

# Time-Muxed Parsing in Marking-based Network Telemetry

Alon Riesenber<sup>\*</sup>, Yonnie Kirzon<sup>\*</sup>, Michael Bunin<sup>\*</sup>,

Elad Galili<sup>\*</sup>, Gidi Navon<sup>•</sup>, Tal Mizrahi<sup>◇\*</sup>

ACM SYSTOR, Haifa, May 2019



# Background

What is network telemetry?



Delay

Packet loss

Queue status

...

Performance measurement + exporting to a remote location

Why do we need telemetry?



Detection

Failures



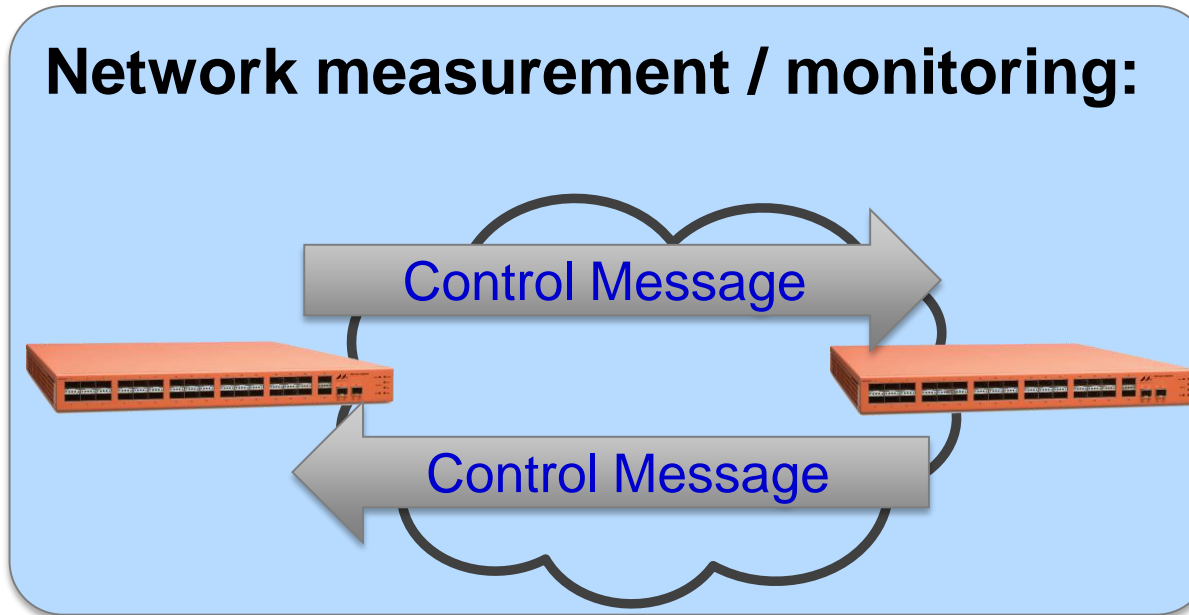
Congestion /  
Bottlenecks

'Elephant'  
flows

...



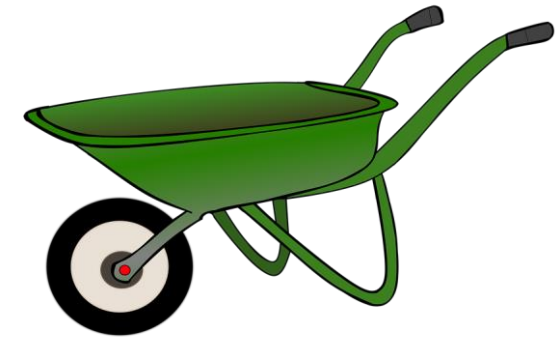
# Operations, Administration, Maintenance (OAM)



**Network Telemetry** 

# Ping / Traceroute

```
C:\Windows\system32\cmd.exe  
C:\>ping www.marvell.com -n 10  
Pinging extranet.marvell.com [10.68.68.50] with 32 bytes of data:  
Reply from 10.68.68.50: bytes=32 time=211ms TTL=57
```



