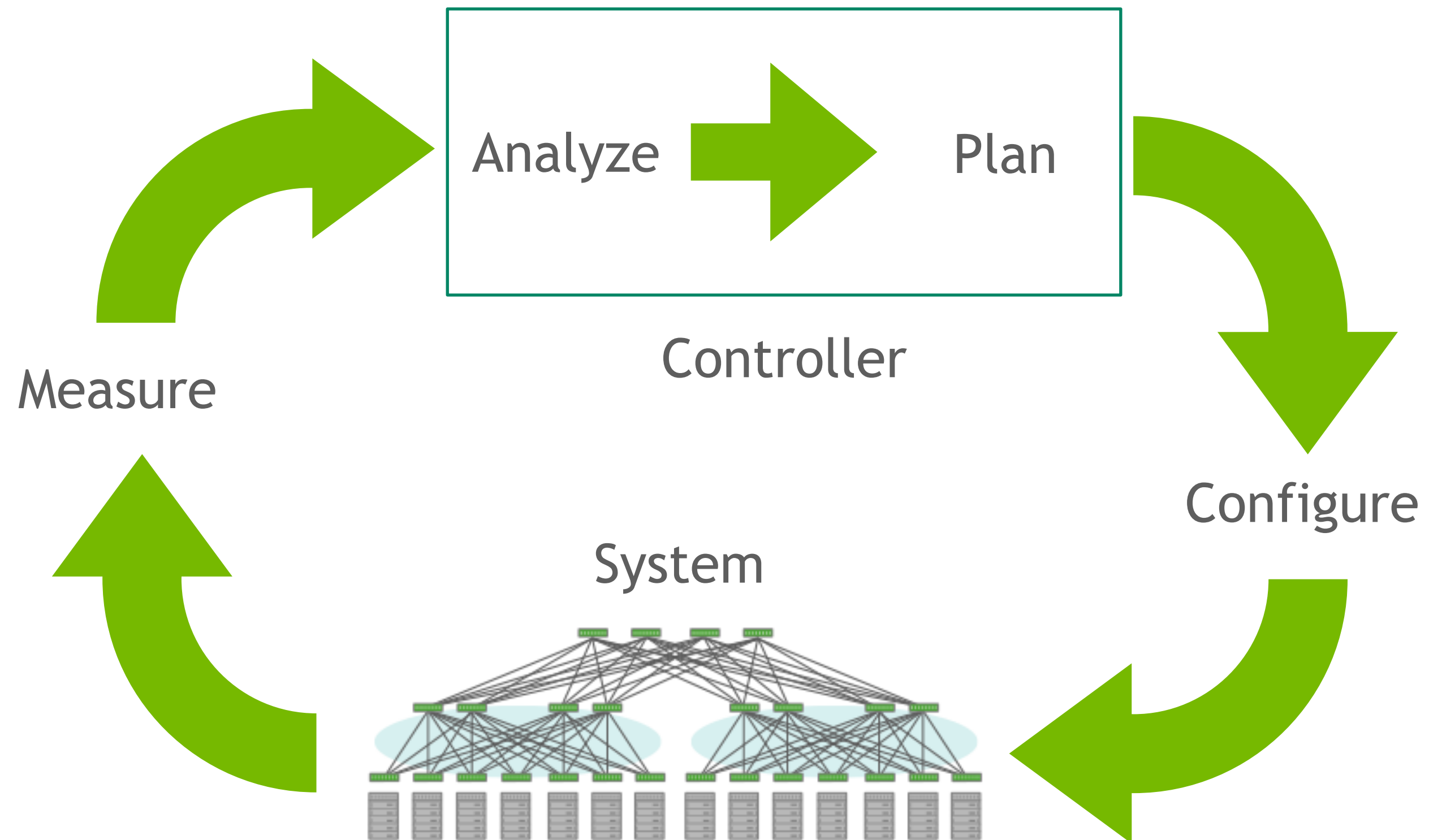


*“If you can’t measure it, you  
can’t improve it”*

— Lord Kelvin

WHY MEASURE?

# FEEDBACK CONTROL



# STABILITY REQUIREMENTS

## ► Observable


- Ability to monitor all important system states
- Example: To safely drive a car you need to see the road ahead, mirrors, and speedometer

## ► Controllable

- Ability to influence all important system states
- Example: To control a car's speed and direction you need to have a steering wheel, accelerator, and brakes

## ► Responsive

- Feedback loop must be fast enough to track changes
- Example: To avoid a car accident you need to react quickly if the car ahead brakes suddenly.

