

BFD
Internet-Draft
Intended status: Informational
Expires: April 29, 2021

S. Pallagatti, Ed.
VMware
G. Mirsky, Ed.
ZTE Corp.
S. Paragiri
Individual Contributor
V. Govindan
M. Mudigonda
Cisco
October 26, 2020

BFD for VXLAN
draft-ietf-bfd-vxlan-16

Abstract

This document describes the use of the Bidirectional Forwarding Detection (BFD) protocol in point-to-point Virtual eXtensible Local Area Network (VXLAN) tunnels used to form an overlay network.

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 <https://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 29, 2021.

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 (<https://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. Conventions Used in this Document	3
2.1. Acronyms	3
2.2. Requirements Language	4
3. Deployment	4
4. Use of the Management VNI	5
5. BFD Packet Transmission over VXLAN Tunnel	6
6. Reception of BFD Packet from VXLAN Tunnel	8
7. Echo BFD	8
8. IANA Considerations	8
9. Security Considerations	9
10. Contributors	9
11. Acknowledgments	9
12. References	10
12.1. Normative References	10
12.2. Informational References	10
Authors' Addresses	11

1. Introduction

"Virtual eXtensible Local Area Network" (VXLAN) [RFC7348] provides an encapsulation scheme that allows building an overlay network by decoupling the address space of the attached virtual hosts from that of the network.

One use of VXLAN is in data centers interconnecting virtual machines (VMs) of a tenant. VXLAN addresses requirements of the Layer 2 and Layer 3 data center network infrastructure in the presence of VMs in a multi-tenant environment by providing a Layer 2 overlay scheme on a Layer 3 network [RFC7348]. Another use is as an encapsulation for Ethernet VPN [RFC8365].

This document is written assuming the use of VXLAN for virtualized hosts and refers to VMs and VXLAN Tunnel End Points (VTEPs) in hypervisors. However, the concepts are equally applicable to non-virtualized hosts attached to VTEPs in switches.

In the absence of a router in the overlay, a VM can communicate with another VM only if they are on the same VXLAN segment. VMs are unaware of VXLAN tunnels as a VXLAN tunnel is terminated on a VTEP.

VTEPs are responsible for encapsulating and decapsulating frames exchanged among VMs.

The ability to monitor path continuity, i.e., perform proactive continuity check (CC) for point-to-point (p2p) VXLAN tunnels, is important. The asynchronous mode of BFD, as defined in [RFC5880], is used to monitor a p2p VXLAN tunnel.

In the case where a Multicast Service Node (MSN) (as described in Section 3.3 of [RFC8293]) participates in VXLAN, the mechanisms described in this document apply and can, therefore, be used to test the continuity of the path between the source NVE and the MSN.

This document describes the use of Bidirectional Forwarding Detection (BFD) protocol to enable monitoring continuity of the path between VXLAN VTEPs that are performing as Network Virtualization Endpoints, and/or between the source NVE and a replicator MSN using a Management VNI (Section 4). All other uses of the specification to test toward other VXLAN endpoints are out of the scope.

2. Conventions Used in this Document

2.1. Acronyms

BFD Bidirectional Forwarding Detection

CC Continuity Check

p2p Point-to-point

MSN Multicast Service Node

NVE Network Virtualization Endpoint

VFI Virtual Forwarding Instance

VM Virtual Machine

VNI VXLAN Network Identifier (or VXLAN Segment ID)

VTEP VXLAN Tunnel End Point

VXLAN Virtual eXtensible Local Area Network

2.2. Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in BCP 14 [RFC2119] [RFC8174] when, and only when, they appear in all capitals, as shown here.

3. Deployment

Figure 1 illustrates the scenario with two servers, each of them hosting two VMs. The servers host VTEPs that terminate two VXLAN tunnels with VXLAN Network Identifier (VNI) number 100 and 200 respectively. Separate BFD sessions can be established between the VTEPs (IP1 and IP2) for monitoring each of the VXLAN tunnels (VNI 100 and 200). Using a BFD session to monitor a set of VXLAN VNIs between the same pair of VTEPs might help to detect and localize problems caused by misconfiguration. An implementation that supports this specification MUST be able to control the number of BFD sessions that can be created between the same pair of VTEPs. This method is applicable whether the VTEP is a virtual or physical device.

