

Workgroup: Network Working Group
Internet-Draft: draft-ietf-nvo3-vxlan-gpe-13
Published: 5 November 2023
Intended Status: Informational
Expires: 8 May 2024
Authors: F. Maino, Ed. L. Kreeger, Ed. U. Elzur, Ed.
 Cisco Systems Arrcus Intel
 Generic Protocol Extension for VXLAN (VXLAN-GPE)

Abstract

This document describes extending Virtual eXtensible Local Area Network (VXLAN), via changes to the VXLAN header, with four new capabilities: support for multi-protocol encapsulation, support for operations, administration and maintenance (OAM) signaling, support for ingress-replicated BUM Traffic (i.e. Broadcast, Unknown unicast, or Multicast), and explicit versioning. New protocol capabilities can be introduced via shim headers.

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 [[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 <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 8 May 2024.

Copyright Notice

Copyright (c) 2023 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 Revised BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Revised BSD License.

Table of Contents

- [1. Introduction](#)
- [2. VXLAN Without Protocol Extension](#)
- [3. Generic Protocol Extension for VXLAN \(VXLAN-GPE\)](#)
 - [3.1. VXLAN-GPE Header](#)
 - [3.2. Multi Protocol Support](#)
 - [3.3. Replicated BUM Traffic](#)
 - [3.4. OAM Support](#)
 - [3.5. Version Bits](#)
- [4. Outer Encapsulations](#)
 - [4.1. Inner VLAN Tag Handling](#)
 - [4.2. Fragmentation Considerations](#)
- [5. Implementation and Deployment Considerations](#)
 - [5.1. Applicability Statement](#)
 - [5.2. Congestion Control Functionality](#)
 - [5.3. UDP Checksum](#)
 - [5.3.1. UDP Zero Checksum Handling with IPv6](#)
 - [5.4. DSCP, ECN, TTL, and 802.1Q](#)
- [6. Backward Compatibility](#)
 - [6.1. VXLAN VTEP to VXLAN-GPE VTEP](#)
 - [6.2. VXLAN-GPE VTEP to VXLAN VTEP](#)
 - [6.3. VXLAN-GPE UDP Ports](#)
 - [6.4. VXLAN-GPE and Encapsulated IP Header Fields](#)
- [7. VXLAN-GPE Examples](#)
- [8. Security Considerations](#)
- [9. Contributors](#)
- [10. Acknowledgments](#)
- [11. IANA Considerations](#)
 - [11.1. UDP Port](#)
 - [11.2. VXLAN-GPE Next Protocol](#)
 - [11.3. VXLAN-GPE Flag and Reserved Bits](#)
- [12. References](#)
 - [12.1. Normative References](#)
 - [12.2. Informative References](#)
- [Authors' Addresses](#)

1. Introduction

Virtual eXtensible Local Area Network VXLAN [[RFC7348](#)] defines an encapsulation format that encapsulates Ethernet frames in an outer

UDP/IP transport. As data centers evolve, the need to carry other protocols encapsulated in an IP packet is required, as well as the need to provide increased visibility and diagnostic capabilities within the overlay. The VXLAN header does not specify the protocol being encapsulated and therefore is currently limited to encapsulating only Ethernet frame payload, nor does it provide the ability to define OAM protocols. In addition, [RFC6335] requires that new transports not use transport layer port numbers to identify tunnel payload, rather it encourages encapsulations to use their own identifiers for this purpose. VXLAN-GPE is intended to extend the existing VXLAN protocol to provide protocol typing, OAM, and versioning capabilities.

The Version and OAM bits are introduced in [Section 3](#), and the choice of location for these fields is driven by minimizing the impact on existing deployed hardware.

In order to facilitate deployments of VXLAN-GPE with hardware currently deployed to support VXLAN, changes from legacy VXLAN have been kept to a minimum. [Section 6](#) provides a detailed discussion about how VXLAN-GPE addresses the requirement for backward compatibility with VXLAN.

The capabilities of the VXLAN-GPE protocol can be extended by defining next protocol "shim" headers that are used to implement new data plane functions. As an example In-situ Operations, Administration, and Maintenance (IOAM) metadata functionalities can be added as specified in [[I-D.brockners-ippm-ioam-vxlan-gpe](#)].

2. VXLAN Without Protocol Extension

As described in [Section 1](#), the VXLAN header has no protocol identifier that indicates the type of payload being carried. Because of this, VXLAN is limited to carrying ethernet payloads.

The VXLAN header [RFC7348] contains a single flag 'I' that, when set, indicates the presence of a VXLAN Network Identifier (VNI).