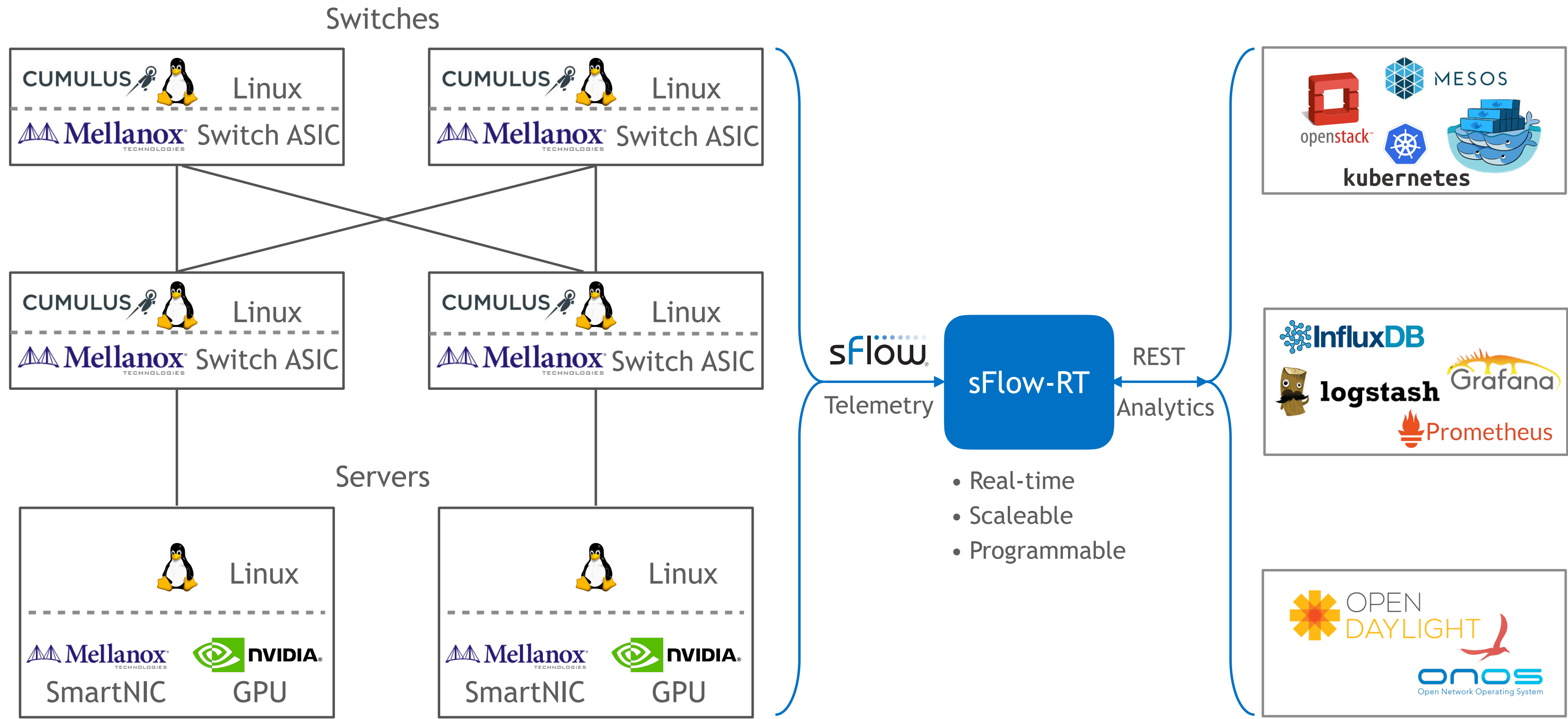


*“In God we trust. All others  
bring data.”*

— Dr. Edwards Deming

WHAT IS SFLOW?

# NETWORK AND SYSTEM OBSERVABILITY



Embedded monitoring of all switches, all servers, all applications, all the time

Consistent measurements shared between multiple management tools

# STANDARD COUNTERS

- ▶ Network (maintained by hardware in network devices)
  - MIB-2 ifTable: ifInOctets, ifInUcastPkts, ifInMulticastPkts, ifInBroadcastPkts, ifInDiscards, ifInErrors, ifUnknownProtos, ifOutOctets, ifOutUcastPkts, ifOutMulticastPkts, ifOutBroadcastPkts, ifOutDiscards, ifOutErrors
- ▶ Host (maintained by operating system kernel)
  - GPU: device\_count, processes, gpu\_time, mem\_time, mem\_total, mem\_free, ecc\_errors, energy, temperature, fan\_speed
  - CPU: load\_one, load\_five, load\_fifteen, proc\_run, proc\_total, cpu\_num, cpu\_speed, uptime, cpu\_user, cpu\_nice, cpu\_system, cpu\_idle, cpu\_wio, cpu\_intr, cpu\_sintr, interrupts, contexts
  - Memory: mem\_total, mem\_free, mem\_shared, mem\_buffers, mem\_cached, swap\_total, swap\_free, page\_in, page\_out, swap\_in, swap\_out
  - Disk IO: disk\_total, disk\_free, part\_max\_used, reads, bytes\_read, read\_time, writes, bytes\_written, write\_time
- ▶ Application (maintained by application)
  - HTTP: method\_option\_count, method\_get\_count, method\_head\_count, method\_post\_count, method\_put\_count, method\_delete\_count, method\_trace\_count, method\_connect\_count, method\_other\_count, status\_1xx\_count, status\_2xx\_count, status\_3xx\_count, status\_4xx\_count, status\_5xx\_count, status\_other\_count

