

# Enhancing perfSONAR Measurement Capabilities using P4 Programmable Data Planes

Ali Mazloun

University of South Carolina  
Columbia, SC, USA  
amazloun@email.sc.edu

Elie Kfoury

University of South Carolina  
Columbia, SC, USA  
ekfoury@email.sc.edu

Jose Gomez

University of South Carolina  
Columbia, SC, USA  
gomezgaj@email.sc.edu

Jorge Crichigno

University of South Carolina  
Columbia, SC, USA  
jcrichigno@cec.sc.edu

## ABSTRACT

As a key element in the Science Demilitarized Zone (Science DMZ) framework, perfSONAR is a crucial tool for monitoring and troubleshooting network performance. This paper presents a system that leverages the flexibility and granularity of P4 programmable data planes to enhance perfSONAR measurements. P4 is a programming language used to define how network packets are processed in the data plane. The proposed system employs a P4 programmable data plane that operates passively on real traffic. This approach ensures real-time per-flow monitoring and detailed insights, assisting administrators in understanding network traffic patterns. Moreover, the system detects microbursts and seamlessly integrates with a regular perfSONAR node. Experimental evaluations show that the proposed system provides full visibility of the real traffic, reports problems undetectable by perfSONAR, and notifies administrators if an observed low performance is caused by the endpoints of the connection (i.e., the network is not the bottleneck).

## KEYWORDS

perfSONAR, Science DMZ, P4, data transfer node, congestion control, Bandwidth-Delay Product (BDP), router's buffer size.

## 1 INTRODUCTION

In an era where scientific and engineering endeavors are generating unprecedented volumes of data, efficient and reliable network infrastructure has become paramount. From the expansive data produced by cutting-edge facilities such as the Large Hadron Collider [1] to the intricate genetic insights derived from portable DNA sequencers [2], the need for seamless data transfer across networks is more pronounced than ever. While conventional enterprise networks can handle operational data, they face significant challenges when dealing with the influx of terabyte and petabyte-scale data generated by these scientific instruments.

In response to these challenges, the Science DMZ has emerged [3]. Aimed to facilitate large data transfers, the Science DMZ provides a dedicated network segment, designed explicitly for the high-speed transmission of substantial scientific data. At its core, the Science DMZ is supported by four key pillars: specialized data transfer nodes

(DTNs) used for high-rate data movement, high-throughput paths characterized by frictionless data flow, performance measurement mechanisms for real-time network assessment, and security protocols optimized for the unique demands of high-performance scientific environments.

While Science DMZs play an essential role in tackling data transfers, a critical dimension consists of advanced monitoring and analytical capabilities to visualize and troubleshoot performance issues. At the core of the Science DMZ, perfSONAR [4] emerges as the tool to effectively implement and execute the crucial tasks of monitoring network performance and troubleshooting. Developed by a collaborative effort among research and education institutions, perfSONAR offers a comprehensive framework for measuring, analyzing, and diagnosing network issues. It enables network administrators, researchers, and operators to gain valuable insights into network performance, ensuring the efficient and reliable operation of their infrastructure. Although perfSONAR is a key tool used to diagnose and troubleshoot network issues, it presents shortcomings such as time accuracy, packet level analysis, and custom network measurements.

This paper proposes a system that enhances perfSONAR measurements by integrating passive P4 programmable data planes. P4 is a programming language that defines how network devices process network data packets [5]. It enables the customization of data plane behavior, allowing network operators to specify how packets are parsed, processed, and forwarded, providing high flexibility and programmability in network operations.

The system operates passively on a copy of the real traffic directed to a P4 programmable data plane. This approach minimizes network interference and latency while offering real-time performance insights, making it particularly well-suited for critical and high-performance networks. With this approach, the system implements per-flow monitoring and provides detailed reports for active flows. The reports include individual flow throughput, round-trip time (RTT), packet loss, queue occupancy, and more. This granularity helps network administrators fine-tune network configurations and in troubleshooting problems.

Furthermore, the proposed system introduces microburst detection, a metric critical for network stability. Microbursts are rapid packet influxes that can impair network performance. The system's nanosecond-level reporting of microbursts empowers administrators to pinpoint and address performance bottlenecks. Besides, the system seamlessly integrates with perfSONAR, aligning with the

pSConfig command and the default archiver for streamlined reporting and data storage. Administrators can adjust reporting rates and set alert thresholds, enhancing coordination between the system and perfSONAR. These functionalities are implemented as an extension of the pSConfig command.

## 1.1 Contributions

This paper leverages the flexibility of P4 programmable data planes to enhance the network measurements produced by a perfSONAR node. The main contributions of this paper can be framed as follows:

- Enhancing monitoring without actively engaging with the network. Passive P4 programmable data planes operate over a copy of the network traffic, minimizing interference with regular network operations.
- Introducing real-time performance reporting and per-flow detailed statistics. The system enhances perfSONAR’s monitoring abilities on a per-packet basis to improve visibility and accurate network analysis.
- Detecting and reporting microbursts. The system offers the capability to pinpoint and address issues arising from microbursts, a challenge often overlooked by conventional monitoring tools.
- Seamless integration with perfSONAR. The system incorporates pSConfig model-based configuration and utilizes the perfSONAR archiver for data storage, facilitating control over reporting rates and alert thresholds throughout a perfSONAR node.

## 1.2 Paper Organization

This paper is organized as follows. Section 2 provides background on Science DMZ, perfSONAR, and P4 programmable data planes. Section 3 describes the proposed system. Section 4 describes the implementation of the system. Section 5 illustrates the performed experiments and use cases. Section 6 presents the related works. Section 7 concludes the paper.

