

## Global Foundation Services



DATA CENTERS



NETWORK  
S



SERVERS



ENERG  
Y

C#

SOFTWARE



SECURIT  
Y

# Large-Scale Passive Monitoring using SDN

Mohan Nanduri  
[mnanduri@microsoft.com](mailto:mnanduri@microsoft.com)

Justin Scott  
[juscott@microsoft.com](mailto:juscott@microsoft.com)

Onur Karaagaoglu  
[onurka@microsoft.com](mailto:onurka@microsoft.com)

# What is this about?

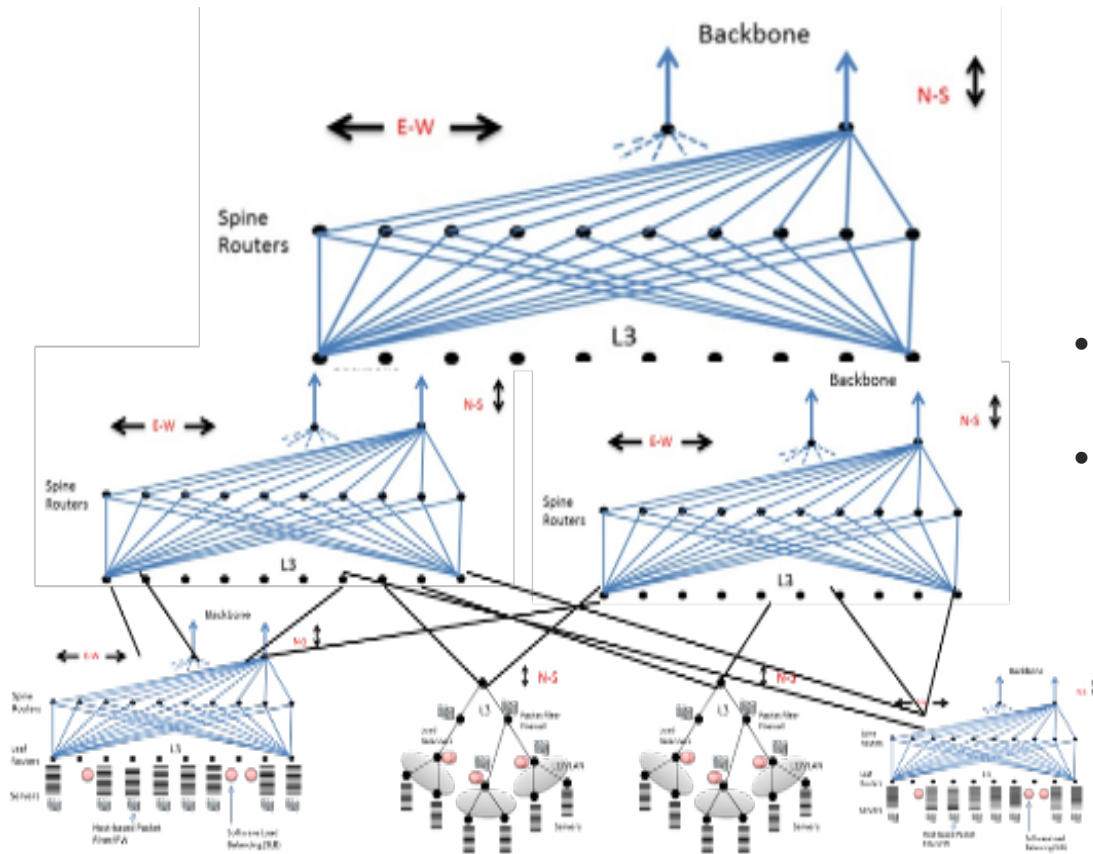
A different way to aggregate data from an optical  
TAP/SPAN

Architecture based on Openflow and commodity  
merchant silicon

# Why?

Exonerate network quickly, reduce MTTR and MTTM  
Lack of ability to data mine telemetry data at the network layer  
Reduce operations engagement to troubleshoot service  
Multi-tenant offering  
Solution that's on-demand and always available  
Provides large-scale tap aggregation leveraging commodity hardware

# Hyper Scale



- Thousands of 10G links per Data Center
- Cost makes it a non-starter with commercial solutions

# Prior Attempts

## Capture-Net

- Off the shelf aggregation gear, was too expensive at scale
- Resulted in lots of gear gathering dust
- Operations not mature enough to back such a solution

## PMA/PUMA – “Passive Measurement Architecture”

- Lower cost than Capture-net
- Designed for a specific environment and not intended to scale
- Extremely feature rich

**Stuck with shuffling sniffers around**

*NOT*

~~THE~~

Just took a step back  
END

# What Features Make Up a Packet Broker?

Terminates taps  
Match on a 5-tuple  
Duplication  
Packets unaltered  
Low Latency  
Statistics  
80%



