

Telemetry Report Format Specification

(working draft)

The P4.org Applications Working Group
Contributions from *Barefoot Networks*, *VMware*, *Xilinx*

2018-3-14

Contents

1. Introduction	2
1.1. Scope	2
2. Key Concepts	3
2.1. Telemetry Report Definition	3
2.2. Telemetry Report Associations	3
2.3. Telemetry Report Events	4
2.4. Telemetry Modes	4
2.4.1. Postcard mode	4
2.4.2. In-band (In-situ) Telemetry mode	4
2.4.3. Using Different Telemetry Modes for Different Telemetry Categories	7
2.5. Correlation of Telemetry Reports	7
3. Telemetry Report Formats	7
3.1. Outer Encapsulation	7
3.1.1. UDP header (8 octets)	8
3.2. Telemetry Report Header (16+ octets)	8
3.3. Telemetry Opaque Data Header	9
3.4. Flow Identification	10
3.5. Embedded Telemetry Metadata	10
3.6. Parsing Considerations	10
A. Acknowledgements	14
B. Change log	14

1. Introduction

Traditional network monitoring has relied on statistics and probe packets such as ICMP echo requests/replies. Recent innovations provide greater insight into network behavior by generating detailed reports of telemetry metadata such as paths, queue occupancy, latency experienced by data packets, and timestamps that can be used to determine hop-by-hop and end-to-end delay. Generation of telemetry reports can be triggered by various events in categories such as flow monitoring, queue congestion, and packet drops. For further information regarding the motivation and usage of detailed telemetry information can be found in the IETF draft for In-situ OAM ¹.

Specifications are being defined for embedding telemetry metadata within data packets, such as INT ² and IOAM ³. This allows for telemetry metadata to be collected as packets traverse a network. When the packets reach the edge of the network, the telemetry metadata is removed and telemetry reports are generated.

This specification defines packet formats for telemetry reports from data plane network devices (e.g. switches) to a distributed telemetry monitoring system. The packet formats use headers that describe the contents of telemetry reports, along with existing (non-telemetry specific) packet headers that can be used to categorize flows.

1.1. Scope

The scope of this specification is interoperability between network devices that generate telemetry reports based on what they see in the data plane, and the initial preprocessors within distributed telemetry monitoring systems that receive the telemetry reports. This specification is applicable when telemetry reports are generated by network devices at the edges of a network, with source and transit network devices embedding telemetry metadata in data packets according to specifications

¹Requirements for In-situ OAM, [draft-brockners-inband-oam-requirements-03](#), March 2017.

²[In-band Network Telemetry \(INT\)](#), October 2017.

³Data Fields for In-situ OAM, [draft-ietf-ippm-ioam-data-01](#), October 2017.

such as INT ² and IOAM ³. This specification is also applicable when each network device directly generates telemetry reports (including transit network devices in the middle of the network), without affecting data packet formats between successive network devices.

Telemetry report encapsulation formats are defined that allow for the inclusion of additional telemetry metadata, beyond the (optional) telemetry metadata embedded between other packet headers as defined in INT and IOAM. The embedded telemetry metadata is included as is in telemetry reports, so the packet formats defined in INT and IOAM also define some aspects of the telemetry report format. See Section 3.5 for further discussion.

This specification does not address any of the following, which are considered out of scope:

- Configuration of network devices so that they can determine when to generate telemetry reports, and what information to include in those reports, such as SAI DTel ⁴.
- Events that trigger generation of telemetry reports.
- Selection of particular destinations within distributed telemetry monitoring systems, to which telemetry reports will be sent.
- Export format for flow statistics or summarized flow records such as IPFIX ⁵.

2. Key Concepts

2.1. Telemetry Report Definition

We define a telemetry report as a message that a network device sends to the monitoring system. A telemetry report carries a snapshot of the original data packet (mostly the inner + outer headers), which triggered the reporting, together with additional telemetry metadata collected from the reporting network device, and possibly from its upstream network devices (in case of in-band mechanism like INT or IOAM). The report message is encapsulated by IP+UDP, hence it can be forwarded from the reporting network device through the data network, and to the destination monitoring system.

The following sections will cover the details on the report generation, report format and encapsulation.