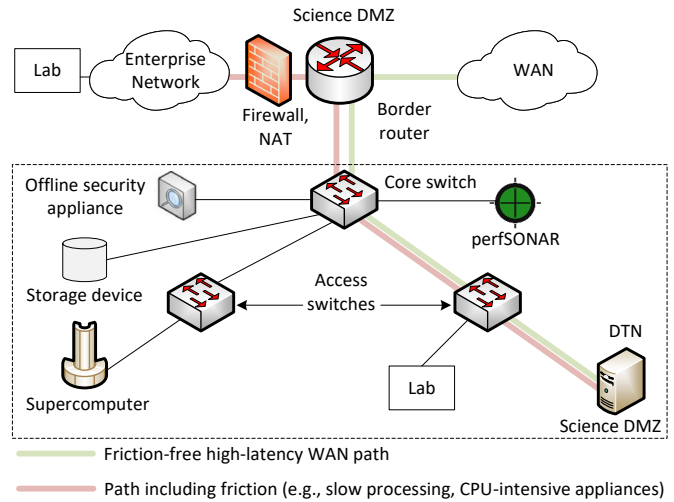
## 2 BACKGROUND

### 2.1 Science DMZ

The Science DMZ is a network architecture specifically designed to address the unique data transfer and performance requirements of large-scale scientific research and data-intensive applications [5]. It emerged as a response to the escalating demand for efficient and high-speed data transfer capabilities within research institutions and facilities, where conventional network infrastructures present limitations to accommodate the massive volumes of data generated by modern scientific experiments and simulations.

Typically, enterprise networks are optimized for general-purpose data communication. However, they are inadequate for the seamless and rapid transmission of vast datasets inherent to scientific endeavors. The need for specialized infrastructure to accommodate the data challenges of scientific research led to the development of the Science DMZ framework.

The Science DMZ architecture is characterized by the following elements [6]:



**Figure 1: A Science DMZ co-located with an enterprise network.**

- **Dedicated, high-Performance Data Transfer Nodes (DTNs):** Science DMZs include specialized DTNs optimized for high-speed data movement. These nodes are equipped with hardware and software configurations that enable efficient and rapid data transfer across Wide Area Networks (WANs).
- **High-Throughput, friction-free Paths:** Science DMZs establish dedicated, high-throughput network paths that minimize bottlenecks and ensure friction-free data movement between endpoints, including instruments, storage systems, and computing resources.
- **Performance Monitoring and Measurement devices:** Comprehensive performance monitoring and measurement mechanisms are integral to Science DMZs. These tools continuously assess network metrics such as bandwidth, latency, jitter, and packet loss to optimize data transfer and troubleshoot potential issues.
- **Specialized Security Policies:** Science DMZs implement security policies and configurations optimized for scientific applications while facilitating data movement.

Fig. 1 illustrates a basic architecture of a Science DMZ co-located with an enterprise network. Notice the absence of a firewall or any stateful inline security appliance in the friction-free path. As data sets traverse from the WAN to a DTN, they can either be locally stored within the DTN or written in a storage device. DTNs have a dual-homed configuration featuring a secondary interface linked to the storage device. This strategy enables the DTN to concurrently receive data from the WAN and transfer it to the storage device, avoiding redundant copying. Within the Science DMZ, users in an internal laboratory can have friction-free access to the data residing on the storage device, experiencing uninterrupted availability. In contrast, users located in a laboratory within the enterprise network, protected by a firewall, experience reasonable performance when accessing stored data within the Science DMZ. This differentiation results from the reduced latency between the Science DMZ and enterprise users, which diminishes the performance impact caused by

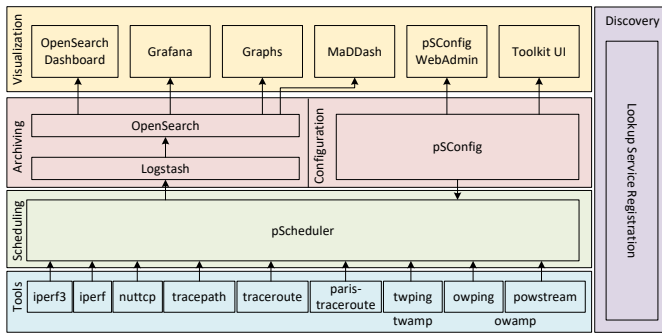


Figure 2: perfSONAR architecture [4].

retransmissions resulting from security appliances. Unlike the gradual recovery observed in high-latency WAN situations, transport protocols facilitate applications to recover from packet losses when operating under low-latency conditions [3]. The main challenge lies in providing long-distance data transfers with a friction-free service.

## 2.2 perfSONAR

PerfSONAR is a monitoring framework designed to measure the performance of network infrastructures. It plays a crucial role in enhancing network operations, troubleshooting issues, and optimizing data transfer. Developed through collaborative efforts within the research and education community, perfSONAR addresses the challenges posed by the ever-increasing demand for high-speed and reliable data transfer across distributed networks. The rapid growth of data-intensive applications, such as scientific research, requires efficient and dependable network connectivity. Traditional monitoring tools often are limited in providing accurate and real-time insights into network performance. PerfSONAR emerged as a response to the need for a comprehensive and standardized solution to measure, diagnose, and enhance network operations in research, education, and scientific environments.

Fig. 2 shows the components of perfSONAR architecture. PerfSONAR comprises a modular architecture integrating open-source tools (e.g., iperf3, ping, traceroute) and perfSONAR-specific modules (e.g., pScheduler, pSConfig, MaDDash, and others). The implementation of perfSONAR consists of distributed software agents deployed across network nodes. These nodes collaborate to collect, analyze, and report network performance metrics. The framework is designed to be vendor-neutral and adaptable to different network environments, making it suitable for a wide range of applications. These perfSONAR nodes can coordinate measurement tests using standard protocols to gather performance data, such as throughput, latency, packet loss, and jitter. The collected data is aggregated and presented in a user-friendly format.

Network administrators can configure perfSONAR nodes to measure and monitor network performance continuously. These nodes exchange measurement data to provide insights into various aspects of network behavior, such as congestion, latency variation, and link utilization. Administrators can access the collected data through web-based interfaces, enabling them to make informed decisions to optimize network performance.

## 2.3 P4 programmable data planes

P4 is a programming language that defines how network devices process network data packets. It enables the customization of data plane behavior, allowing network operators to specify how packets are parsed, processed, and forwarded, providing high flexibility and programmability in network operations.

Traditional network devices operate on fixed, proprietary hardware and software architectures, which limits the type of processing that can be performed by a switch or router. The surge in diverse applications, from cloud computing to IoT, demanded a network infrastructure that can dynamically adapt and scale. P4 emerged in response to this demand, enabling the creation of programmable data planes that can be tailored to the unique needs of various use cases [5, 7–9].

Network devices equipped with P4-compatible hardware can be programmed to implement specific packet processing functions. P4 language, a high-level programming language, is used to define packet parsing, processing, and forwarding instructions. These instructions are compiled into target-specific configurations that are loaded onto the programmable devices.

