

Internet Engineering Task Force
Internet-Draft
Intended status: Standards Track
Expires: April 20, 2016

S. Pallagatti, Ed.
B. Saji
S. Paragiri
Juniper Networks
V. Govindan
M. Mudigonda
Cisco
G. Mirsky
Ericsson
October 18, 2015

BFD for VXLAN
draft-spallagatti-bfd-vxlan-02

Abstract

This document describes use of Bidirectional Forwarding Detection (BFD) protocol for VXLAN .

Requirements Language

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](#)].

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). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

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."

This Internet-Draft will expire on April 20, 2016.

Copyright Notice

Copyright (c) 2015 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](http://trustee.ietf.org/license-info) 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	2
2.	Use cases	3
3.	Deployment	4
4.	BFD Packet Transmission	5
4.1.	BFD Packet Encapsulation	5
5.	Reception of BFD packet	6
5.1.	Demultiplexing of the BFD packet	6
6.	Use of reserved VNI	6
7.	Echo BFD	6
8.	IANA Considerations	7
9.	Security Considerations	7
10.	Contributors	7
11.	Acknowledgements	7
12.	Normative References	7
	Authors' Addresses	8

[1.](#) Introduction

"Virtual eXtensible Local Area Network (VXLAN)" has been defined in [[RFC7348](#)] that provides an encapsulation scheme which allows VM's to communicate in data center network.

VXLAN is typically deployed in data centers interconnecting virtualized hosts, which may be spread across multiple racks. The individual racks may be part of a different Layer 3 network or they could be in a single Layer 2 network. The VXLAN segments/overlay networks are overlaid on top of these Layer 2 or Layer 3 networks.

A virtual machine (VM) can communicate with a VM in other host only if they are on same VXLAN. VM's are unaware of VXLAN tunnels as VXLAN tunnel is terminated on VXLAN Tunnel End Point(VTEP) (hypervisor/TOR). VTEPs (hypervisor/TOR) are responsible for encapsulating and decapsulating frames exchanged among VM's.

Since underlay is a L3 network, continuity check for these tunnels becomes important. BFD as defined in [[RFC5880](#)] can be used to

monitor the VXLAN tunnels. Use of [[I-D.ietf-bfd-multipoint](#)] is for future study.

This draft addresses requirements outlined in [[I-D.ashwood-nvo3-operational-requirement](#)]. Specifically with reference to the OAM model to Figure 3 of [[I-D.ashwood-nvo3-operational-requirement](#)], this draft outlines proposal to implement the OAM mechanism between the NV Edges using BFD.

2. Use cases

Main use case of BFD for VXLAN is for tunnel continuity check. BFD packets between VTEPs will exercise the VXLAN path in underlay/overlay ensuring the VXLAN path reachability. BFD failure detection can be used for maintenance. There are other use cases such as

Layer 2 VM's:

Most deployments will have VM's with only L2 capabilities that may not support L3. BFD being a L3 protocol can be used as tunnel CC mechanism, where BFD will start and terminate at the Network Virtualization (NV) Edge (VTEPs).

It is possible to aggregate the CC sessions for multiple tenants by running a BFD session between the VTEPs over VxLAN tunnel. In rest of this document terms NV Edge and VTEP are used interchangeably.

Fault localization:

It is also possible that VM's are L3 aware and can possibly host a BFD session. In these cases BFD sessions can be established among VM's for CC. In addition BFD sessions can be established among VTEPs for tunnel CC. Having a hierarchical OAM model helps localize faults though requires additional consideration.

Service node reachability:

Service node is responsible for sending BUM traffic. In case of service node tunnel terminates at VTEP and it might not even host VM's. BFD session between TOR/hypervisor and service node can be used to monitor service node reachability.