2.2. Telemetry Report Associations

There are many reasons why users may want telemetry reports to be generated. This specification currently considers three categories for telemetry report generation:

Tracked Flows

Telemetry reports are generated matching certain flow definitions. A telemetry specific access control list (called a flow watchlist in this specification) determines which data packets to monitor by matching packet header fields and optionally identification of the ingress interface. (Note that the telemetry specific watchlist is not performing any access control. It only makes decisions related to monitoring actions.) The expectation is that telemetry reports can be generated for those packets that match the flow watchlist. The telemetry reports include information about the path that packets traverse as well as other telemetry metadata such as hop latency and queue occupancy.

Dropped Packets

Telemetry reports are generated for all dropped packets matching a telemetry specific access control list (called a drop watchlist in this specification). This provides visibility into the impact of packet drops on user traffic.

²In-band Network Telemetry (INT), October 2017.

³Data Fields for In-situ OAM, [draft-ietf-ippm-ioam-data-01](#), October 2017.

⁴SAI Data Plane Telemetry Proposal, November 2017.

⁵Specification of the IP Flow Information Export (IPFIX) Protocol for the Exchange of Flow Information, [RFC 7011](#), September 2013.

Congested Queues

Telemetry reports are generated for traffic entering a specific queue during a period of queue congestion. This provides visibility into the traffic causing and prolonging queue congestion, for example a few large elephant flows that overwhelm a queue, as well as the victim traffic (mice flows) getting hurt by the congestion. This also enables the detection and “re-play” of a short microburst, caused by a large number of mice flows arriving at the queue at the same time.

Each telemetry report may be associated with one or more of these categories. This is indicated in the telemetry report by defining association bits, one for each category, as will be shown in Section 3. New categories (and corresponding association bits) may be added to future versions of this specification.

Network devices will need to be configured so that they can determine when to generate telemetry reports, and what information to include in those reports. Such configuration is considered to be beyond the scope of this specification. See ⁴ for one API proposal to enable data plane telemetry capabilities in network devices across all three categories.

2.3. Telemetry Report Events

Telemetry reports are typically triggered by packet processing at a network device. However, even when processed packets match a watchlist for a telemetry report category, it is not necessary for each inspected packet to trigger generation of a telemetry report. Network devices may apply filters to determine when significant events occur that should be reported. This is called event detection in this specification. For example, a network device may trigger telemetry report generation whenever a packet matching a tracked application flow is received or transmitted on a different path than previous packets, or if a significant change in latency is experienced at one particular hop.

Determination of which packets trigger reports, in other words the specific conditions and logic to determine the events of interest, is left open for implementations to differentiate themselves, and is considered to be beyond the scope of this specification.

2.4. Telemetry Modes

There are two different modes which differ with regard to the locations from which telemetry reports are generated.

2.4.1. Postcard mode

In the postcard mode, each network device generates its own telemetry reports, as shown in Figure 1. The distributed telemetry monitoring system will receive reports from different network devices, each describing the telemetry metadata (such as switch IDs, port IDs, latency) for one hop. There is no change to data packets traversing the network. When using postcard mode, the telemetry metadata precedes the original packet headers within the telemetry report.

2.4.2. In-band (In-situ) Telemetry mode

In the other telemetry mode, telemetry metadata is embedded in between the original headers of data packets as they traverse the network, as shown in Figure 2. This may be done using any of the telemetry data plane specifications such as INT or IOAM. When a packet enters the network, the source switch may insert a telemetry instruction header, thereby instructing downstream switches to add the desired telemetry metadata. At each hop, the transit switch inserts its telemetry metadata. The sink switch extracts the telemetry instruction header before progressing the original packet. Depending on the result of event detection, the sink switch may generate a telemetry report containing all of the telemetry metadata from all hops across the network.

⁴[SAI Data Plane Telemetry Proposal](#), November 2017.