# SFLOW'S SCALABLE PUSH PROTOCOL

## ► Simple

- standard structures containing densely packed blocks of counters
- extensible (tag, length, value)
- RFC 1832: XDR encoded (big endian, quad-aligned, binary) - simple to encode/decode
- unicast UDP transport

## ► Minimal configuration

- collector address
- polling interval

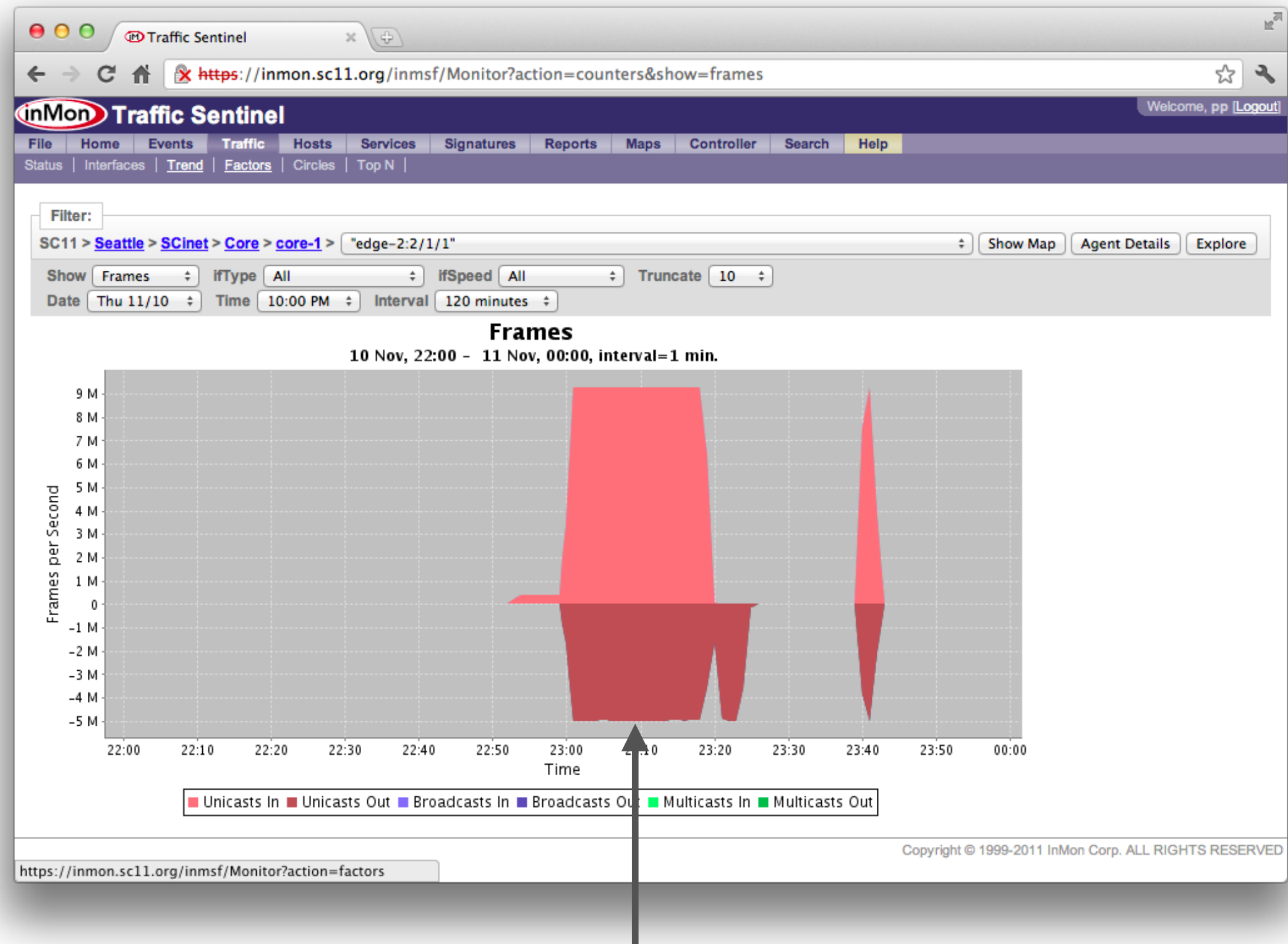
## ► Cloud friendly

- flat, two tier architecture: many embedded agents → central “smart” collector
- sFlow agents automatically start sending metrics on startup, automatically discovered
- eliminates complexity of maintaining polling daemons (and associated configurations)

# COUNTERS AREN'T ENOUGH

Don't tell you why there is a spike

- Counters tell you there is a problem, but not why.
- Counters summarize performance by dropping high cardinality attributes:
  - ip addresses
  - protocols
  - ports
- Need to be able to efficiently disaggregate counter by attributes in order to understand root cause of performance problems
- How do you get this data when there are millions of packets per second?