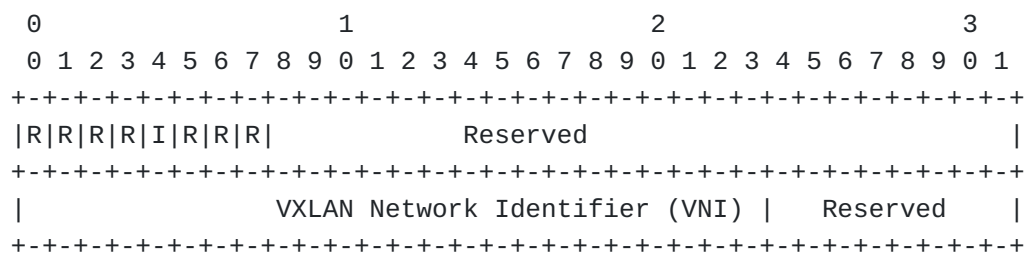


Figure 1: VXLAN Header

3. Generic Protocol Extension for VXLAN (VXLAN-GPE)

3.1. VXLAN-GPE Header

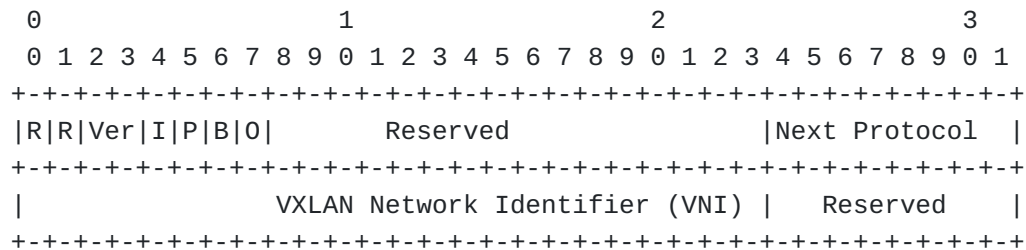


Figure 2: VXLAN-GPE Header

Flags (8 bits): The first 8 bits of the header are the flag field. The bits designated "R" above are reserved flags. These MUST be set to zero on transmission and ignored on receipt.

Version (Ver): Indicates VXLAN-GPE protocol version. The initial version is 0. If a receiver does not support the version indicated it MUST drop the packet.

Instance Bit (I bit): The I bit MUST be set to indicate a valid VNI.

Next Protocol Bit (P bit): The P bit is set to indicate that the Next Protocol field is present.

BUM Traffic Bit (B bit): The B bit is set to indicate that this is ingress-replicated BUM Traffic (ie, Broadcast, Unknown unicast, or Multicast).

OAM Flag Bit (O bit): The O bit is set to indicate that the packet is an OAM packet.

Next Protocol: This 8 bit field indicates the protocol header immediately following the VXLAN-GPE header.

VNI: This 24 bit field identifies the VXLAN overlay network the inner packet belongs to. Inner packets belonging to different VNIs cannot communicate with each other (unless explicitly allowed by policy).

Reserved: Reserved fields MUST be set to zero on transmission and ignored on receipt.

3.2. Multi Protocol Support

This draft defines the following two changes to the VXLAN header in order to support multi-protocol encapsulation:

P Bit: Flag bit 5 is defined as the Next Protocol bit. The P bit MUST be set to 1 to indicate the presence of the 8 bit next protocol field.

When UDP dest port=4790, P = 0 the "Next Protocol" field must be set to zero and the payload MUST be ETHERNET(L2) as defined by [\[RFC7348\]](#).

Flag bit 5 was chosen as the P bit because this flag bit is currently reserved in VXLAN.

Next Protocol Field: When the P-bit is set to 1, the lower 8 bits of the first word are used to carry a Next Protocol. This 'Next Protocol' field contains the protocol of the encapsulated payload packet. Values are tracked by the IANA "LISP-GPE Next Protocol" registry, defined in Section 6.1 of [\[RFC9305\]](#), that VXLAN-GPE shares with the LISP-GPE protocol. At the time this document is edited the IANA "LISP-GPE Next Protocol" registry specifies the following Next Protocol values:

0x00 : Reserved

0x01 : IPv4

0x02 : IPv6

0x03 : Ethernet

0x04 : Network Service Header [\[RFC8300\]](#)

0x05 to 0x7D: Unassigned

0x7E, 0x7F: Experimentation and testing

0x80 to 0xFD: Unassigned (shim headers)

0xFE, 0xFF: Experimentation and testing (shim headers)

Next protocol values 0x7E, 0x7F and 0xFE, 0xFF are assigned for experimentation and testing as per [\[RFC3692\]](#).

Next protocol values from 0x80 to 0xFD are assigned to protocols encoded as generic shim headers. All shim protocols MUST use the header structure in [Figure 3](#), which includes a 'Next Protocol' field. When shim headers are used with other protocols identified by

Next Protocol values from 0x0 to 0x7F, all the shim headers MUST come first.

Shim headers can be used to incrementally deploy new GPE features, keeping the processing of shim headers known to a given VTEP implementation in the 'fast' path (typically an Application-Specific Integrated Circuit (ASIC)) while punting the processing of the remaining new GPE features to the 'slow' path.

Shim protocols MUST have the first 32 bits defined as:

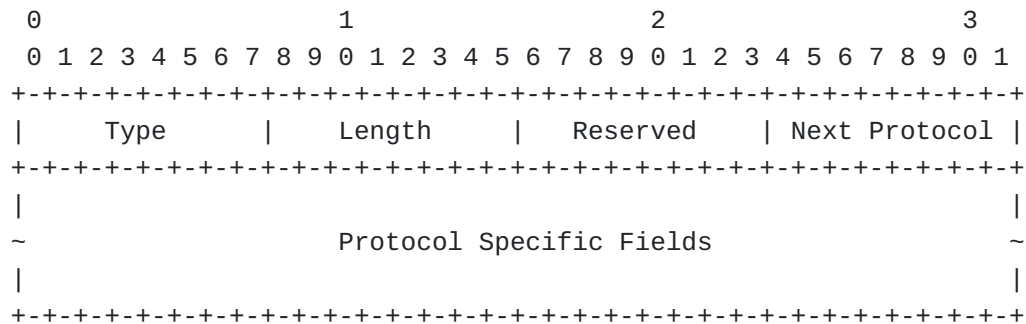


Figure 3: Shim Header

Where:

Type: This field identifies the different messages of this protocol.

