

Network Working Group
Internet Draft
Intended status: Informational
Expires: September 27, 2020

L. Dunbar
Futurewei
B. Sarikaya
Denpel Informatique
B.Khasnabish
Independent
T. Herbert
Intel
S. Dikshit
Aruba-HPE
March 27, 2020

Virtual Machine Mobility Solutions for L2 and L3 Overlay Networks
draft-ietf-nvo3-vmm-10

Abstract

This document describes virtual machine mobility solutions commonly used in data centers built with overlay-based network. This document is intended for describing the solutions and the impact of moving VMs (or applications) from one Rack to another connected by the Overlay networks.

For layer 2, it is based on using an NVA (Network Virtualization Authority) - NVE (Network Virtualization Edge) protocol to update ARP (Address Resolution Protocol) table or neighbor cache entries after a VM (virtual machine) moves from an Old NVE to a New NVE. For Layer 3, it is based on address and connection migration after the move.

Status of this Memo

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

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#). This document may not be modified, and derivative works of it may not be created, except to publish it as an RFC and to translate it into languages other than English.

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 September 26, 2020.

Copyright Notice

Copyright (c) 2020 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.

Table of Contents

1.	Introduction.....	3
2.	Conventions used in this document.....	4
3.	Requirements.....	5
4.	Overview of the VM Mobility Solutions.....	6
4.1.	Inter-VNs communication.....	6

4.2. VM Migration in Layer 2 Network.....	6
4.3. VM Migration in Layer-3 Network.....	8
4.4. Address and Connection Management in VM Migration.....	9
5. Handling Packets in Flight.....	10
6. Moving Local State of VM.....	10
7. Handling of Hot, Warm and Cold VM Mobility.....	11
8. Other Options.....	11
9. VM Lifecycle Management.....	12
10. Security Considerations.....	12
11. IANA Considerations.....	13
12. Acknowledgments.....	13
13. Change Log.....	13
14. References.....	13
14.1. Normative References.....	14
14.2. Informative References.....	15

1. Introduction

This document describes the overlay-based data center networks solutions in supporting multitenancy and VM (Virtual Machine) mobility. Being able to move VMs dynamically, from one server to another, makes it possible for dynamic load balancing or work distribution. Therefore, dynamic VM Mobility is highly desirable for large scale multi-tenant DCs.

This document is strictly within the DCVPN, as defined by the NV03 Framework [[RFC 7365](#)]. The intent is to describe Layer 2 and Layer 3 Network behavior when VMs are moved from one NVE to another.

This document assumes that the VMs move is initiated by VM management system, i.e. planed move. How and when to move VM are out of the scope of this document. [RFC7666](#) already has the description of the MIB for VMs controlled by Hypervisor. The impact of VM mobility on higher layer protocols and applications is outside its scope.

Many large DCs (Data Centers), especially Cloud DCs, host tasks (or workloads) for multiple tenants. A tenant can be a department of one organization or an organization. There are communications among tasks belonging to one tenant and communications among tasks belonging to different tenants or with external entities.

Server Virtualization, which is being used in almost all of today's data centers, enables many VMs to run on a single physical computer or server sharing the processor/memory/storage. Network connectivity among VMs is provided by the network virtualization edge (NVE) [[RFC8014](#)]. It is highly desirable [[RFC7364](#)] to allow VMs to be moved dynamically (live, hot, or cold move) from one

server to another for dynamic load balancing or optimized work distribution.

There are many challenges and requirements related to VM mobility in large data centers, including dynamic attaching/detaching VMs to/from Virtual Network Edges (VNEs). In addition, retaining IP addresses after a move is a key requirement [[RFC7364](#)]. Such a requirement is needed in order to maintain existing transport connections.

In traditional Layer-3 based networks, retaining IP addresses after a move is generally not recommended because the frequent move will cause fragmented IP addresses, which introduces complexity in IP address management.