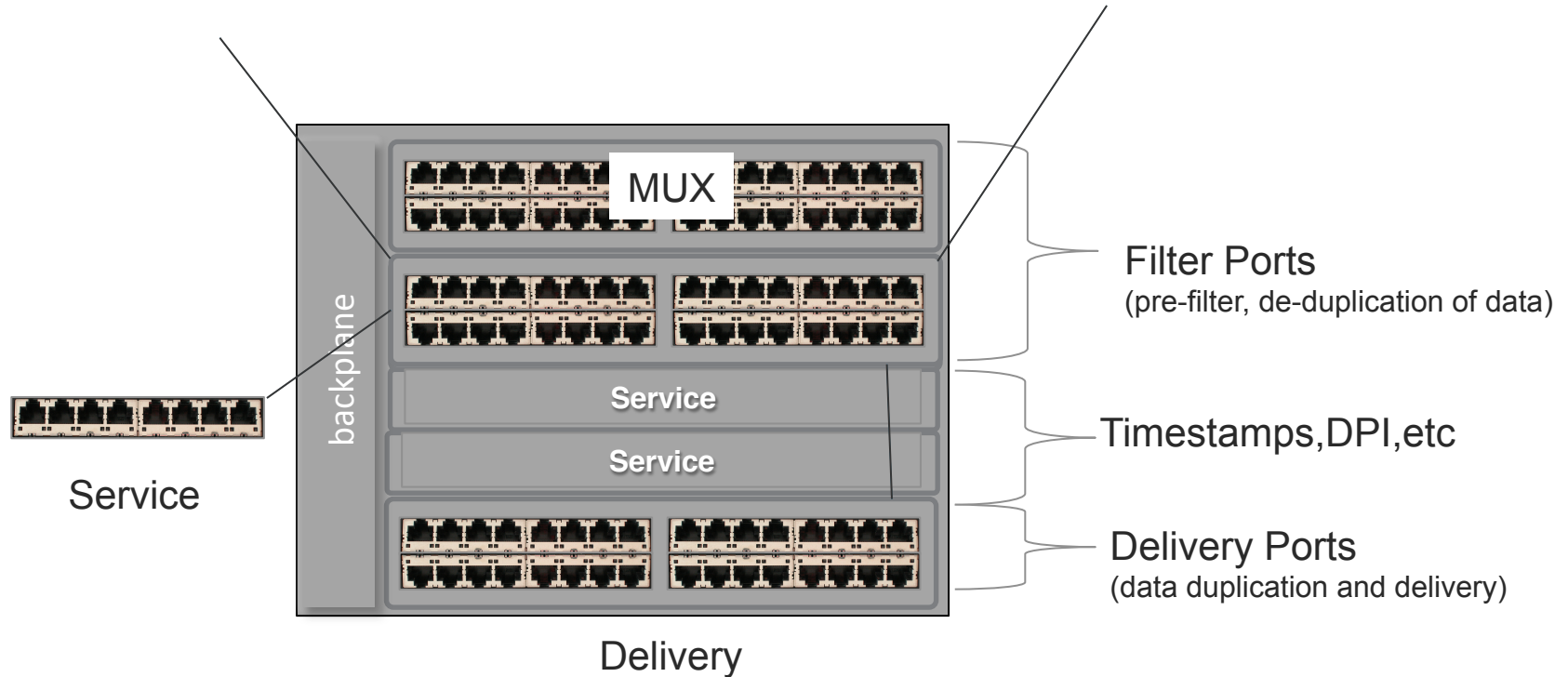
Merchant  
Silicon

Layer 7 packet inspection  
Time stamps  
Frame Slicing  
Microburst detection  
20%