Length: This field indicates the length, in in 4-octet units, of this protocol message not including the first 4 octets.

Reserved: The use of this field is reserved to the protocol defined in this message.

Next Protocol Field: This field contains the protocol of the encapsulated payload. The values are tracked in the IANA "LISP-GPE Next Protocol" registry, defined in Section 6.1 of [[RFC9305](#)].

3.3. Replicated BUM Traffic

Flag bit 6 is defined as the B bit. When the B bit is set to 1, the packet is marked as an ingress-replicated BUM Traffic (i.e. Broadcast, Unknown unicast, or Multicast) to help egress VTEP to differentiate between known and unknown unicast. The details of using the B bit are out of scope for this document, but please see [[RFC8365](#)] for an example in the EVPN context. As with the P-bit, bit 6 is currently a reserved flag in VXLAN.

3.4. OAM Support

Flag bit 7 is defined as the O bit. When the O bit is set to 1, the packet is an OAM packet and OAM processing MUST occur. Other header fields including Next Protocol MUST adhere to the definitions in [Section 3](#). The OAM protocol details are out of scope for this document. As with the P-bit, bit 7 is currently a reserved flag in VXLAN.

3.5. Version Bits

VXLAN-GPE bits 2 and 3 are defined as version bits. These bits are reserved in VXLAN. The version field is used to ensure backward compatibility going forward with future VXLAN-GPE updates.

The initial version for VXLAN-GPE is 0.

4. Outer Encapsulations

In addition to the VXLAN-GPE header, the packet is further encapsulated in UDP and IP. Data centers based on Ethernet, will then send this IP packet over Ethernet.

Outer UDP Header:

Destination UDP Port: IANA has assigned the value 4790 for the VXLAN-GPE UDP port. This well-known destination port is used when sending VXLAN-GPE encapsulated packets.

Source UDP Port: The source UDP port is used as entropy for devices forwarding encapsulated packets across the underlay (ECMP for IP routers, or load splitting for link aggregation by bridges). Tenant traffic flows should all use the same source UDP port to lower the chances of packet reordering by the underlay for a given flow. It is recommended for VTEPs to generate this port number using a hash of the inner packet headers. Implementations MAY use the entire 16 bit source UDP port for entropy.

UDP Checksum: see [Section 5.3](#) for considerations related to UDP Checksum processing.

Outer IP Header:

This is the header used by the underlay network to deliver packets between VTEPs. The destination IP address can be a unicast or a multicast IP address. The source IP address must be the source VTEP IP address which can be used to return tenant packets to the tenant system source address within the inner packet header.

When the outer IP header is IPv4, VTEPs MUST set the DF bit.

Outer Ethernet Header:

Most data centers networks are built on Ethernet. Assuming the outer IP packet is being sent across Ethernet, there will be an Ethernet header used to deliver the IP packet to the next hop, which could be the destination VTEP or be a router used to forward the IP packet towards the destination VTEP. If VLANs are in use within the data center, then this Ethernet header would also contain a VLAN tag.

The following figures show the entire stack of protocol headers that would be seen on an Ethernet link carrying encapsulated packets from a VTEP across the underlay network for both IPv4 and IPv6 based underlay networks.


```

0                               1                               2                               3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1

```

Outer Ethernet Header:

```

+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                               Outer Destination MAC Address                               |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Outer Destination MAC Address | Outer Source MAC Address |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                               Outer Source MAC Address                               |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Opt Ethertype = C-Tag 802.1Q | Outer VLAN Tag |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Ethertype = 0x0800 |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

Outer IPv4 Header:

```

+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|Version| IHL |Type of Service| Total Length |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Identification |Flags| Fragment Offset |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Time to Live |Protocol=17(UDP)| Header Checksum |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                               Outer Source IPv4 Address                               |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                               Outer Destination IPv4 Address                               |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

Outer UDP Header:

```

+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Source Port | Dest Port = 4790 |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| UDP Length | UDP Checksum |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

VXLAN-GPE Header:

```

+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|R|R|Ver|I|P|R|O| Reserved |Next Protocol |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| VXLAN Network Identifier (VNI) | Reserved |
+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

Payload:

```

+-+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
| Depends on VXLAN-GPE Next Protocol field above. |

```

```
|    Note that if the payload is Ethernet, then the original    |
|    Ethernet Frame's FCS is not included.                      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
```

Frame Check Sequence:

```
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|    New FCS (Frame Check Sequence) for Outer Ethernet Frame    |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
```

Figure 4: Outer Headers for VXLAN-GPE over IPv4

0	1	2	3
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1			

Outer Ethernet Header:

```

+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|               Outer Destination MAC Address               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Outer Destination MAC Address | Outer Source MAC Address   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|               Outer Source MAC Address                   |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Opt Ethertype = C-Tag 802.1Q | Outer VLAN Tag             |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Ethertype = 0x86DD          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Outer IPv6 Header:

```

+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|Version| Traffic Class |               Flow Label               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|      Payload Length      | NxtHdr=17(UDP)| Hop Limit |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|
+
|
+
|               Outer Source IPv6 Address
+
|
+
|
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|
+
|
+
|               Outer Destination IPv6 Address
+
|
+
|
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Outer UDP Header:

```

+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|      Source Port      |      Dest Port = 4790      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|      UDP Length      |      UDP Checksum          |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

VXLAN-GPE Header:

```

+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|R|R|Ver|I|P|R|O|      Reserved      |Next Protocol |

```


Figure 5: Outer Headers for VXLAN-GPE over IPv6

4.1. Inner VLAN Tag Handling

If the inner packet (as indicated by the VXLAN-GPE Next Protocol field) is an Ethernet frame, it is recommended that it does not contain a VLAN tag. In the most common scenarios, the tenant VLAN tag is translated into a VXLAN Network Identifier. In these scenarios, VTEPs should never send an inner Ethernet frame with a VLAN tag, and a VTEP performing decapsulation should discard any inner frames received with a VLAN tag. However, if the VTEPs are specifically configured to support it for a specific VXLAN Network Identifier, a VTEP may support transparent transport of the inner VLAN tag between all tenant systems on that VNI. The VTEP never looks at the value of the inner VLAN tag, but simply passes it across the underlay.

4.2. Fragmentation Considerations

VTEPs MUST never fragment an encapsulated VXLAN-GPE packet, and when the outer IP header is IPv4, VTEPs MUST set the DF bit in the outer IPv4 header. It is recommended that the underlay network be configured to carry an MTU at least large enough to accommodate the added encapsulation headers. It is recommended that VTEPs perform Path MTU discovery [[RFC1191](#)] [[RFC1981](#)] to determine if the underlay network can carry the encapsulated payload packet.

5. Implementation and Deployment Considerations

5.1. Applicability Statement

VXLAN-GPE conforms, as an UDP-based encapsulation protocol, to the UDP usage guidelines as specified in [[RFC8085](#)]. The applicability of these guidelines are dependent on the underlay IP network and the nature of the encapsulated payload.

[[RFC8085](#)] outlines two applicability scenarios for UDP applications, 1) general Internet and 2) controlled environment. The controlled environment means a single administrative domain or adjacent set of cooperating domains. A network in a controlled environment can be managed to operate under certain conditions whereas, in the general Internet, this cannot be done. Hence requirements for a tunnel protocol operating under a controlled environment can be less restrictive than the requirements of the general Internet.

VXLAN-GPE is intended to be deployed in a data center network environment operated by a single operator or adjacent set of cooperating network operators that fits with the definition of controlled environments in [[RFC8085](#)].

For the purpose of this document, a Traffic-Managed Controlled Environment (TMCE), outlined in [\[RFC8086\]](#), is defined as an IP network that is traffic-engineered and/or otherwise managed (e.g., via use of traffic rate limiters) to avoid congestion. Significant portions of text in this Section are based on [\[RFC8086\]](#).

It is the responsibility of the network operators to ensure that the guidelines/requirements in this section are followed as applicable to their VXLAN-GPE deployments

5.2. Congestion Control Functionality

VXLAN-GPE does not natively provide congestion control functionality and relies on the payload protocol traffic for congestion control. As such VXLAN-GPE MUST be used with congestion controlled traffic or within a network that is traffic managed to avoid congestion (TMCE). An operator of a traffic managed network (TMCE) may avoid congestion by careful provisioning of their networks, rate-limiting of user data traffic and traffic engineering according to path capacity.

Keeping in mind the recommendation above, new encapsulated payloads, when registered, MUST be accompanied by a set of guidelines derived from Section 5 of [\[RFC9300\]](#). Such new protocols should be designed for explicit congestion signals to propagate consistently from lower-layer protocols into IP. Then, the IP internetwork layer can act as a portability layer to carry congestion notifications from non-IP-aware congested nodes up to the transport layer (L4). By following the guidelines in [\[I-D.ietf-tsvwg-ecn-encap-guidelines\]](#), subnetwork designers can enable a Layer 2 protocol to participate in congestion control without dropping packets, via propagation of Explicit Congestion Notification (ECN) data [\[RFC3168\]](#) to receivers.

5.3. UDP Checksum

For IP payloads, Section 5.3 of [\[RFC9300\]](#) specifies how to handle UDP checksums, encouraging implementors to consider UDP checksum usage guidelines in Section 3.4 of [\[RFC8085\]](#) when it is desirable to protect UDP and LISP headers against corruption.

