

Network Working Group  
Internet-Draft  
Intended status: Standards Track  
Expires: May 3, 2018

F. Brockners  
S. Bhandari  
V. Govindan  
C. Pignataro  
Cisco  
H. Gredler  
RtBrick Inc.  
J. Leddy  
Comcast  
S. Youell  
JMPC  
T. Mizrahi  
Marvell  
D. Mozes  
Mellanox Technologies Ltd.  
P. Lapukhov  
Facebook  
R. Chang  
Barefoot Networks  
October 30, 2017

VXLAN-GPE Encapsulation for In-situ OAM Data  
draft-brockners-ioam-vxlan-gpe-00

Abstract

In-situ Operations, Administration, and Maintenance (IOAM) records operational and telemetry information in the packet while the packet traverses a path between two points in the network. This document outlines how IOAM data fields are encapsulated in VXLAN-GPE.

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 May 3, 2018.

## Copyright Notice

Copyright (c) 2017 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 . . . . .	2
2. Conventions . . . . .	3
2.1. Requirement Language . . . . .	3
2.2. Abbreviations . . . . .	3
3. IOAM Data Field Encapsulation in VXLAN-GPE . . . . .	3
3.1. IOAM Trace Data in VXLAN-GPE . . . . .	3
3.2. IOAM POT Data in VXLAN-GPE . . . . .	7
3.3. IOAM Edge-to-Edge Data in VXLAN-GPE . . . . .	8
4. Discussion of the encapsulation approach . . . . .	9
5. IANA Considerations . . . . .	10
6. Security Considerations . . . . .	10
7. Acknowledgements . . . . .	11
8. References . . . . .	11
8.1. Normative References . . . . .	11
8.2. Informative References . . . . .	12
Authors' Addresses . . . . .	12

## 1. Introduction

In-situ OAM (IOAM) records OAM information within the packet while the packet traverses a particular network domain. The term "in-situ" refers to the fact that the IOAM data fields are added to the data packets rather than being sent within packets specifically dedicated to OAM. This document defines how IOAM data fields are transported as part of the VXLAN-GPE [I-D.ietf-nvo3-vxlan-gpe] encapsulation. The IOAM data fields are defined in [I-D.ietf-ippm-ioam-data]. An implementation of IOAM which leverages VXLAN-GPE to carry the IOAM data is available from the FD.io open source software project [FD.io].

## 2. Conventions

### 2.1. Requirement 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].

### 2.2. Abbreviations

Abbreviations used in this document:

IOAM: In-situ Operations, Administration, and Maintenance

MTU: Maximum Transmit Unit

OAM: Operations, Administration, and Maintenance

POT: Proof of Transit

SFC: Service Function Chain

VXLAN-GPE: Virtual eXtensible Local Area Network, Generic Protocol Extension

## 3. IOAM Data Field Encapsulation in VXLAN-GPE

For encapsulating IOAM data fields into VXLAN-GPE [I-D.ietf-nvo3-vxlan-gpe] the different IOAM data fields are added as options within new IOAM protocol headers in VXLAN-GPE. In an administrative domain where IOAM is used, insertion of the IOAM protocol header(s) in VXLAN GPE is enabled at the VXLAN-GPE tunnel endpoints which also serve as IOAM encapsulating/decapsulating nodes by means of configuration. The VXLAN-GPE header is defined in [I-D.ietf-nvo3-vxlan-gpe]. IOAM specific fields for VXLAN-GPE are defined in this document.

### 3.1. IOAM Trace Data in VXLAN-GPE

IOAM tracing data represents data that is inserted at nodes that a packet traverses. To allow for optimal implementations in both software as well as hardware forwarders, two different ways to encapsulate IOAM data are defined: "Pre-allocated" and "Incremental". See [I-D.ietf-ippm-ioam-data] for details on IOAM tracing and the pre-allocated and incremental IOAM trace options.

The packet formats of the pre-allocated IOAM trace and incremental IOAM trace when encapsulated in VXLAN-GPE are defined as below.

