
M. Toy
Comcast
A. Isaac
Bloomberg
V. Manral
Hewlett-Packard
L. Dunbar
Huawei

Expires: January 2014

July 11, 2013

Use Cases for DC Network Virtualization Overlays

[draft-ietf-nvo3-use-case-02](#)

Abstract

This document describes the DC Network Virtualization (NV03) use cases that may be potentially deployed in various data centers and apply to different applications.

Status of this Memo

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

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

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

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/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 January, 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..... | 3 |
| 1.1. | Contributors..... | 4 |
| 1.2. | Terminology..... | 4 |
| 2. | Basic Virtual Networks in a Data Center..... | 4 |
| 3. | Interconnecting DC Virtual Network and External Networks..... | 6 |
| 3.1. | DC Virtual Network Access via Internet..... | 6 |
| 3.2. | DC VN and Enterprise Sites interconnected via SP WAN..... | 7 |
| 4. | DC Applications Using NV03..... | 8 |
| 4.1. | Supporting Multi Technologies and Applications in a DC.... | 9 |
| 4.2. | Tenant Network with Multi-Subnets or across multi DCs.... | 9 |
| 4.3. | Virtual Data Center (vDC)..... | 11 |
| 5. | OAM Considerations..... | 12 |
| 6. | Summary..... | 13 |
| 7. | Security Considerations..... | 14 |
| 8. | IANA Considerations..... | 14 |
| 9. | Acknowledgements..... | 14 |
| 10. | References..... | 14 |
| 10.1. | Normative References..... | 14 |
| 10.2. | Informative References..... | 15 |
| | Authors' Addresses..... | 15 |

1. Introduction

Server Virtualization has changed IT industry in terms of efficiency, cost, and the speed in providing a new applications and/or services. However the problems in today's data center networks hinder the support of cloud applications and multi tenant networks [[NV03PRBM](#)]. The goal of DC Network Virtualization Overlays, i.e. NV03, is to decouple the communication among tenant systems from DC physical networks and to allow one physical network infrastructure to provide: 1) traffic isolation among tenant virtual networks over the same physical network; 2) independent address space in each virtual network and address isolation from the infrastructure's; 3) Flexible VM placement and move from one server to another without VM address and configuration change. These characteristics will help address the issues in today's cloud applications [[NV03PRBM](#)].

Although NV03 enables a true network virtualization environment, the NV03 solution has to address the communication between a virtual network and a physical network. This is because 1) many DCs that need to provide network virtualization are currently running over physical networks, the migration will be in steps; 2) a lot of DC applications are served to Internet users which run directly on physical networks; 3) some applications are CPU bound like Big Data analytics and may not need the virtualization capability.

This document is to describe general NV03 use cases that apply to various data centers. Three types of the use cases described here are:

- o Basic virtual networks in DC. A virtual network connects many tenant systems in a Data Center site (or more) and forms one L2 or L3 communication domain. Many virtual networks are over same DC physical network. The case may be used for DC internal applications that constitute the DC East-West traffic.
- o DC virtual network access from external. A DC provider offers a secure DC service to an enterprise customer and/or Internet users. An enterprise customer may use a traditional VPN provided by a carrier or an IPsec tunnel over Internet connecting to a virtual network within a provider DC site. This mainly constitutes DC North-South traffic.
- o DC applications or services that may use NV03. Three scenarios are described: 1) use NV03 and other network technologies to build a tenant network; 2) construct several virtual networks as a tenant network; 3) apply NV03 to a virtual DC (vDC) service.

The document uses the architecture reference model defined in [[NV03FRWK](#)] to describe the use cases.

[1.1.](#) Contributors

Vinay Bannai
PayPal
2211 N. First St,
San Jose, CA 95131
Phone: +1-408-967-7784
Email: vbannai@paypal.com

Ram Krishnan
Brocade Communications
San Jose, CA 95134
Phone: +1-408-406-7890
Email: ramk@brocade.com

[1.2.](#) Terminology

This document uses the terminologies defined in [[NV03FRWK](#)], [[RFC4364](#)]. Some additional terms used in the document are listed here.

CPE: Customer Premise Equipment

DMZ: Demilitarized Zone. A computer or small subnetwork that sits between a trusted internal network, such as a corporate private LAN, and an un-trusted external network, such as the public Internet.