In view of many VM mobility schemes that exist today, there is a desire to document comprehensive VM mobility solutions that cover both IPv4 and IPv6. The large Data Center networks can be organized as one large Layer-2 network geographically distributed in several buildings/cities or Layer-3 networks with large number of host routes that cannot be aggregated as the result of frequent moves from one location to another without changing their IP addresses. The connectivity between Layer 2 boundaries can be achieved by the network virtualization edge (NVE) functioning as Layer 3 gateway routing across bridging domain such as in Warehouse Scale Computers (WSC).

[2. Conventions used in this document](#)

This document uses the terminology defined in [[RFC7364](#)]. In addition, we make the following definitions:

VM: Virtual Machine

Tasks: Task is a program instantiated or running on a virtual machine or container. Tasks in virtual machines or containers can be migrated from one server to another. We use task, workload and virtual machine interchangeably in this document.

Hot VM Mobility: A given VM could be moved from one server to another in running state.

Warm VM Mobility: In case of warm VM mobility, the VM states are mirrored to the secondary server (or domain) at a predefined (configurable) regular intervals. This reduces the overheads and complexity, but this may also lead to a situation when both servers may not contain the exact same data (state information)

Cold VM Mobility: A given VM could be moved from one server to another in stopped or suspended state.

Old NVE: refers to the old NVE where packets were forwarded to before migration.

New NVE: refers to the new NVE after migration.

Packets in flight: refers to the packets received by the Old NVE sent by the correspondents that have old ARP or neighbor cache entry before VM or task migration.

Users of VMs in diskless systems or systems not using configuration files are called end user clients.

Cloud DC: Third party data centers that host applications, tasks or workloads owned by different organizations or tenants.

3. Requirements

This section states requirements on data center network virtual machine mobility.

Data center network should support both IPv4 and IPv6 VM mobility.

Virtual machine (VM) mobility should not require changing VMs' IP addresses after the move.

There is "Hot Migration" with transport service continuing, and "Cold Migration" with transport service restarted, i.e. the task running is stopped on the Old NVE, moved to the New NVE and the task is restarted. Not all DCs support "Hot Migration. DCs that only support Cold Migration should make their customers aware of the potential service interruption during the Cold Migration.

VM mobility solutions/procedures should minimize triangular routing except for handling packets in flight.

VM mobility solutions/procedures should not need to use tunneling except for handling packets in flight.

4. Overview of the VM Mobility Solutions

Layer 2 and Layer 3 mobility solutions are described respectively in the following sections.

4.1. Inter-VNs communication

Inter VNs (Virtual Networks) communication refers to communication among tenants (or hosts) belonging to different VNs. Those tenants can be attached to the NVEs co-located in the same Data Center or in different Data centers. This document assumes that the inter-VNs communication is via the NV03 Gateway as described in [RFC8014](#) (NV03 Architecture). [RFC 8014](#) ([Section 5.3](#)) describes the NV03 Gateway function which is to relay traffic onto and off of a virtual network, i.e. among different VNs.

After a VM is moved to a new NVE, the VM's corresponding Gateway may need to change as well. If such a change is not possible, then the path to the external entity need to be hair-pinned to the NV03 Gateway used prior to the VM move.

4.2. VM Migration in Layer 2 Network

In a Layer-2 based approach, VM moving to another NVE does not change its IP address. But this VM is now under a new NVE, previously communicating NVEs may continue sending their packets to the Old NVE. Therefore, Address Resolution Protocol (ARP) cache in IPv4 [[RFC0826](#)] or neighbor cache in IPv6 [[RFC4861](#)] in the NVEs that have attached VMs communicating with the VM being moved need to be updated promptly. If the VM being moved has communication with external entities, the NV03 gateway needs to be notified of the new NVE where the VM is moved to.