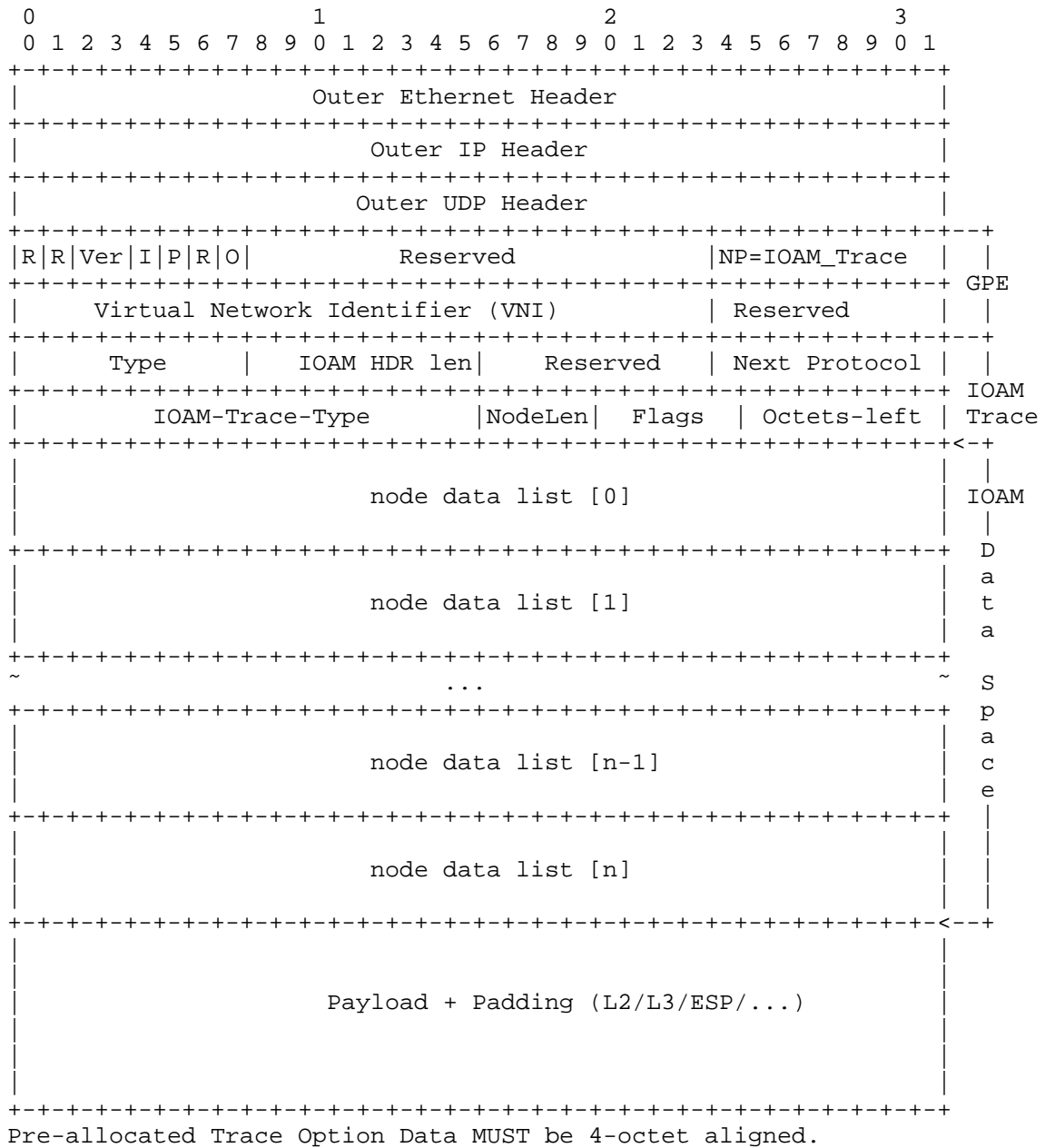


Figure 1: IOAM Pre-allocated Trace Option Format over VXLAN-GPE

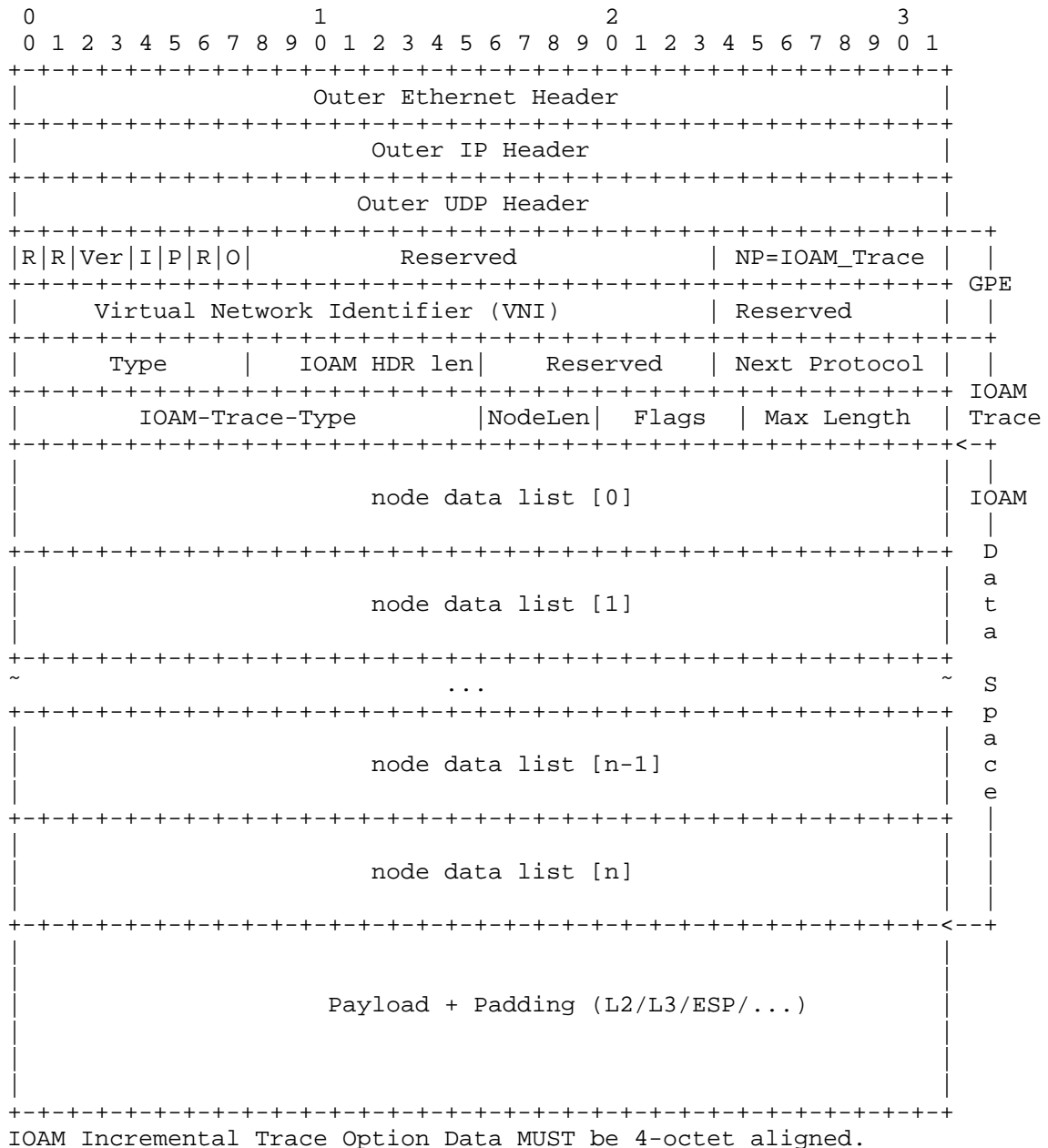


Figure 2: IOAM Incremental Trace Option Format over VXLAN-GPE

The IOAM Trace header consists of 8 octets, as illustrated in Figure 1 and Figure 2. The format of the first 4 octets (Figure 3)

is specific to VXLAN-GPE, and is defined in this document. The format of the next 4 octets (trace option header) is defined in [I-D.ietf-ippm-ioam-data], and is described here for the sake of clarity.