DNS: Domain Name Service

NAT: Network Address Translation

VIRB: Virtual Integrated Routing/Bridging

Note that a virtual network in this document is an overlay virtual network instance.

[2.](#) Basic Virtual Networks in a Data Center

A virtual network may exist within a DC. The network enables a communication among Tenant Systems (TSs) that are in a Closed User Group (CUG). A TS may be a physical server/device or a virtual machine (VM) on a server. The network virtual edge (NVE) may co-

exist with Tenant Systems, i.e. on a same end-device, or exist on a different device, e.g. a top of rack switch (ToR). A virtual network has a unique virtual network identifier (may be local or global unique) for an NVE to properly differentiate it from other virtual networks.

The TSs attached to the same NVE may belong to the same or different virtual network. The multiple CUGs can be constructed in a way so that the policies are enforced when the TSs in one CUG communicate with the TSs in other CUGs. An NVE provides the reachability for Tenant Systems in a CUG, and may also have the policies and provide the reachability for Tenant Systems in different CUGs (See [section 4.2](#)). Furthermore in a DC operators may construct many tenant networks that have no communication in between at all. In this case, each tenant network may use its own address space. One tenant network may have one or more virtual networks.

A Tenant System may also be configured with multiple addresses and participate in multiple virtual networks, i.e. use different address in different virtual networks. For examples, a TS may be a NAT GW; or a firewall for multiple CUGs.

Network Virtualization Overlay in this context means that a virtual network is implemented in overlay, i.e. traffic from an NVE to another is sent via a tunnel.[NV03FMWK] This architecture decouples tenant system address scheme and configuration from the infrastructure's, which brings a great flexibility for VM placement and mobility. This also makes the transit nodes in the infrastructure not aware of the existence of the virtual networks. One tunnel may carry the traffic belonging to different virtual networks; a virtual network identifier is used for traffic demultiplexing.

A virtual network may be an L2 or L3 domain. The TSs attached to an NVE may belong to different virtual networks that may be in L2 or L3. A virtual network may carry unicast traffic and/or broadcast/multicast/unknown traffic from/to tenant systems. There are several ways to transport BUM traffic.[[NV03MCAST](#)]

It is worth to mention two distinct cases here. The first is that TSs and NVE are co-located on a same end device, which means that the NVE can be made aware of the TS state at any time via internal API. The second is that TSs and NVE are remotely connected, i.e. connected via a switched network or point-to-point link. In this case, a protocol is necessary for NVE to know TS state.

One virtual network may connect many TSes that attach to many different NVEs. TS dynamic placement and mobility results in frequent changes in the TS and NVE bindings. The TS reachability update mechanism need be fast enough to not cause any service interruption. The capability of supporting many TSs in a virtual network and many more virtual networks in a DC is critical for NV03 solution.

If a virtual network spans across multiple DC sites, one design is to allow the network seamlessly to span across the sites without DC gateway routers' termination. In this case, the tunnel between a pair of NVEs may in turn be tunneled over other intermediate tunnels over the Internet or other WANs, or the intra DC and inter DC tunnels are stitched together to form an end-to-end virtual network across DCs.

3. Interconnecting DC Virtual Network and External Networks

For customers (an enterprise or individuals) who utilize the DC provider's compute and storage resources to run their applications, they need to access their systems hosted in a DC through Internet or Service Providers' WANs. A DC provider may construct a virtual network that connect all the resources designated for a customer and allow the customer to access their resources via a virtual gateway (vGW). This, in turn, becomes the case of interconnecting a DC virtual network and the network at customer site(s) via Internet or WANs. Two cases are described here.

3.1. DC Virtual Network Access via Internet

A customer can connect to a DC virtual network via Internet in a secure way. Figure 1 illustrates this case. A virtual network is configured on NVE1 and NVE2 and two NVEs are connected via an L3 tunnel in the Data Center. A set of tenant systems are attached to NVE1 on a server. The NVE2 resides on a DC Gateway device. NVE2 terminates the tunnel and uses the VNID on the packet to pass the packet to the corresponding vGW entity on the DC GW. A customer can access their systems, i.e. TS1 or TSn, in the DC via Internet by using IPsec tunnel [[RFC4301](#)]. The IPsec tunnel is configured between the vGW and the customer gateway at customer site. Either static route or BGP may be used for peer routes. The vGW provides IPsec functionality such as authentication scheme and encryption. Note that: 1) some vGW functions such as firewall and load balancer may also be performed by locally attached network appliance devices; 2) The virtual network in DC may use different address space than external users, then vGW need to provide the NAT function; 3) more than one IPsec tunnels can be configured for the redundancy; 4) vGW

may be implemented on a server or VM. In this case, IP tunnels or IPsec tunnels may be used over DC infrastructure.