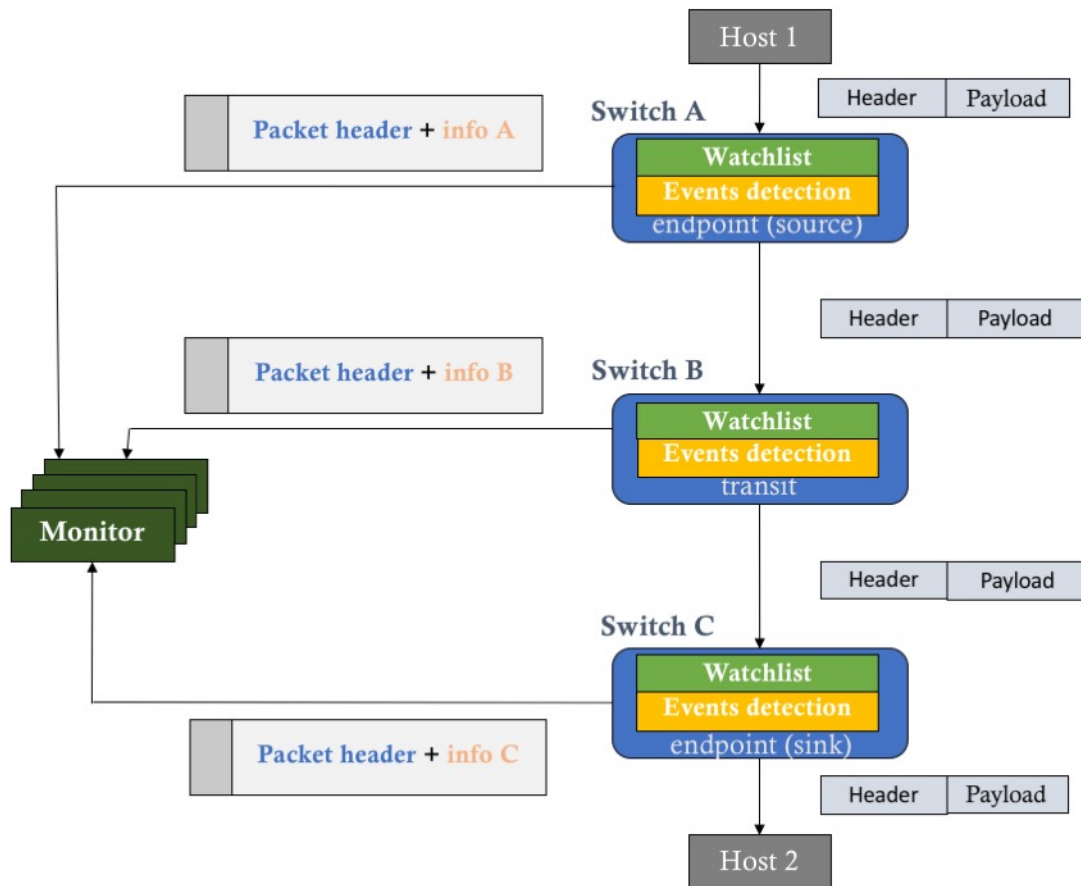


Figure 1. Telemetry Architecture with reports generated by all switches, aka postcard.