When deployed, P4 programmable data planes execute the custom-defined packet processing rules. This enables the network to dynamically adapt to changing traffic patterns, optimize resource utilization, and enhance overall performance. P4 enables the implementation of innovative algorithms and protocols, making it possible to support new applications and services seamlessly.

Although the latest version of perfSONAR (version 5) aggregates the tests into a single value, the authors are aware that perfSONAR can be as granular as the tool it uses (e.g., iPerf, ping). The latest perfSONAR processes the measurements collected by the Tools layer using Logstash [4]. By default, Logstash is programmed to aggregate the result of a perfSONAR test. For throughput tests, Logstash only reports the average value. For RTT tests, Logstash reports the minimum, the maximum, and the mean RTT. Because the default configuration does not reflect the maximum granularity supported by perfSONAR, Logstash’s filters can be adjusted. However, this minor granularity addition only applies to the performed active tests and does not allow perfSONAR to have per-flow visibility. Having per-flow visibility is one of the main goals that perfSONAR is aiming to achieve which cannot be done without operating passively on the real traffic [10].

## 3 PROPOSED SYSTEM

### 3.1 Overview

The proposed system utilizes a P4 programmable switch to passively monitor live traffic. The system is depicted in Fig. 3. A legacy switch (core switch) connects the Science DMZ to the Internet. A TAP duplicates the traffic arriving at and departing from the legacy switch to the programmable switch.

The data plane of the programmable switch collects per-flow measurements where flows are characterized by their 5-tuple (source IP, destination IP, source port, destination port, and protocol). The collected measurements are extracted by the switch’s control plane for additional processing before being reported to a perfSONAR archiver. The rate of measurement extraction by the control plane is a variable that a perfSONAR node can configure. The proposed system

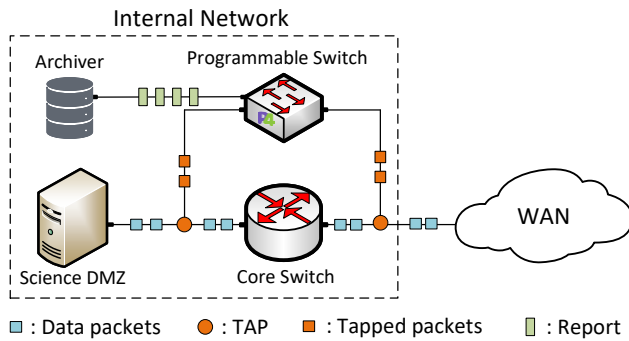


Figure 3: Proposed system overview.

also allows the administrator to define an alerting threshold for each monitored metric. If the threshold is exceeded, the switch alerts the administrator and increases the rate of measurement collection in order to get higher visibility. The number of reports per second, the alerting threshold, and the increase in the reporting rate when the threshold is exceeded are all customizable and can be configured from the configuration layer of the perfSONAR node (i.e., using pSConfig).

### 3.2 Architecture

As depicted in Fig. 4, the proposed system has four main components. The first component is the data plane of the programmable switch. The data plane is responsible for monitoring the traffic. The second component is the control plane of the switch. The control plane is responsible for extracting and processing the measurements from the data plane before sending them to the perfSONAR archiver, which represents the third component in the proposed system. Besides storing the collected data, the archiver allows the real-time visualization of the metrics extracted from the data plane. The last component is the configuration daemon. The daemon runs on a perfSONAR node and configures the control plane of the programmable switch.

As indicated by Fig. 5(a), the configuration layer of perfSONAR is responsible for configuring the time intervals during run time. The four time intervals,  $t_N$ ,  $t_P$ ,  $t_R$ , and  $t_Q$ , represent the rate at which the control plane extracts the number of bytes, packet losses, RTT, and queue occupancy measurements, respectively from the data plane. This layer can also configure alerting thresholds ( $a_N$ ,  $a_P$ ,  $a_R$ , and  $a_Q$ ) for the four metrics. If one of the alerting thresholds is exceeded, the control plane notifies the administrator and increases the rate of collection to a value defined by the administrator.

The interaction between the data plane and the control plane is shown in Fig. 5(b). The data plane receives packets from  $n$  different flows representing the real traffic. For each of the  $n$  flows, the data plane continuously monitors four metrics: the number of bytes, the round-trip time (RTT), the retransmission count (packet loss), and the queuing delay (queue occupancy). The control plane utilizes the APIs provided by the switch manufacturer to access the measurements maintained by the data plane at run-time. The control plane extracts and reports the measurements of all the flows simultaneously, as shown in Fig. 5(c).

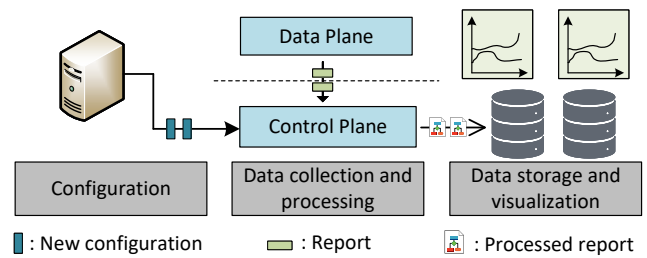


Figure 4: Proposed system architecture.

### 3.3 Features

**3.3.1 Passive Measurement.** Unlike active measurement, passive measurement does not add overhead or consume network resources. Through tapping, the programmable switch operates on a copy of the actual network traffic. The switch does not actively participate in the network, and consequently, there is no risk of increasing network latency, bottlenecks, or affecting the overall performance. This makes passive monitoring suitable for critical and high-performance networks where minimal interference is crucial [11].

**3.3.2 Per-flow Monitoring.** The proposed system significantly enhances the monitoring capabilities of perfSONAR. Typically, perfSONAR runs periodic active measurements to monitor the network. This provides a periodic estimation of the state of the network (e.g., available bandwidth); however, it lacks the ability to report the performance of the real traffic. The proposed system enhances the visibility of perfSONAR by reporting the network performance in real-time. Besides, the system provides per-flow detailed reports. The data plane updates the statistics of each flow on a per-packet basis. The statistics are continuously updated and maintained by dedicated stateful registers. Because providing a per-packet report overwhelms the data collector (i.e., perfSONAR archiver), the administrator has the flexibility to configure the reporting interval, which can be in the order of milliseconds.

