

Network Working Group
Internet-Draft
Intended status: Experimental
Expires: May 3, 2017

F. Brockners
S. Bhandari
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, 2016

Data Formats for In-situ OAM
draft-brockners-inband-oam-data-02

Abstract

In-situ Operations, Administration, and Maintenance (OAM) records operational and telemetry information in the packet while the packet traverses a path between two points in the network. This document discusses the data types and data formats for in-situ OAM data records. In-situ OAM data records can be embedded into a variety of transports such as NSH, Segment Routing, VXLAN-GPE, native IPv6 (via extension header), or IPv4. In-situ OAM is to complement current out-of-band OAM mechanisms based on ICMP or other types of probe packets.

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

Internet-Draft

In-situ OAM Data Formats

October 2016

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, 2017.

Copyright Notice

Copyright (c) 2016 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
3.	In-situ OAM Data Types and Data Format	4
3.1.	In-situ OAM Tracing Options	4
3.1.1.	Pre-allocated Trace Option	6
3.1.2.	Incremental Trace Option	9
3.1.3.	In-situ OAM node data element format	11
3.1.4.	Examples of In-situ OAM node data	14
3.2.	In-situ OAM Proof of Transit Option	16
3.3.	In-situ OAM Edge-to-Edge Option	18
4.	In-situ OAM Data Export	18
5.	IANA Considerations	19
6.	Manageability Considerations	19
7.	Security Considerations	19
8.	Acknowledgements	19
9.	References	19
9.1.	Normative References	19
9.2.	Informative References	20
	Authors' Addresses	20

[1.](#) Introduction

This document defines data record types for "in-situ" Operations, Administration, and Maintenance (OAM). In-situ OAM records OAM information within the packet while the packet traverses a particular network domain. The term "in-situ" refers to the fact that the OAM

data is added to the data packets rather than is being sent within packets specifically dedicated to OAM. A discussion of the motivation and requirements for in-situ OAM can be found in [[I-D.brockners-inband-oam-requirements](#)]. In-situ OAM is to complement "out-of-band" or "active" mechanisms such as ping or traceroute, or more recent active probing mechanisms as described in [[I-D.lapukhov-dataplane-probe](#)]. In-situ OAM mechanisms can be leveraged where current out-of-band mechanisms do not apply or do not offer the desired results, such as proving that a certain set of traffic takes a pre-defined path, SLA verification for the live data traffic, detailed statistics on traffic distribution paths in networks that distribute traffic across multiple paths, or scenarios where probe traffic is potentially handled differently from regular data traffic by the network devices.

This document defines the data types and data formats for in-situ OAM data records. The in-situ OAM data records can be transported by a variety of transport protocols, including NSH, Segment Routing, VXLAN-GPE, IPv6, IPv4. Encapsulation details for these different transport protocols are outside the scope of this document.

[2.](#) Conventions

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](#)].

Abbreviations used in this document:

MTU:	Maximum Transmit Unit
NSH:	Network Service Header
OAM:	Operations, Administration, and Maintenance
SFC:	Service Function Chain

SID: Segment Identifier

SR: Segment Routing

TLV: Type-Length-Value

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

[3.](#) In-situ OAM Data Types and Data Format

This section defines in-situ OAM data types and data formats of the data records required for in-situ OAM. The different uses of in-situ OAM require the definition of different types of data. The in-situ OAM data format for the data being carried corresponds to the three main categories of in-situ OAM data defined in [\[I-D.brockners-inband-oam-requirements\]](#), which are : edge-to-edge, per node, and for selected nodes only.

Transport options for in-situ OAM data are found in [\[I-D.brockners-inband-oam-transport\]](#). In-situ OAM data is defined as options in Type-Length-Value (TLV) format. The TLV format for each of the three different types of in-situ OAM data is defined in this document.

In-situ OAM is expected to be deployed in a specific domain rather than on the overall Internet. The part of the network which employs in situ OAM is referred to as the "in-situ OAM-domain". In-situ OAM data is added to a packet upon entering the in-situ OAM-domain and is removed from the packet when exiting the domain. Within the in-situ OAM-domain, the in-situ OAM data may be updated by network nodes that the packet traverses. The device which adds in-situ OAM data container to the packet to capture in-situ OAM data is called the "in-situ OAM encapsulating node", whereas the device which removes the in-situ OAM data container is referred to as the "in-situ OAM decapsulating node". Nodes within the domain which are aware of in-situ OAM data and read and/or write or process the in-situ OAM data are called "in-situ OAM transit nodes". Note that not every node in

an in-situ OAM domain needs to be an in-situ OAM transit node. For example, a Segment Routing deployment might require the segment routing path to be verified. In that case, only the SR nodes would also be in-situ OAM transit nodes rather than all nodes.