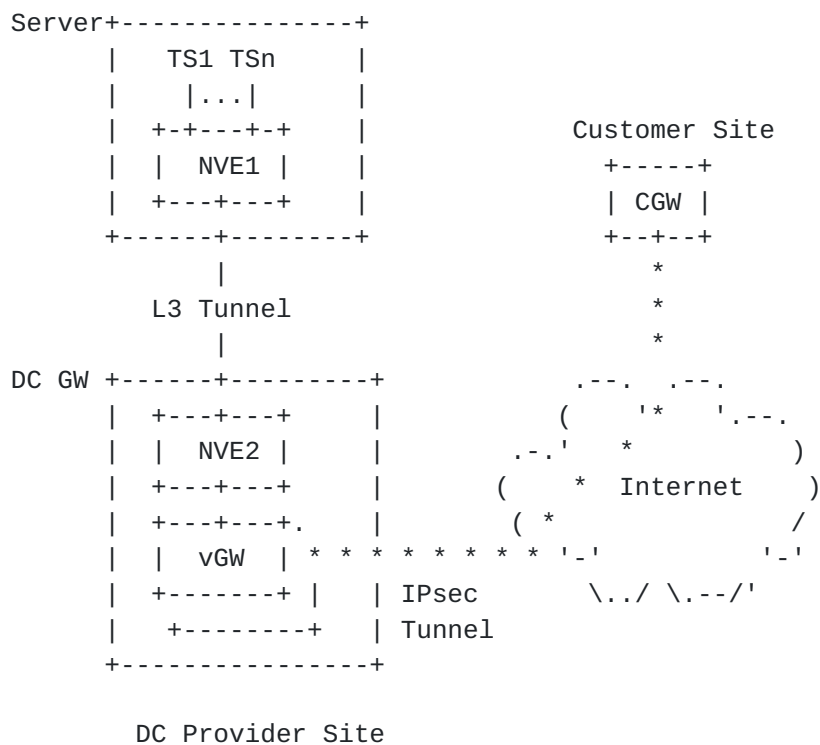


Figure 1 DC Virtual Network Access via Internet

3.2. DC VN and Enterprise Sites interconnected via SP WAN

An enterprise company may lease the VM and storage resources hosted in the 3rd party DC to run its applications. For example, the 3rd party company may run its web applications at 3 party sites but run backend applications in own DCs. The Web applications and backend applications need to communicate privately. The 3 party DC may construct one or more virtual networks to connect all VMs and storage running the Enterprise Web applications. The company may buy a p2p private tunnel such as VPWS from a SP to interconnect its site and the virtual network at the 3rd party site. A protocol is necessary for exchanging the reachability between two peering points and the traffic are carried over the tunnel. If an enterprise has multiple sites, it may buy multiple p2p tunnels to form a mesh interconnection among the sites and the 3 party site. This requires each site peering with all other sites for route distribution.

Another way to achieve multi-site interconnection is to use Service Provider (SP) VPN services, in which each site only peers with SP PE site. A DC Provider and VPN SP may build a DC virtual network (VN) and VPN independently. The VPN interconnects several enterprise sites and the DC virtual network at DC site, i.e. VPN site. The DC VN and SP VPN interconnect via a local link or a tunnel. The control plan interconnection options are described in [RFC4364](#) [[RFC4364](#)]. In Option A with VRF-LITE [[VRF-LITE](#)], both DC GW and SP PE maintain a routing/forwarding table, and perform the table lookup in forwarding. In Option B, DC GW and SP PE do not maintain the forwarding table, it only maintains the VN and VPN identifier mapping, and swap the identifier on the packet in the forwarding process. Both option A and B requires tunnel termination. In option C, DC GW and SP PE use the same identifier for VN and VPN, and just perform the tunnel stitching, i.e. change the tunnel end points. Each option has pros/cons (see [RFC4364](#)) and has been deployed in SP networks depending on the applications. The BGP protocols may be used in these options for route distribution. Note that if the provider DC is the SP Data Center, the DC GW and PE in this case may be on one device.