# Old-School Passive Monitoring



Counters

Per port

Queue State

Per flow

Latency

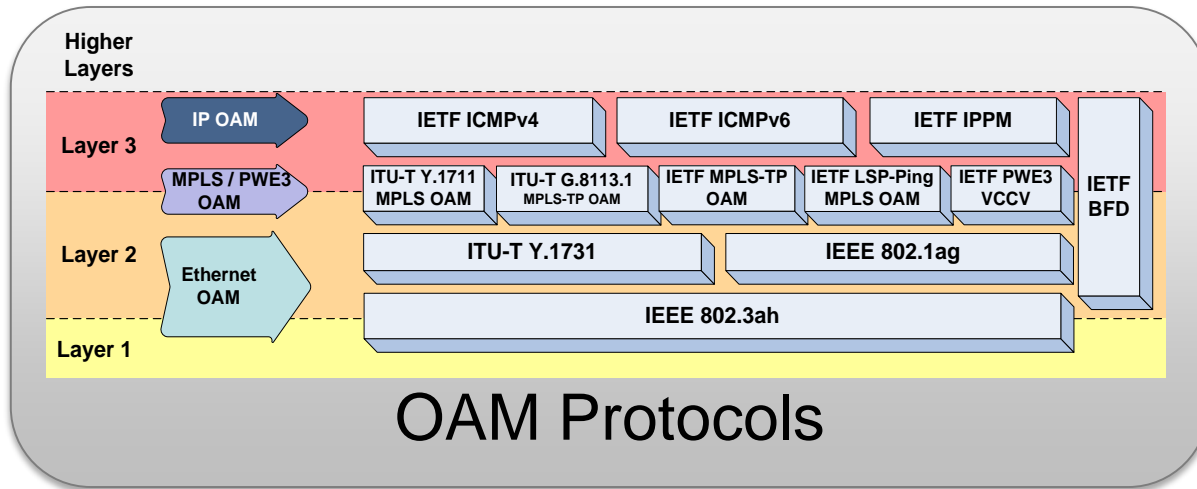
Per queue

⋮

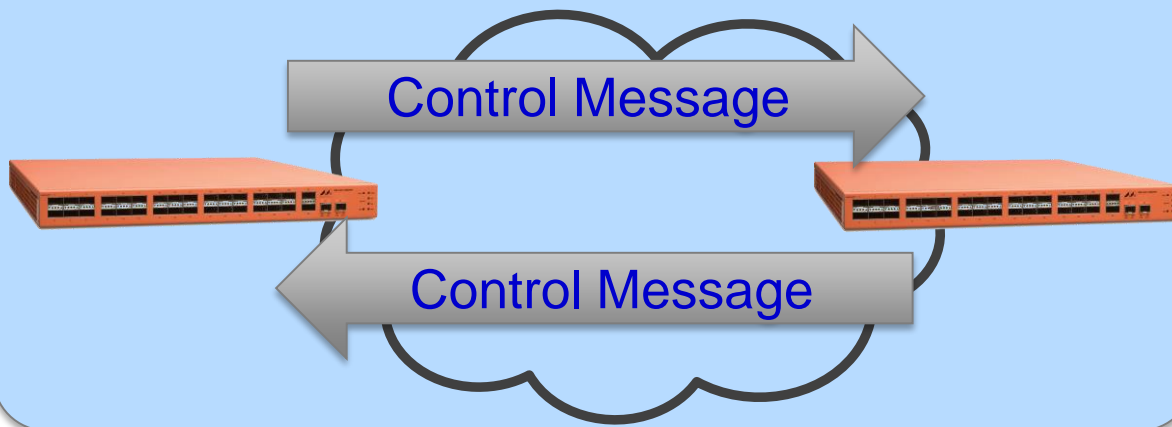
⋮



# Carrier Network OAM



## Active measurement / monitoring:

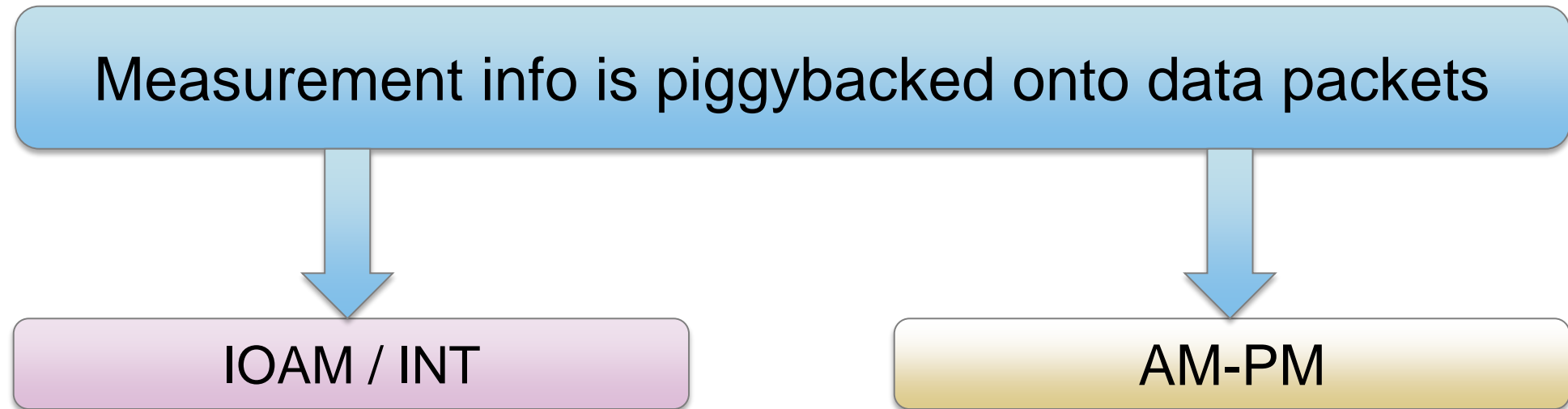


# Fate Sharing



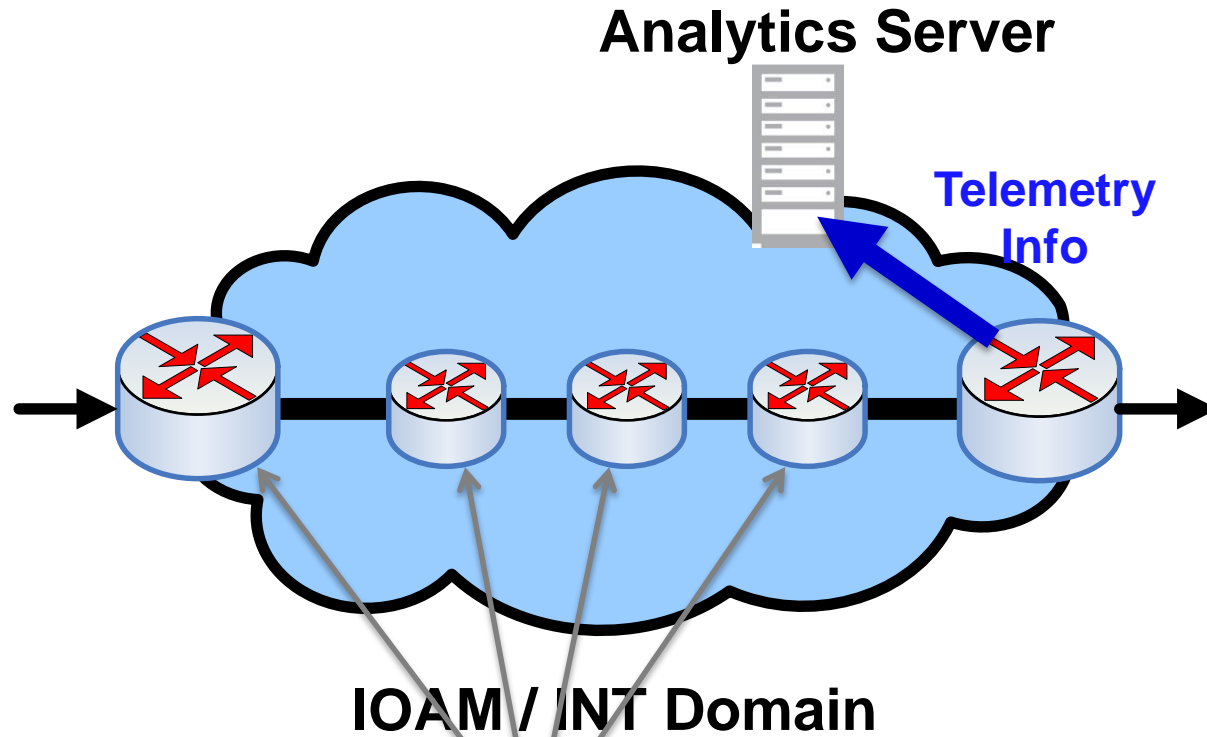
<http://www.speedtest.net>