In order to protect the integrity of VXLAN-GPE headers, options, and payload (for example to avoid mis-delivery of payload to different tenant systems in case of data corruption), outer UDP checksum SHOULD be used with VXLAN-GPE when transported over IPv4. The UDP checksum provides a statistical guarantee that a payload was not corrupted in transit. These integrity checks are not strong from a coding or cryptographic perspective and are not designed to detect physical-layer errors or malicious modification of the datagram (see Section 3.4 of [\[RFC8085\]](#)). In deployments where such a risk exists,

an operator SHOULD use additional data integrity mechanisms such as offered by IPSec.

An operator MAY choose to disable UDP checksum and use zero checksum if VXLAN-GPE packet integrity is provided by other data integrity mechanisms, such as IPsec or additional checksums, or if one of the conditions in [Section 5.3.1](#) (a, b, or c) are met.

5.3.1. UDP Zero Checksum Handling with IPv6

By default, UDP checksum MUST be used when VXLAN-GPE is transported over IPv6. A tunnel endpoint MAY be configured for use with a zero UDP checksum if additional requirements described in this section are met.

When VXLAN-GPE is used over IPv6, UDP checksum is used to protect IPv6 headers, UDP headers and VXLAN-GPE headers and payload from potential data corruption. As such by default VXLAN-GPE MUST use UDP checksum when transported over IPv6. An operator MAY choose to configure to operate with zero UDP checksum if operating in a traffic managed controlled environment as stated in [Section 5.1](#) if one of the following conditions are met:

- a. It is known that the packet corruption is exceptionally unlikely (perhaps based on knowledge of equipment types in their underlay network), and the operator is willing to take a risk of undetected packet corruption
- b. It is determined through observational measurements (perhaps through historic or current traffic flows that use non zero checksum) that the level of packet corruption is tolerably low and where the operator is willing to take the risk of undetected corruption
- c. VXLAN-GPE payloads are carrying applications that are tolerant of misdelivered or corrupted packets (perhaps through higher-layer checksum validation and/or reliability through retransmission)

In addition VXLAN-GPE tunnel implementations using Zero UDP checksum MUST meet the following requirements:

1. Use of UDP checksum over IPv6 MUST be the default configuration for all VXLAN-GPE tunnels
2. If VXLAN-GPE is used with zero UDP checksum over IPv6 then such VTEP implementation MUST meet all the requirements specified in section 4 of [\[RFC6936\]](#) and requirements 1 as specified in section 5 of [\[RFC6936\]](#)

3. The VTEP that decapsulates the packet SHOULD check the source and destination IPv6 addresses are valid for the VXLAN-GPE tunnel that is configured to receive Zero UDP checksum and discard other packets for which such check fails
4. The VTEP that encapsulates the packet MAY use different IPv6 source addresses for each VXLAN-GPE tunnel that uses Zero UDP checksum mode in order to strengthen the decapsulator's check of the IPv6 source address (i.e the same IPv6 source address is not to be used with more than one IPv6 destination address, irrespective of whether that destination address is a unicast or multicast address). When this is not possible, it is RECOMMENDED to use each source address for as few VXLAN-GPE tunnels that use zero UDP checksum as is feasible
5. Measures SHOULD be taken to prevent VXLAN-GPE traffic over IPv6 with zero UDP checksum from escaping into the general Internet. Examples of such measures include employing packet filters at the gateways or edge of a VXLAN-GPE network, and/or keeping logical or physical separation of VXLAN network from networks carrying General Internet

The above requirements do not change either the requirements specified in [[RFC2460](#)] as modified by [[RFC6935](#)] or the requirements specified in [[RFC6936](#)].

The requirement to check the source IPv6 address in addition to the destination IPv6 address, plus the recommendation against reuse of source IPv6 addresses among VXLAN-GPE tunnels collectively provide some mitigation for the absence of UDP checksum coverage of the IPv6 header. A traffic-managed controlled environment that satisfies at least one of three conditions listed at the beginning of this section provides additional assurance.

5.4. DSCP, ECN, TTL, and 802.1Q