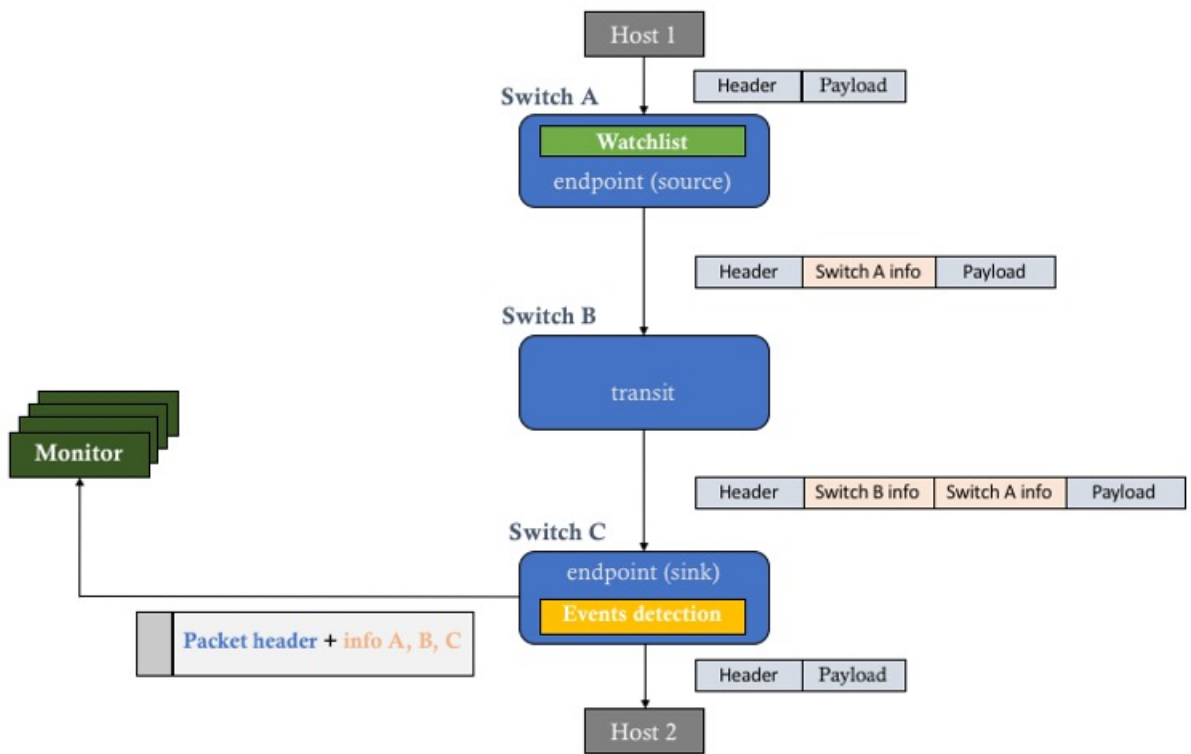


Figure 2. Telemetry Architecture with reports generated by sink switches.

In order to reduce complexity at the sink switch, some telemetry reports may include embedded telemetry metadata intermingled with the original packet headers. This simplifies generation of telemetry reports due to receipt of data packets with embedded telemetry metadata. The telemetry data plane specification such as INT or IOAM specifies the format for this portion of the telemetry metadata. This approach reduces data plane complexity, allowing for all telemetry report processing and generation to be done in the data plane itself without any need to punt to the control plane for further processing.

The sink switch has the option to add its local telemetry metadata either in the telemetry report header defined in this specification, or in the embedded telemetry metadata intermingled with the original packet headers.

2.4.3. Using Different Telemetry Modes for Different Telemetry Categories

Even when in-band (in-situ) telemetry mode is used for the category of tracked flows, it is possible to use the postcard telemetry mode for other categories such as dropped packets and congested queues. The latter categories are often monitored as per switch, per port, or per queue local events, suggesting that telemetry reports should be generated directly from the affected switch(es).

2.5. Correlation of Telemetry Reports

Telemetry reports for a specific application flow matching a flow watchlist may be received from multiple network devices. In case of postcard mode, each hop will generate a separate telemetry report. Even when telemetry metadata is embedded in the data plane according to a specification such as INT or IOAM, telemetry reports for one flow may still be generated by multiple network devices in case of path change or in case of dropped packets.

The distributed telemetry monitoring system may want to correlate these telemetry reports, based on the original packet header fields included in each telemetry report. The telemetry reports include one association bit for each telemetry report category, providing hints to the distributed telemetry monitoring system that it can use to assist with telemetry report correlation. In particular, the distributed telemetry monitoring system may want to apply certain types of telemetry report correlation only when the corresponding bits are set.

The mechanisms for correlation are left to each implementation, and are considered to be beyond the scope of this specification.

3. Telemetry Report Formats

This section specifies the packet formats for telemetry reports.

3.1. Outer Encapsulation

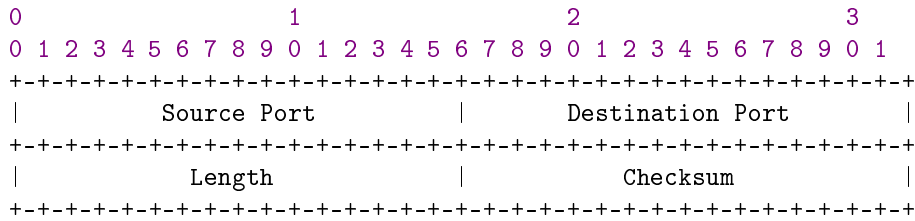
Telemetry reports are defined using a UDP-based encapsulation. Various outer encapsulations may be used to transport the UDP packets. Typically this would simply be an Ethernet header, followed by an IPv4 or IPv6 header, followed by the UDP header. This specification does not preclude the use of different transport encapsulations.

The source IP address identifies the network device that generates the telemetry report.