Requires  
specialized  
Hardware

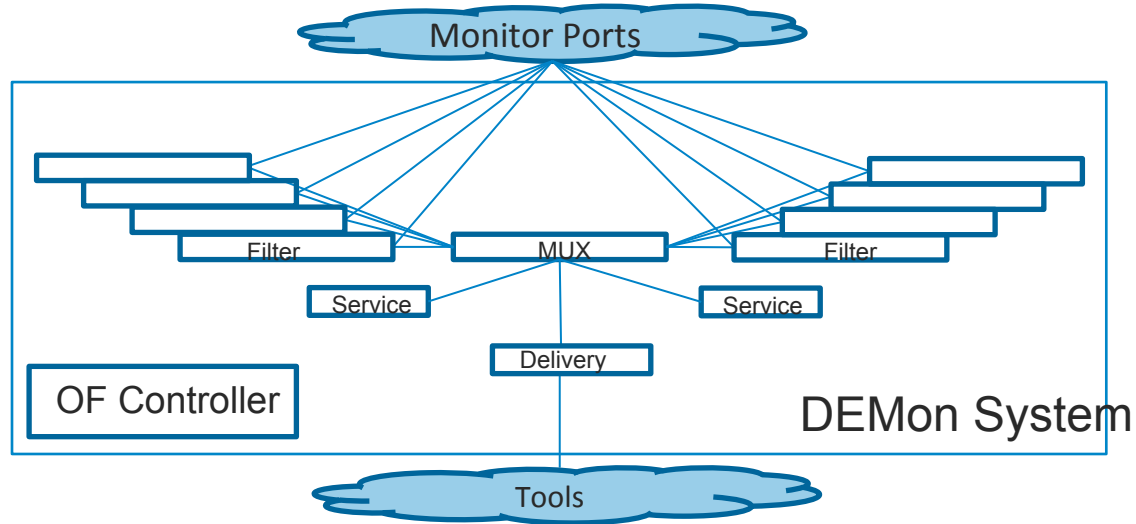
# Reverse Engineering Packet Broker



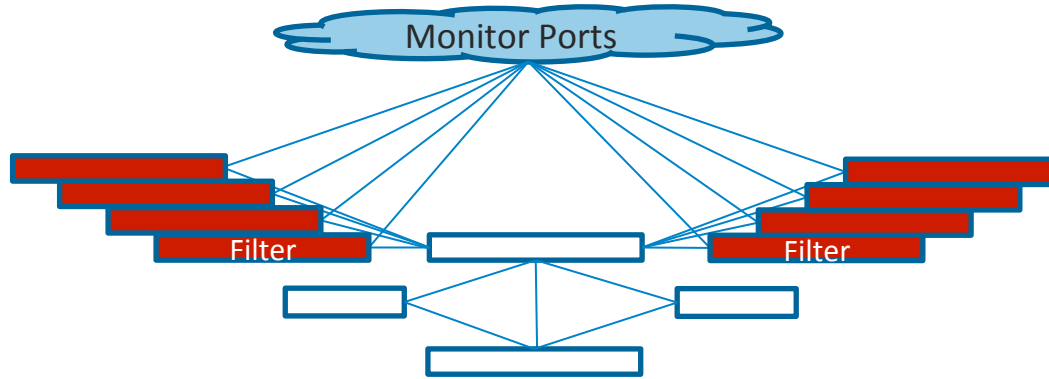


Architecture

# Architecture Overview

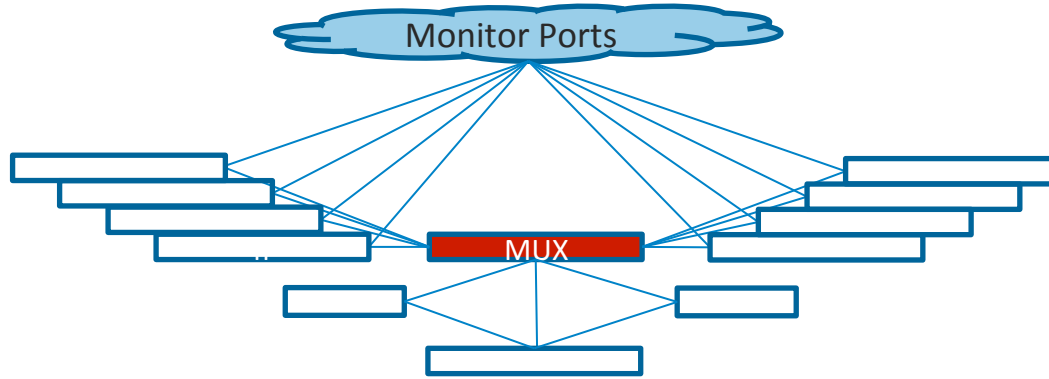


# Filter Layer



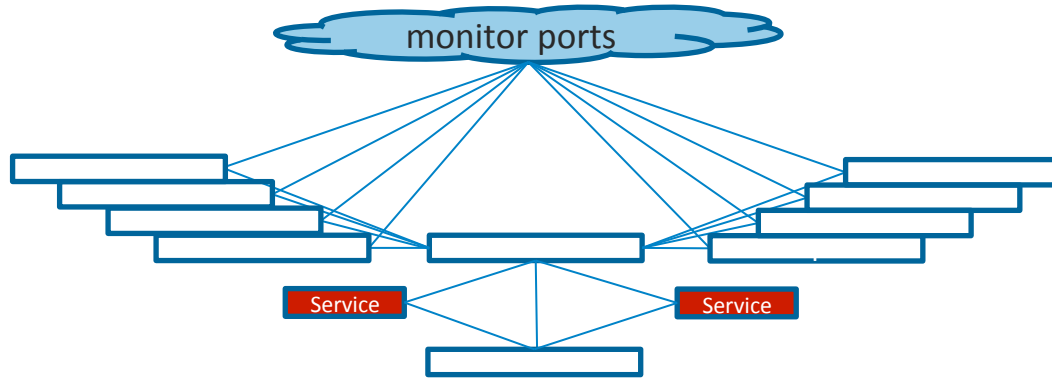
- Terminates all monitor ports
- Drops all traffic by default
- De-duplication of data if needed
- Aggressive sFlow exports

# MUX Layer



- Aggregates all filter switches in a data center
- Directs traffic to either service nodes or delivery interfaces
- Enables service chaining per policy

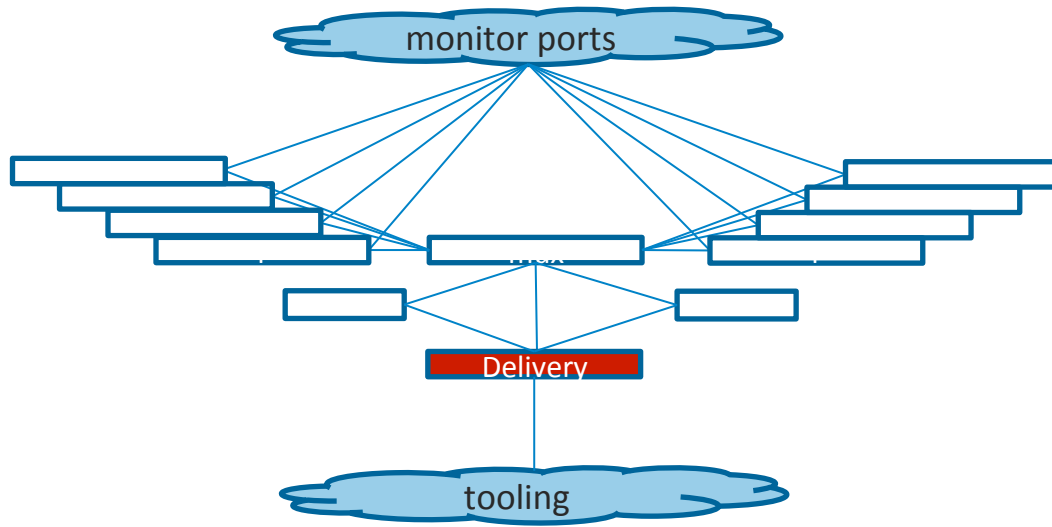
# Service Layer



Aggregated by MUX layer  
Flows are sent through services nodes to perform extended functions  
Resulting flows are sent back to the delivery switch and then to tools