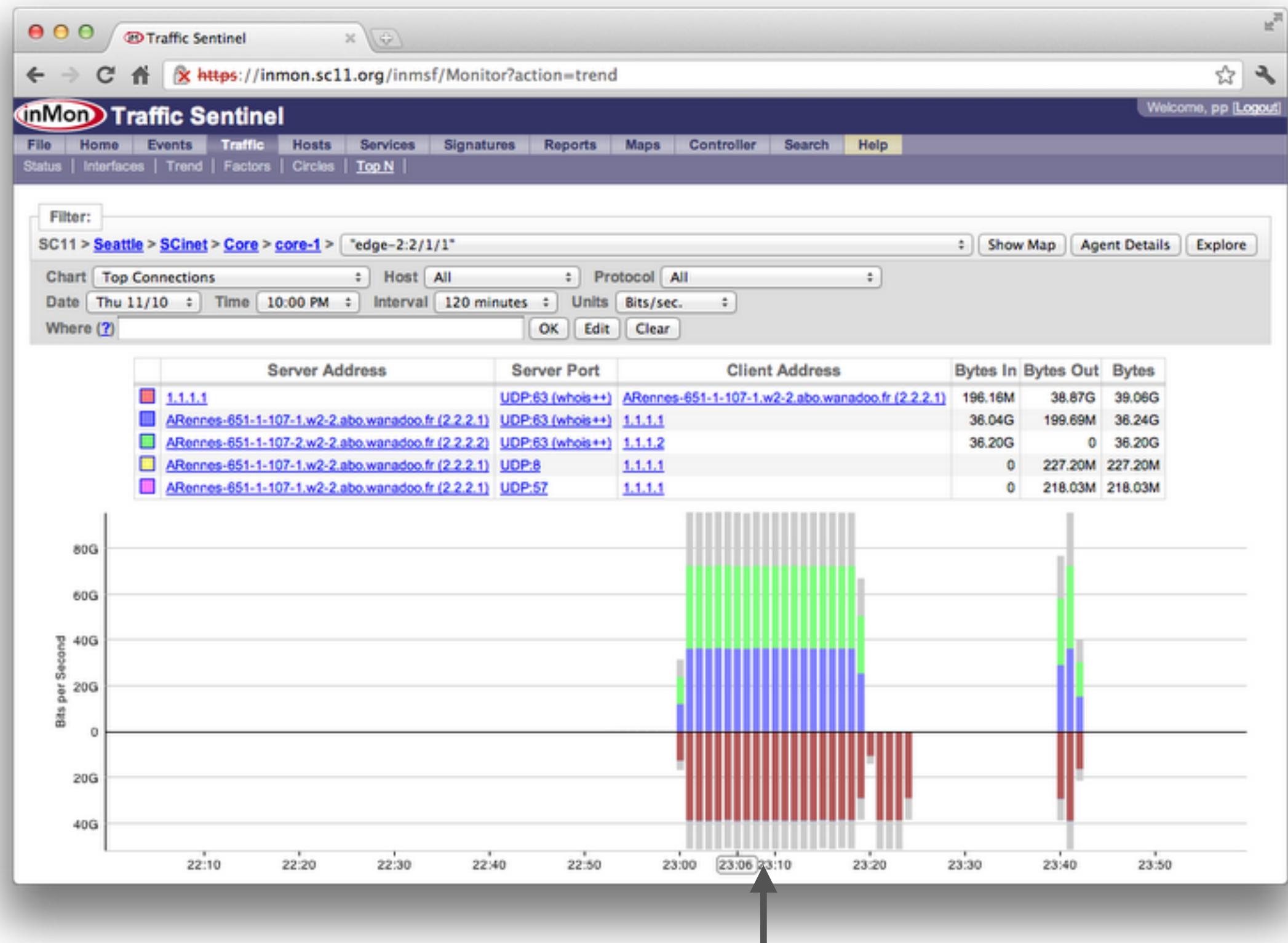
Why the spike in traffic?