In IPv4, the VM immediately after the move should send a gratuitous ARP request message containing its IPv4 and Layer 2 MAC address in its new NVE. Upon receiving this message, the New NVE can update its ARP cache. The New NVE should send a notification

of the newly attached VM to the central directory [[RFC7067](#)] embedded in the NVA to update the mapping of the IPv4 address & MAC address of the moving VM along with the new NVE address. An NVE-to-NVA protocol is used for this purpose [[RFC8014](#)]. The old NVE, upon a VM is moved away, should send an ARP scan to all its attached VMs to refresh its ARP Cache.

Reverse ARP (RARP) which enables the host to discover its IPv4 address when it boots from a local server [[RFC0903](#)], is not used by VMs if the VM already knows its IPv4 address (most common scenario). Next, we describe a case where RARP is used.

There are some vendor deployments (diskless systems or systems without configuration files) wherein the VM's user, i.e. end-user client asks for the same MAC address upon migration. This can be achieved by the clients sending RARP request message which carries the MAC address looking for an IP address allocation. The server, in this case the new NVE needs to communicate with NVA, just like in the gratuitous ARP case to ensure that the same IPv4 address is assigned to the VM. NVA uses the MAC address as the key in the search of ARP cache to find the IP address and informs this to the new NVE which in turn sends RARP reply message. This completes IP address assignment to the migrating VM.

Other NVEs that have attached VMs or the NV03 Gateway that have external entities communicating with this VM may still have the old ARP entry. To avoid old ARP entries being used by other NVEs, the old NVE upon discovering a VM is detached should send a notification to all other NVEs and its NV03 Gateway to time out the ARP cache for the VM [[RFC8171](#)]. When an NVE (including the old NVE) receives packet or ARP request destined towards a VM (its MAC or IP address) that is not in the NVE's ARP cache, the NVE should send query to NVA's Directory Service to get the associated NVE address for the VM. This is how the Old NVE tunneling these in-flight packets to the New NVE to avoid packets loss.

When VM address is IPv6, the operation is similar:

In IPv6, after the move, the VM immediately sends an unsolicited neighbor advertisement message containing its IPv6 address and Layer-2 MAC address to its new NVE. This message is sent to the IPv6 Solicited Node Multicast Address corresponding to the target address which is the VM's IPv6 address. The NVE receiving this message should send request to update VM's neighbor cache entry in the central directory of the NVA. The NVA's neighbor cache entry should include IPv6 address of the VM, MAC address of the VM and

the NVE IPv6 address. An NVE-to-NVA protocol is used for this purpose [[RFC8014](#)].

To avoid other NVEs communicating with this VM using the old neighbor cache entry, the old NVE upon discovering a VM being moved or VM management system which initiates the VM move should send a notification to all NVEs to timeout the ND cache for the VM being moved. When a ND cache entry for those VMs times out, their corresponding NVEs should send query to the NVA for an update.

4.3. VM Migration in Layer-3 Network

Traditional Layer-3 based data center networks usually have all hosts (tasks) within one subnet attached to one NVE. By this design, the NVE becomes the default route for all hosts (tasks) within the subnet. But this design requires IP address of a host (task) to change after the move to comply with the prefixes of the IP address under the new NVE.

A VM migration in Layer 3 Network solution is to allow IP addresses staying the same after moving to different locations. The Identifier Locator Addressing or ILA [[I-D.herbert-intarea-ila](#)] is one of such solutions.

Because broadcasting is not available in Layer-3 based networks, multicast of neighbor solicitations in IPv6 and ARP for IPv4 would need to be emulated. Scalability of the multicast (such as IPv6 ND and IPv4 ARP) can become problematic because the hosts belonging to one subnet (or one VLAN) can span across many NVEs. Sending broadcast traffic to all NVEs can cause unnecessary traffic in the DCN if the hosts belonging to one subnet are only attached to a very small number of NVEs. It is preferable to have a directory [[RFC7067](#)] or NVA to manage the updates to an NVE of the potential other NVEs a specific subnet may be attached and get periodic reports from an NVE of all the subnets being attached/detached, as described by [RFC8171](#).