Some Applications:  
Deeper (layer 7) filtering  
Time stamping  
Microburst detection  
Traffic Ratio's (SYN/SYNC ACK)  
Frame slicing (64, 128 and 256 byte)  
Payload removal for compliance  
Rate limiting

# Delivery Layer



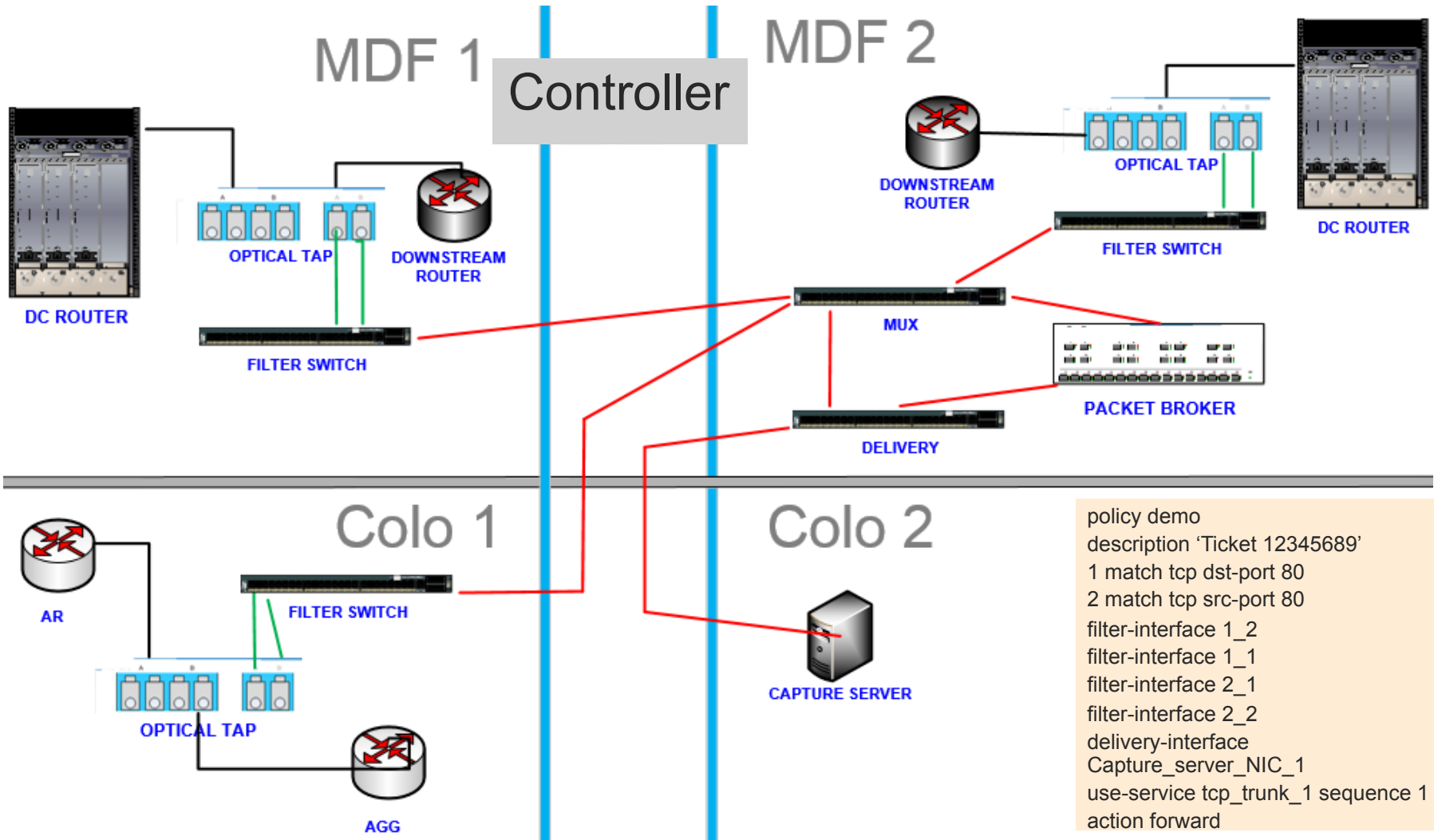
1:N and N:1 delivery to duplication of data  
Delivery to local or tunnel traffic to remote tools