# Piggybacked Measurement





# Piggybacked Metadata – IOAM / INT



Switches push local metadata into header: delay, queue state, ...



IOAM In situ OAM  
INT In-band Network Telemetry

Per-packet metadata 😊  
Per-packet overhead ☹️

# AM-PM: Alternate Marking – Performance Measurement

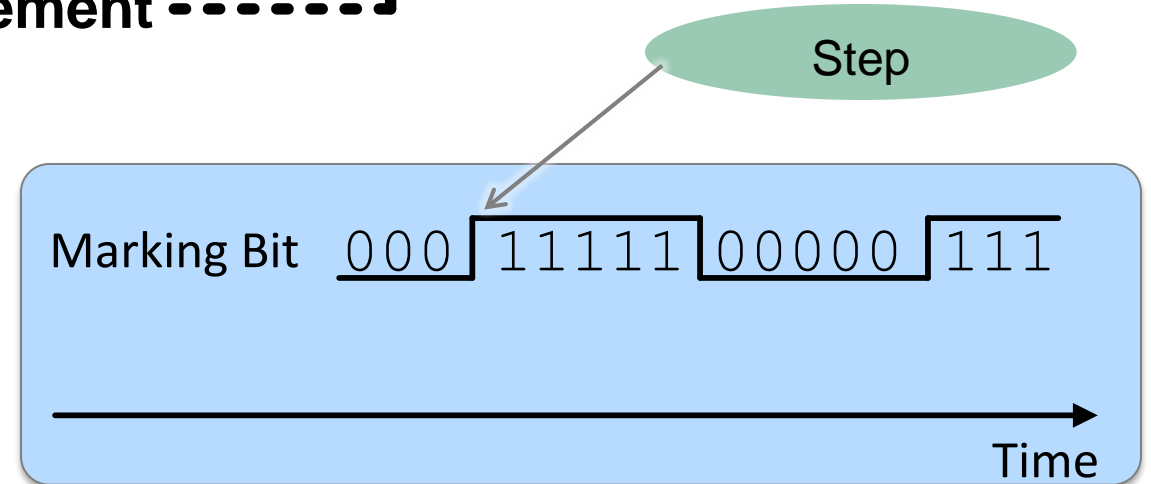
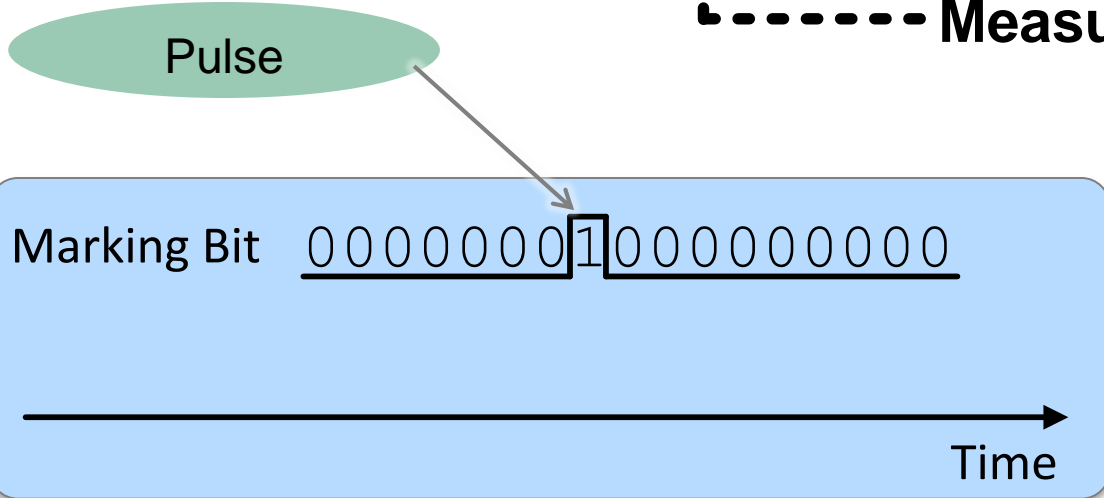
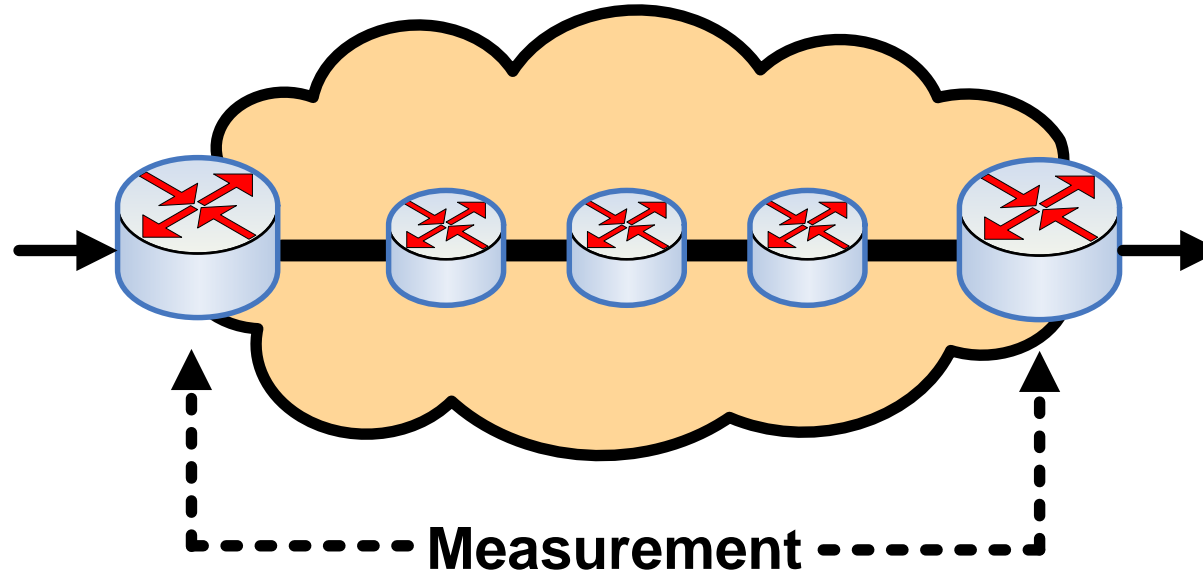
RFC 8321