The proposed system also provides a detailed report for each terminated long flow. The report uses the nanosecond granularity of the switch to provide the flow's start and end times. The report also includes the total number of packets sent, the total number of bytes, the average throughput, and the number and percentage of retransmissions. This report helps administrators understand the type of traffic they are interacting with and consequently helps them tune their network.

**3.3.3 Microburst Detection.** In addition to enhancing the visibility of perfSONAR, the proposed system reports microbursts, a metric not supported by any tool used by perfSONAR. Microbursts occur when a massive amount of packets arrives in a very short duration. The duration of a microburst can be lower than a hundred microseconds [12], making it undetectable without per-packet visibility [13]. Yet, microbursts can significantly reduce the performance of networks. The burst of packets fills a router's buffer, increasing the flows' RTT. If the burst is big enough, it might bloat the buffer, leading the router to drop all the packets arriving during the burst. Because of the nature of some TCP congestion control algorithms, packet losses can lead to significant performance degradation. The

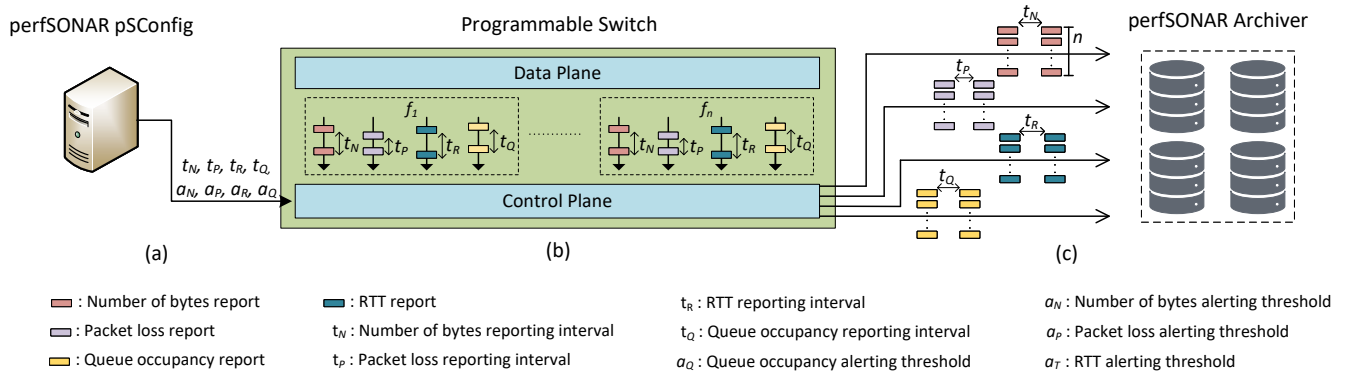


Figure 5: The interaction between the different components of the proposed system.

proposed system is able to report the starting time and the duration for any microburst with nanosecond granularity. To do so, the switch monitors the queue occupancy on a per-packet basis. The sudden increase in queue occupancy is reported as a microburst.

**3.3.4 Identifying if the Connection is Limited by the Network or by the Sender/Receiver.** In a network, connections can encounter various levels of restrictions, stemming from different elements of the communication process: the sender, the receiver, or the network infrastructure itself [14]. The proposed system is capable of determining whether these limitations arise from the network or from the sender/receiver.

Specifically, the system can identify whether a connection’s performance is restricted by the network or by the sender/receiver. If the network is the source of the limitation and the underlying cause isn’t related to queuing issues (e.g., small buffer size, microbursts), it might be necessary to conduct active measurements to precisely pinpoint the root of the problem.

However, when a connection’s restriction originates from the sender or receiver, engaging in active measurements is not advisable. This is due to two key reasons. Firstly, introducing active measurements in such cases would introduce additional overhead to the network, potentially exacerbating the issue. Secondly, these measurements would likely fail to accurately identify the underlying problem when it’s sender/receiver-related. DTN nodes do not interfere with perfSONAR tests, and consequently, perfSONAR measurements cannot pinpoint the problems caused by the DTNs themselves.

In essence, the system’s ability to differentiate between network-related and sender/receiver-related limitations offers valuable insights into the appropriate course of action. Active measurements

can be selectively applied based on the nature of the limitation, ensuring efficient troubleshooting.

**3.3.5 Seamless Integration with perfSONAR.** To support seamless integration with perfSONAR, the proposed system had to satisfy two conditions: 1) the programmable switch should be configurable through the configuration layer of perfSONAR (i.e., through pSConfig model), and 2) the switch should use perfSONAR archiver to store the collected data. Regarding the first requirement, new functionality has been added to the pSConfig model, which is responsible for configuring the switch’s control plane at run time. The added functionality allows a perfSONAR node to control the number of reports per second sent by the control plane to the archiver. Consider Fig. 6. *config-P4* is the added command to pSConfig through which the node configures the control plane. The *--metric* parameter specifies to which metric the configuration is to be applied. The *--samples\_per\_second* specifies the number of samples to be reported by the control plane to the archiver. The first line sets the rate of throughput reports to 1 per second, and the second line sets the number of RTT reports to 2 per second. The configuration will be applied to all metrics if the administrator does not use the *--metric* parameter. The *config-P4* command also allows the administrator to define a threshold for each metric that will trigger an alert if its value is exceeded. A triggered alert will notify the administrator and increase the reporting rate. In the third line, the rate of queue occupancy reports will be set to 10 reports per second if the queue occupancy exceeds 30%.

## 4 IMPLEMENTATION

To support per-flow monitoring, the data plane of the switch is programmed to group packets into flows using the hash of the 5-tuple. After that, the data plane detects long flows using count-min sketches (CMS) [15]. After detecting a long flow, the data plane reports the ID of the flow (i.e., the hash of the 5-tuple), its source and destination IP, and its reversed ID. The reversed ID is calculated by reversing the hashing order of the source and destination fields of the 5-tuple (the source and destination IP and the source and destination ports). The reversed ID is used to identify the flow to which an acknowledgment packet (ACK) belongs, which is essential for calculating RTT.

```
1. psconfig config-P4 --metric throughput --samples_per_second 1
2. psconfig config-P4 --metric RTT --samples_per_second 2
3. psconfig config-P4 --metric queue_occupancy --alert --threshold 30
   --samples_per_second 10
```

Figure 6: Configuration example.