3. Deployment

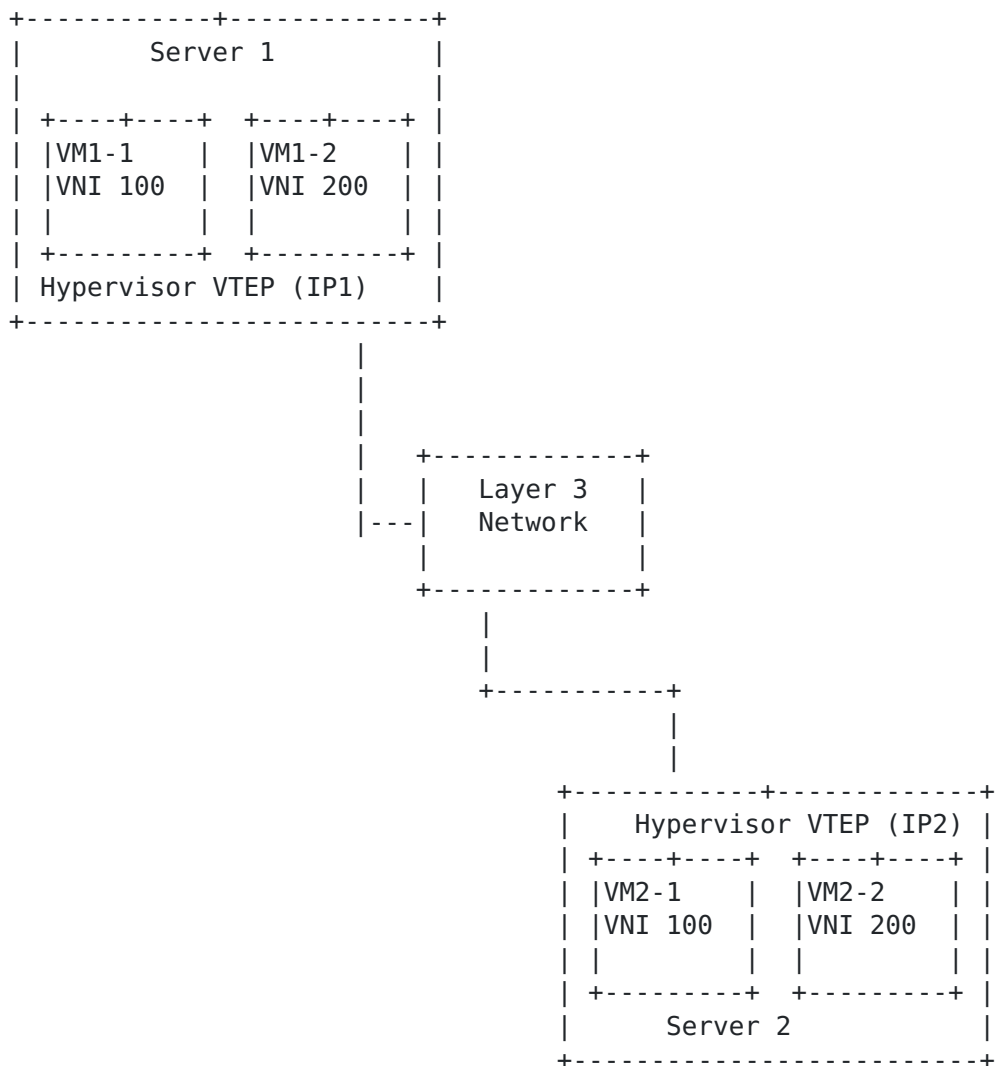


Figure 1

Figure 1 illustrates the scenario where we have two servers, each of them hosting two VMs. These VTEPs terminate two VXLAN tunnels with VNI number 100 and 200 between them. Separate BFD sessions can be established between the VTEPs (IP1 and IP2) for monitoring each of the VXLAN tunnels (VNI 100 and 200). No BFD packets intended to Hypervisor VTEP should be forwarded to a VM as VM may drop BFD packets leading to false negative. This method is applicable whether VTEP is a software or a physical device.

4. BFD Packet Transmission

BFD packet MUST be encapsulated and sent to remote VTEP as explained in [Section 4.1](#). Implementations SHOULD ensure that the BFD packets follow the same lookup path of VXLAN packets within the sender system.

4.1. BFD Packet Encapsulation

VXLAN packet format has been defined in [Section 5 of \[RFC7348\]](#). The Outer IP/UDP and VXLAN headers MUST be encoded by the sender as per [\[RFC7348\]](#).

If VTEP is equipped with Generic Protocol Extension (GPE) header capabilities and decides to use GPE instead of VXLAN then GPE header MUST be encoded as per Section 3.3 of [\[I-D.quinn-vxlan-gpe\]](#). Next Protocol Field in GPE header MUST be set to IPv4 or IPv6.

Details of how VTEP decides to use VXLAN or GPE header are outside the scope of this document.

The BFD packet MUST be carried inside the inner MAC frame of the VxLAN packet. The inner MAC frame carrying the BFD payload has the following format:

Ethernet Header:

Destination MAC: This MUST be a well-known MAC [TBD] OR the MAC address of the destination VTEP. The details of how the destination MAC address is obtained are outside the scope of this document.

Source MAC: MAC address of the originating VTEP

IP header:

Source IP: IP address of the originating VTEP.

Destination IP: IP address of the terminating VTEP.

TTL: This MUST be set to 1. This is to ensure that the BFD packet is not routed within the L3 underlay network.

[Ed.Note]:Use of inner source and destination IP addresses needs more discussion by the WG.

The fields of the UDP header and the BFD control packet are encoded as specified in [RFC 5881](#) for p2p VXLAN tunnels.

5. Reception of BFD packet

Once a packet is received, VTEP MUST validate the packet as described in [Section 4.1 of \[RFC7348\]](#). If the Destination MAC of the inner MAC frame matches the well-known MAC or the MAC address of the VTEP the packet MUST be processed further.