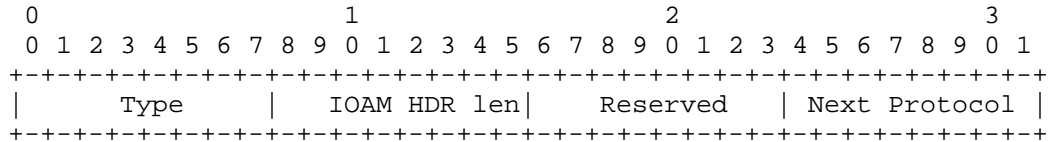


Figure 3: Trace Shim Header for VXLAN-GPE

The fields of the trace shim header are as follows:

Type: 8-bit unsigned integer defining IOAM header type  
 IOAM\_TRACE\_Preallocated or IOAM\_Trace\_Incremental are defined here.

IOAM HDR len: 8-bit unsigned integer. Length of the IOAM HDR in 4-octet units.

Reserved: 8-bit reserved field MUST be set to zero.

Next Protocol: 8-bit unsigned integer that determines the type of header following IOAM protocol. The value is from the IANA registry setup for VXLAN GPE Next Protocol defined in [I-D.ietf-nvo3-vxlan-gpe].

The fields of the trace option header [I-D.ietf-ippm-ioam-data] are as follows:

IOAM-Trace-Type: 16-bit identifier of IOAM Trace Type as defined in [I-D.ietf-ippm-ioam-data] IOAM-Trace-Types.

Node Data Length: 4-bit unsigned integer as defined in [I-D.ietf-ippm-ioam-data].

Flags: 5-bit field as defined in [I-D.ietf-ippm-ioam-data].

Octets-left: 7-bit unsigned integer as defined in [I-D.ietf-ippm-ioam-data].

Maximum-length: 7-bit unsigned integer as defined in [I-D.ietf-ippm-ioam-data].

Node data List [n]: Variable-length field as defined in [I-D.ietf-ippm-ioam-data].

### 3.2. IOAM POT Data in VXLAN-GPE

IOAM proof of transit (POT, see also [I-D.brockners-proof-of-transit]) offers a means to verify that a packet has traversed a defined set of nodes. IOAM POT data fields are encapsulated in VXLAN-GPE using a dedicated VXLAN-GPE protocol header:

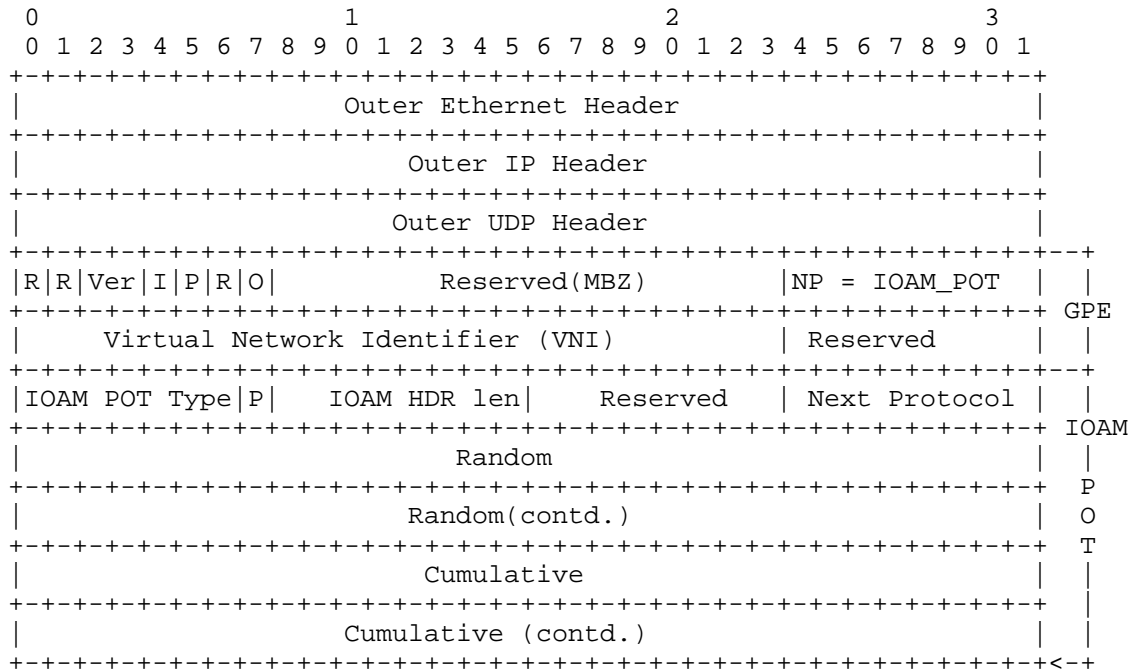


Figure 4: IOAM POT Header Following the VXLAN-GPE Header

The IOAM POT Shim Header (Figure 5), which is defined in this document, is a 4-octet header, that includes the following fields:

**IOAM POT Type:** 7-bit identifier of a particular POT variant that specifies the POT data that is to be included as defined in [I-D.ietf-ippm-ioam-data].

**Profile to use (P):** 1-bit as defined in [I-D.ietf-ippm-ioam-data] IOAM POT Option.

**IOAM HDR len:** 8-bit unsigned integer. Length of the IOAM HDR in 4-octet units.

**Reserved:** 8-bit reserved field MUST be set to zero.

Next Protocol: 8-bit unsigned integer that determines the type of header following IOAM protocol. The value is from the IANA registry setup for VXLAN GPE Next Protocol defined in [I-D.ietf-nvo3-vxlan-gpe].

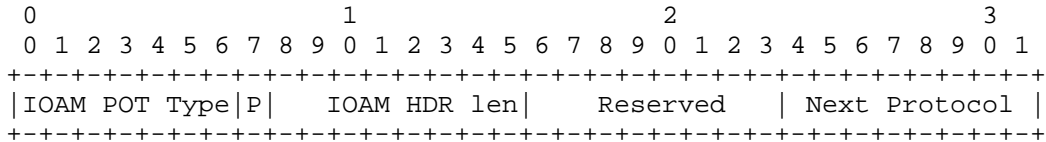


Figure 5: POT Shim Header for VXLAN-GPE

The rest of the fields in the POT option [I-D.ietf-ippm-ioam-data] are as follows:

Random: 64-bit Per-packet random number.

Cumulative: 64-bit Cumulative value that is updated by the Service Functions.

### 3.3. IOAM Edge-to-Edge Data in VXLAN-GPE

The IOAM edge-to-edge option is to carry data that is added by the IOAM encapsulating node and interpreted by the IOAM decapsulating node. IOAM specific fields to encapsulate IOAM Edge-to-Edge data fields are defined as follows:

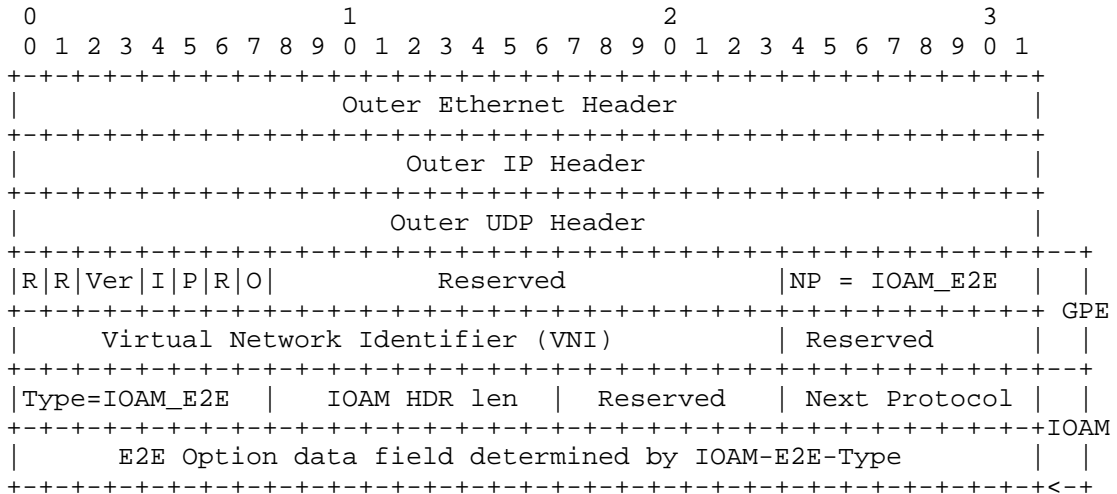


Figure 6: IOAM Edge-to-Edge over a VXLAN-GPE Header



The IOAM E2E Shim Header, which is defined in this document, is a 4-octet header, that includes the following fields:

Type: 8-bit identifier of a particular E2E variant that specifies the E2E data that is to be included as defined in [I-D.ietf-ippm-ioam-data].

IOAM HDR len: 8-bit unsigned integer. Length of the IOAM HDR in 4-octet units.

Reserved: 8-bit reserved field MUST be set to zero.

Next Protocol: 8-bit unsigned integer that determines the type of header following IOAM protocol. The value is from the IANA registry setup for VXLAN GPE Next Protocol defined in [I-D.ietf-nvo3-vxlan-gpe].

```

      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
+-----+-----+-----+-----+-----+-----+-----+-----+
|Type=IOAM_E2E |   IOAM HDR len   |   Reserved   | Next Protocol |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

Figure 7: E2E Shim Header for VXLAN-GPE

The rest of the E2E option [I-D.ietf-ippm-ioam-data] consists of:

E2E Option data field: Variable length field as defined in [I-D.ietf-ippm-ioam-data] IOAM E2E Option.

#### 4. Discussion of the encapsulation approach

This section is to support the working group discussion in selecting the most appropriate approach for encapsulating IOAM data fields in VXLAN-GPE.

An encapsulation of IOAM data fields in VXLAN-GPE should be friendly to an implementation in both hardware as well as software forwarders. Hardware forwarders benefit from an encapsulation that minimizes iterative look-ups of fields within the packet: Any operation which looks up the value of a field within the packet, based on which another lookup is performed, consumes additional gates and time in an implementation - both of which are desired to be kept to a minimum. This means that flat TLV structures are to be preferred over nested TLV structures. IOAM data fields are grouped into three option categories: Trace, proof-of-transit, and edge-to-edge. Each of these three options defines a TLV structure. A hardware-friendly encapsulation approach avoids grouping these three option categories

into yet another TLV structure, but would rather carry the options as a serial sequence.

Two approaches for encapsulating IOAM data fields in VXLAN-GPE could be considered:

1. Use a single GPE protocol type for all IOAM types: IOAM would receive a single GPE protocol type code point. A "sub-type" (e.g. 4 bit wide) would then specify what IOAM options type(s) (trace, proof-of-transit, edge-to-edge) are carried. In case there is a need for additional IOAM options, changes would be contained within the single GPE protocol type for IOAM.
2. Use one GPE protocol type per IOAM options type: Each IOAM data field option (trace, proof-of-transit, and edge-to-edge) would be specified by its own "next protocol", i.e. each IOAM options type becomes its own GPE protocol type with a dedicated code point. This implies that in case additional IOAM option types would be added in the future, additional GPE protocol type code points would need to be allocated.

The second option has been chosen here, because it avoids the additional layer of TLV nesting that the use of a single GPE protocol type for all IOAM option types would result in.

## 5. IANA Considerations

IANA is requested to allocate protocol numbers for the following VXLAN-GPE "Next Protocols" related to IOAM:

Next Protocol	Description	Reference
x	IOAM_Trace	This document
y	IOAM_POT	This document
z	IOAM_E2E	This document

## 6. Security Considerations

The security considerations of VXLAN-GPE are discussed in [I-D.ietf-nvo3-vxlan-gpe], and the security considerations of IOAM in general are discussed in [I-D.ietf-ippm-ioam-data].