(100Gbit link carrying 14,000,000 packets/second)

# RANDOM SAMPLING

Scaleable method of adding details

- Random sampling is lightweight
- Unbiased results with known accuracy
- Sampled packet header (128 bytes)
- Forwarding state associated with sampled packet (e.g. ingress/egress port, VLAN, next hop, CIDR etc.)
- Linux kernel instrumentation randomly samples packets
- Offload to hardware instrumentation on network switch scales functionality to Terabit/s speeds
- Identify top sources, destinations, connections, protocols



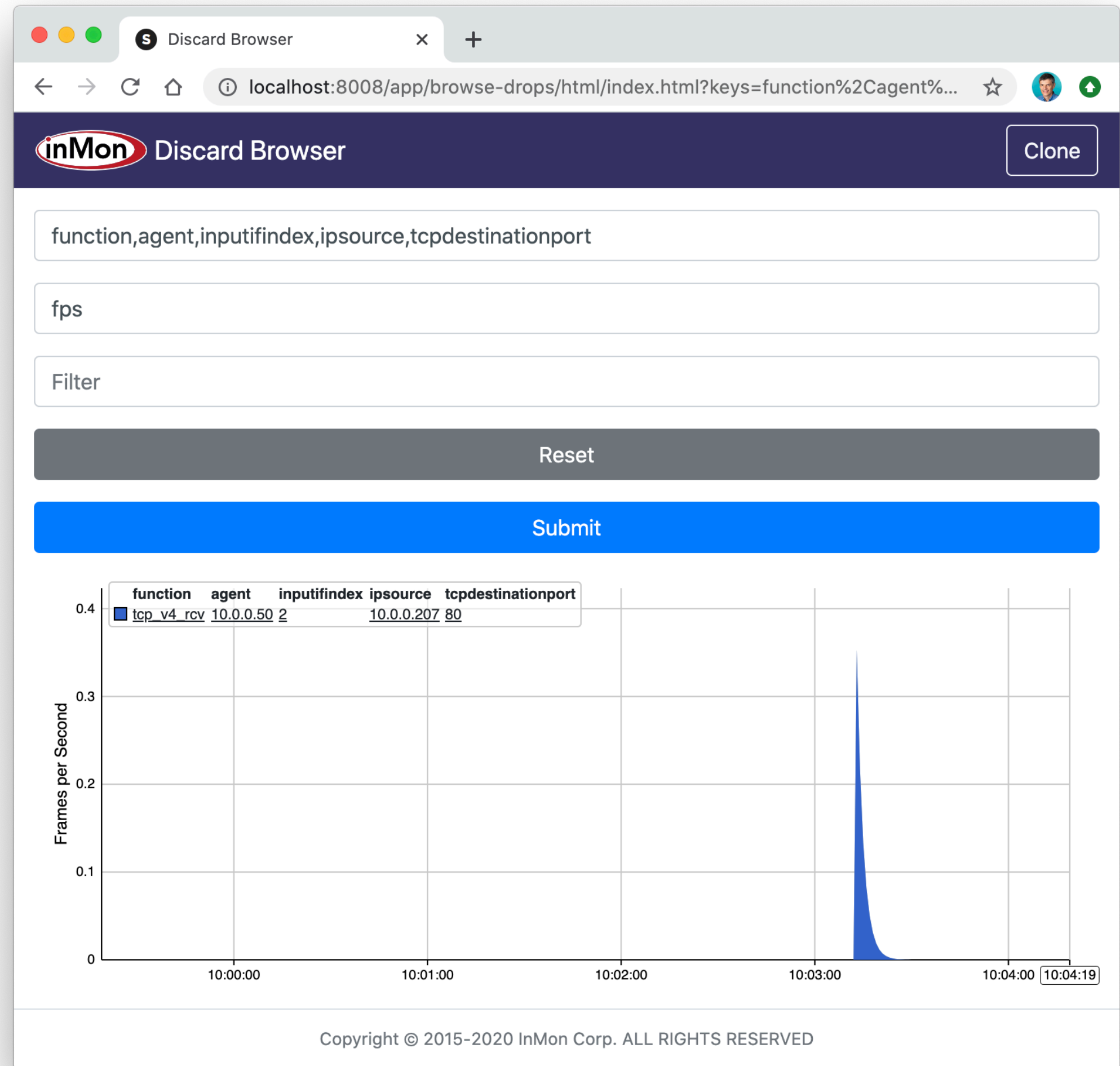
Break out traffic by client, server and port

(graph based on samples from 100Gbit link carrying 14,000,000 packets/second)