## 4.1 Throughput Monitoring

For each flow, the programmable switch calculates the number of bytes, the queuing delay, the packet losses, and the RTT. The calculated measurements are stored inside stateful registers which are indexed by the flow ID. The most straightforward metric to calculate in the data plane is the number of bytes. The programmable switch utilizes the total length field of the IPv4 header to get the packet's length in bytes. After that, the switch updates the register cell that maintains the number of bytes of the current flow. The control plane uses the APIs provided by the manufacturer of the switch to read the register values. The throughput is then calculated in the control plane by dividing the number of bits over the reporting duration.

## 4.2 Queue Occupancy Monitoring

Furthermore, although monitoring queue occupancy on a fine-grain resolution is not supported by perfSONAR due to software limitations, it can be easily performed by a programmable switch. The programmable switch calculates the time packets spent inside the core switch to calculate its queue occupancy. The TAP duplicates each packet twice. The first duplication occurs at the ingress port of the core switch, and the second duplication occurs at the egress port of the switch. The programmable switch calculates the time difference between the arrival of the two copies to calculate the queuing delay. The queuing delay of the flow to which the packet belongs is stored in a stateful register. The control plane samples the queuing delay and calculates the queue occupancy. Queue occupancy equals the queuing delay divided by the buffer size of the core switch.

However, because the duration of microbursts can be in the order of tens of microseconds, the sampling approach might not detect them. For this, microburst detection should be fully implemented in the data plane. The data plane measures the queue occupancy on a per-packet basis. Any sudden increase can be instantaneously detected by the data plane. The data plane can then report the event to the control plane.

## 4.3 RTT and Packet loss Monitoring

The last two metrics calculated by the data plane are the RTT and the packet losses. The proposed system adopts RTT and packet loss calculation algorithm from [16]. RTT represents the time between a packet being sent and its corresponding acknowledgment received. As the sender transmits data packets to the receiver, the receiver acknowledges the receipt of these packets by sending acknowledgment packets back to the sender. The acknowledgment number in these packets represents the next expected sequence number, indicating that all bytes up to that number have been successfully received. If the sender receives an acknowledgment packet with an acknowledgment number that is less than the last transmitted sequence number, it signifies that the receiver has not yet confirmed the reception of the data. This inconsistency indicates a requirement for retransmission, which, in turn, implies a packet loss scenario.

Any packet with a non-zero payload might be acknowledged by a future ACK number. The sequence number of the packets is used to calculate the expected future ACK number (eACK). When an ACK packet is received, the ACK number is extracted and compared against the previously stored eACKs. If a match occurs, the

---

### Algorithm 1: RTT and packet loss calculation

---

```

Data plane():
  hdr ← pkt.extract(headers);
  flow_ID ← hash(5_tuple);
  pkt_type ← get_type(hdr.tcp.flag, hdr.ip.total_len);
  if pkt_type == Seq then
    prev_seq_no ← prev_seq_register[flow_ID];
    if hdr.tcp.seq_no < prev_seq_no then
      | pkt_loss_register[flow_ID] += 1;
    else
      | prev_seq_register[flow_ID] ← hdr.tcp.seq_no;
  rev_flow_ID ← hash(rev_5_tuple);
  exp_ack ← hdr.tcp.seq_no + (hdr.ip.total_len - 4 *
    | hdr.ip.ihl - 4 * hdr.tcp.data_offset);
  pkt_sig ← (rev_flow_ID, exp_ack);
  eack_register[pkt_sig] ← time.now();
  else if pkt_type == ACK then
    pkt_sig ← (flow_ID, hdr.tcp.ack_no);
    arrival_time ← eack_register[pkt_sig];
    if arrival_time != 0 then
      | rtt ← time.now() - arrival_time;
      | rtt_register[flow_ID] ← rtt;

```

---

switch decrements the current timestamp from the timestamp of the matched eACK. This time difference is the RTT of the flow.

The pseudocode of the RTT and packet loss calculation program is summarized in Algorithm 1. The data plane first parses the packets' headers and obtains the *flow\_ID*. It then identifies the packet type (i.e., whether it is a data packet or an ACK packet). The packet type is obtained by inspecting its TCP flag and total length. For data packets (referred to by Seq in the pseudocode), the data plane compares the current packet's sequence number with the previous packet's sequence number (*prev\_seq\_no*). If the current sequence number is smaller, the switch increments the count of packet losses because retransmission has occurred. Otherwise, the switch updates the value of the previous sequence number to the current sequence number. After that, the data plane calculates the reversed *flow\_ID* by hashing the 5-tuple in a reversed order. Then, the data plane calculates the eACK number of the current packet. The *flow\_ID* and the eACK are used to create the packet's signature. The arrival time of the packet is then stored in a register indexed by its signature (*eack\_register* in the pseudocode).

For ACK packets, the data plane creates a signature using the *flow\_ID* and the ACK number. After that, the created signature is used to index *eack\_register*. The extracted value represents the timestamp of the data packet that corresponds to the current ACK packet. The RTT is then calculated by subtracting the retrieved timestamp from the current timestamp. The final step is to store the RTT in the *rtt\_register* at the index *flow\_ID*. The control plane obtains the per-flow packet losses and per-flow RTT by reading the *pkt\_loss\_register* and the *rtt\_register*. The control plane can also report the percentage of packet losses by dividing the number of packet losses of a flow by the number of packets sent by the same flow.

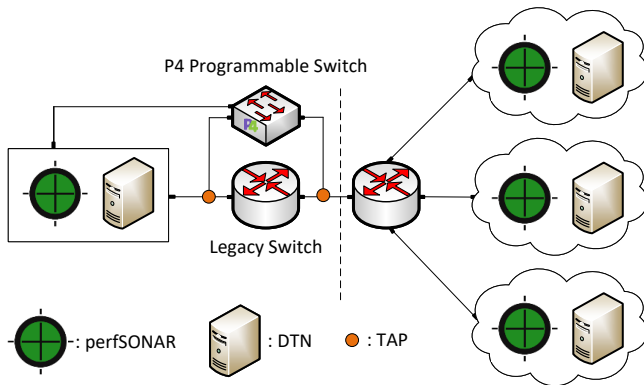


Figure 7: Topology.