Fioccola, G., Capello, A., Cociglio, M., Castaldelli, L., Chen, M., Zheng, L., Mirsky, G., and T. Mizrahi, “Alternate Marking method for passive and hybrid performance monitoring”, [RFC 8321](#), 2018.

draft-mizrahi-ippm-multiplexed-alternate-marking  
(internet draft)

T. Mizrahi, C. Arad, G. Fioccola, M. Cociglio, M. Chen, L. Zheng, and G. Mirsky. “Compact Alternate Marking Methods for Passive Performance Monitoring”, [draft-mizrahi-ippm-compact-alternate-marking](#), work in progress, IETF, 2018.

# AM-PM: What Can We Do with ONE Bit Per Packet?



# AM-PM: **Pulse** Marking – Delay Measurement



Servers



Servers



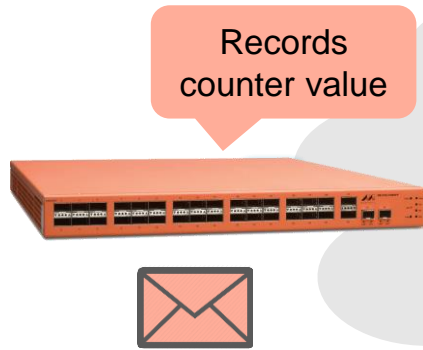
Analytics Server

Time Sent: March 8th, 16:02, 123400789 nsec (UTC)  
Time Received: March 8th, 16:02, 123500789 nsec (UTC)  
**Network Delay: 100  $\mu$ sec**

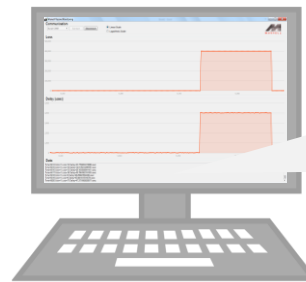
# AM-PM: **Pulse** Marking – Loss Measurement