[3.1.](#) In-situ OAM Tracing Options

"In-situ OAM tracing data" is expected to be collected at every node that a packet traverses, i.e., in a typical deployment all nodes in an in-situ OAM-domain would participate in in-situ OAM and thus be in-situ OAM transit nodes, in-situ OAM encapsulating or in-situ OAM decapsulating nodes. The maximum network diameter of the in-situ OAM domain is assumed to be known.

To optimize hardware and software implementations tracing is defined as two separate options. Any deployment MAY choose to configure and support one or both of the following options. An implementation of the transport protocol that carries these in-situ OAM data MAY choose to support only one of the options. In the event that both options

are utilized at the same time, the Incremental Trace Option MUST be placed before the Pre-allocated Trace Option.

Pre-allocated Trace Option: This trace option is defined as a container of node-data elements with pre-allocated space for each node to populate its information. This option is useful for software implementations where it is efficient to allocate the space once and index into the array to populate the data during transit. The in-situ OAM encapsulating node allocates the option header and sets the fields in the option header. The in situ OAM encapsulating node allocates an array which is to store operational data retrieved from every node while the packet traverses the domain. In-situ OAM transit nodes update the content of the array. A pointer which is part of the in-situ OAM trace data points to the next empty slot in the array, which is where the next in-situ OAM transit node fills in its data.

Incremental Trace Option: This trace options is defined as a container of node-data elements where each node allocates and pushes its node data immediately following the option header. The number of node-data recorded and maximum number of node data that can be recorded are written into the option header. This format

of trace recording is useful for some of the hardware implementations as this eliminates the need for the transit network elements to read the full array in the option and allows for arbitrarily long packets as the MTU allows. The in-situ OAM encapsulating node allocates the option header. The in-situ OAM encapsulating node based on operational state and configuration sets the fields in the header to control how large the node data list can grow. In-situ OAM transit nodes pushes its node data to the node data list and increments the number of node data records in the header.

Every node data entry is to hold information for a particular in situ OAM transit node that is traversed by a packet. The in-situ OAM decapsulating node removes the in-situ OAM data and process and/or export the metadata. In-situ OAM data uses its own name-space for information such as node identifier or interface identifier. This allows for a domain-specific definition and interpretation. For example: In one case an interface-id could point to a physical interface (e.g., to understand which physical interface of an aggregated link is used when receiving or transmitting a packet) whereas in another case it could refer to a logical interface (e.g., in case of tunnels).

The following in-situ OAM data is defined for in-situ OAM tracing:

- o Identification of the in-situ OAM node. An in-situ OAM node identifier can match to a device identifier or a particular control point or subsystem within a device.
- o Identification of the interface that a packet was received on.
- o Identification of the interface that a packet was sent out on.
- o Time of day when the packet was processed by the node. Different definitions of processing time are feasible and expected, though it is important that all devices of an in-situ OAM domain follow the same definition.
- o Generic data: Format-free information where syntax and semantic of the information is defined by the operator in a specific

deployment. For a specific deployment, all in-situ OAM nodes should interpret the generic data the same way. Examples for generic in-situ OAM data include geo-location information (location of the node at the time the packet was processed), buffer queue fill level or cache fill level at the time the packet was processed, or even a battery charge level.