# SDN Controller

OpenFlow 1.0 based

Discovers topology via LLDP

Roles of each layer are assigned and discovered automatically



```

policy demo
description 'Ticket 12345689'
1 match tcp dst-port 80
2 match tcp src-port 80
filter-interface 1_2
filter-interface 1_1
filter-interface 2_1
filter-interface 2_2
delivery-interface
Capture_server_NIC_1
use-service tcp_trunk_1 sequence 1
action forward

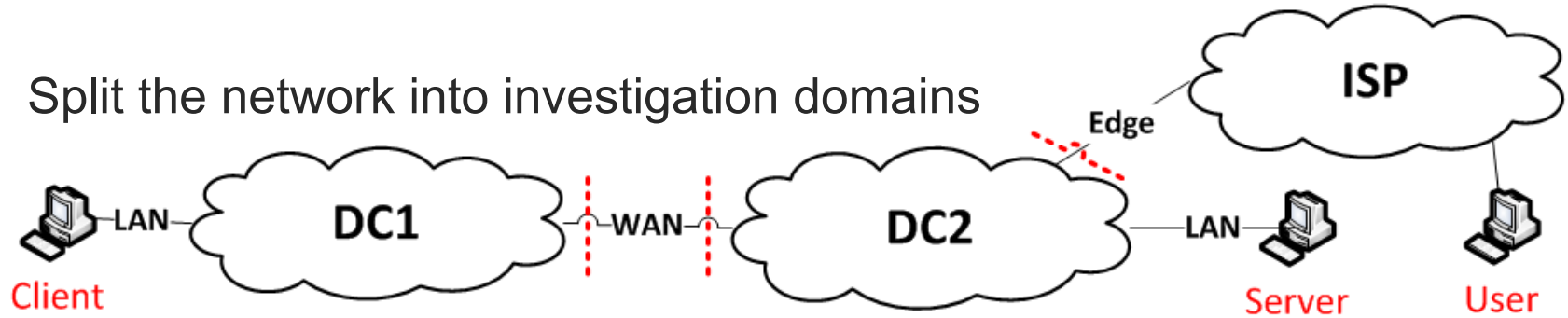
```



# Use Cases and Examples

# Reactive

Split the network into investigation domains



Quickly exonerate or implicate a network segment

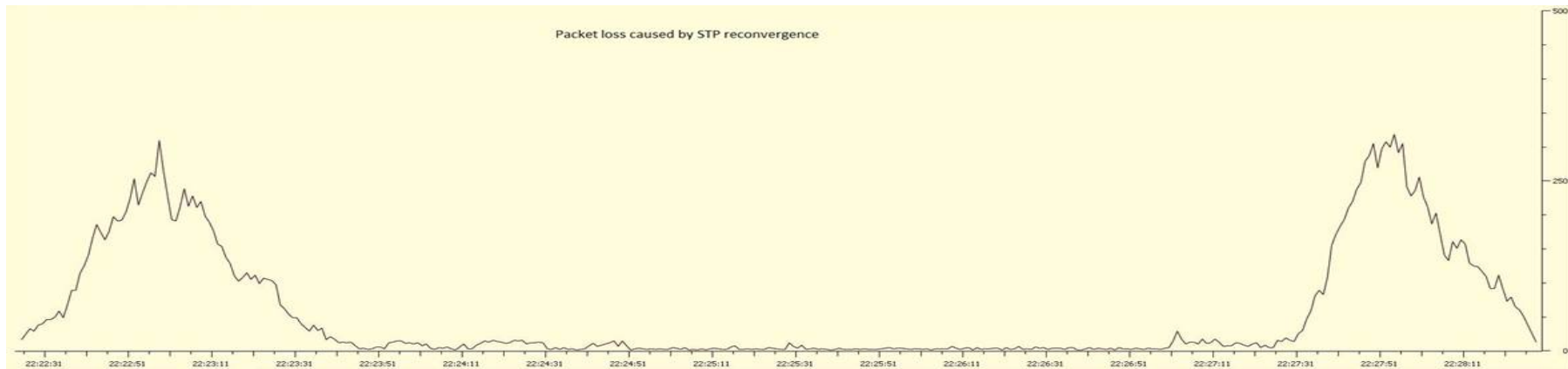
Verify TCP intelligent network appliance are operating as expected

# Proactive Monitoring

Relying sole on SNMP polling and syslog's gives you false confidence

Performance data can be gleaned from exposing TCP telemetry data

Ability to detect re-transmissions (TCP-SACK)



# IPv6

Users were unable to connect to the service intermittently via IPv6

## Repro facts:

3-way TCP connection setup's up

9-way SSL handshake fails

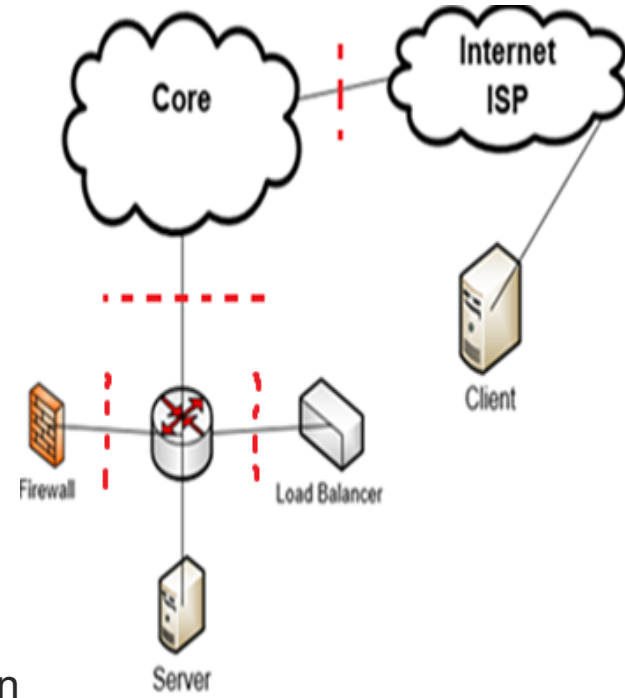
Ack for client hello was not making it back to load balancer