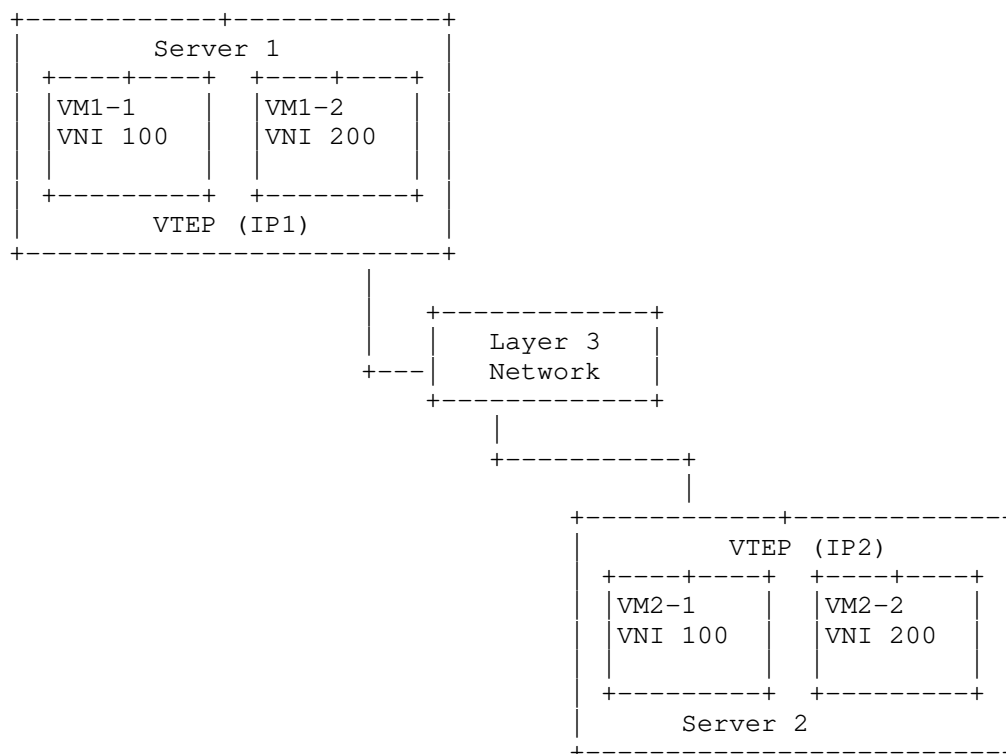


Figure 1: Reference VXLAN Domain

At the same time, a service layer BFD session may be used between the tenants of VTEPs IP1 and IP2 to provide end-to-end fault management (this use case is outside the scope of this document). In such a case, for VTEPs, the BFD Control packets of that session are indistinguishable from data packets.

For BFD Control packets encapsulated in VXLAN (Figure 2), the inner destination IP address SHOULD be set to one of the loopback addresses from 127/8 range for IPv4 or to one of IPv4-mapped IPv6 loopback addresses from `::ffff:127.0.0.0/104` range for IPv6.

4. Use of the Management VNI

In most cases, a single BFD session is sufficient for the given VTEP to monitor the reachability of a remote VTEP, regardless of the number of VNIs. BFD control messages MUST be sent using the Management VNI which acts as the as control and management channel between VTEPs. An implementation MAY support operating BFD on

another (non-Management) VNI although the implications of this are outside the scope of this document. The selection of the VNI number of the Management VNI MUST be controlled through a management plane. An implementation MAY use VNI number 1 as the default value for the Management VNI. All VXLAN packets received on the Management VNI MUST be processed locally and MUST NOT be forwarded to a tenant.

5. BFD Packet Transmission over VXLAN Tunnel

BFD packets MUST be encapsulated and sent to a remote VTEP as explained in this section. Implementations SHOULD ensure that the BFD packets follow the same forwarding path as VXLAN data packets within the sender system.

BFD packets are encapsulated in VXLAN as described below. The VXLAN packet format is defined in Section 5 of [RFC7348]. The value in the VNI field of the VXLAN header MUST be set to the value selected as the Management VNI. The Outer IP/UDP and VXLAN headers MUST be encoded by the sender as defined in [RFC7348].