Servers



Servers



Analytics Server

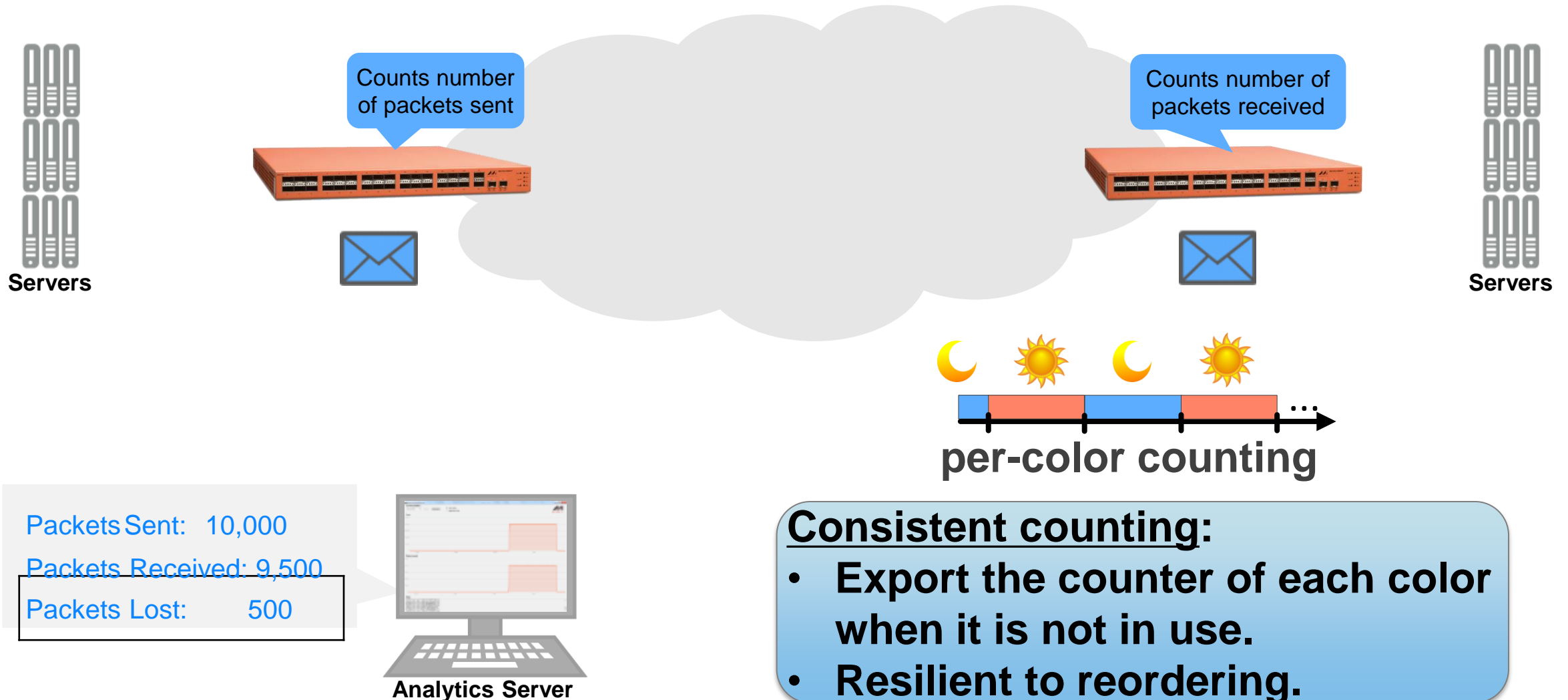
Counter: 2100

Counter: 2000

Packets lost: 100

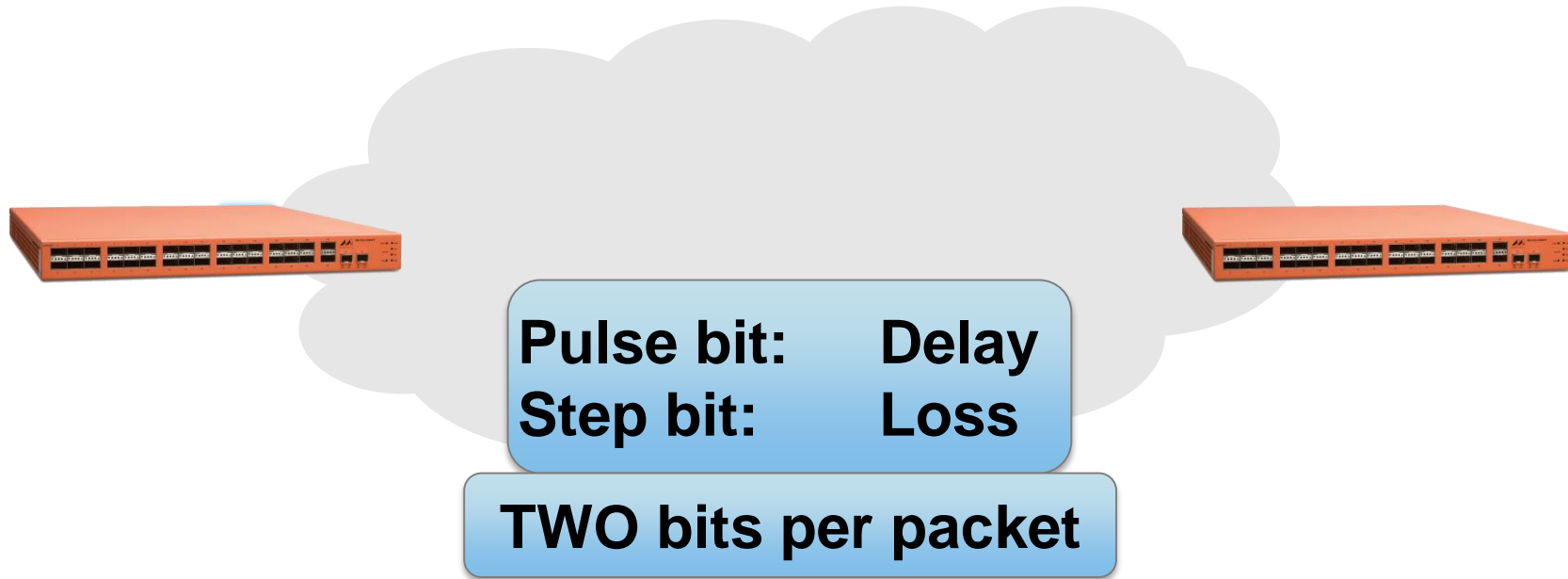
**Out of order?**

# AM-PM: **Alternate** Marking – Loss Measurement



- Consistent counting:**
- Export the counter of each color when it is not in use.
  - Resilient to reordering.