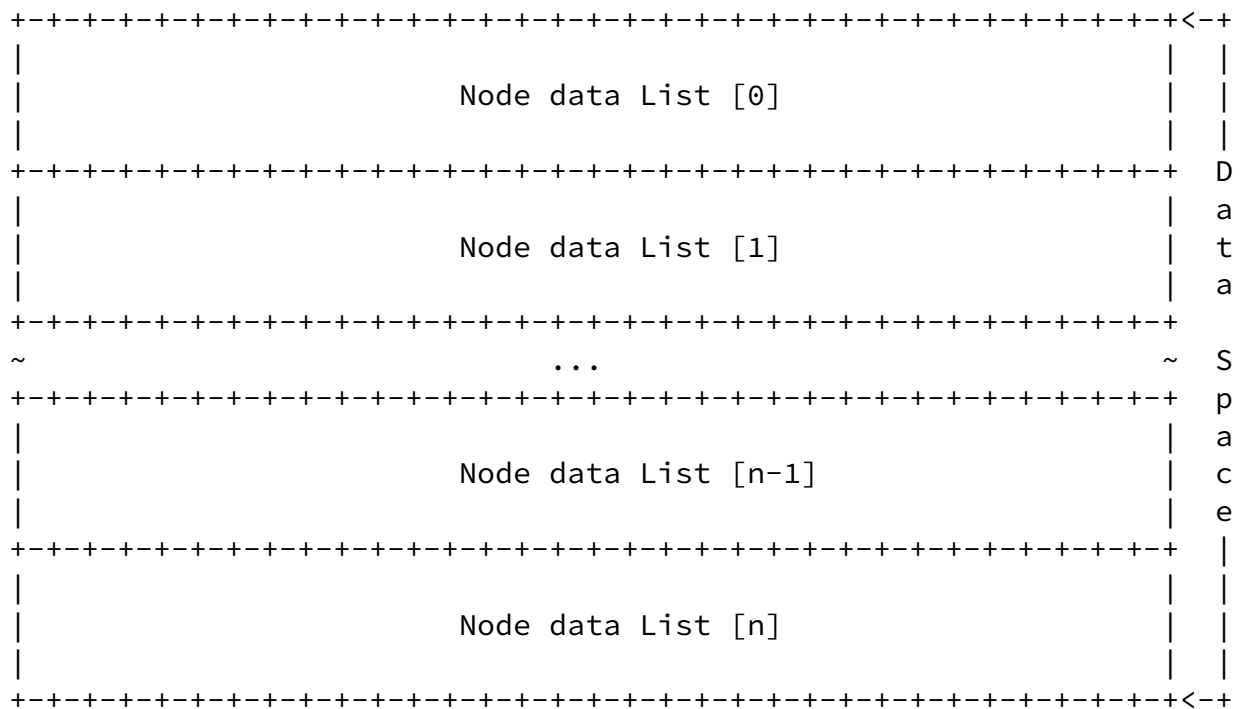
- o A mechanism to detect whether in-situ OAM trace data was added at every hop or whether certain hops in the domain weren't in-situ OAM transit nodes.

The "Node data List" array in the packet is populated iteratively as the packet traverses the network, starting with the last entry of the array, i.e., "Node data List [n]" is the first entry to be populated, "Node data List [n-1]" is the second one, etc.

[3.1.1.1.](#) Pre-allocated Trace Option

In-situ OAM Pre-allocated Trace Option:

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
+-+-+-----	+-+-+-----	+-+-+-----	+-+-+-----
Option Type	Opt Data Len	Elements-left	Reserved1
+-+-+-----	+-+-+-----	+-+-+-----	+-+-+-----
	OAM-Trace-Type	Reserved2	



Option Type: 8-bit identifier of the type of option. Option number is defined based on the encapsulation protocol.

Opt Data Len: 8-bit unsigned integer. Length of the Option Data field of this option, in octets.

Elements-left: 8-bit unsigned integer. A pointer that indicates the next data recording point in the data space of the packet in octets. It is the index into the "Node data List" array shown above.

Reserved1: 8-bit unused field in this document and MUST be set to zero.

OAM-trace-type: 16-bit identifier of a particular trace element variant.

The trace type value is a bit field. The following bit fields are

defined in this document, with details on each field described in the [Section 3.1.3](#). The order of packing the trace data in each Node-data element follows the bit order for setting each trace data element.

- Bit 0 When set indicates presence of Hop_Lim and node_id in the Node data.
- Bit 1 When set indicates presence of ingress_if_id and egress_if_id in the Node data.
- Bit 2 When set indicates presence of timestamp seconds in the Node data
- Bit 3 When set indicates presence of timestamp nanoseconds in the Node data.
- Bit 4 When set indicates presence of transit delay in the Node data.
- Bit 5 When set indicates presence of app_data in the Node data.
- Bit 6 When set indicates presence of queue depth in the Node data.
- Bit 7 - 14 Undefined in this document.
- Bit 15 When set indicates wide data format for all the node data elements that are present. When unset indicates short data format for all the node data elements that are present.