## Solution:

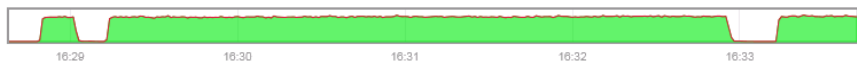
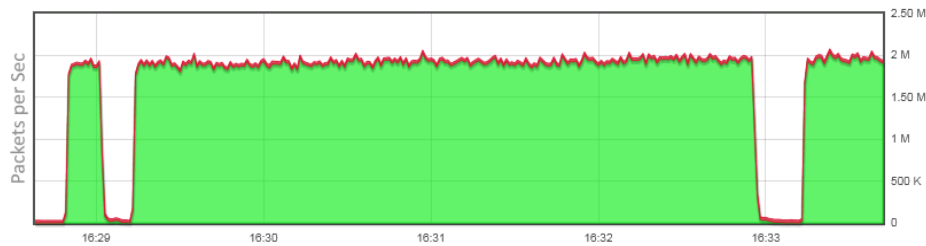
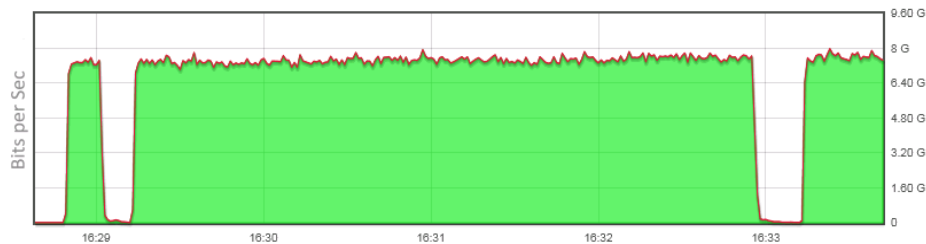
Implicates or exonerates Layer7 devices that are commonly finger pointed

## Root cause:

Race condition - If the client hello was received on the load balancer before the backend connection was made it would trigger the bug



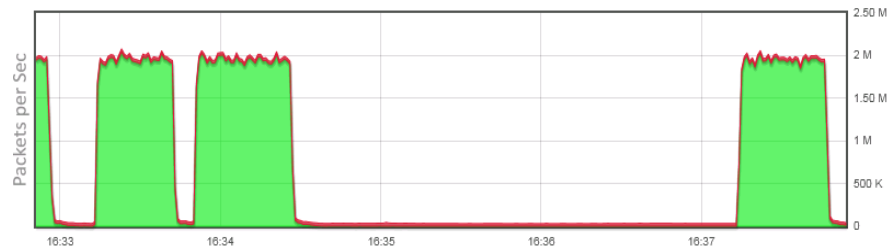
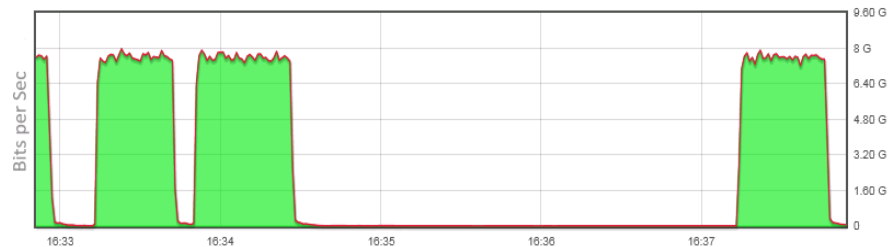
# DDOS



Fit All | Selected: 16:28:38 - 16:33:42 | <-5min 16:28:38 +5min> | current | live/pause

live

Passed Dropped



Fit All | Selected: 16:32:51 - 16:37:54 | <-5min 16:32:51 +5min> | current | live/pause