When encapsulating IP (including over Ethernet) packets, [[RFC2983](#)] provides guidance for mapping packets that contain Differentiated Services Code Point (DSCP) information between inner and outer IP headers. The Pipe model typically fits better with network virtualization. The DSCP value on the tunnel header is set based on a policy (which may be a fixed value, one based on the inner traffic class, or some other mechanism for grouping traffic). Some aspects of the Uniform model (which treats the inner and outer DSCP value as a single field by copying on ingress and egress) may also apply, such as the ability to remark the inner header on tunnel egress based on transit marking. However, the Uniform model is not conceptually consistent with network virtualization, which seeks to

provide strong isolation between encapsulated traffic and the physical network.

[[RFC6040](#)] describes the mechanism for exposing ECN capabilities on IP tunnels and propagating congestion markers to the inner packets. This behavior MUST be followed for IP packets encapsulated in VXLAN-GPE.

Though the Uniform model or the Pipe model could be used for TTL (or Hop Limit in the case of IPv6) handling when tunneling IP packets, the Pipe model is more aligned with network virtualization.

[[RFC2003](#)] provides guidance on handling TTL between inner IP headers and outer IP tunnels; this model is more aligned with the Pipe model and is recommended for use with VXLAN-GPE for network-virtualization applications.

When a VXLAN-GPE VTEP performs Ethernet encapsulation, the inner 802.1Q 3-bit Priority Code Point ('PCP') field [IEEE.802.1Q_2014] MAY be mapped from the encapsulated frame to the DSCP codepoint of the Differentiated Services ('DS') field defined in [[RFC2474](#)].

When a VXLAN-GPE VTEP performs Ethernet encapsulation, the inner-header 802.1Q VLAN Identifier (VID) [IEEE.802.1Q_2014] MAY be mapped to, or used to determine, the 'VXLAN Network Identifier' (VNI) field.

Refer to [Section 8](#) for considerations about the use of integrity protection for deployments, such as the public Internet, concerned with on-path attackers.

6. Backward Compatibility

6.1. VXLAN VTEP to VXLAN-GPE VTEP

A VXLAN VTEP conforms to VXLAN frame format and uses UDP destination port 4789 when sending traffic to VXLAN-GPE VTEP. As per VXLAN, reserved bits 5 and 7, VXLAN-GPE P and O-bits respectively must be set to zero. The remaining reserved bits must be zero, including the VXLAN-GPE version field, bits 2 and 3. The encapsulated payload MUST be Ethernet.

6.2. VXLAN-GPE VTEP to VXLAN VTEP

A VXLAN-GPE VTEP MUST NOT encapsulate non-Ethernet frames to a VXLAN VTEP. When encapsulating Ethernet frames to a VXLAN VTEP, the VXLAN-GPE VTEP MUST conform to VXLAN frame format and hence will set the P bit to 0, the Next Protocol to 0 and use UDP destination port 4789. A VXLAN-GPE VTEP MUST also set O = 0 and Ver = 0 when encapsulating Ethernet frames to VXLAN VTEP. The receiving VXLAN VTEP will treat this packet as a VXLAN packet.

A method for determining the capabilities of a VXLAN VTEP (GPE or non-GPE) is out of the scope of this draft.

6.3. VXLAN-GPE UDP Ports

VXLAN-GPE uses a IANA assigned UDP destination port, 4790, when sending traffic to VXLAN-GPE VTEPs.

6.4. VXLAN-GPE and Encapsulated IP Header Fields

When encapsulating IP (including over Ethernet) packets [[RFC2983](#)] provides guidance for mapping DSCP between inner and outer IP headers. The Pipe model typically fits better Network virtualization. The DSCP value on the tunnel header is set based on a policy (which may be a fixed value, one based on the inner traffic class, or some other mechanism for grouping traffic). Some aspects of the Uniform model (which treats the inner and outer DSCP value as a single field by copying on ingress and egress) may also apply, such as the ability to remark the inner header on tunnel egress based on transit marking. However, the Uniform model is not conceptually consistent with network virtualization, which seeks to provide strong isolation between encapsulated traffic and the physical network.

[[RFC6040](#)] describes the mechanism for exposing ECN capabilities on IP tunnels and propagating congestion markers to the inner packets. This behavior MUST be followed for IP packets encapsulated in VXLAN-GPE.