This configuration allows the enterprise networks communicating to the tenant systems attached to the VN in a provider DC without interfering with DC provider underlying physical networks and other virtual networks in the DC. The enterprise may use its own address space on the tenant systems in the VN. The DC provider can manage which VM and storage attachment to the VN. The enterprise customer manages what applications to run on the VMs in the VN. See [Section 4](#) for more.

The interesting feature in this use case is that the VN and compute resource are managed by the DC provider. The DC operator can place them at any server without notifying the enterprise and WAN SP because the DC physical network is completely isolated from the carrier and enterprise network. Furthermore, the DC operator may move the VMs assigned to the enterprise from one sever to another in the DC without the enterprise customer awareness, i.e. no impact on the enterprise 'live' applications running these resources. Such advanced features bring DC providers great benefits in serving cloud applications but also add some requirements for NV03 [[NV03PRBM](#)].

4. DC Applications Using NV03

NV03 brings DC operators the flexibility in designing and deploying different applications in an end-to-end virtualization overlay environment, where the operators no longer need to worry about the constraints of the DC physical network configuration when creating

VMs and configuring a virtual network. DC provider may use NV03 in various ways and also use it in the conjunction with physical networks in DC for many reasons. This section just highlights some use cases.

4.1. Supporting Multi Technologies and Applications in a DC

Most likely servers deployed in a large data center are rolled in at different times and may have different capacities/features. Some servers may be virtualized, some may not; some may be equipped with virtual switches, some may not. For the ones equipped with hypervisor based virtual switches, some may support VxLAN [[VXLAN](#)] encapsulation, some may support NVGRE encapsulation [[NVGRE](#)], and some may not support any types of encapsulation. To construct a tenant network among these servers and the ToR switches, it may construct one virtual network and one traditional VLAN network; or two virtual networks that one uses VxLAN encapsulation and another uses NVGRE.

In these cases, a gateway device or virtual GW is used to participate in multiple virtual networks. It performs the packet encapsulation/decapsulation and may also perform address mapping or translation, and etc.

A data center may be also constructed with multi-tier zones. Each zone has different access permissions and run different applications. For example, the three-tier zone design has a front zone (Web tier) with Web applications, a mid zone (application tier) with service applications such as payment and booking, and a back zone (database tier) with Data. External users are only able to communicate with the Web application in the front zone. In this case, the communication between the zones MUST pass through the security GW/firewall. One virtual network may be configured in each zone and a GW is used to interconnect two virtual networks. If individual zones use the different implementations, the GW needs to support these implementations as well.

4.2. Tenant Network with Multi-Subnets or across multi DCs

A tenant network may contain multiple subnets. The DC physical network needs support the connectivity for many tenant networks. The inter-subnets policies may be placed at some designated gateway devices only. Such design requires the inter-subnet traffic to be sent to one of the gateways first for the policy checking, which may cause traffic hairpin at the gateway in a DC. It is desirable that an NVE can hold some policies and be able to forward inter-subnet traffic directly. To reduce NVE burden, the hybrid design may be

deployed, i.e. an NVE can perform forwarding for the selected inter-subnets and the designated GW performs for the rest. For example, each NVE performs inter-subnet forwarding for a tenant, and the designated GW is used for inter-subnet traffic from/to the different tenant networks.

A tenant network may span across multiple Data Centers in distance. DC operators may configure an L2 VN within each DC and an L3 VN between DCs for a tenant network. For this configuration, the virtual L2/L3 gateway can be implemented on DC GW device. Figure 2 illustrates this configuration.

Figure 2 depicts two DC sites. The site A constructs one L2 VN, say L2VNa, on NVE1, NVE2, and NVE3. NVE1 and NVE2 reside on the servers which host multiple tenant systems. NVE3 resides on the DC GW device. The site Z has similar configuration with L2VNz on NVE3, NVE4, and NVE6. One L3 VN, say L3VNx, is configured on the NVE5 at site A and the NVE6 at site Z. An internal Virtual Interface of Routing and Bridging (VIRB) is used between L2VNI and L3VNI on NVE5 and NVE6, respectively. The L2VNI is the MAC/NVE mapping table and the L3VNI is the IP prefix/NVE mapping table. A packet to the NVE5 from L2VNa will be decapsulated and converted into an IP packet and then encapsulated and sent to the site Z. The policies can be checked at VIRB.