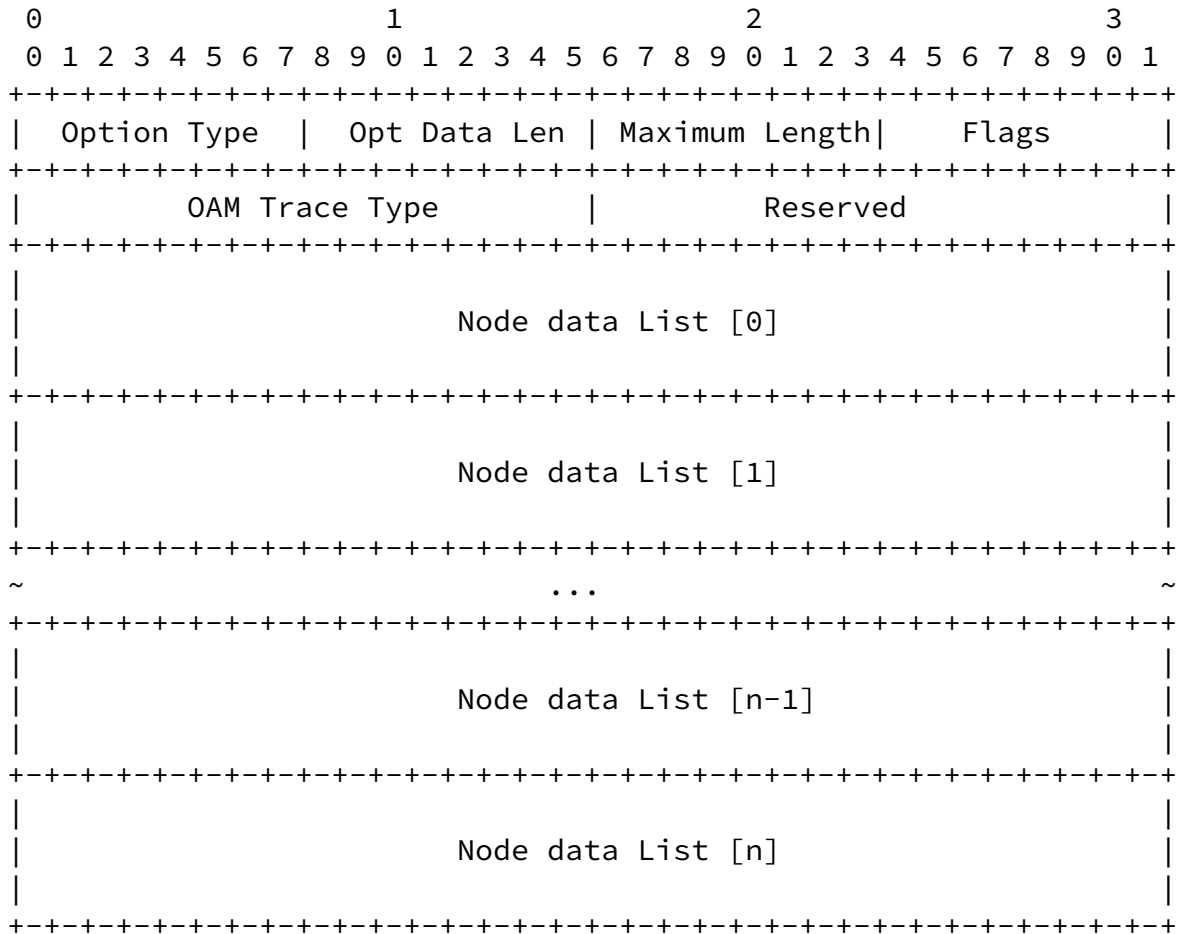
[Section 3.1.4](#) describes the format of a number of trace types.

Reserved2: 16-bit unused field in this document and MUST be set to zero.

Node data List [n]: Variable-length field. The format of which is determined by the OAM Type representing the n-th Node data in the Node data List. The Node data List is encoded starting from the last Node data of the path. The first element of the node data list (Node data List [0]) contains the last node of the path while the last node data of the Node data List (Node data List[n]) contains the first Node data of the path traced. The index contained in "Elements-left" identifies the current active Node data to be populated.