Though Uniform or Pipe models could be used for TTL (or Hop Limit in case of IPv6) handling when tunneling IP packets, Pipe model is more aligned with network virtualization. [[RFC2003](#)] provides guidance on handling TTL between inner IP header and outer IP tunnels; this model is more aligned with the Pipe model and is recommended for use with VXLAN-GPE for network virtualization applications.

When a VXLAN-GPE router performs Ethernet encapsulation, the inner 802.1Q 3-bit priority code point (PCP) field MAY be mapped from the encapsulated frame to the DSCP codepoint of the DS field defined in [[RFC2474](#)].

When a VXLAN-GPE router performs Ethernet encapsulation, the inner header 802.1Q VLAN Identifier (VID) MAY be mapped to, or used to determine the VXLAN Network Identifier (VNI) field.

7. VXLAN-GPE Examples

This section provides three examples of protocols encapsulated using the Generic Protocol Extension for VXLAN described in this document.

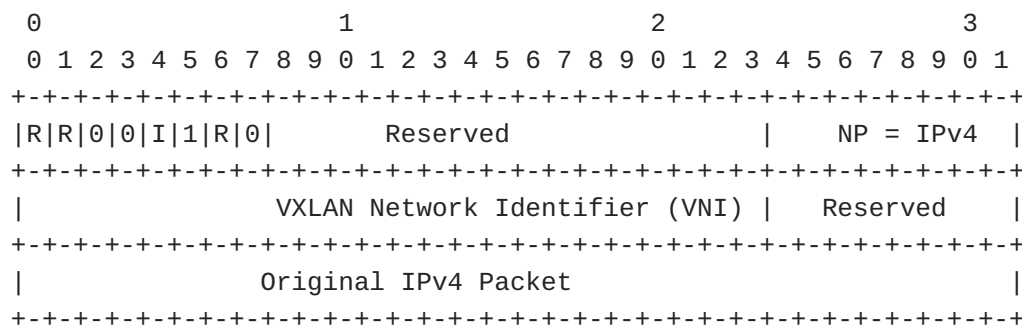


Figure 6: IPv4 and VXLAN-GPE

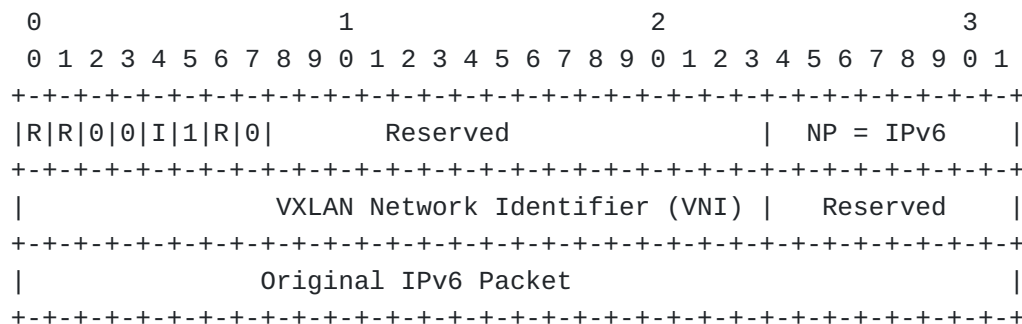


Figure 7: IPv6 and VXLAN-GPE

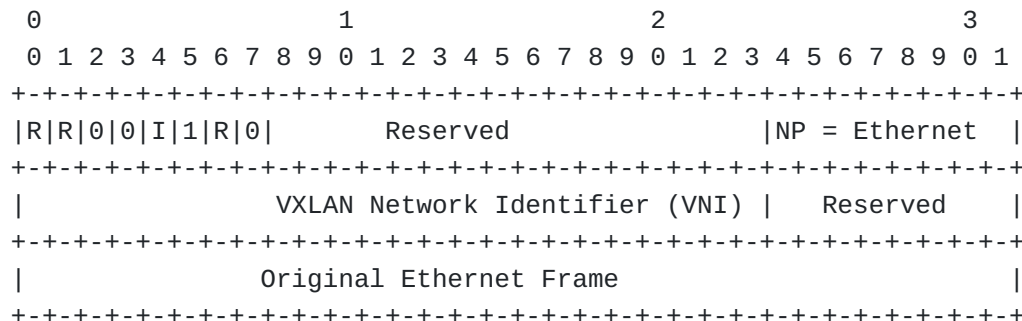


Figure 8: Ethernet and VXLAN-GPE

8. Security Considerations

The security considerations for VXLAN applies to VXLAN-GPE, see [\[RFC7348\]](#).

As is the case for many encapsulations that use optional extensions, VXLAN-GPE is subject to on-path adversaries that can make arbitrary modifications to the packet (including the P-bit) to change or remove any part of the payload, or claim to encapsulate any protocol payload type. Typical integrity protection mechanisms (such as IPsec) SHOULD be used in combination with VXLAN-GPE by those protocol extensions that want to protect against on-path attackers.

With VXLAN-GPE, issues such as data plane spoofing, flooding, and traffic redirection may depend on the particular protocol payload encapsulated.

Operators have to make an assessment based on their network environment and determine the risks that are applicable to their specific environment and use appropriate mitigation approaches, as applicable.

9. Contributors

Paul Quinn
Cisco Systems
paulq@cisco.com

Rajeev Manur
Broadcom
rmanur@broadcom.com

Michael Smith
Cisco Systems
michsmit@cisco.com

Darrel Lewis
Cisco Systems
darlewis@cisco.com

Puneet Agarwal
Innovium, Inc
puneet@acm.org

Lucy Yong
Huawei USA
lucy.yong@huawei.com

Xiaohu Xu
Alibaba Inc
xiaohu.xxh@alibaba-inc.com

Pankaj Garg
Microsoft
pankajg@microsoft.com

David Melman
Marvell
davidme@marvell.com

Jennifer Lemon
Broadcom Limited
jennifer.lemon@broadcom.com

10. Acknowledgments

A special thank you goes to Dino Farinacci for his guidance and detailed review. Thanks to Tom Herbert for the suggestion to assign codepoints for experimentations and testing.

11. IANA Considerations

11.1. UDP Port

UDP 4790 port has been assigned by IANA for VXLAN-GPE in the "Service Name and Transport Protocol Port Number" registry.

11.2. VXLAN-GPE Next Protocol

VXLAN-GPE 'Next Protocol' values are tracked by the IANA "LISP-GPE Next Protocol" registry, defined in Section 6.1 of [[RFC9305](#)] that VXLAN-GPE shares with LISP-GPE. New values are assigned under the Specification Required policy [[RFC8126](#)]. The protocols that are being assigned values do not themselves need to be IETF Standards Track protocols.

At the time this document is edited the IANA "LISP-GPE Next Protocol" registry specifies the following Next Protocol values:

Next Protocol	Description	Reference
0x00	Reserved	RFC 9305
0x01	IPv4	RFC 9305
0x02	IPv6	RFC 9305
0x03	Ethernet	RFC 9305
0x04	NSH	RFC 9305
0x05..0x7D	Unassigned	
0x7E, 0x7F	Experimentation and testing	RFC 9305
0x80..0xFD	Unassigned (shim headers)	
0xFE, 0xFF	Experimentation and testing (shim headers)	RFC 9305

Table 1

11.3. VXLAN-GPE Flag and Reserved Bits

There are ten flag bits at the beginning of the VXLAN-GPE header, followed by 16 reserved bits and an 8-bit reserved field at the end of the header. New bits are assigned via Standards Action [[RFC5226](#)].

Bits 0-1 - Reserve6

Bits 2-3 - Version

Bit 4 - Instance ID (I bit)

Bit 5 - Next Protocol (P bit)

Bit 6 - Reserved

Bit 7 - OAM (O bit)

Bit 8-23 - Reserved

Bits 24-31 in the 2nd Word -- Reserved

Reserved bits/fields MUST be set to 0 by the sender and ignored by the receiver.

12. References

12.1. Normative References

- [RFC1191] Mogul, J. and S. Deering, "Path MTU discovery", RFC 1191, DOI 10.17487/RFC1191, November 1990, <<https://www.rfc-editor.org/info/rfc1191>>.
- [RFC1981] McCann, J., Deering, S., and J. Mogul, "Path MTU Discovery for IP version 6", RFC 1981, DOI 10.17487/RFC1981, August 1996, <<https://www.rfc-editor.org/info/rfc1981>>.