#### 4.4 Detecting Flows not Constrained by the Network

Besides calculating different metrics, the data plane can identify if a flow is limited by the network or by the sender/receiver. To determine whether a flow’s limitations stem from network factors or sender/receiver constraints, the programmable switch employs a method that involves monitoring the fluctuation of the flight size (the count of transmitted bytes awaiting acknowledgment) in relation to observed packet losses. If the flight size remains stable despite the absence of packet losses (implying a lack of increment), the limiting factor is likely the sender or the receiver. Conversely, when the flight size expands in conjunction with detected packet losses, the network is likely the influencing factor behind the limitations. This approach is drawn from insights presented in the work by Ghasemi et al. [14].

### 5 EXPERIMENTATION

It’s important to emphasize that this paper does not seek to replace perfSONAR; instead, its objective is to complement and enhance its capabilities. The system’s intention is to augment perfSONAR by providing finer detail to certain metrics (such as throughput, RTT, and packet loss) and by introducing support for novel measurements (including queue occupancy and microbursts). This section will delve into an evaluation of the integration of a P4 programmable switch with perfSONAR. First, the section will highlight the improved granularity that supplements the existing perfSONAR metrics (throughput, RTT, and packet loss). Following that, the section will show the significance of the newly added metrics in enhancing the functionality of perfSONAR.

#### 5.1 Topology Setup

In reference to Fig. 7, the experimental topology encompasses an internal network and three external networks. Each of these networks is equipped with a DTN for data exchange and a perfSONAR node dedicated to monitoring network performance. The interconnections between networks involve two legacy switches. Notably, the link interconnecting these switches acts as a performance bottleneck, operating at a throughput of 5 Gigabits per second (Gb/s), while all other links provide a capacity of 10 Gb/s.

To perform passive measurements, a TAP captures packets arriving at the ingress and egress ports of the legacy switch directly linked to the internal network. These captured packets are then directed to a P4 programmable switch. The programmable switch’s control plane is connected to the local perfSONAR node, leveraging its archiving capabilities to store collected measurements. The data generation process is facilitated by the iPerf3 tool [17], which generates controlled traffic patterns for testing and evaluation purposes. Additionally, to visually present the collected data in real-time, the Grafana visualization platform [18] is employed.

In all tests, there will be an exchange of traffic from the internal DTN to all external DTNs. Unless specified otherwise, the interval used by the programmable switch to report the measurements is one second, and Grafana will group the reported measurements by their destination IP address (i.e., destination DTN). The RTTs between the local DTN and the three external DTNs are 50, 75, and 100 milliseconds (ms), respectively. All perfSONAR tests are initiated from the internal DTN.

#### 5.2 Per-flow Monitoring

This section is dedicated to demonstrating the real-time reporting capabilities of the proposed system, showcasing its ability to provide instantaneous updates on key metrics such as throughput, queue occupancy, RTT, and packet losses for individual flows. To illustrate this functionality, an experiment was conducted, involving the introduction of a third data transfer alongside two pre-existing data transfers. The visual representation of this experiment can be observed in Fig. 8.

Within the depicted figure, the upper left graph presents per-flow throughput, while the bottom left graph presents per-flow RTT. Moving to the upper right graph, queue occupancy is presented, while the bottom right graph presents per-flow packet losses. Through these four graphs, a clear illustration of TCP behavior emerges when a novel flow is introduced into the network scenario. These visualizations provide valuable insights into how the proposed system responds to changing conditions, enabling swift analysis of network dynamics as new flows are incorporated.

Initially, TCP endeavors to distribute the available bandwidth equitably among all active connections. This equitable distribution becomes evident through the observations in the first graph, where the throughput of the three flows converges to an approximate parity (around 1.7 Gbps for each flow).

Additionally, when a fresh TCP connection is established, or following a period of inactivity, the sender’s objective is to ascertain the network’s existing bandwidth allocation and congestion status. To achieve this, the sender initiates a burst of packet transmission. This burst aids the sender in promptly identifying the feasible data transmission volume without inducing congestion. This behavior is illustrated through the analysis of queue occupancy and the count of packet losses. As a new flow becomes part of the network, the queue rapidly fills up, manifesting as a sudden surge in the queue occupancy graph. Notably, during the burst phase, the queue’s capacity is often exceeded, leading to significant packet losses, as depicted in the packet loss percentage graph.

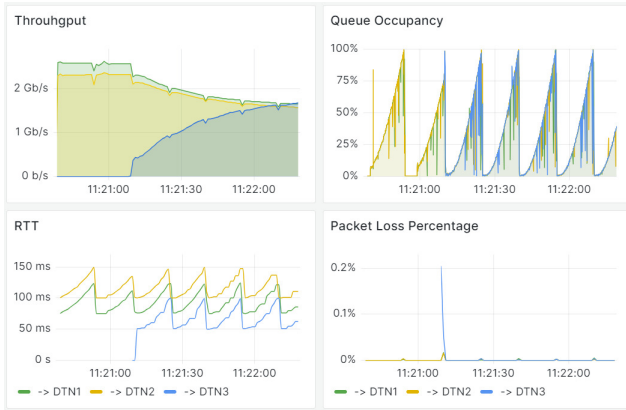


Figure 8: Per-flow measurements.

### 5.3 Additional Traffic Statistics

In addition to the measurements collected by the data plane, the control plane can provide invaluable insights into the network’s traffic dynamics. This is achieved by leveraging access to the measurements of all active flows. Consequently, the control plane gains the capability to extract numerous metrics such as total link utilization, the count of active flows, aggregate packets sent, aggregate bytes transmitted, and more.

Moreover, the control plane performs computations that surpass the data plane’s computational and resource constraints. For instance, the control plane can calculate the fairness among TCP flows. Notably, this involves determining the Jain’s fairness index, a quantifiable fairness metric, by employing the following equation [19]:

$$F = \frac{\left(\sum_{i=1}^N x_i\right)^2}{N \cdot \sum_{i=1}^N x_i^2} \quad (1)$$

The components of the equation are as follows:

- $F$  is the Jain’s fairness index.
- $N$  represents the number of flows.
- $x_i$  corresponds to the resource allocation (e.g., bandwidth) for the  $i$ th flow.