The Destination IP address identifies a location in the distributed telemetry monitoring system that will receive the telemetry report.

In case of IPv4, as is the case for any other IP packet, either the Don't Fragment (DF) bit must be set, or the IPv4 ID field must be set so that the value does not repeat within the maximum datagram lifetime for a given source address/destination address/protocol tuple.

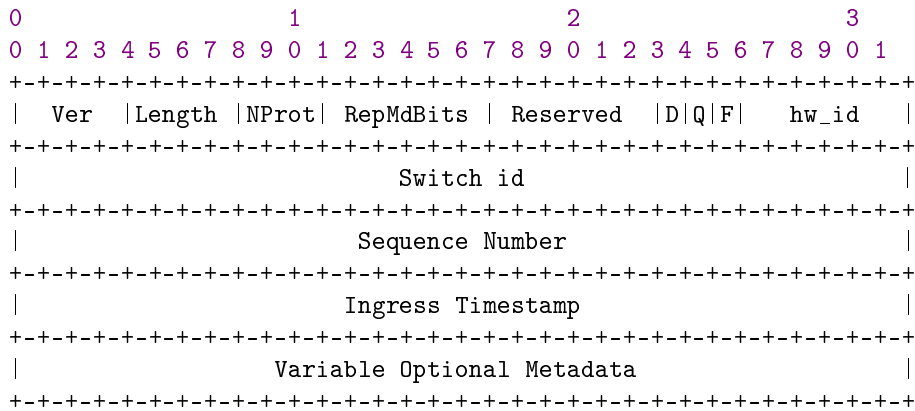
3.1.1. UDP header (8 octets)



The Source Port may optionally be used to carry flow entropy, for example based on a hash of the inner 5-tuple. Otherwise, it should be set to 0.

The Destination Port is user configurable. The expectation is that the same Destination Port value will be used for all telemetry reports in a particular deployment.

3.2. Telemetry Report Header (16+ octets)



Ver: Version

This specification defines **version 1**.

Length

Indicates the length of the telemetry report header in multiples of 4 octets, including the variable optional metadata.

NProt: Next Protocol

- 0: Ethernet
- 1: IPv4
- 2: IPv6
- 3: Telemetry opaque data

RepMdBits: Report Metadata Bits

Bitmap that indicates which optional metadata is present in the telemetry report header. Each bit represents 4 octets of optional metadata.

- bit 0 (MSB): Ingress port id (16 bits) + Egress port id (16 bits)
- bit 1: Hop latency
- bit 2: Queue id (8 bits) + Queue occupancy (24 bits)
- bit 3: Egress Timestamp (32 bits)
- bit 4: Queue id (8 bits) + Drop reason (8 bits) + Padding (16 bits)
- bit 5: Egress port tx utilization

This specification defines the following metadata:

Drop reason

An enumeration that indicates the reason why a packet was dropped, for example as defined in github.com/p4lang/switch.

See the INT specification ² for definitions of the remaining metadata.

D: Dropped

Indicates that at least one packet matching a drop watchlist was dropped.

Q: Congested Queue Association

Indicates the presence of congestion on a monitored queue.

F: Tracked Flow Association

Indicates that this telemetry report is for a tracked flow, i.e. the packet matched a flow watchlist somewhere (in case of INT or IOAM) or locally (in case of postcard). The report might include INT or IOAM metadata beyond the inner ethernet header. Other telemetry reports are likely to be received for the same tracked flow, from the same network device and (in case of drop reports, postcard or path changes) from other network devices.

hw_id

Identifies the hardware subsystem within the network device that generated this report. For example, in a chassis with multiple linecards this could identify a specific linecard, or a subsystem within a linecard. The hw_id is unique within the scope of a Switch id.

Switch id

The unique ID of a switch (generally administratively assigned). Switch IDs must be unique within a management domain.