live

Passed Dropped

5 Minutes later

# DDOS: Packet Capture

4 0.000003000 ).000001000 [redacted] [redacted] NTP

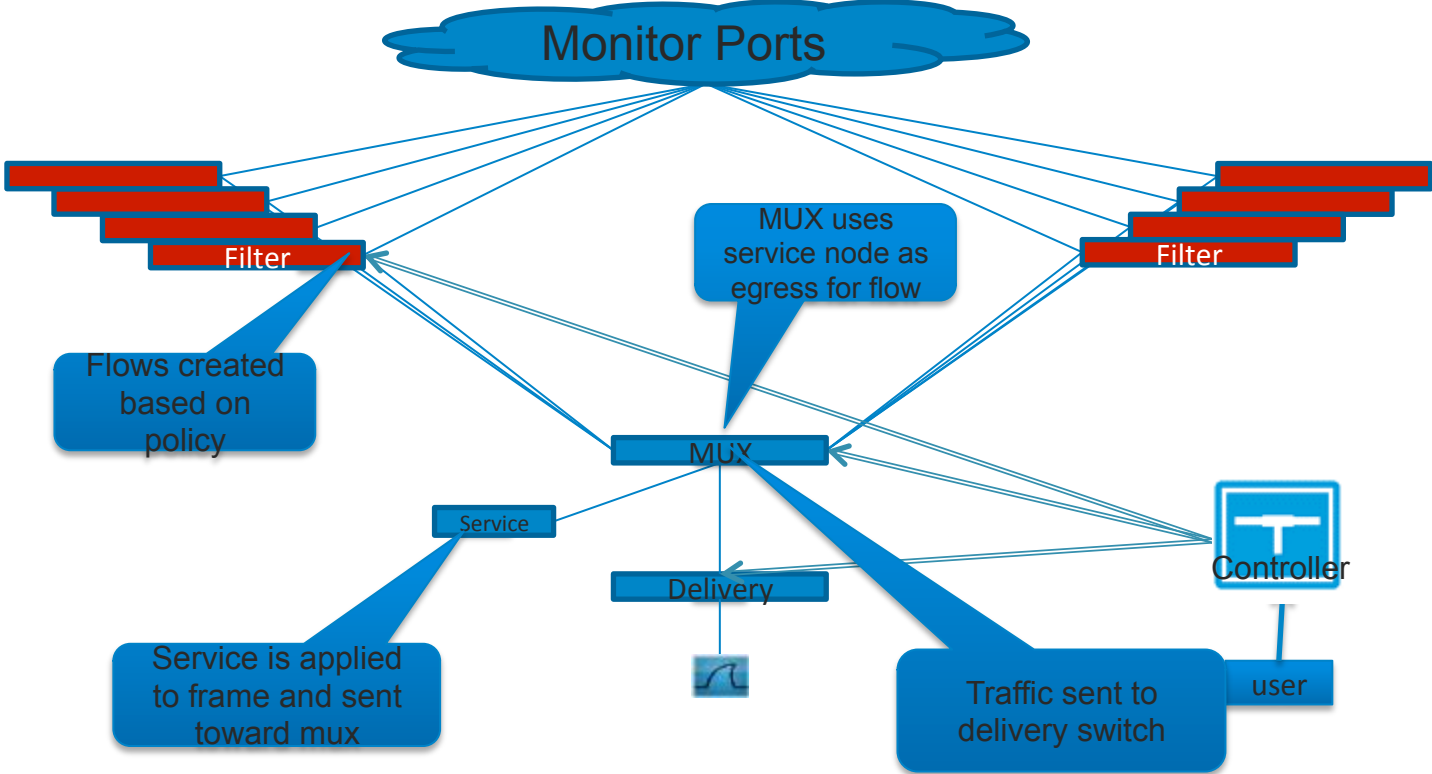
<

- ⊕ Frame 4: 482 bytes on wire (3856 bits), 482 bytes captured (3856 bits) on interface 0
- ⊕ Ethernet II, Src: JuniperN\_bb:17:c2 (00:1f:12:bb:17:c2), Dst: Cisco\_fd:1a:3c (00:05:73:fd:1a:3c)
- ⊕ Internet Protocol Version 4, Src: [redacted] ([redacted]), Dst: [redacted] ([redacted])
- ⊕ User Datagram Protocol, Src Port: ntp (123), Dst Port: http (80)
- ⊖ Network Time Protocol (NTP Version 2, private)
  - ⊕ Flags: 0xd7
  - ⊕ Auth, sequence: 69
    - Implementation: XNTPD (3)

Request code: MON\_GETLIST\_1 (42)

```
0000 00 05 73 fd 1a 3c 00 1f 12 bb 17 c2 08 00 45 00  ..S..<.. .....E.
0010 01 d4 fc 7d 00 00 3a 11 36 5e c2 2c c0 11 a8 3d  ...}...: 6^.,...=
0020 21 c2 00 7b 00 50 01 c0 49 ec d7 45 03 2a 00 06  !..{.P.. I..E.
0030 00 48 00 00 00 01 00 00 1e 3a 00 00 00 00 00 00  .H..... :.....
0040 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00  .....
```

# sFLOW



# Caveats and Cost



# Caveats

TCP/IP fields with MPLS encapsulated packets cannot be matched

Lack of IPv6 source and destination matching in OF 1.0

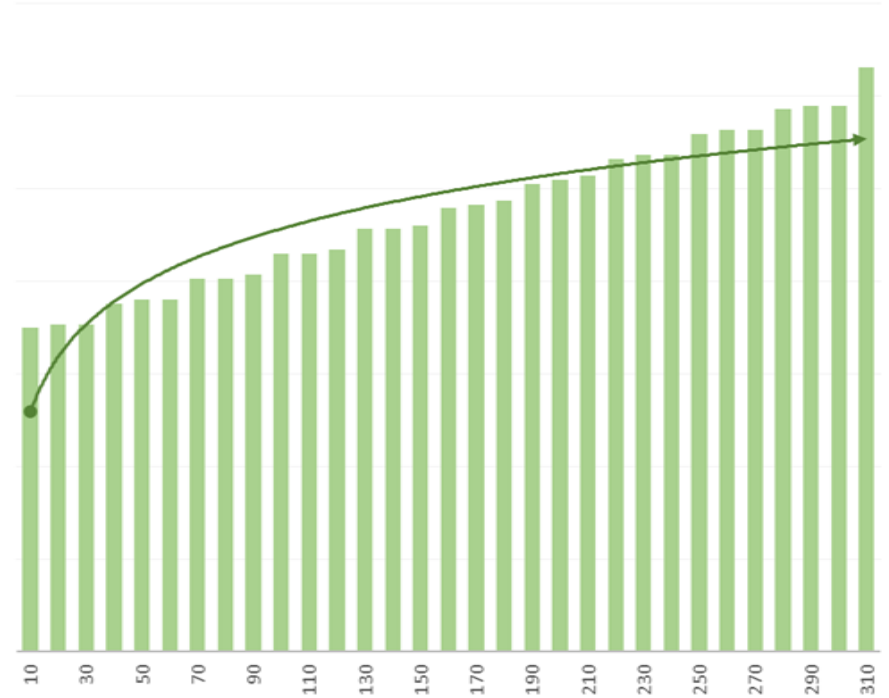
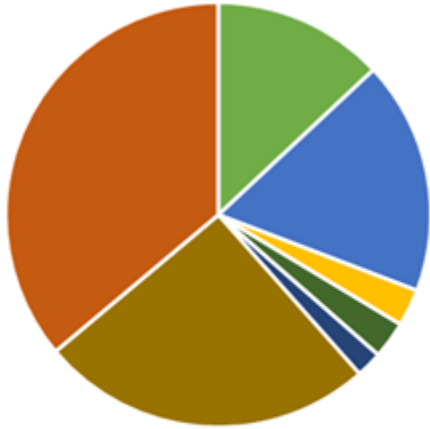
Limited number of flow rules due to TCAM limitation