- [RFC2003] Perkins, C., "IP Encapsulation within IP", RFC 2003, DOI 10.17487/RFC2003, October 1996, <<https://www.rfc-editor.org/info/rfc2003>>.
- [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>>.
- [RFC2460] Deering, S. and R. Hinden, "Internet Protocol, Version 6 (IPv6) Specification", RFC 2460, DOI 10.17487/RFC2460, December 1998, <<https://www.rfc-editor.org/info/rfc2460>>.
- [RFC2474] Nichols, K., Blake, S., Baker, F., and D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers", RFC 2474, DOI 10.17487/RFC2474, December 1998, <<https://www.rfc-editor.org/info/rfc2474>>.
- [RFC2983] Black, D., "Differentiated Services and Tunnels", RFC 2983, DOI 10.17487/RFC2983, October 2000, <<https://www.rfc-editor.org/info/rfc2983>>.
- [RFC3168] Ramakrishnan, K., Floyd, S., and D. Black, "The Addition of Explicit Congestion Notification (ECN) to IP", RFC 3168, DOI 10.17487/RFC3168, September 2001, <<https://www.rfc-editor.org/info/rfc3168>>.
- [RFC3692] Narten, T., "Assigning Experimental and Testing Numbers Considered Useful", BCP 82, RFC 3692, DOI 10.17487/RFC3692, January 2004, <<https://www.rfc-editor.org/info/rfc3692>>.

- [RFC5226] Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA Considerations Section in RFCs", RFC 5226, DOI 10.17487/RFC5226, May 2008, <<https://www.rfc-editor.org/info/rfc5226>>.
- [RFC6040] Briscoe, B., "Tunnelling of Explicit Congestion Notification", RFC 6040, DOI 10.17487/RFC6040, November 2010, <<https://www.rfc-editor.org/info/rfc6040>>.
- [RFC6335] Cotton, M., Eggert, L., Touch, J., Westerlund, M., and S. Cheshire, "Internet Assigned Numbers Authority (IANA) Procedures for the Management of the Service Name and Transport Protocol Port Number Registry", BCP 165, RFC 6335, DOI 10.17487/RFC6335, August 2011, <<https://www.rfc-editor.org/info/rfc6335>>.
- [RFC6935] Eubanks, M., Chimento, P., and M. Westerlund, "IPv6 and UDP Checksums for Tunneled Packets", RFC 6935, DOI 10.17487/RFC6935, April 2013, <<https://www.rfc-editor.org/info/rfc6935>>.
- [RFC6936] Fairhurst, G. and M. Westerlund, "Applicability Statement for the Use of IPv6 UDP Datagrams with Zero Checksums", RFC 6936, DOI 10.17487/RFC6936, April 2013, <<https://www.rfc-editor.org/info/rfc6936>>.
- [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>>.
- [RFC8085] Eggert, L., Fairhurst, G., and G. Shepherd, "UDP Usage Guidelines", BCP 145, RFC 8085, DOI 10.17487/RFC8085, March 2017, <<https://www.rfc-editor.org/info/rfc8085>>.
- [RFC8086] Yong, L., Ed., Crabbe, E., Xu, X., and T. Herbert, "GRE-in-UDP Encapsulation", RFC 8086, DOI 10.17487/RFC8086, March 2017, <<https://www.rfc-editor.org/info/rfc8086>>.
- [RFC8126] Cotton, M., Leiba, B., and T. Narten, "Guidelines for Writing an IANA Considerations Section in RFCs", BCP 26, RFC 8126, DOI 10.17487/RFC8126, June 2017, <<https://www.rfc-editor.org/info/rfc8126>>.
- [RFC8300] Quinn, P., Ed., Elzur, U., Ed., and C. Pignataro, Ed., "Network Service Header (NSH)", RFC 8300, DOI 10.17487/

RFC8300, January 2018, <<https://www.rfc-editor.org/info/rfc8300>>.

[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>>.

[RFC9300] Farinacci, D., Fuller, V., Meyer, D., Lewis, D., and A. Cabellos, Ed., "The Locator/ID Separation Protocol (LISP)", RFC 9300, DOI 10.17487/RFC9300, October 2022, <<https://www.rfc-editor.org/info/rfc9300>>.

[RFC9305] Maino, F., Ed., Lemon, J., Agarwal, P., Lewis, D., and M. Smith, "Locator/ID Separation Protocol (LISP) Generic Protocol Extension", RFC 9305, DOI 10.17487/RFC9305, October 2022, <<https://www.rfc-editor.org/info/rfc9305>>.

12.2. Informative References

[I-D.brockners-ippm-ioam-vxlan-gpe]

Brockners, F., Bhandari, S., Govindan, V. P., Pignataro, C., Gredler, H., Leddy, J., Youell, S., Mizrahi, T., Kfir, A., Gafni, B., Lapukhov, P., and M. Spiegel, "VXLAN-GPE Encapsulation for In-situ OAM Data", Work in Progress, Internet-Draft, draft-brockners-ippm-ioam-vxlan-gpe-03, 4 November 2019, <<https://datatracker.ietf.org/doc/html/draft-brockners-ippm-ioam-vxlan-gpe-03>>.

[I-D.ietf-tsvwg-ecn-encap-guidelines]

Briscoe, B. and J. Kaippallimalil, "Guidelines for Adding Congestion Notification to Protocols that Encapsulate IP", Work in Progress, Internet-Draft, draft-ietf-tsvwg-ecn-encap-guidelines-20, 18 October 2023, <<https://datatracker.ietf.org/doc/html/draft-ietf-tsvwg-ecn-encap-guidelines-20>>.

Authors' Addresses

Fabio Maino (Editor)
Cisco Systems

Email: fmaino@cisco.com

Larry Kreeger (editor)
Arrcus

Email: lkreeger@gmail.com

Uri Elzur (editor)
Intel

Email: uri.elzur@intel.com