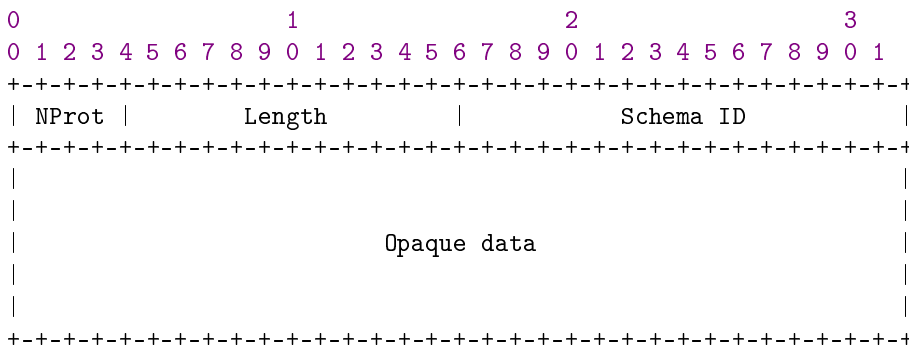
Sequence Number

Reflects the sequence of reports from a specific combination of (Switch id, hw_id) to a particular telemetry report destination. This can be used to detect loss of telemetry reports before they reach their intended destination.

Ingress Timestamp

The device local time when the packet was first received on the ingress physical or logical port, in nanoseconds.

3.3. Telemetry Opaque Data Header



NProt: Next Protocol

Same definition as NProt in the Telemetry Report Header defined in Section 3.2.

Length

Indicates the length of the telemetry opaque data header in multiples of 4 octets.

Schema ID

²[In-band Network Telemetry \(INT\)](#), October 2017.

2-octet unsigned integer identifying the schema of the opaque data.

Opaque data

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

3.4. Flow Identification

There is no explicit metadata defined for flow identification. The expectation is that packet headers will be included in the telemetry report, allowing the distributed telemetry monitoring system to categorize and identify flows in any manner that it desires.

3.5. Embedded Telemetry Metadata

There may still be further telemetry metadata embedded within the payload after the Telemetry Report header. For example, this is typically the case when there is telemetry metadata from hops prior to the network device generating the report. The telemetry metadata will typically be encoded using a defined data plane format such as INT and IOAM.

A network device generating a telemetry report may include its local telemetry metadata in any of the following:

- the embedded telemetry metadata,
- the Telemetry Report header in the same telemetry report as the embedded telemetry metadata from previous hops, or
- the Telemetry Report header in a separate telemetry report from the embedded telemetry metadata from previous hops. Note that in this case the ingress timestamp will be the same in the Telemetry Report Header in both telemetry reports.

If the Tracked Flow Association bit is set to 0, then there will not be any embedded telemetry metadata in the report.

If the Tracked Flow Association bit is set to 1, there may or may not be any embedded telemetry metadata in the report. See Section 3.6 for parsing considerations.

3.6. Parsing Considerations

When a telemetry report is received by the distributed telemetry monitoring system, it must parse the packet to retrieve the telemetry metadata and to identify the flow. Figure 3 shows which headers will be present at the beginning of the packet, assuming a simple Ethernet/IP transport of the telemetry report packet.

The packet format after this point can vary depending on the format of the original packet, and whether embedded telemetry metadata is present. The following figures show a few examples of the remaining packet format. These examples are not intended to be complete or exclusive.

Figure 4 shows the remaining packet format when the original packet is a simple flat packet and there is no embedded telemetry metadata.

Figure 5 shows the remaining packet format when the original packet is a simple flat packet and there is embedded INT over TCP/UDP/ICMP telemetry metadata.

Even when using INT over TCP/UDP/ICMP, the original packet may be an encapsulated packet such as a VXLAN packet. When processing a telemetry report for an encapsulated packet, the distributed telemetry monitoring system may desire to categorize flows based on inner headers. In this case, it should parse the telemetry report all the way down past any embedded telemetry metadata (if present), even when the Telemetry Report Header includes optional metadata such as drop reason. It may also want to process the embedded telemetry metadata, for example to

