

Technical Report of NetView

Category	Primary Metadata	Description
Switch Metadata	switch_id	The unique id of a switch within a management domain.
	control_plane_state_version	Whenever a control plane state changes (e.g., IP FIB update), the control plane can also update this version number in the data plane.
Ingress Metadata	ingress_port_id	The port on which the packet was received.
	ingress_timestamp	The device local time when the packet was received on the ingress port.
	ingress_port_rx_pkt_count	Total number of packets received on the ingress port where the packet was received.
	ingress_port_rx_byte_count	Total number of bytes received on the ingress port where the packet was received.
	ingress_port_rx_drop_count	Total number of packet drops occurred on the ingress port where the packet was received.
	ingress_port_rx_utilization	Current utilization of the ingress port where the packet was received.
Egress Metadata	egress_port_id	The id of the port forwarding the packet.
	egress_timestamp	Device local time when the packet leaves the egress port.
	egress_port_tx_pkt_count	Total number of packets forwarded through the egress port.
	egress_port_tx_byte_count	Total number of bytes forwarded through the egress port where the packet was forwarded.
	egress_port_tx_drop_count	Total number of packet drops occurred on the egress port where the packet was forwarded.
	egress_port_tx_utilization	Current utilization of the egress port via which the packet was sent out.
Buffer Metadata	queue_id	The id of the queue the device used to serve the packet.
	instantaneous_queue_length	The instantaneous length (in bytes, cells, or packets) of the queue the packet has observed in the device while being forwarded.
	average_queue_length	The average length (in bytes, cells, or packets) of the queue via which the packet was served.
	congestion_status	The ratio of the current queue length to the configured maximum queue limit.
	queue_drop_count	Total number of packets dropped by the queue.

Figure 1. Primary metadata of NetView.

I. METADATA AND API

To simplify telemetry policy enforcement for operators, we provide a suite of expressive and convenient telemetry APIs. Then, we present several non-exhaustive telemetry applications to show the usefulness of NetView.

1) *Supported Metadata*: NetView supports two kinds of metadata: primary metadata and derived metadata. (1) Primary metadata, as shown in Figure 1, is directly supported by data plane [1], comprising switch metadata (e.g., switch_id), ingress metadata (e.g., ingress_port_rx_pkt_count), buffer metadata (e.g., instantaneous_queue_length), and egress metadata (e.g., egress_port_tx_utilization). In this paper, we take all of the primary metadata into consideration, towards fine-grained visibility of the network. (2) Derived metadata

Category	Derived Metadata	Description
Path-level Metadata	path_length	The number of hops in the given path.
	path_trace	Each hop information in the given path in sequence, i.e., "node:port, node:port...".
	path_rtt	Round trip time in the given path.
	path_utilization	Link utilization of each link in the given path.
	path_throughput	Instantaneous throughput of each link in the given path.
Node-level Metadata	switch_workload	Total traffic on all in/out ports of the switch.
	forward_packet_count	Total number of packets forwarded in the switch.
	drop_packet_count	Total number of packets dropped in the switch.
	alive_port_count	Total number of alive ports in the switch.
Probe-level Metadata	tx_probe_count	Total number of sent probes.
	tx_probe_timestamp	The timestamp of the first/last sent probe.
	rx_probe_count	Total number of received probes.
	rx_probe_timestamp	The timestamp of the first/last received probe.

Figure 2. Derived metadata of NetView.

is high-level empirical metadata, which can be calculated by primary metadata. As is shown in Figure 2, it contains path-level metadata (e.g., path_utilization), node-level metadata (e.g., switch_workload), and probe-level metadata (e.g., rx_probe_timestamp). Particularly, we can tackle path-level metadata by approximating it into primary metadata at the node. As for the delay, it consists of four parts: transmission delay, propagation delay, queueing delay, and processing delay. The queueing delay and processing delay represents network performance deviation among the above four parts. As for the throughput, the path throughput is decided by the minimal port throughput along the path. Therefore, without measuring delay (throughput) along the path, instead, NetView measures queueing delay (port throughput) of each node, and use node-level queueing delay (port throughput) to compose path-level information.

2) *Network Telemetry API*: NetView provides declarative APIs that express queries for a wide range of telemetry tasks.

- PathQuery("SrcNode:SrcPort", "DstNode:DstPort") means that query the information of (a) path(s) between the port of source node and the port of destination node.
- NodeQuery("Node:Port") means that query the information of (a) given node(s).
- ProbeQuery("TerminalIP") means that query the information of probes from a single terminal.
- Where("Scope") means that specify the metadata scope.
- Period("DesiredTime") means that set a timer to query network periodically.
- Return("Metadata") means that return a derived/primary metadata as the telemetry result.

Application	Query	Description
End-to-end Throughput Measurement	Q = PathQuery("10.1.0.1:1", "10.2.0.1:1") .Where("path_throughput < 1Gbps") .Period("2s") .Return("path_throughput ")	Measure the path throughput between the port 1 of switch 1 and the port 1 of switch 2, which conforms to the limitation.
Queue Length Measurement	Q = NodeQuery("*,*") .Period("5ms") .Return("instantaneous_queue_length")	Measure the instantaneous length of the queue in all of switches every 5ms.
Node Black Hole Discovery	Q = NodeQuery("*,*") .Where("switch_workload == 0") .Period("10s") .Return("switch_id")	Locate switch id of the node black hole by discovering the switch not transmitting packets.
Random Packet Drop Detection	Q = NodeQuery("10.0.*,*") .Where("ingress_port_rx_drop_count > 0" "egress_port_tx_drop_count > 0") .Period("1ms") .Return("switch_id", "ingress_port_id", "egress_port_id")	Detect ports where packet drop happens in the first pod.

Figure 3. Sample applications of NetView.

3) *Network-Telemetry-Based Applications*: Operators can write queries to express various telemetry applications, including the query for the entire data center (global query) or just for parts of it (non-global query), *e.g.*, query for an end-to-end path or query for switches in a pod. As Figure 3 shows, we introduce non-exhaustive telemetry applications as examples. (a) Prove Network Innocence. Operators can let NetView periodically measure end-to-end throughput, identifying whether end hosts or the network should be blamed for service interruption. (b) Network Planning. With NetView, Operators can accumulate a large volume of telemetry data (*e.g.*, queue length) for the long term network capacity planning and topology augmentation. (c) Failure Location. Operators can query parts of the network prone to failure or even the entire data center network to detect sporadic failure in the first place. (d) Event Tracking. Numerous network events are of interest to network operators. For example, Network operators always want to learn where and when packets are dropped for an application flow in time. With NetView, network operators can easily get a vital clue (*e.g.*, random packet drop count) about non-reproducible events.

II. THE CHOICE OF k_{Fwd}

It is worth noting that the mainstream commercial data center topology is quite regular (*e.g.*, Fat-Tree, VL2, Leaf-Spine), making sure that the probe can return to the vantage server within 5 hops (three-tier topology) or 3 hops (two-tier topology) no matter where the probe is now in the data center. Thus, as long as k_{Fwd} is larger than 9 (three-tier topology) or 5 (two-tier topology), the probe can surely return to the vantage server after traveling in the data center.

REFERENCES

- [1] Ed Doe Hugh Holbrook Anoop Ghan-wani Dan Daly Mukesh Hira Changhoon Kim, Parag Bhide and Bruce Davie. In-band network telemetry. Website, 2016. <https://p4.org/assets/INT-current-spec.pdf>.