Policy will not load balance traffic amongst ECMP links

Not all switch vendors OF implementation is the same

Commercial controller support is splintering

# Cost breakdown

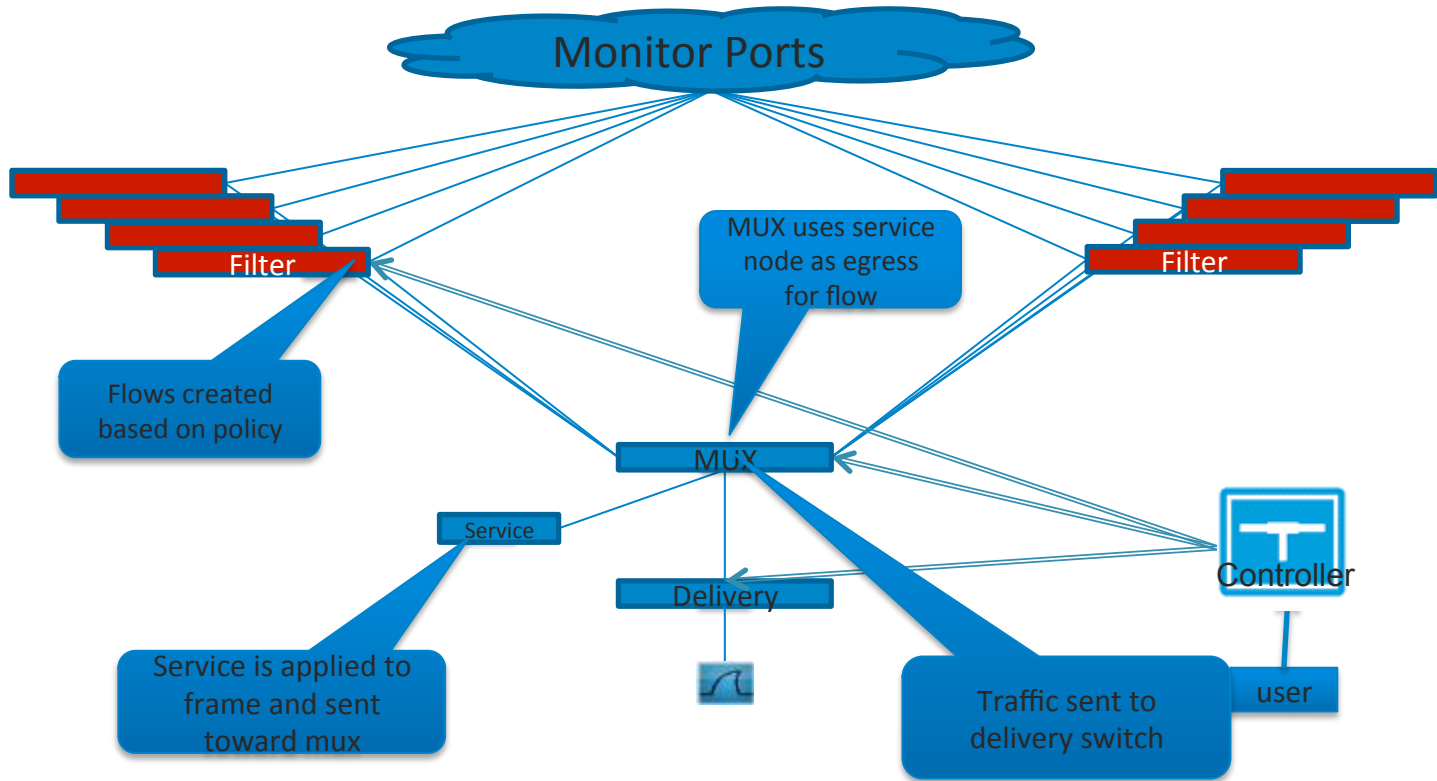


Number of links

Questions?

# Backup

# Creating a Service Chain



# Sniffer Features

Terminate taps

Match on 5-tuple

Duplication of packets

Low latency

Layer 7 packet inspection

Time stamping

Frame slicing