[3.1.2.](#) Incremental Trace Option

In-situ OAM Incremental Trace Option:



Option Type: 8-bit identifier of the type of option. Option number is defined based on the encapsulation protocol.

Opt Data Len: 8-bit unsigned integer. Length of the Option Data field of this option, in octets.

Maximum Length: 8-bit unsigned integer. This field specifies the maximum length of the node data list in octets. Given that the sender knows the minimum path MTU, the sender can set the maximum of node data bytes allowed before exceeding the MTU. Thus, a simple comparison between "Opt data Len" and "Max Length" allows to decide whether or not data could be added.

Flags 8-bit field. Following flags are defined:

- 1 "Overflow" (0-bit) (least significant bit). This bit is set by the network element if the number of records on the packet is at the maximum limit as specified by the packet: i.e., the packet is already "full" of telemetry information. This is useful for transit nodes to ignore further processing of the option. If inserting a new node data record would cause "Opt Data Len" to exceed "Max Length", no record is added and the overflow "0-bit" must be set to "1" in the header.

OAM-trace-type: 16-bit identifier of a particular trace element variant.

The trace type value is a bit field. The following bit fields are defined in this document, with details on each field described in the [Section 3.1.3](#). The order of packing the trace data in each Node-data element follows the bit order for setting each trace data element.

- | | |
|-------|---|
| Bit 0 | When set indicates presence of Hop_Lim and node_id in the Node data. |
| Bit 1 | When set indicates presence of ingress_if_id and egress_if_id in the Node data. |
| Bit 2 | When set indicates presence of timestamp seconds in the Node data |
| Bit 3 | When set indicates presence of timestamp nanoseconds in the Node data. |
| Bit 4 | When set indicates presence of transit delay in the Node data. |
| Bit 5 | When set indicates presence of app_data in the Node data. |
| Bit 6 | When set indicates presence of queue depth in the Node data. |
| Bit 7 | When set indicates presence of variable length Opaque |

State Snapshot field.

Bit 8-14 Undefined in this draft.

Bit 15 When set indicates wide data format for all the node data elements that are present. When unset indicates short data format for all the node data elements that are present.

Brockners, et al.

Expires May 3, 2017

[Page 10]

Internet-Draft

In-situ OAM Data Formats

October 2016

[Section 3.1.4](#) describes the format of a number of trace types.

Reserved: 2 bytes unused field in this document and MUST be set to zero.

Node data List [n]: Variable-length field. The format of which is determined by the OAM Type representing the n-th Node data in the Node data List. The Node data List is encoded starting from the last Node data of the path. The first element of the node data list (Node data List [0]) contains the last node of the path while the last node data of the Node data List (Node data List[n]) contains the first Node data of the path traced.

[3.1.3](#). In-situ OAM node data element format

The in-situ OAM node data elements are defined in 2 formats - short and wide that is selected by bit 15 in the OAM-trace-type field. All the data records MUST be 4-byte aligned in both the formats.

Data type and format for each of the data records in short format is shown below:

Hop_Lim and node_id: 4-octet field defined as follows:

```

 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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Hop_Lim | node_id |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
```

Hop_Lim: 1-octet unsigned integer. It is set to the Hop Limit value in the packet at the node that records this data. Hop Limit information is used to identify the location of the node

in the communication path. This is copied from the lower layer for e.g. TTL value in IPv4 header or hop limit field from IPv6 header of the packet.

node_id: 3-octet unsigned integer. Node identifier field to uniquely identify a node within in-situ OAM domain. The procedure to allocate, manage and map the node_ids is beyond the scope of this document.

ingress_if_id and egress_if_id: 4-octet field defined as follows:

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|           ingress_if_id           |           egress_if_id           |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

ingress_if_id: 2-octet unsigned integer. Interface identifier to record the ingress interface the packet was received on.

egress_if_id: 2-octet unsigned integer. Interface identifier to record the egress interface the packet is forwarded out of.

timestamp seconds: 4-octet unsigned integer. Absolute timestamp in seconds that specifies the time at which the packet was received by the node. The format of this field is identical to the most significant 32 bits of 64 least significant bits of the [IEEE1588v2]. This truncated format consists of a 32-bit seconds field. As defined in [IEEE1588v2], the timestamp specifies the number of seconds elapsed since 1 January 1970 00:00:00 according to the International Atomic Time (TAI).