Hot VM Migration in Layer 3 involves coordination among many entities, such as VM management system and NVA. Cold task migration, which is a common practice in many data centers, involves the following steps:

- Stop running the task.
- Package the runtime state of the job.

- Send the runtime state of the task to the New NVE where the task is to run.
- Instantiate the task's state on the new machine.
- Start the tasks for the task continuing from the point at which it was stopped.

[RFC7666](#) has the more detailed description of the State Machine of VMs controlled by Hypervisor

4.4. Address and Connection Management in VM Migration

Since the VM attached to the New NVE needs to be assigned with the same address as VM attached to the Old NVE, extra processing or configuration is needed, such as:

- Configure IPv4/v6 address on the target VM/NVE.
- Suspend use of the address on the old NVE. This includes the old NVE sending query to NVA upon receiving packets destined towards the VM being moved away. If there is no response from NVA for the new NVE for the VM, the old NVE can only drop the packets. Referring to the VM State Machine described in [RFC7666](#).
- Trigger NVA to push the new NVE-VM mapping to other NVEs which have the attached VMs communicating with the VM being moved.

Connection management for the applications running on the VM being moved involves reestablishing existing TCP connections in the new place.

The simplest course of action is to drop all TCP connections to the applications running on the VM during a migration. If the migrations are relatively rare events in a data center, impact is relatively small when TCP connections are automatically closed in the network stack during a migration event. If the applications running are known to handle this gracefully (i.e. reopen dropped connections) then this approach may be viable.

More involved approach to connection migration entails a proxy to the application (or the application itself) to pause the connection, package connection state and send to target, instantiate connection state in the peer stack, and restarting the connection. From the time the connection is paused to the time it is running again in the new stack, packets received for the connection could be silently dropped. For some period of time,

the old stack will need to keep a record of the migrated connection. If it receives a packet, it can either silently drop the packet or forward it to the new location, as described in [Section 5](#).

5. Handling Packets in Flight

The Old NVE may receive packets from the VM's ongoing communications. These packets should not be lost; they should be sent to the New NVE to be delivered to the VM. The steps involved in handling packets in flight are as follows:

Preparation Step: It takes some time, possibly a few seconds for a VM to move from its Old NVE to a New NVE. During this period, a tunnel needs to be established so that the Old NVE can forward packets to the New NVE. Old NVE gets New NVE address from its NVA assuming that the NVA gets the notification when a VM is moved from one NVE to another. It is out of the scope of this document on which entity manages the VM move and how NVA gets notified of the move. The Old NVE can store the New NVE address for the VM with a timer. When the timer expired, the entry for the New NVE for the VM can be deleted.

Tunnel Establishment - IPv6: Inflight packets are tunneled to the New NVE using the encapsulation protocol such as VXLAN in IPv6.

Tunnel Establishment - IPv4: Inflight packets are tunneled to the New NVE using the encapsulation protocol such as VXLAN in IPv4.

Tunneling Packets - IPv6: IPv6 packets received for the migrating VM are encapsulated in an IPv6 header at the Old NVE. New NVE decapsulates the packet and sends IPv6 packet to the migrating VM.

Tunneling Packets - IPv4: IPv4 packets received for the migrating VM are encapsulated in an IPv4 header at the Old NVE. New NVE decapsulates the packet and sends IPv4 packet to the migrating VM.

Stop Tunneling Packets: When the Timer for storing the New NVE address for the VM expires. The Timer should be long enough for all other NVEs that need to communicate with the VM to get their NVE-VM cache entries updated.

6. Moving Local State of VM

In addition to the VM mobility related signaling (VM Mobility Registration Request/Reply), the VM state needs to be transferred to the New NVE. The state includes its memory and file system if

the VM cannot access the memory and the file system after moving to the New NVE.