Note that the L2VNa, L2VNz, and L3VNx in Figure 2 are overlay virtual networks.

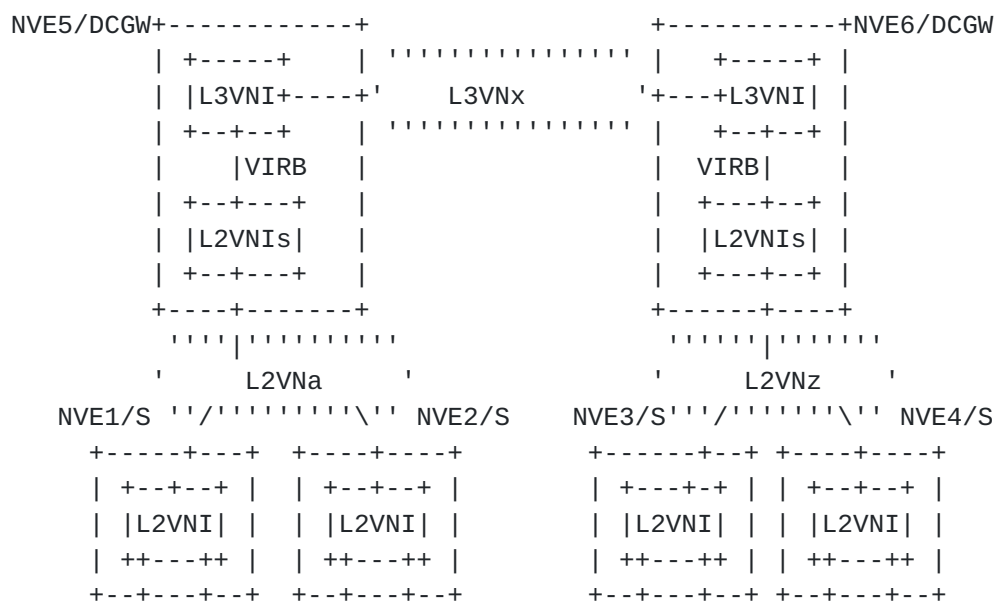




Figure 2 Tenant Virtual Network with Bridging/Routing

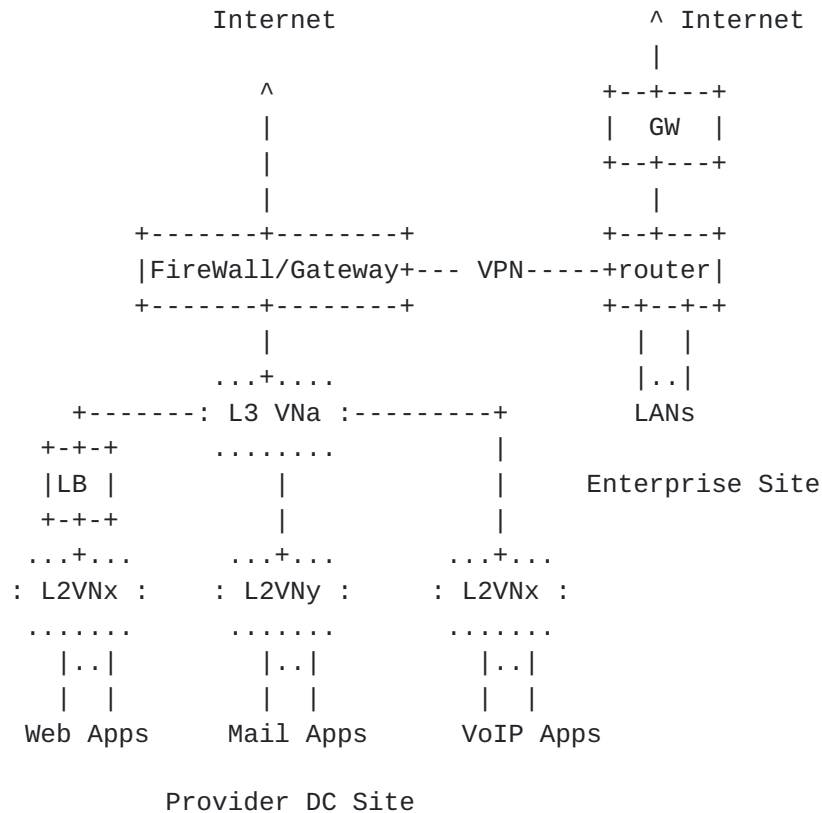
4.3. Virtual Data Center (vDC)