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|           timestamp seconds           |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

timestamp nanoseconds: 4-octet unsigned integer in the range 0 to 10^9-1 . This timestamp specifies the fractional part of the wall clock time at which the packet was received by the node in units of nanoseconds. It is nanoseconds that are recorded in 32 least significant bits of absolute time as per [IEEE1588v2]. This

fields allows for delay computation between any two nodes in the network when the nodes are time synchronized.

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                                     timestamp nanoseconds                                     |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

transit delay: 4-octet unsigned integer in the range 0 to $2^{30}-1$.
It is the time in nanoseconds packet spent in transiting a node.
This can serve to give an indication of queuing delay at the node.
If the transit delay exceeds $2^{30}-1$ nanoseconds then the top bit
'0' is set to indicate overflow.

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|0|                                     transit delay                                     |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

app_data: 4-octet placeholder which can be used by the node to add application specific data

queue depth: 4-octet unsigned integer field. This field indicates the length of the egress interface queue of the interface where the packet is forwarded out of.

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                                     queue depth                                     |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Data type and format for each of the elements in wide format follows when Most Significant Bit (MSB) i.e., bit 15 of OAM-Trace-Type is set:

Hop_Lim and node_id: 8-octet field defined as follows:

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Hop_Lim |                                     node_id                                     | ~

```

```

+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
~                               node_id (contd)                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Hop_Lim: 1-octet unsigned integer. It is set to the Hop Limit value in the packet at the node that records this data. Hop Limit information is used to identify the location of the node in the communication path. This is copied from the lower layer for e.g. TTL value in IPv4 header or hop limit field from IPv6 header of the packet.

node_id: 7-octet unsigned integer. Node identifier field to uniquely identify a node within in-situ OAM domain. The procedure to allocate, manage and map the node_ids is beyond the scope of this document.

ingress_if_id and egress_if_id: 8-octet field defined as follows:

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               ingress_if_id                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               egress_if_id                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

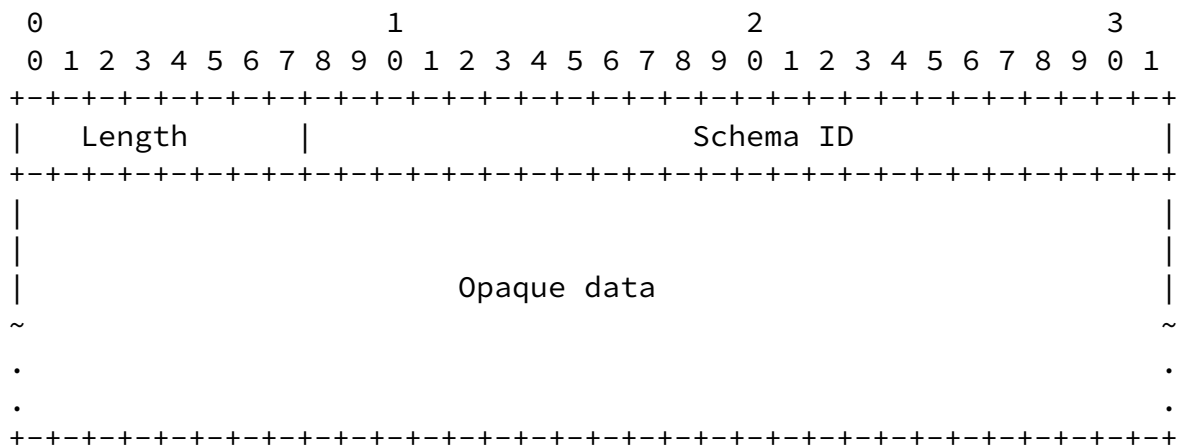
```

ingress_if_id: 4-octet unsigned integer. Interface identifier to record the ingress interface the packet was received on.

egress_if_id: 4-octet unsigned integer. Interface identifier to record the egress interface the packet is forwarded out of.

app_data: 8-octet placeholder which can be used by the node to add application specific data.

Opaque State Snapshot: Variable length field. It allows the network element to store arbitrary state in the node data record, without a pre-defined schema. The schema needs to be made known to the analyzer by some out-of-band means. The 24-bit "Schema Id" field in the record is supposed to let the analyzer know which particular schema to use, and it is expected to be configured on the network element by the operator. This ID is expected to be configured on the device by the network operator.



Length: 1-octet unsigned integer. It is the length of the Opaque data field that follows Schema Id. It MUST always be a multiple of 4.

Schema ID: 3-octet unsigned integer identifying the schema of Opaque data.

Opaque data: Variable length field. This field is interpreted as specified by the schema identified by the Schema ID.

The fields - timestamp seconds, timestamp nanoseconds and transit delay have the same format as defined in short format.

3.1.4. Examples of In-situ OAM node data

An entry in the "Node data List" array can have different formats, following the needs of the deployment. Some deployments might only be interested in recording the node identifiers, whereas others might be interested in recording node identifier and timestamp. The section defines different formats that an entry in "Node data List" can take.

0x002B: In-situ OAM-trace-type is 0x2B then the format of node data is:

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
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---