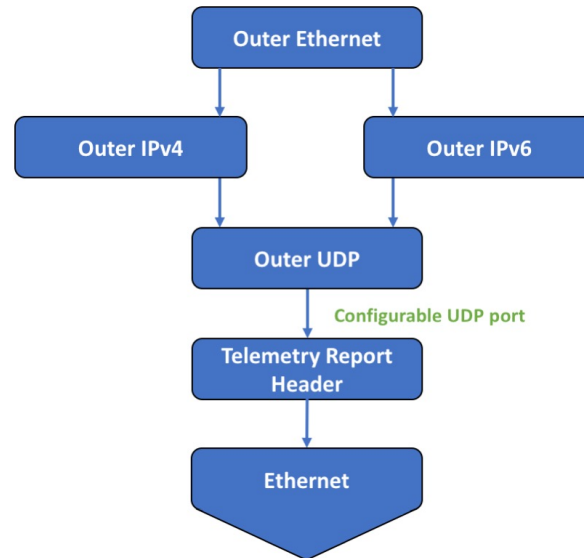


Figure 3. Telemetry Report Outer Encapsulation Format

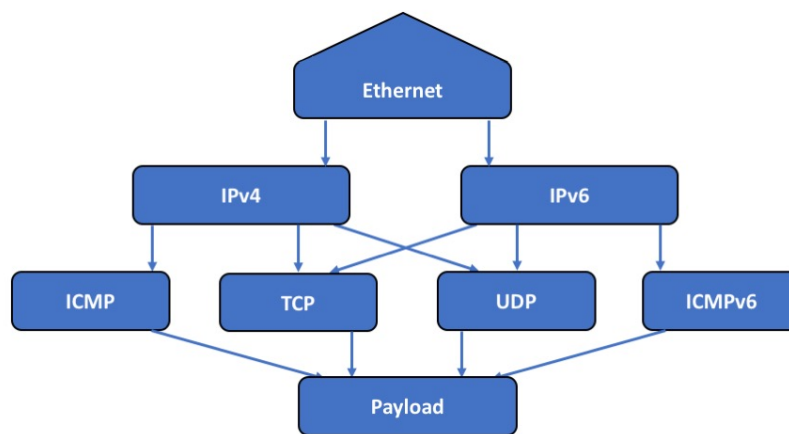


Figure 4. Remaining Packet Format - Flat Packet

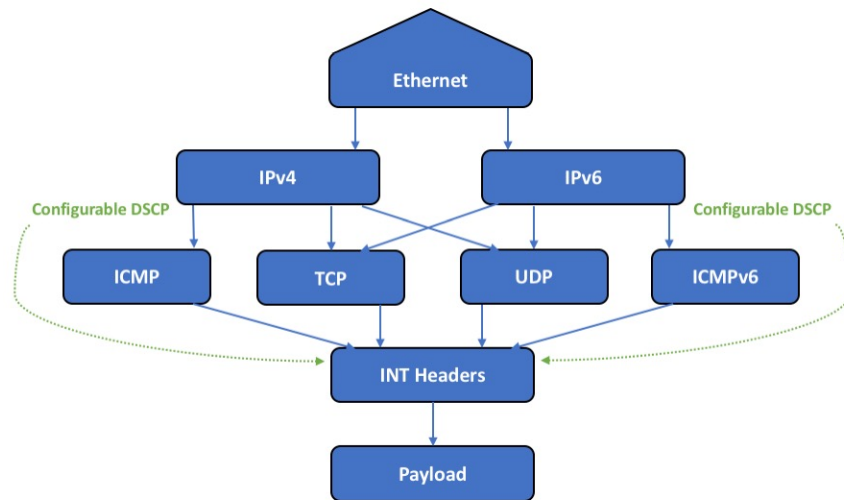


Figure 5. Remaining Packet Format - Flat Packet with INT over TCP/UDP/ICMP

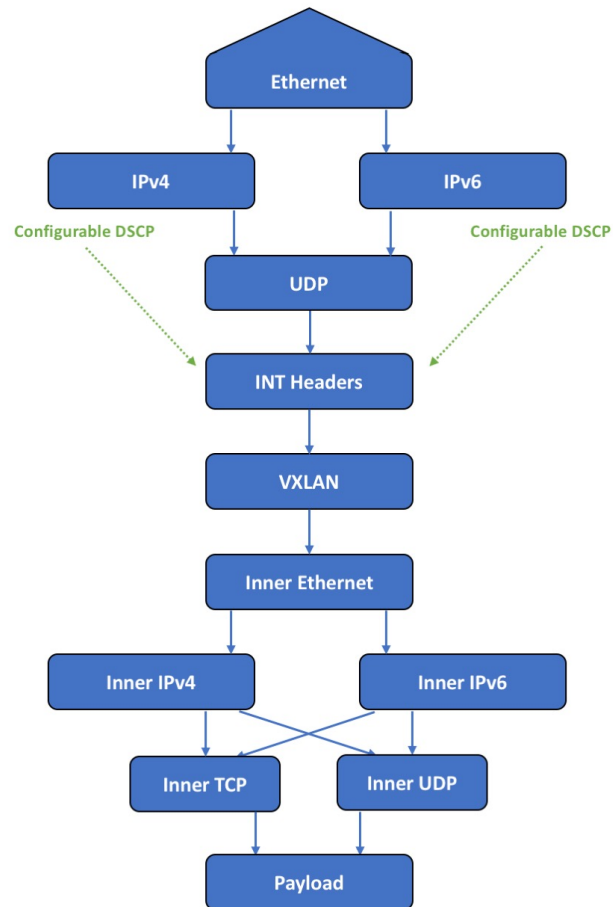


Figure 6. Remaining Packet Format - VXLAN Packet with INT over TCP/UDP

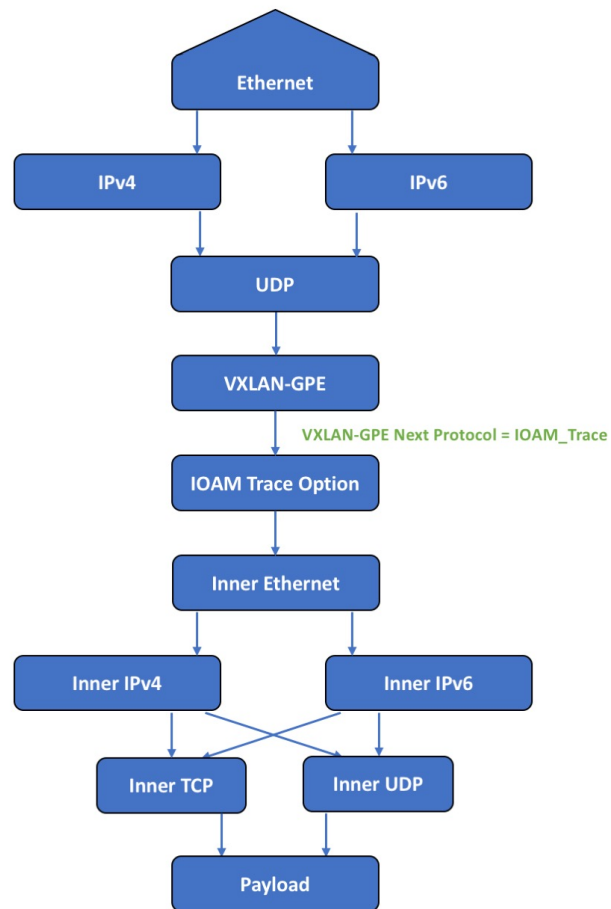


Figure 7. Remaining Packet Format - VXLAN Packet with IOAM Trace

recognize the case where a path change directs traffic to a congested switch where packets are being dropped.

Figure 6 shows the remaining packet format when the original packet is a VXLAN packet and there is embedded INT over TCP/UDP/ICMP telemetry metadata.

Figure 7 shows the remaining packet format when the original packet is a VXLAN packet and there is embedded IOAM Trace telemetry metadata. See IOAM drafts ³⁶ for further details.

A. Acknowledgements

We thank the following individuals for their contributions to this specification.

- Gordon Brebner
- Mukesh Hira
- Jeongkeun Lee
- Mickey Spiegel

B. Change log

- 2017-11-10
 - Initial release
- 2018-2-14
 - Promote *Switch id* to fixed portion of the Telemetry Report Header
 - * The Switch id is always present.
 - Flexible format allowing for arbitrary combinations of optional telemetry metadata in the Telemetry Report Header
 - * Replaces previous Telemetry Drop Header and Telemetry Switch Local Report Header
 - * Adds a 4 bit length field indicating the Telemetry Report Header length in multiples of 4 octets
 - * Adds a bitmap indicating which optional metadata is present in the Telemetry Report Header
 - * Rearranges fields in the first 32 bits of the Telemetry Report Header in order to achieve proper alignment, and to place reserved bits between the report metadata bitmap and the association bits so that either one can expand as necessary
- 2018-3-14
 - Add telemetry opaque data header

³Data Fields for In-situ OAM, [draft-ietf-ippm-ioam-data-01](#), October 2017.

⁶VXLAN-GPE Encapsulation for In-situ OAM Data, [draft-brockners-ioam-vxlan-gpe-00](#), October 2017.