Enterprise DC's today may deploy routers, switches, and network appliance devices to construct its internal network, DMZ, and external network access and have many servers and storage running various applications. A DC Provider may offer a virtual DC service to enterprise customers. A vDC provides the same capability as a physical DC. A customer manages what and how applications to run in the vDC. Instead of using many hardware devices to do it, with the network virtualization overlay technology, DC operators may build such vDCs on top of a common DC infrastructure for many such customers and run network service application per vDC. The network service applications may include firewall, DNS, load balancer, gateway, etc. The network virtualization overlay further enables potential for vDC mobility when a customer moves to different locations because vDC configuration is decouple from the infrastructure network.

Figure 3 below illustrates one scenario. For the simple illustration, it only shows the L3 VN or L2 VN as virtual routers or switches. In this case, DC operators create several L2 VNs (L2VN_x, L2VN_y, L2VN_z) in Figure 3 to group the tenant systems together per application basis, create one L3 VN, e.g. VN_a for the internal routing. A net device (may be a VM or server) runs firewall/gateway applications and connects to the L3VN_a and Internet. A load balancer (LB) is used in L2 VN_x. A VPWS p2p tunnel is also built between the gateway and enterprise router. Enterprise customer runs Web/Mail/Voice applications at the provider DC site; lets the users at Enterprise site to access the applications via the VPN tunnel and Internet via a gateway at the Enterprise site; let Internet users access the applications via the gateway in the provider DC.

The customer decides which applications are accessed by intranet only and which by both intranet and extranet and configures the proper security policy and gateway function. Furthermore a customer may want multi-zones in a vDC for the security and/or set different QoS levels for the different applications.

This use case requires the NV03 solution to provide the DC operator an easy way to create a VN and NVEs for any design and to quickly assign TSs to VNIs on a NVE they attach to, easily to set up virtual topology and place or configure policies on an NVE or VMs that run net services, and support VM mobility. Furthermore a DC operator and/or customer should be able to view the tenant network topology and configure the tenant network functions. DC provider may further let a tenant to manage the vDC itself.



firewall/gateway and Load Balancer (LB) may run on a server or VMs

Figure 3 Virtual Data Center by Using NV03

5. OAM Considerations

NV03 brings the ability for a DC provider to segregate tenant traffic. A DC provider needs to manage and maintain NV03 instances. Similarly, the tenant needs to be informed about underlying network

failures impacting tenant applications or the tenant network is able to detect both overlay and underlay network failures and builds some resiliency mechanisms.

Various OAM and SOAM tools and procedures are defined in [IEEE 802.1ag], [ITU-T Y.1731], [[RFC4378](#)], [[RFC5880](#)], [ITU-T Y.1564] for L2 and L3 networks, and for user, including continuity check, loopback, link trace, testing, alarms such as AIS/RDI, and on-demand and periodic measurements. These procedures may apply to tenant overlay networks and tenants not only for proactive maintenance, but also to ensure support of Service Level Agreements (SLAs).

As the tunnel traverses different networks, OAM messages need to be translated at the edge of each network to ensure end-to-end OAM.

6. Summary

The document describes some general potential use cases of NV03 in DCs. The combination of these cases should give operators flexibility and capability to design more sophisticated cases for various purposes.

DC services may vary from infrastructure as a service (IaaS), platform as a service (PaaS), to software as a service (SaaS), in which the network virtualization overlay is just a portion of an application service. NV03 decouples the service construction/configurations from the DC network infrastructure configuration, and helps deployment of higher level services over the application.