The mechanism of transferring VM States and file system is out of the scope of this document. Referring to [RFC7666](#) for detailed information.

7. Handling of Hot, Warm and Cold VM Mobility

Both Cold and Warm VM mobility (or migration) refers to the VM being completely shut down at the Old NVE before restarted at the New NVE. Therefore, all transport services to the VM are restarted.

In this document, all VM mobility is initiated by VM Management System. The Cold VM mobility only exchange the needed states between the Old NVE and the New NVE after the VM attached to the Old NVE is completely shut down. There is time delay before the new VM is launched. The cold mobility option can be used for non-critical applications and services that can tolerate interrupted TCP connections.

The Warm VM mobility refers to having the backup entities receive backup information at more frequent intervals, so that it can take less time to launch the VM under the new NVE and other NVEs that communicate with the VM can be notified prior to the VM move. The duration of the interval determines the effectiveness (or benefit) of Warm VM mobility. The larger the duration, the less effective the Warm VM mobility option becomes.

For Hot VM Mobility, once a VM moves to a New NVE, the VM IP address does not change and the VM should be able to continue to receive packets to its address(es). The VM needs to send a gratuitous Address Resolution message or unsolicited Neighbor Advertisement message upstream after each move.

Upon starting at the New NVE, the VM should send an ARP or Neighbor Discovery message. Cold VM mobility also allows the Old NVE and all communicating NVEs to time out ARP/neighbor cache entries of the VM. It is necessary for the NVA to push the updated ARP/neighbor cache entry to NVEs or for NVEs to pull the updated ARP/neighbor cache entry from NVA.

8. Other Options

VM Hot mobility is to enable uninterrupted running of the application or workload instantiated on the VM when the VM running

conditions changes, such as utilization overload, hardware running condition changes, or others.

There is also a Hot Standby option to prevent unexpected failure conditions, where there are VMs in both primary and secondary NVEs. They have identical information and can provide services simultaneously as in load-share mode of operation. If the VM in the primary NVE fails, there is no need to actively move the VM to the secondary NVE because the VM in the secondary NVE already contain identical information. The Hot Standby option is the costliest mechanism, and hence this option is utilized only for mission-critical applications and services. In Hot Standby option, regarding TCP connections, one option is to start with and maintain TCP connections to two different VMs at the same time. The least loaded VM responds first and pickup providing service while the sender (origin) still continues to receive Ack from the heavily loaded (secondary) VM and chooses not to use the service of the secondary responding VM. If the situation (loading condition of the primary responding VM) changes the secondary responding VM may start providing service to the sender (origin).

9. VM Lifecycle Management

The VM lifecycle management is a complicated task, which is beyond the scope of this document. Not only it involves monitoring server utilization, balanced distribution of workload, etc., but also needs to manage seamlessly VM migration from one server to another.

10. Security Considerations

Security threats for the data and control plane for overlay networks are discussed in [[RFC8014](#)]. ARP (IPv4) and ND (IPv6) are not secure, especially if we accept gratuitous versions in multi-tenant environment.

In Layer-3 based overlay data center networks, the problem of address spoofing may arise. An NVE may have untrusted VMs attached. This usually happens in cases like the VMs running third party applications. Those untrusted VMs can send falsified ARP (IPv4) and ND (IPv6) messages, causing NVE, NV03 Gateway, and NVA to be overwhelmed and not able to perform legitimate functions. The attacker can intercept, modify, or even stop data in-transit ARP/ND messages intended for other VNs and initiate DDOS attacks to other VMs attached to the same NVE.

The locator-identifier mechanism given as an example (ILA) doesn't include secure binding. It doesn't discuss how to securely bind the new locator to the identifier.

This requires VM management system to apply stronger security mechanisms when add a VM to an NVE. VM Management system is out of scope of this document.

11. IANA Considerations

This document makes no request to IANA.

12. Acknowledgments

The authors are grateful to Bob Briscoe, David Black, Dave R. Worley, Qiang Zu, Andrew Malis for helpful comments.