# AM-PM: **Double** Marking



# AM-PM: **Multiplexed** Marking



Servers



Servers

**Pulse:**                      **Delay**  
**Step:**                        **Loss**

**ONE** bit per packet

**Accurate loss and  
delay measurement!**



# Design and Implementation of AM-PM

## Match-Action Lookup

TCAM / Exact match / P4

## Time-as-a-match

TimeFlip

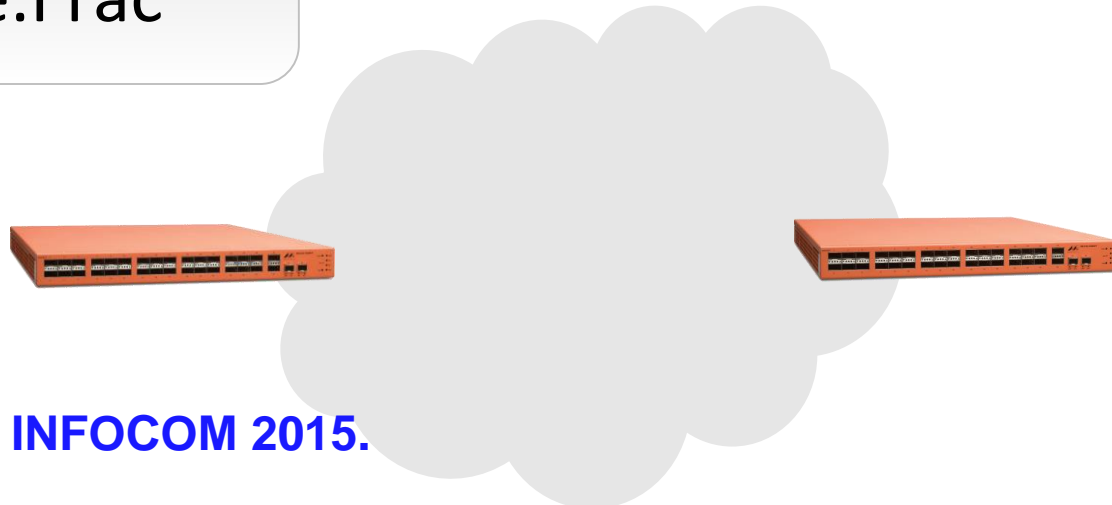
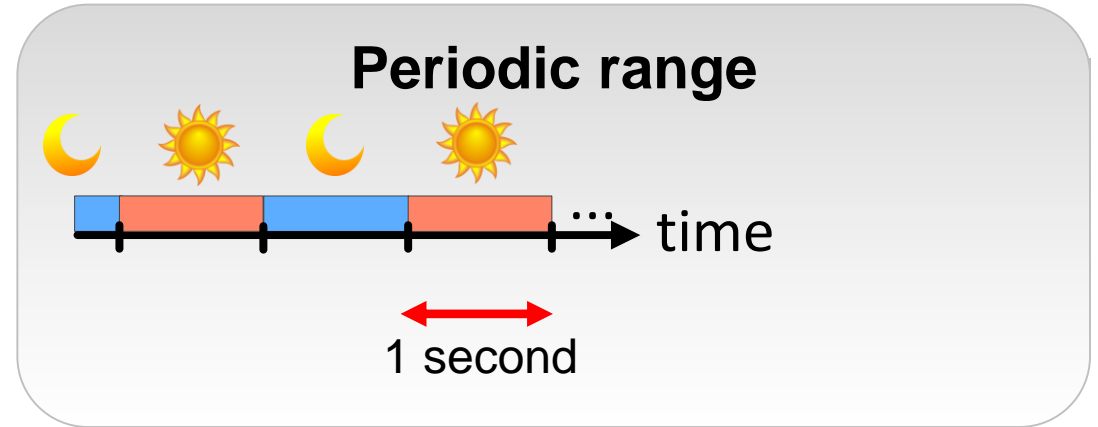
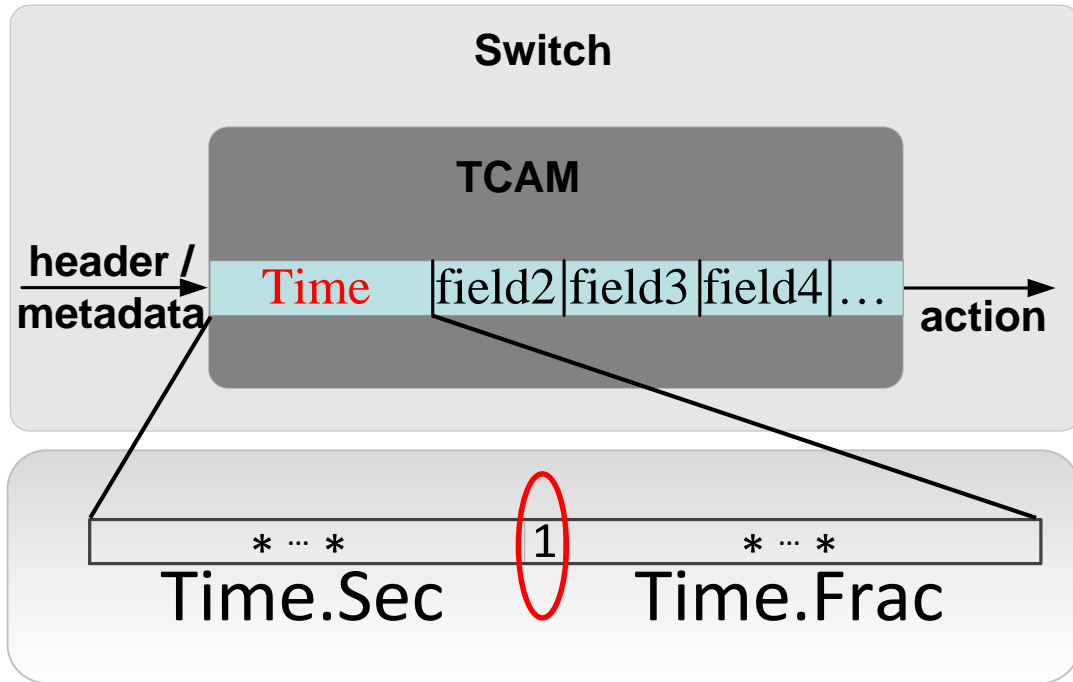


## State

Detect first packet  
(pulse/step)



# Time-as-a-match: TimeFlip [MRM]



[MRM] Mizrahi, Rottenstreich, Moses, INFOCOM 2015.