NV03's underlying network provides the tunneling between NVEs so that two NVEs appear as one hop to each other. Many tunneling technologies can serve this function. The tunneling may in turn be tunneled over other intermediate tunnels over the Internet or other WANs. It is also possible that intra DC and inter DC tunnels are stitched together to form an end-to-end tunnel between two NVEs.

A DC virtual network may be accessed by external users in a secure way. Many existing technologies can help achieve this.

NV03 implementations may vary. Some DC operators prefer to use centralized controller to manage tenant system reachability in a tenant network, other prefer to use distributed protocols to advertise the tenant system location, i.e. associated NVEs. For the migration and special requirement, the different solutions may apply to one tenant network in a DC. When a tenant network spans across multiple DCs and WANs, each network administration domain may use

different methods to distribute the tenant system locations. Both control plane and data plane interworking are necessary.

7. Security Considerations

Security is a concern. DC operators need to provide a tenant a secured virtual network, which means one tenant's traffic isolated from the other tenant's traffic and non-tenant's traffic; they also need to prevent DC underlying network from any tenant application attacking through the tenant virtual network or one tenant application attacking another tenant application via DC networks. For example, a tenant application attempts to generate a large volume of traffic to overload DC underlying network. The NV03 solution has to address these issues.

8. IANA Considerations

This document does not request any action from IANA.

9. Acknowledgements

Authors like to thank Sue Hares, Young Lee, David Black, Pedro Marques, Mike McBride, David McDysan, Randy Bush, and Uma Chunduri for the review, comments, and suggestions.

10. References

10.1. Normative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997
- [RFC4364] Rosen, E. and Y. Rekhter, "BGP/MPLS IP Virtual Private Networks (VPNs)", [RFC 4364](#), February 2006.
- [IEEE 802.1ag] "Virtual Bridged Local Area Networks - Amendment 5: Connectivity Fault Management", December 2007.
- [ITU-T G.8013/Y.1731] OAM Functions and Mechanisms for Ethernet based Networks, 2011.
- [ITU-T Y.1564] "Ethernet service activation test methodology", 2011.
- [RFC4378] Allan, D., Nadeau, T., "A Framework for Multi-Protocol Label Switching (MPLS) Operations and Management (OAM)", [RFC4378](#), February 2006

- [RFC4301] Kent, S., "Security Architecture for the Internet Protocol", [rfc4301](#), December 2005
- [RFC5880] Katz, D. and Ward, D., "Bidirectional Forwarding Detection (BFD)", [rfc5880](#), June 2010.

10.2. Informative References

- [NVGRE] Sridharan, M., "NVGRE: Network Virtualization using Generic Routing Encapsulation", [draft-sridharan-virtualization-nvgre-02](#), work in progress.
- [NV03PRBM] Narten, T., et al "Problem Statement: Overlays for Network Virtualization", [draft-ietf-nvo3-overlay-problem-statement-03](#), work in progress.
- [NV03FRWK] Lasserre, M., Motin, T., and et al, "Framework for DC Network Virtualization", [draft-ietf-nvo3-framework-03](#), work in progress.
- [NV03MCAST] Ghanwani, A., "Multicast Issues in Networks Using NV03", [draft-ghanwani-nvo3-mcast-issues-00](#), work in progress.
- [VRF-LITE] Cisco, "Configuring VRF-lite", <http://www.cisco.com>
- [VXLAN] Mahalingam, M., Dutt, D., etc "VXLAN: A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks", [draft-mahalingam-dutt-dcops-vxlan-03.txt](#), work in progress.

Authors' Addresses

Lucy Yong
Huawei Technologies,
5340 Legacy Dr.
Plano, TX 75025

Phone: +1-469-277-5837
Email: lucy.yong@huawei.com

Mehmet Toy
Comcast
1800 Bishops Gate Blvd.,
Mount Laurel, NJ 08054

Phone : +1-856-792-2801
E-mail : mehmet_toy@cable.comcast.com

Aldrin Isaac
Bloomberg
E-mail: aldrin.isaac@gmail.com

Vishwas Manral
Hewlett-Packard Corp.
3000 Hanover Street, Building 20C
Palo Alto, CA 95014

Phone: 650-857-5501
Email: vishwas.manral@hp.com

Linda Dunbar
Huawei Technologies,
5340 Legacy Dr.
Plano, TX 75025 US

Phone: +1-469-277-5840
Email: linda.dunbar@huawei.com