The UDP destination port and the TTL of the inner MAC frame MUST be validated to determine if the received packet can be processed by BFD. BFD packet with inner MAC set to VTEP or well-known MAC address MUST not be forwarded to VM's.

To ensure BFD detects the proper configuration of VXLAN Network Identifier(VNI) in a remote VTEP, a lookup SHOULD be performed with the MAC-DA and VNI as key in the Virtual Forwarding Instance(VFI) table of the originating/ terminating VTEP in order to exercise the VFI associated with the VNI.

5.1. Demultiplexing of the BFD packet

Demultiplexing of IP BFD packet has been defined in [Section 3 of \[RFC5881\]](#). Since multiple BFD sessions may be running between two VTEPs, there needs to be a mechanism for demultiplexing received BFD packets to the proper session. The procedure for demultiplexing packets with Your Discriminator = 0 is different from [\[RFC5880\]](#). For such packets, the BFD session MUST be identified using the inner headers, i.e. the source IP and the destination IP present in the IP header carried by the payload of the VXLAN encapsulated packet. The VNI of the packet SHOULD be used to derive interface related information for demultiplexing the packet. If BFD packet is received with non-zero your discriminator then BFD session should be demultiplexed only with your discriminator as the key.

6. Use of reserved VNI

BFD session MAY be established for the reserved VNI 0. One way to aggregate BFD sessions between VTEP's is to establish a BFD session with VNI 0. A VTEP MAY also use VNI 0 to establish a BFD session with a service node.

7. Echo BFD

Support for echo BFD is outside the scope of this document.

8. IANA Considerations

The well-known MAC to be used for the Destination MAC address of the inner MAC frame needs to be defined

9. Security Considerations

Document recommends setting of inner IP TTL to 1 which could lead to DDoS attack, implementation MUST have throttling in place. Throttling MAY be relaxed for BFD packeted based on port number.

Other than inner IP TTL set to 1 this specification does not raise any additional security issues beyond those of the specifications referred to in the list of normative references.

10. Contributors

Reshad Rahman
rrahman@cisco.com
Cisco

11. Acknowledgements

Authors would like to thank Jeff Hass of Juniper Networks for his reviews and feedback on this material.

Authors would also like to thank Nobo Akiya, Marc Binderberger and Shahram Davari for the extensive review.

12. Normative References

- [I-D.ashwood-nvo3-operational-requirement]
Ashwood-Smith, P., Iyengar, R., Tsou, T., Sajassi, A., Boucadair, M., Jacquenet, C., and M. Daikoku, "NV03 Operational Requirements", [draft-ashwood-nvo3-operational-requirement-03](#) (work in progress), July 2013.
- [I-D.ietf-bfd-multipoint]
Katz, D., Ward, D., and J. Networks, "BFD for Multipoint Networks", [draft-ietf-bfd-multipoint-07](#) (work in progress), August 2015.
- [I-D.ietf-bfd-seamless-base]
Akiya, N., Pignataro, C., Ward, D., Bhatia, M., and J. Networks, "Seamless Bidirectional Forwarding Detection (S-BFD)", [draft-ietf-bfd-seamless-base-05](#) (work in progress), June 2015.

[I-D.quinn-vxlan-gpe]

Quinn, P., Manur, R., Kreeger, L., Lewis, D., Maino, F., Smith, M., Agarwal, P., Yong, L., Xu, X., Elzur, U., Garg, P., and D. Melman, "Generic Protocol Extension for VXLAN", [draft-quinn-vxlan-gpe-04](#) (work in progress), February 2015.

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), DOI 10.17487/RFC2119, March 1997, <<http://www.rfc-editor.org/info/rfc2119>>.

[RFC5880] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD)", [RFC 5880](#), DOI 10.17487/RFC5880, June 2010, <<http://www.rfc-editor.org/info/rfc5880>>.

[RFC5881] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD) for IPv4 and IPv6 (Single Hop)", [RFC 5881](#), DOI 10.17487/RFC5881, June 2010, <<http://www.rfc-editor.org/info/rfc5881>>.

[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, <<http://www.rfc-editor.org/info/rfc7348>>.

Authors' Addresses

Santosh Pallagatti (editor)
Juniper Networks
Embassy Business Park
Bangalore, KA 560093
India

Email: santoshpk@juniper.net

Basil Saji
Juniper Networks
Embassy Business Park
Bangalore, KA 560093
India

Email: sbasil@juniper.net

Sudarsan Paragiri
Juniper Networks
1194 N. Mathilda Ave.
Sunnyvale, California 94089-1206
USA

Email: sparagiri@juniper.net

Vengada Prasad Govindan
Cisco

Email: venggovi@cisco.com

Mallik Mudigonda
Cisco

Email: mmudigon@cisco.com

Greg Mirsky
Ericsson

Email: gregory.mirsky@ericsson.com