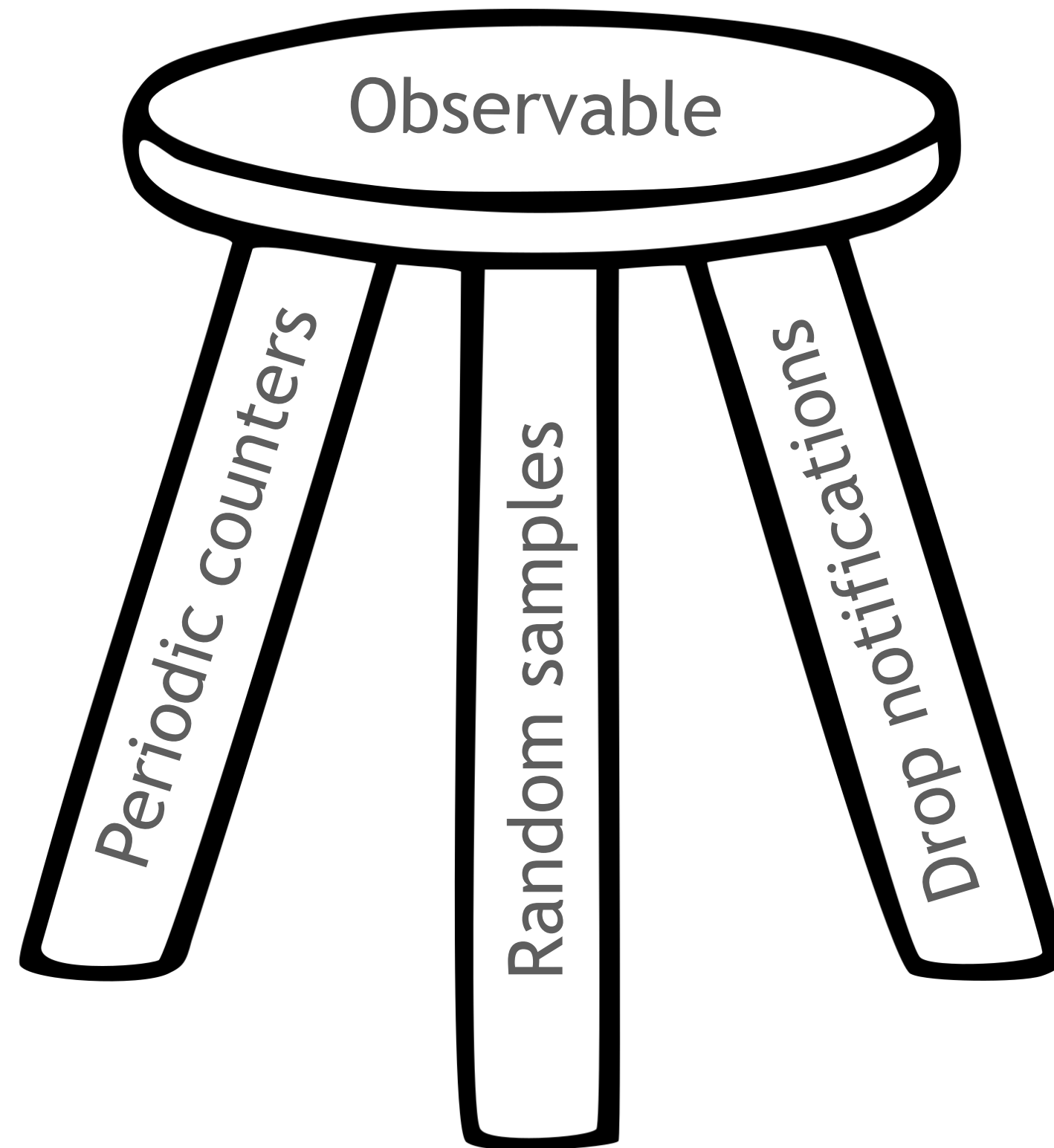
# NEW! PACKET DROP EVENTS

## Details of rare but critical drop events

- Discard counters tell you that packets are being dropped, but not root cause
- Packet discards are rare and unlikely to be sampled
- Packet discards can severely impact performance and availability of critical services
- Packet discard notifications in sFlow:
  - Discarded packet header (128 bytes)
  - Reason for dropping packet
- Linux kernel netlink drop\_monitor API reports each dropped packet and the reason it was dropped
- Offload to hardware instrumentation on a network switch scales functionality to Terabit/s network speeds
- Identify sources, destinations, protocols, locations of discards



# MEASUREMENTS STREAMS COMBINE FOR OBSERVABILITY



# SFLOW IS AN INDUSTRY STANDARD

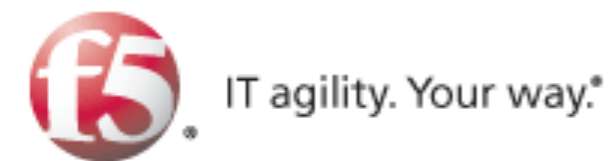
Broad multi-vendor support ensures observability of physical network