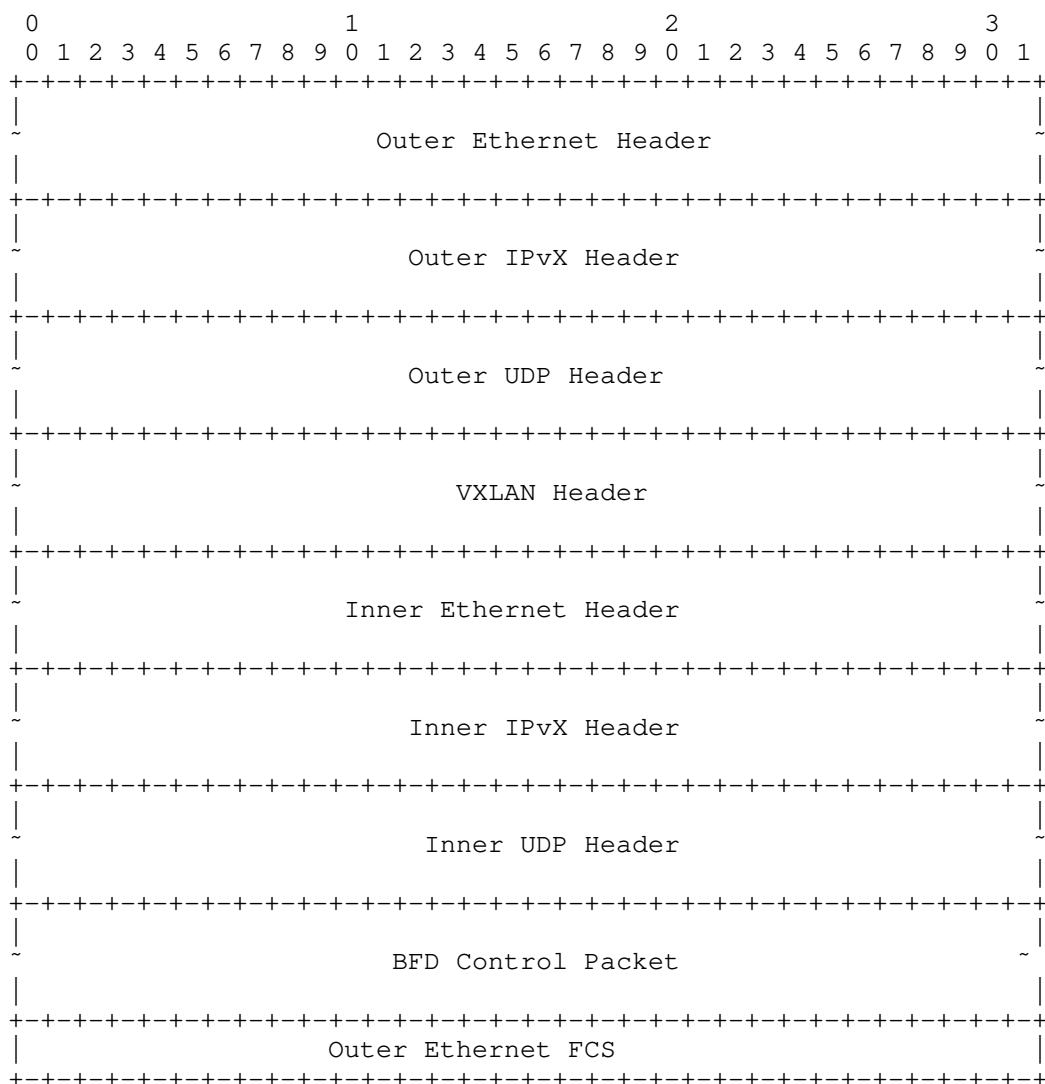


Figure 2: VXLAN Encapsulation of BFD Control Packet

The BFD packet MUST be carried inside the inner Ethernet frame of the VXLAN packet. The choice of Destination MAC and Destination IP addresses for the inner Ethernet frame MUST ensure that the BFD Control packet is not forwarded to a tenant but is processed locally at the remote VTEP. The inner Ethernet frame carrying the BFD Control packet- has the following format:

Ethernet Header:

Destination MAC: A Management VNI, which does not have any tenants, will have no dedicated MAC address for decapsulated traffic. The value (TBD1) SHOULD be used in this field.

Source MAC: MAC address associated with the originating VTEP.

Ethertype: is set to 0x0800 if the inner IP header is IPv4, and is set to 0x86DD if the inner IP header is IPv6.