By applying this index, the control plane effectively assesses the equity among different flows, allowing for a comprehensive evaluation of traffic distribution within the network. This index quantifies the fairness of resource allocation among multiple flows, with a value closer to 1 indicating a more equitable distribution of resources. There are situations when one host is bottlenecked by a bandwidth smaller than the others, which reduces fairness. Monitoring the fairness index enhances the overall monitoring and management of network resources, resulting in more optimized and balanced performance.

Illustrated in Fig. 9, the link utilization and fairness of traffic are depicted within the time interval showcased in Fig. 8. Despite the link being fully utilized, the fairness index reveals a notable departure from equitable resource distribution for approximately 20 seconds. This departure implies that although the link was saturated,

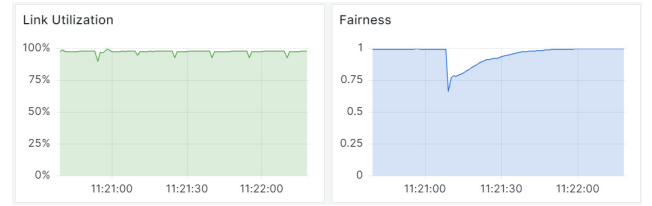


Figure 9: Additional traffic statistics calculated by the control plane.

the available bandwidth wasn’t uniformly shared among the flows during this duration. The fairness index is

To identify the underlying cause of this observed disparity, the fairness graph in Fig. 9 should be cross-referenced with the throughput graph in Fig. 8. This analysis shows that the period of unfairness aligns with the time taken by the three TCP flows to converge subsequent to the introduction of the third flow.

### 5.4 Use Cases

**5.4.1 Detecting small-sized buffers.** In this experiment, the average RTT for the flows is set to 100 ms. As per the established guideline, the buffer size should be one bandwidth-delay product (BDP) [5]. The BDP is calculated as the product of the bandwidth and the RTT, resulting in  $5Gb/s \times 100ms = 500$  Megabits, which translates to 62.5 Megabytes. Accordingly, the buffer size was configured to be  $BDP/2$  or precisely 31.25 Megabytes, representing a small-sized buffer.

The proposed system detects microbursts and shows their impact on queue occupancy, percentage of packet losses, and on throughput. If the queue fails to absorb the microburst, leading to significant packet losses and degradation in the throughput, then the size of the buffer should be reevaluated. The experimental results are illustrated in Fig. 10. Across the three flows, the percentage of lost packets notably escalated, surpassing 0.5% for the first two flows and exceeding 1% for the third. It took approximately 40 seconds for the throughput of these flows to recover. The programmable switch’s data demonstrates that the buffer was unable to efficiently absorb the microburst, signaling a requirement for buffer size adjustment. By tracking the count and magnitude of microbursts, administrators can ascertain an appropriate buffer size that effectively mitigates microburst impact without incurring undue queuing delays.

**5.4.2 Determining if a connection is limited by sender/receiver or by the network.** In this experiment, three tests are performed. In the test involving DTN1, the network is set to be the bottleneck by introducing 0.01% packet losses. In the test involving DTN2, the receiver is set to be the bottleneck by reducing its TCP buffer size. In the test involving DTN3, the sender is set to be the bottleneck by setting its sending rate to 500 Mb/s which is less than the bottleneck link (5 Gb/s).

The results of the experiment are depicted in Fig. 11. The measurements are grouped into flows using the hash of the 5-tuple. The throughput of Flow1, or the flow from the internal DTN to DTN1, is fluctuating because of the induced packet losses. The programmable switch reported that the flow is limited by the network. For Flow2 and Flow3, the throughput is steady at around 250 Mb/s and 500





Figure 10: Microburst detection.

Mb/s, respectively. The programmable switch reported that both flows are limited by either the sender or the receiver but not by the network.

## 6 RELATED WORK

NetSage [20] is a data analysis and visualization platform that combines data from various sources into a single unified view. It is open to the public and accessible via <https://portal.netsage.global/grafana>, and its software is open source. NetSage uses active and passive measurements to provide performance visualizations. It can collect data from routers, switches, active testing sites, and science data archives. NetSage analysis longitudinal data to understand longer-term trends and behaviors.

OSG Network Monitoring Platform [21] is a comprehensive network monitoring platform that supports data collection, processing, storage, and visualization. This platform collects measurements from the tests performed by perfSONAR. This platform stores data in a short-term store location (last six months), in a long-term store location, and in a backup location. Further, the platform has a centralized configuration system that configures the tools to be used in the tests, the participating nodes in the tests, and the test schedule for the entire infrastructure.

Bezerra et al. [22] integrates In-band Network Telemetry (INT) [23] with AmLight to provide sub-second network monitoring and performance evaluation metrics. Programmable switches were utilized to collect measurements at the data plane level. The collected measurement is then directed to AmLight Collector, which is responsible for receiving, parsing, processing, and generating operational reports. Through INT, AmLight was capable of providing instantaneous bandwidth utilization, monitoring packet drop probability and jitter, and tracing packets.

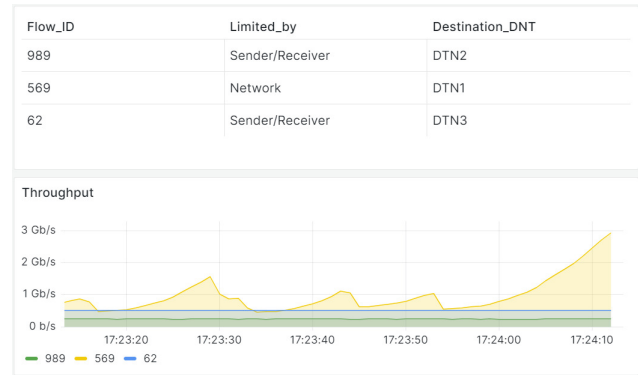


Figure 11: Reporting the flows limited by network and those limited by the sender/receiver.

The proposed system in this paper utilizes programmable switches to enhance the monitoring and troubleshooting functionalities of perfSONAR. The traces produced by perfSONAR are used to understand longitudinal traffic behavior (NetSage and OSG Network Monitoring Platform), detect anomalies in the traffic [24], and train network-wide anomaly detection systems [25]. By providing higher visibility and more detailed reports, the proposed system can enhance the performance of all the platforms that utilize perfSONAR traces. Furthermore, the proposed system allows perfSONAR to detect degradation problems not caused by the network (e.g., the connection is limited by the network).