ARISTA

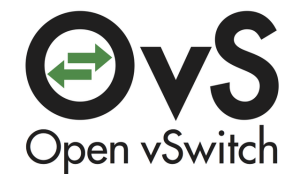


DELL EMC

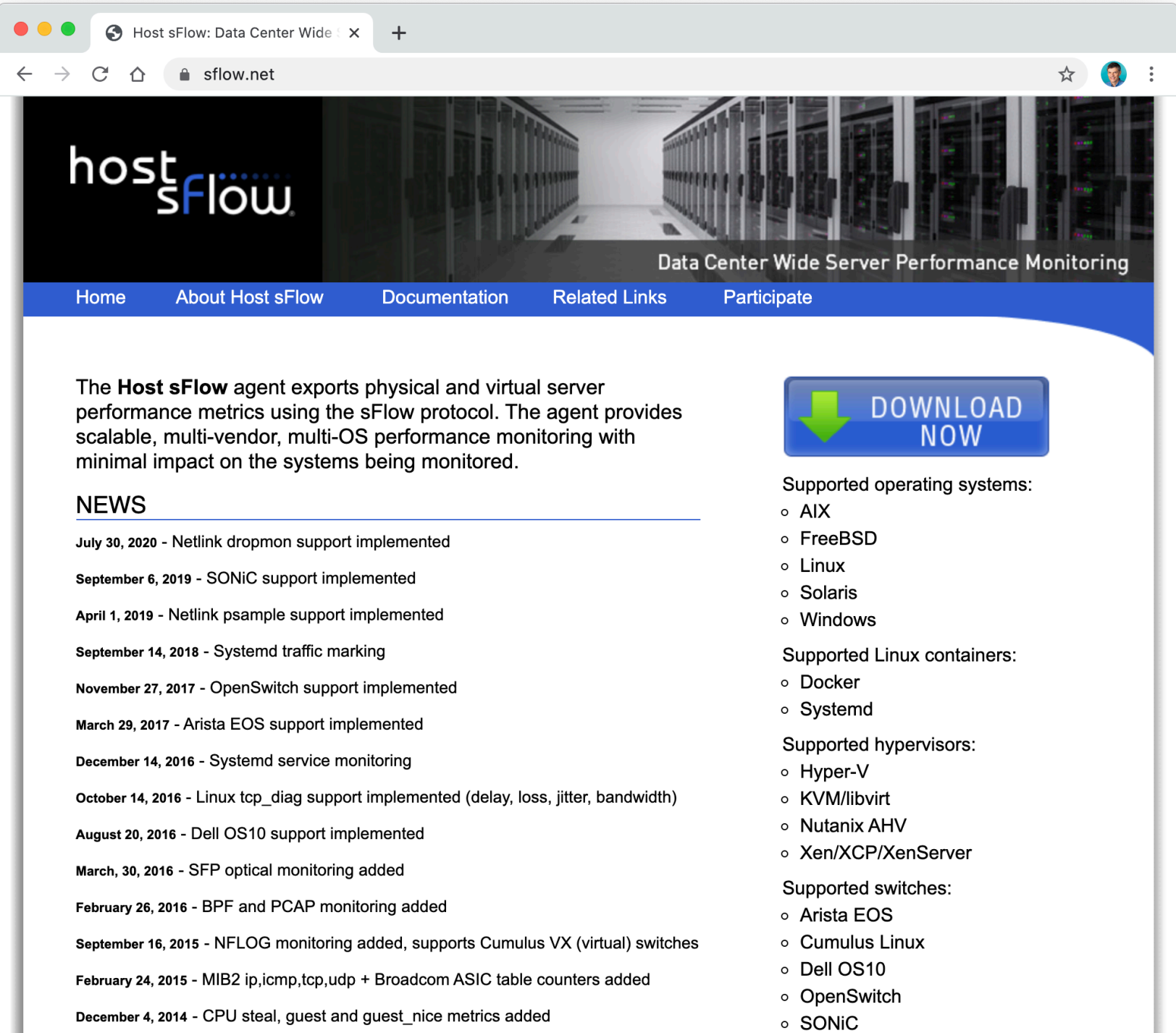


HITACHI  
Inspire the Next

Innovium

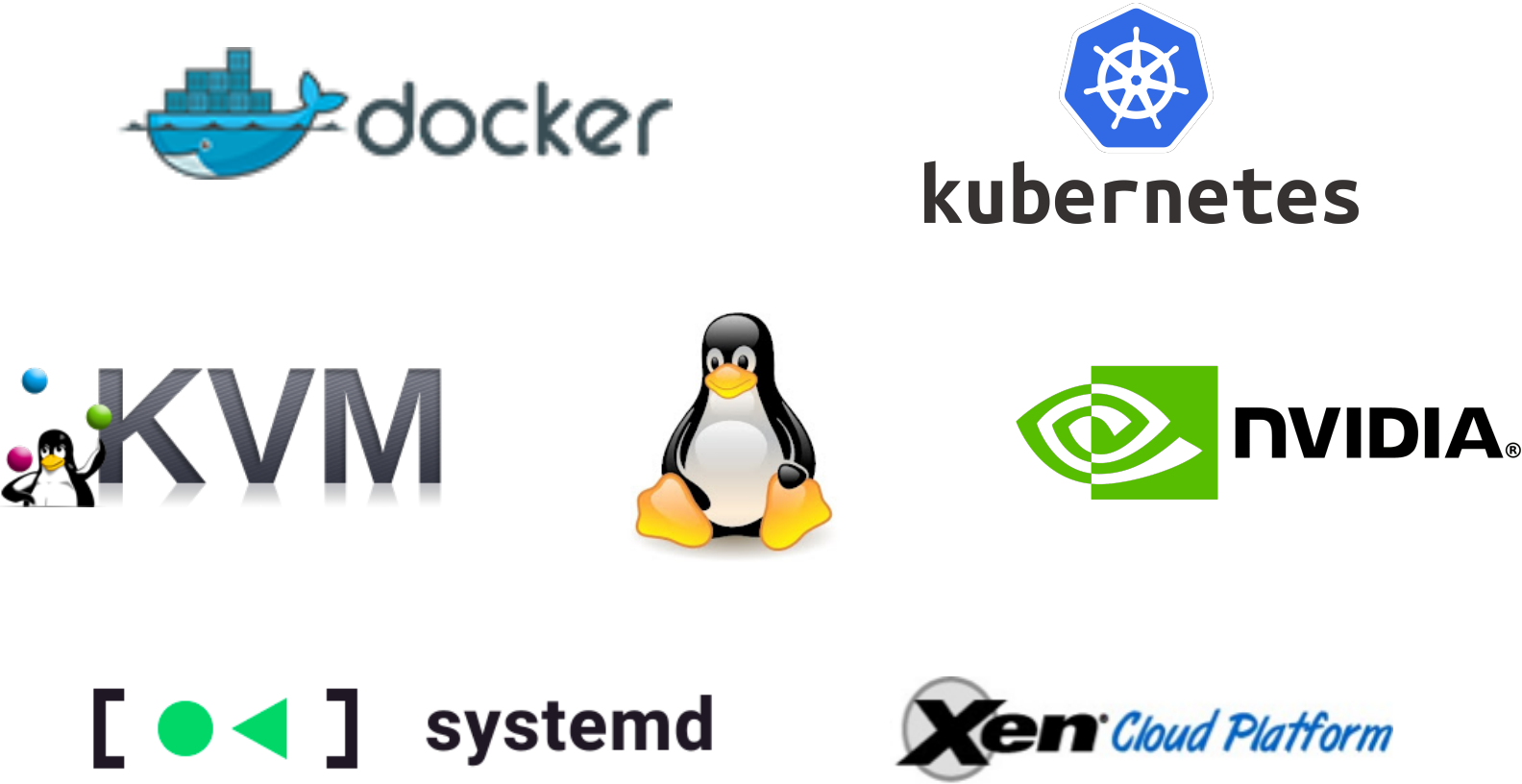


ZTE中兴



# HOST SFLOW

Extends observability into software network edge





# GPU CLUSTER DEMO

# Configure sFlow agents

Edit `/etc/hsflowd.conf`

Host sFlow agent pre-installed on Cumulus Linux

Download server package from <https://sflow.net>

Same software and configuration on switches and servers  
minimises operational complexity

```
# Minimal configuration
sflow {
    collector { ip = 10.31.234.46 }
}
```

# sFlow-RT analyzer

Pre-packaged Docker image

Real-time sFlow analytics as a microservice:

- sFlow telemetry in on UDP port 6343
- Analytics out through REST API on port 8008

Download packages from <https://sflow-rt.com>

```
docker run -d \  
-p 6343:6343/udp \  
-p 8008:8008 \  
sflow/prometheus
```

www.youtube.com/watch?v=XdMkcK\_AoQc&fea

YouTube

Search

SIGN IN

**inMon** sFlow-RT: Real-time GPU and network telemetry

**GPU Utilization**

Percent

Buffers  
Memory  
Execution

**Network Traffic**

Source Destination  
192.168.1.46 192.168.2.47  
192.168.2.47 192.168.1.46

Bits per Second

**Network Drops**

Reason Agent Source Destination  
acl 10.31.234.137 192.168.1.46 192.168.2.47

Frames per Second

9:29

9:36 / 11:54

Copyright © 2020 InMon Corp. ALL RIGHTS RESERVED

Using Advanced Telemetry to Correlate GPU and Network Performance Issues

2 views • Nov 11, 2020

0 0

SHARE

SAVE

...

**Peter Phaal**  
23 subscribers

SUBSCRIBE

Discussion of telemetry and requirements for network automation. Overview of sFlow measurement architecture. Discussion of recently added packet drop monitoring functionality. Live demonstration of GPU compute cluster analytics.

SHOW MORE

Talk and demo available on YouTube

[https://youtu.be/XdMkcK\\_AoQc](https://youtu.be/XdMkcK_AoQc)