IP header:

Destination IP: IP address MUST NOT be of one of tenant's IP addresses. The IP address SHOULD be selected from the range 127/8 for IPv4, for IPv6 - from the range ::ffff:127.0.0.0/104. Alternatively, the destination IP address MAY be set to VTEP's IP address.

Source IP: IP address of the originating VTEP.

TTL or Hop Limit: MUST be set to 255 in accordance with [RFC5881].

The fields of the UDP header and the BFD Control packet are encoded as specified in [RFC5881].

6. Reception of BFD Packet from VXLAN Tunnel

Once a packet is received, the VTEP MUST validate the packet. If the packet is received on the management VNI and is identified as BFD control packet addressed to the VTEP, and then the packet can be processed further. Processing of BFD control packets received on non-management VNI is outside the scope of this specification.

The received packet's inner IP payload is then validated according to Sections 4 and 5 in [RFC5881].

7. Echo BFD

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

8. IANA Considerations

IANA is requested to assign a single MAC address to the value TBD1 from the "IANA Unicast 48-bit MAC Address" registry from the "Unassigned (small allocations)" block. The Usage field will be "BFD for VXLAN" with a Reference field of this document.

9. Security Considerations

Security issues discussed in [RFC5880], [RFC5881], and [RFC7348] apply to this document.

This document recommends using an address from the Internal host loopback addresses 127/8 range for IPv4 or an IP4-mapped IPv6 loopback address from ::ffff:127.0.0.0/104 range for IPv6 as the destination IP address in the inner IP header. Using such an address prevents the forwarding of the encapsulated BFD control message by a transient node in case the VXLAN tunnel is broken as according to [RFC1812].

A router SHOULD NOT forward, except over a loopback interface, any packet that has a destination address on network 127. A router MAY have a switch that allows the network manager to disable these checks. If such a switch is provided, it MUST default to performing the checks.

The use of IPv4-mapped IPv6 addresses has the same property as using the IPv4 network 127/8, moreover, the IPv4-mapped IPv6 addresses prefix is not advertised in any routing protocol.

If the implementation supports establishing multiple BFD sessions between the same pair of VTEPs, there SHOULD be a mechanism to control the maximum number of such sessions that can be active at the same time.

10. Contributors

Reshad Rahman
rrahman@cisco.com
Cisco

11. Acknowledgments

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

Authors would also like to thank Nobo Akiya, Marc Binderberger, Shahram Davari, Donald E. Eastlake 3rd, Anoop Ghanwani, Dinesh Dutt, Joel Halpern, and Carlos Pignataro for the extensive reviews and the most detailed and constructive comments.

12. References

12.1. Normative References

- [RFC1812] Baker, F., Ed., "Requirements for IP Version 4 Routers", RFC 1812, DOI 10.17487/RFC1812, June 1995, <<https://www.rfc-editor.org/info/rfc1812>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC5880] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD)", RFC 5880, DOI 10.17487/RFC5880, June 2010, <<https://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, <<https://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, <<https://www.rfc-editor.org/info/rfc7348>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in RFC 2119 Key Words", BCP 14, RFC 8174, DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.

12.2. Informational References

- [RFC8293] Ghanwani, A., Dunbar, L., McBride, M., Bannai, V., and R. Krishnan, "A Framework for Multicast in Network Virtualization over Layer 3", RFC 8293, DOI 10.17487/RFC8293, January 2018, <<https://www.rfc-editor.org/info/rfc8293>>.
- [RFC8365] Sajassi, A., Ed., Drake, J., Ed., Bitar, N., Shekhar, R., Uttaro, J., and W. Henderickx, "A Network Virtualization Overlay Solution Using Ethernet VPN (EVPN)", RFC 8365, DOI 10.17487/RFC8365, March 2018, <<https://www.rfc-editor.org/info/rfc8365>>.

Authors' Addresses

Santosh Pallagatti (editor)
VMware

Email: santosh.pallagatti@gmail.com

Greg Mirsky (editor)
ZTE Corp.

Email: gregimirsky@gmail.com

Sudarsan Paragiri
Individual Contributor

Email: sudarsan.225@gmail.com

Vengada Prasad Govindan
Cisco

Email: venggovi@cisco.com

Mallik Mudigonda
Cisco

Email: mmudigon@cisco.com