## 7 CONCLUSION

In summary, the integration of P4 with perfSONAR yields substantial benefits by passively monitoring real traffic in real-time and on a per-flow basis. The real-time passive analysis executed through a P4 programmable switch significantly augments perfSONAR's troubleshooting capabilities. This integration equips perfSONAR with the ability to discern flows constrained by sender/receiver limitations, facilitating informed decisions on the need for active measurements. This synergy not only elevates monitoring precision but also streamlines the diagnostic process for optimized network performance.

## REFERENCES

- [1] European Organization for Nuclear Research, "Towards the future: tackling upcoming data challenges." [Online]. Available: <https://home.cern/about/computing>, Accessed on 08-06-2023.
- [2] Emily Waltz, "Portable DNA Sequencer MinION Helps Build the Internet of Living Things." [Online]. Available: <https://tinyurl.com/fysezkef>, Accessed on 08-06-2023.
- [3] E. Dart, L. Rotman, B. Tierney, M. Hester, and J. Zurawski, "The science DMZ: A network design pattern for data-intensive science," in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, 2013.
- [4] perfSONAR Project, "What is perfSONAR?." [Online]. Available: <https://tinyurl.com/4xef95h9>, Accessed on 08-06-2023.
- [5] J. Crichigno, E. Bou-Harb, and N. Ghani, "A comprehensive tutorial on science DMZ," *IEEE Communications Surveys & Tutorials*, 2018.
- [6] Energy Sciences Network, "Science DMZ." [Online]. Available: <https://fasterdata.es.net/science-dmz/>, Accessed on 08-06-2023.
- [7] A. AlSabeih, J. Khoury, E. Kfoury, J. Crichigno, and E. Bou-Harb, "A survey on security applications of P4 programmable switches and a STRIDE-based vulnerability assessment," *Computer Networks*, 2022.

- [8] J. Gomez, E. F. Kfoury, J. Crichigno, and G. Srivastava, "A survey on TCP enhancements using P4-programmable devices," *Computer Networks*, 2022.
- [9] A. Mazloun, E. Kfoury, J. Gomez, and J. Crichigno, "A survey on rerouting techniques with P4 programmable data plane switches," *Computer Networks*, 2023.
- [10] B. Tierney, J. Metzger, J. Boote, E. Boyd, A. Brown, R. Carlson, M. Zekauskas, J. Zurawski, M. Swany, and M. Grigoriev, "perfonar: Instantiating a global network measurement framework," *SOSP Wksp. Real Overlays and Distrib. Sys.*, vol. 28, 2009.
- [11] B. Erikssoon, P. Barford, and R. Nowak, "Network discovery from passive measurements," in *Proceedings of the ACM SIGCOMM 2008 conference on Data communication*, pp. 291–302, 2008.
- [12] Q. Zhang, V. Liu, H. Zeng, and A. Krishnamurthy, "High-resolution measurement of data center microbursts," in *Proceedings of the 2017 Internet Measurement Conference*, pp. 78–85, 2017.
- [13] X. Chen, S. L. Feibish, Y. Koral, J. Rexford, O. Rottenstreich, S. A. Monetti, and T.-Y. Wang, "Fine-grained queue measurement in the data plane," in *Proceedings of the 15th International Conference on Emerging Networking Experiments And Technologies*, pp. 15–29, 2019.
- [14] M. Ghasemi, T. Benson, and J. Rexford, "Dapper: Data plane performance diagnosis of TCP," in *Proceedings of the Symposium on SDN Research*, pp. 61–74, 2017.
- [15] G. Cormode and S. Muthukrishnan, "An improved data stream summary: the count-min sketch and its applications," *Journal of Algorithms*, vol. 55, no. 1, pp. 58–75, 2005.
- [16] X. Chen, H. Kim, J. M. Aman, W. Chang, M. Lee, and J. Rexford, "Measuring TCP round-trip time in the data plane," in *Proceedings of the Workshop on Secure Programmable Network Infrastructure*, pp. 35–41, 2020.
- [17] "iPerf - the ultimate speed test tool for TCP, UDP and SCTP." [Online]. Available: <https://tinyurl.com/47vzrv8k>, Accessed on 8-14-2023.
- [18] "Grafana platform." [Online]. Available: <https://grafana.com/>, Accessed on 8-14-2023.
- [19] R. Jain, A. Durresti, and G. Babic, "Throughput fairness index: An explanation," in *ATM Forum contribution*, vol. 99, 1999.
- [20] K. Turner, M. Khanal, T. Seto-Mook, A. Gonzalez, J. Leigh, A. Lake, S. S. Baveha, S. Faci, B. Tierney, D. Doyle, *et al.*, "The NetSage Measurement Framework: Design, Development, and Discoveries," in *2020 IEEE/ACM Innovating the Network for Data-Intensive Science (INDIS)*, pp. 45–56, IEEE, 2020.
- [21] M. Babik, S. McKee, B. P. Bockelman, E. M. F. Hernandez, E. Martelli, I. Vukotic, D. Weitzel, and M. Zvada, "Improving WLCG networks through monitoring and analytics," in *EPJ Web of Conferences*, vol. 214, p. 08006, EDP Sciences, 2019.
- [22] J. Bezerra, I. Brito, A. Quintana, J. Ibarra, V. Chergarova, R. Frez, H. Morgan, M. LeClerc, and A. Paneri, "Deploying per-packet telemetry in a long-haul network: the AmLight use case," in *2021 IEEE Workshop on Innovating the Network for Data-Intensive Science (INDIS)*, pp. 44–49, IEEE, 2021.
- [23] E. F. Kfoury, J. Crichigno, and E. Bou-Harb, "An exhaustive survey on P4 programmable data plane switches: Taxonomy, applications, challenges, and future trends," *IEEE access*, 2021.
- [24] P. Calyam, J. Pu, W. Mandrawa, and A. Krishnamurthy, "Ontimedetect: Dynamic network anomaly notification in perfSONAR deployments," in *2010 IEEE International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems*, pp. 328–337, IEEE, 2010.
- [25] Y. Zhang, S. Debroy, and P. Calyam, "Network-wide anomaly event detection and diagnosis with perfSONAR," *IEEE Transactions on Network and Service Management*, vol. 13, no. 3, pp. 666–680, 2016.