13. Change Log

- . submitted version -00 as a working group draft after adoption
- . submitted version -01 with these changes: references are updated,
 - o added packets in flight definition to [Section 2](#)
- . submitted version -02 with updated address.
- . submitted version -03 to fix the nits.
- . submitted version -04 in reference to the WG Last call comments.
- . Submitted version - 05, 06, 07, and 08 to address IETF LC comments from TSV area.

14. References

14.1. Normative References

- [RFC0826] Plummer, D., "An Ethernet Address Resolution Protocol: Or Converting Network Protocol Addresses to 48.bit Ethernet Address for Transmission on Ethernet Hardware", STD 37, [RFC 826](#), DOI 10.17487/RFC0826, November 1982, <<https://www.rfc-editor.org/info/rfc826>>.
- [RFC0903] Finlayson, R., Mann, T., Mogul, J., and M. Theimer, "A Reverse Address Resolution Protocol", STD 38, [RFC 903](#), DOI 10.17487/RFC0903, June 1984, <<https://www.rfc-editor.org/info/rfc903>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [RFC2629] Rose, M., "Writing I-Ds and RFCs using XML", [RFC 2629](#), DOI 10.17487/RFC2629, June 1999, <<https://www.rfc-editor.org/info/rfc2629>>.
- [RFC4861] Narten, T., Nordmark, E., Simpson, W., and H. Soliman, "Neighbor Discovery for IP version 6 (IPv6)", [RFC 4861](#), DOI 10.17487/RFC4861, September 2007, <<https://www.rfc-editor.org/info/rfc4861>>.
- [RFC7067] L. Dunbar, D. Eastlake, R. Perlman, I. Gashinsky, "directory Assistance Problem and High Level Design Proposal", [RFC7067](#), Nov. 2013
- [RFC7348] Mahalingam, M., Dutt, D., Duda, K., Agarwal, P., Kreeger, L., Sridhar, T., Bursell, M., and C. Wright, "Virtual eXtensible Local Area Network (VXLAN): A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks", [RFC 7348](#), DOI 10.17487/RFC7348, August 2014, <<https://www.rfc-editor.org/info/rfc7348>>.

[RFC7364] Narten, T., Ed., Gray, E., Ed., Black, D., Fang, L., Kreeger, L., and M. Napierala, "Problem Statement: Overlays for Network Virtualization", [RFC 7364](#), DOI 10.17487/RFC7364, October 2014, <<https://www.rfc-editor.org/info/rfc7364>>.

[RFC7666] H. Asai, et al, "Management Information Base for Virtual Machines Controlled by a Hypervisor", [RFC7666](#), Oct 2015.

[RFC8014] Black, D., Hudson, J., Kreeger, L., Lasserre, M., and T. Narten, "An Architecture for Data-Center Network Virtualization over Layer 3 (NV03)", [RFC 8014](#), DOI 10.17487/RFC8014, December 2016, <<https://www.rfc-editor.org/info/rfc8014>>.

[RFC8171] D. Eastlake, L. Dunbar, R. Perlman, Y. Li, "Edge Directory Assistance Mechanisms", [RFC 8171](#), June 2017

[14.2. Informative References](#)

[I-D.herbert-intarea-ila] Herbert, T. and P. Lapukhov, "Identifier-locator addressing for IPv6", [draft-herbert-intarea-ila](#) - 04 (work in progress), March 2017.

Authors' Addresses

Linda Dunbar
Futurewei
Email: ldunbar@futurewei.com

Behcet Sarikaya
Denpel Informatique
Email: sarikaya@ieee.org

Bhumip Khasnabish
Independent
Email: vumip1@gmail.com

Tom Herbert
Intel
Email: tom@herbertland.com

Saumya Dikshit
Aruba-HPE
Bangalore, India
Email: saumya.dikshit@hpe.com