# Design and Implementation of AM-PM: Step/Pulse

**Match-Action  
Lookup**

TCAM / Exact match / P4

**Time-as-a-match**

TimeFlip



**State**

Detect first packet  
(pulse/step)

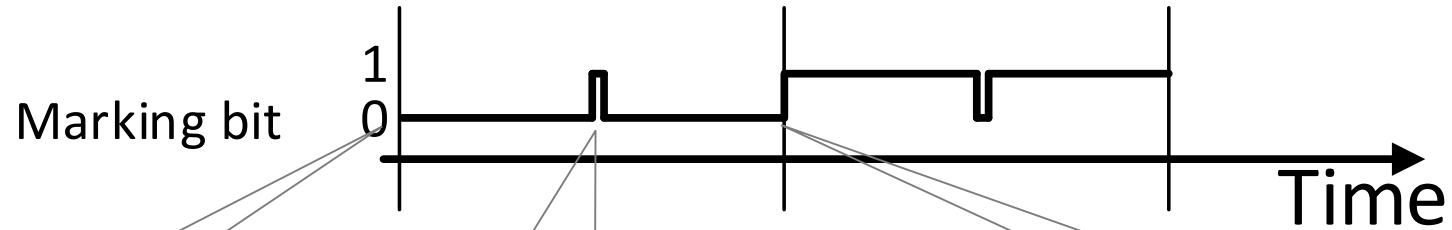
Match	Action
<i>TimeBit = 0</i>	<i>MarkBit = 0, counter0</i>

Match	Action
<i>MarkBit = 0</i>	<i>counter0</i>

Match	Action
<i>TimeBit = 1, Reg = 0</i>	<i>MarkBit = 1, Reg = 1, timestamp</i>
<i>TimeBit = 1, Reg = 1</i>	<i>MarkBit = 0, Reg = 1</i>
<i>TimeBit = 0, Reg = 1</i>	<i>MarkBit = 1, Reg = 0, timestamp</i>

Match	Action
<i>MarkBit = 1</i>	<i>timestamp</i>

# Multiplexed Marking: a Naïve Implementations



Track the value of the marking bit.

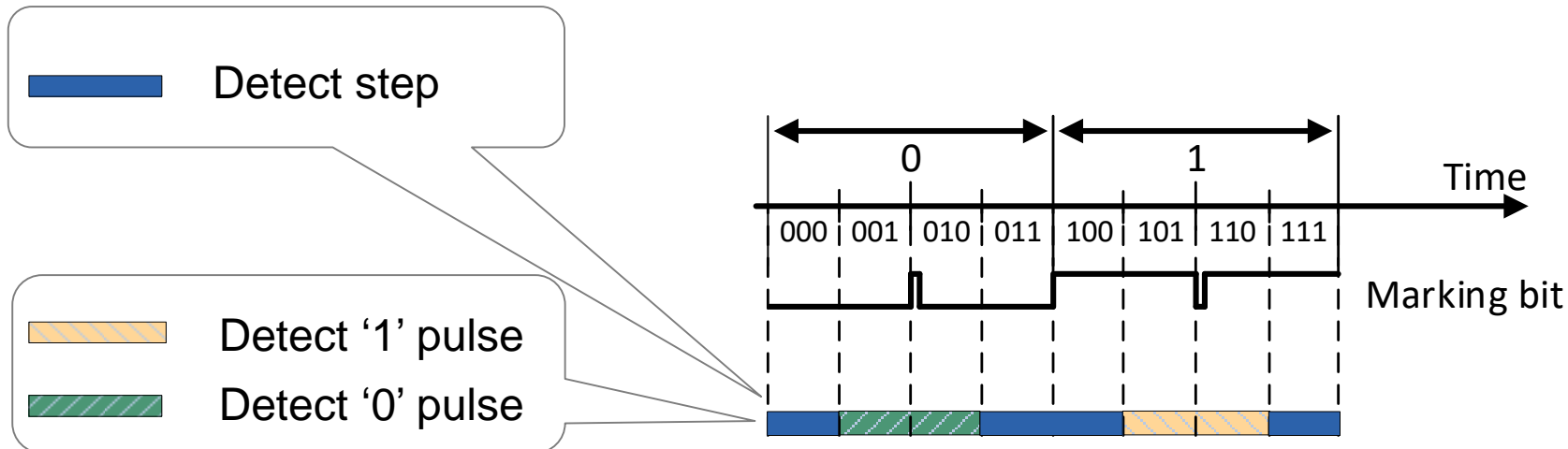
Detect pulse  
When the value changes for **one** packet.

Detect step  
When the value changes for more than **one** packet.

Non-trivial to implement using a match-action abstraction.

# Our Approach: Time-multiplexed Parsing

Header field(s) have a different interpretation in each time slot!



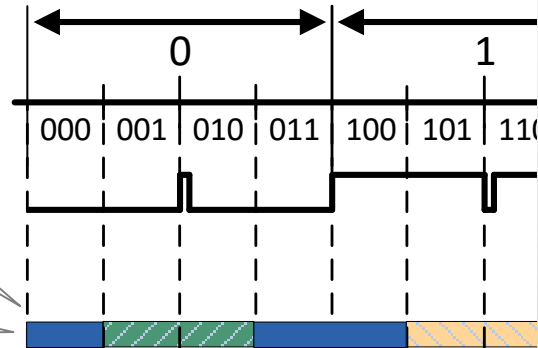
- TimeFlip is used to divide time into time slots.
- The marking bit has a different interpretation in each time slot.
- Requires rough time synchronization, e.g., ~ 1 second.

# Our Approach: Time-multiplexed Parsing

Header field(s) have a different interpretation in each time slot!