```

+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Hop_Lim   |                               node_id                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   ingress_if_id   |                               egress_if_id                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               timestamp nanoseconds                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               app_data                                           |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

0x0003: In-situ OAM-trace-type is 0x0003 then the format is:

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Hop_Lim   |                               node_id                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   ingress_if_id   |                               egress_if_id                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

0x0009: In-situ OAM-trace-type is 0x0009 then the format is:

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Hop_Lim   |                               node_id                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               timestamp nanoseconds                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

0x0021: In-situ OAM-trace-type is 0x0021 then the format is:

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|   Hop_Lim   |                               node_id                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                               app_data                                           |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

0x0029: In-situ OAM-trace-type is 0x0029 then the format is:

```

 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
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|  Hop_Lim      |                      node_id                      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                      timestamp nanoseconds                      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                      app_data                      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

0x104D: In-situ OAM-trace-type is 0x104D then the format is:

```

 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
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|  Hop_Lim      |                      node_id                      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                      node_id(contd)                      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                      timestamp seconds                      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|                      timestamp nanoseconds                      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|  Length      |                      Schema Id                      |
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
|
|
|                      Opaque data                      |
|
~
.
.
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

[3.2.](#) In-situ OAM Proof of Transit Option

In-situ OAM Proof of Transit data is to support the path or service function chain [[RFC7665](#)] verification use cases. Proof-of-transit uses methods like nested hashing or nested encryption of the in-situ OAM data or mechanisms such as Shamir's Secret Sharing Schema (SSSS). While details on how the in-situ OAM data for the proof of transit option is processed at in-situ OAM encapsulating, decapsulating and transit nodes are outside the scope of the document, all of these approaches share the need to uniquely identify a packet as well as iteratively operate on a set of information that is handed from node to node. Correspondingly, two pieces of information are added as in-situ OAM data to the packet:

Profile to use (P): 1-bit. Indicates which POT-profile is used to generate the Cumulative. Any node participating in POT will have a maximum of 2 profiles configured that drive the computation of cumulative. The two profiles are numbered 0, 1. This bit conveys whether profile 0 or profile 1 is used to compute the Cumulative.

Reserved: 7-bit. Reserved for future use.

Random: 64-bit Per packet Random number.

Cumulative: 64-bit Cumulative that is updated at specific nodes by processing per packet Random number field and configured parameters.

Brockners, et al.

Expires May 3, 2017

[Page 17]

Internet-Draft

In-situ OAM Data Formats

October 2016

Note: Larger or smaller sizes of "Random" and "Cumulative" data are feasible and could be required for certain deployments (e.g. in case of space constraints in the transport protocol used). Future versions of this document will address different sizes of data for "proof of transit".

3.3. In-situ OAM Edge-to-Edge Option

The in-situ OAM Edge-to-Edge Option is to carry data which is to be interpreted only by the in-situ OAM encapsulating and in-situ OAM decapsulating node, but not by in-situ OAM transit nodes.

Currently only sequence numbers use the in-situ OAM Edge-to-Edge option. In order to detect packet loss, packet reordering, or packet duplication in an in-situ OAM-domain, sequence numbers can be added to packets of a particular tube (see [\[I-D.hildebrand-spu-d-prototype\]](#)). Each tube leverages a dedicated namespace for its sequence numbers.

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	2	3	4	5	6	7	8	9	
Option Type										Opt Data Len										OAM-E2E-Type										reserved										
E2E Option data format determined by iOAM-E2E-Type																																								~

Option Type: 8-bit identifier of the type of option.