IOAM is considered a "per domain" feature, where one or several operators decide on leveraging and configuring IOAM according to their needs. Still, operators need to properly secure the IOAM

domain to avoid malicious configuration and use, which could include injecting malicious IOAM packets into a domain.

## 7. Acknowledgements

The authors would like to thank Eric Vyncke, Nalini Elkins, Srihari Raghavan, Ranganathan T S, Karthik Babu Harichandra Babu, Akshaya Nadahalli, Stefano Previdi, Hemant Singh, Erik Nordmark, LJ Wobker, and Andrew Yourtchenko for the comments and advice.

## 8. References

### 8.1. Normative References

- [ETYPES] "IANA Ethernet Numbers",  
<<https://www.iana.org/assignments/ethernet-numbers/ethernet-numbers.xhtml>>.
- [I-D.ietf-ippm-ioam-data]  
Brockners, F., Bhandari, S., Pignataro, C., Gredler, H., Leddy, J., Youell, S., Mizrahi, T., Mozes, D., Lapukhov, P., Chang, R., and d. daniel.bernier@bell.ca, "Data Fields for In-situ OAM", draft-ietf-ippm-ioam-data-00 (work in progress), September 2017.
- [I-D.ietf-nvo3-vxlan-gpe]  
Maino, F., Kreeger, L., and U. Elzur, "Generic Protocol Extension for VXLAN", draft-ietf-nvo3-vxlan-gpe-04 (work in progress), April 2017.
- [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>>.
- [RFC2784] Farinacci, D., Li, T., Hanks, S., Meyer, D., and P. Traina, "Generic Routing Encapsulation (GRE)", RFC 2784, DOI 10.17487/RFC2784, March 2000, <<https://www.rfc-editor.org/info/rfc2784>>.
- [RFC3232] Reynolds, J., Ed., "Assigned Numbers: RFC 1700 is Replaced by an On-line Database", RFC 3232, DOI 10.17487/RFC3232, January 2002, <<https://www.rfc-editor.org/info/rfc3232>>.

## 8.2. Informative References

- [FD.io] "Fast Data Project: FD.io", <<https://fd.io/>>.
- [I-D.brockners-proof-of-transit]  
Brockners, F., Bhandari, S., Dara, S., Pignataro, C.,  
Leddy, J., Youell, S., Mozes, D., and T. Mizrahi, "Proof  
of Transit", draft-brockners-proof-of-transit-03 (work in  
progress), March 2017.
- [RFC7665] Halpern, J., Ed. and C. Pignataro, Ed., "Service Function  
Chaining (SFC) Architecture", RFC 7665,  
DOI 10.17487/RFC7665, October 2015, <[https://www.rfc-  
editor.org/info/rfc7665](https://www.rfc-editor.org/info/rfc7665)>.

## Authors' Addresses

Frank Brockners  
Cisco Systems, Inc.  
Hansaallee 249, 3rd Floor  
DUESSELDORF, NORDRHEIN-WESTFALEN 40549  
Germany

Email: [fbrockne@cisco.com](mailto:fbrockne@cisco.com)

Shwetha Bhandari  
Cisco Systems, Inc.  
Cessna Business Park, Sarjapura Marathalli Outer Ring Road  
Bangalore, KARNATAKA 560 087  
India

Email: [shwethab@cisco.com](mailto:shwethab@cisco.com)

Vengada Prasad Govindan  
Cisco Systems, Inc.

Email: [venggovi@cisco.com](mailto:venggovi@cisco.com)

Carlos Pignataro  
Cisco Systems, Inc.  
7200-11 Kit Creek Road  
Research Triangle Park, NC 27709  
United States

Email: [cpignata@cisco.com](mailto:cpignata@cisco.com)

Hannes Gredler  
RtBrick Inc.

Email: hannes@rtbrick.com

John Leddy  
Comcast

Email: John\_Leddy@cable.comcast.com

Stephen Youell  
JP Morgan Chase  
25 Bank Street  
London E14 5JP  
United Kingdom

Email: stephen.youell@jpmorgan.com

Tal Mizrahi  
Marvell  
6 Hamada St.  
Yokneam 20692  
Israel

Email: talmi@marvell.com

David Mozes  
Mellanox Technologies Ltd.

Email: davidm@mellanox.com

Petr Lapukhov  
Facebook  
1 Hacker Way  
Menlo Park, CA 94025  
US

Email: petr@fb.com

Remy Chang  
Barefoot Networks  
2185 Park Boulevard  
Palo Alto, CA 94306  
US