Opt Data Len: 8-bit unsigned integer. Length of the Option Data field of this option, in octets.

OAM-E2E-Type: 8-bit identifier of a particular in situ OAM E2E variant.

0: E2E option data is a 64-bit sequence number added to a specific tube which is used to identify packet loss and reordering for that tube.

Reserved: 8-bit. (Reserved Octet) Reserved octet for future use.

[4.](#) In-situ OAM Data Export

In-situ OAM nodes collect information for packets traversing a domain that supports in-situ OAM. The device at the domain edge (which could also be an end-host) which receives a packet with in-situ OAM information chooses how to process the in-situ OAM data collected

Brockners, et al.

Expires May 3, 2017

[Page 18]

Internet-Draft

In-situ OAM Data Formats

October 2016

within the packet. This decapsulating node can simply discard the information collected, can process the information further, or export the information using e.g., IPFIX.

The discussion of in-situ OAM data processing and export is left for a future version of this document.

[5.](#) IANA Considerations

IANA considerations will be added in a future version of this document.

[6.](#) Manageability Considerations

Manageability considerations will be addressed in a later version of this document..

[7.](#) Security Considerations

Security considerations will be addressed in a later version of this document. For a discussion of security requirements of in-situ OAM, please refer to [[I-D.brockners-inband-oam-requirements](#)].

[8.](#) Acknowledgements

The authors would like to thank Eric Vyncke, Nalini Elkins, Srihari

Raghavan, Ranganathan T S, Karthik Babu Harichandra Babu, Akshaya Nadahalli, LJ Wobker, Erik Nordmark, and Andrew Yourtchenko for the comments and advice. This document leverages and builds on top of several concepts described in [[I-D.kitamura-ipv6-record-route](#)]. The authors would like to acknowledge the work done by the author Hiroshi Kitamura and people involved in writing it.

[9.](#) References

[9.1.](#) Normative References

[I-D.brockners-inband-oam-requirements]
Brockners, F., Bhandari, S., Dara, S., Pignataro, C., Gredler, H., Leddy, J., and S. Youell, "Requirements for In-band OAM", [draft-brockners-inband-oam-requirements-01](#) (work in progress), July 2016.

Brockners, et al.	Expires May 3, 2017	[Page 19]
-------------------	---------------------	-----------

Internet-Draft	In-situ OAM Data Formats	October 2016
----------------	--------------------------	--------------

[IEEE1588v2]
Institute of Electrical and Electronics Engineers, "1588-2008 - IEEE Standard for a Precision Clock Synchronization Protocol for Networked Measurement and Control Systems", IEEE Std 1588-2008, 2008, <<http://standards.ieee.org/findstds/standard/1588-2008.html>>.

[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>>.

[9.2.](#) Informative References

[I-D.brockners-inband-oam-transport]
Brockners, F., Bhandari, S., Pignataro, C., Gredler, H., Leddy, J., and S. Youell, "Encapsulations for In-band OAM Data", [draft-brockners-inband-oam-transport-01](#) (work in

progress), July 2016.

[I-D.hildebrand-spud-prototype]

Hildebrand, J. and B. Trammell, "Substrate Protocol for User Datagrams (SPUD) Prototype", [draft-hildebrand-spud-prototype-03](#) (work in progress), March 2015.

[I-D.kitamura-ipv6-record-route]

Kitamura, H., "Record Route for IPv6 (PR6) Hop-by-Hop Option Extension", [draft-kitamura-ipv6-record-route-00](#) (work in progress), November 2000.

[I-D.lapukhov-dataplane-probe]

Lapukhov, P. and r. remy@barefootnetworks.com, "Data-plane probe for in-band telemetry collection", [draft-lapukhov-dataplane-probe-01](#) (work in progress), June 2016.

[RFC7665] Halpern, J., Ed. and C. Pignataro, Ed., "Service Function Chaining (SFC) Architecture", [RFC 7665](#), DOI 10.17487/RFC7665, October 2015, <<http://www.rfc-editor.org/info/rfc7665>>.

Authors' Addresses

Brockners, et al.

Expires May 3, 2017

[Page 20]

Internet-Draft

In-situ OAM Data Formats

October 2016

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

Email: 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

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

Email: 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

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