— Detect step

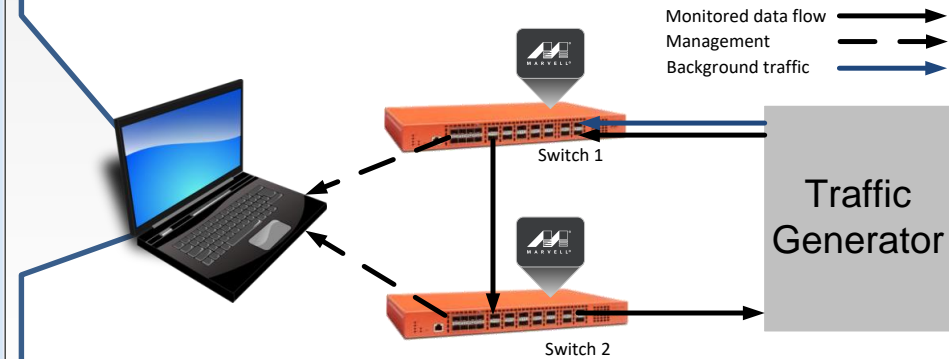
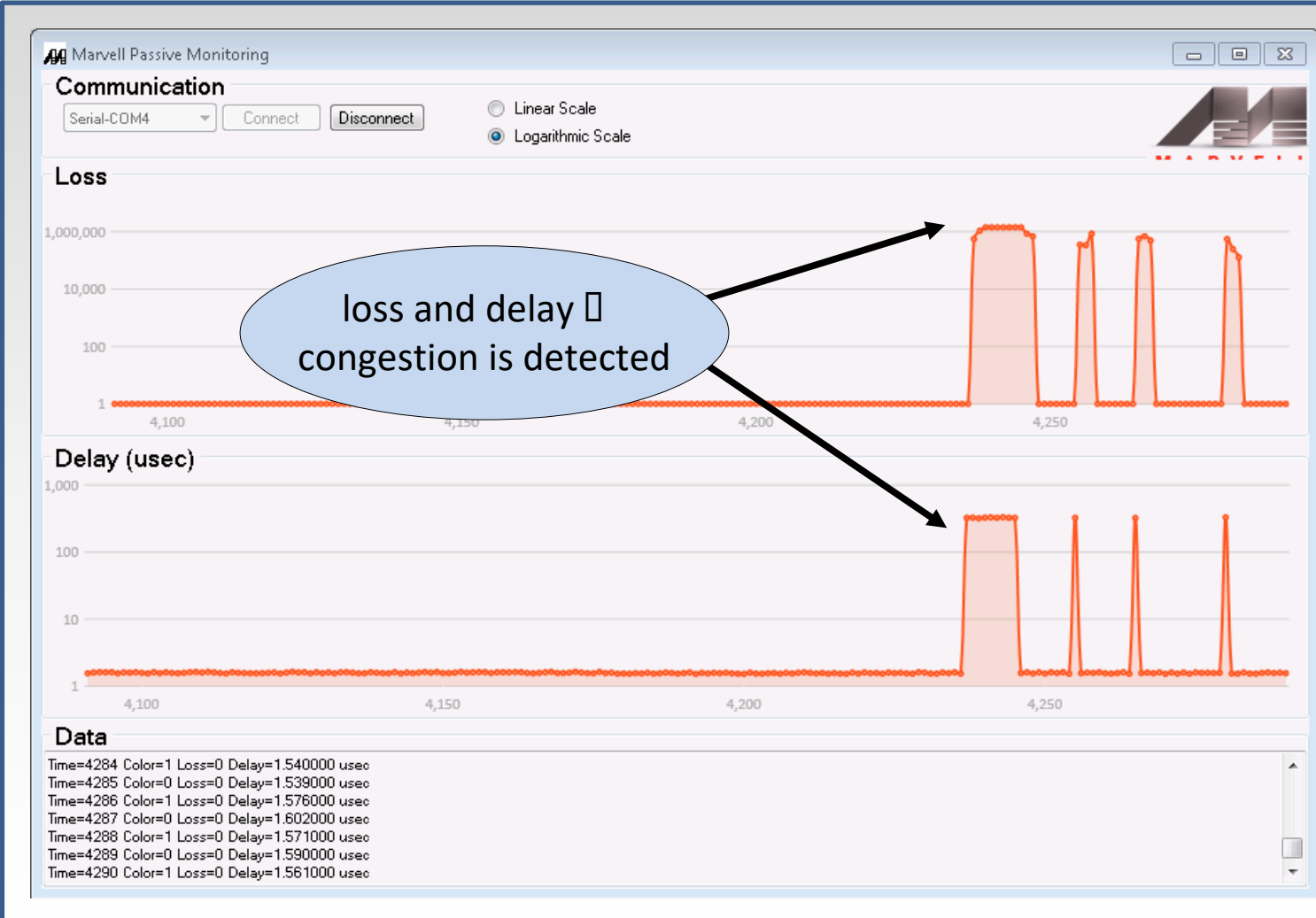
— Detect '1' pulse  
 — Detect '0' pulse



Match	Action
<i>TimeBits = 010. Rea = 0</i>	<i>MarkBit = 1. Rea = 1.</i>
Match	Action
<i>TimeBits = 001, MarkBit = 1</i>	<i>counter0, timestamp</i>
<i>TimeBits = 010, MarkBit = 1</i>	<i>counter0, timestamp</i>
<i>TimeBits = 101, MarkBit = 0</i>	<i>counter1, timestamp</i>
<i>TimeBits = 110, MarkBit = 0</i>	<i>counter1, timestamp</i>
<i>TimeBits = ***, MarkBit = 0</i>	<i>counter0</i>
<i>TimeBits = ***, MarkBit = 1</i>	<i>counter1</i>

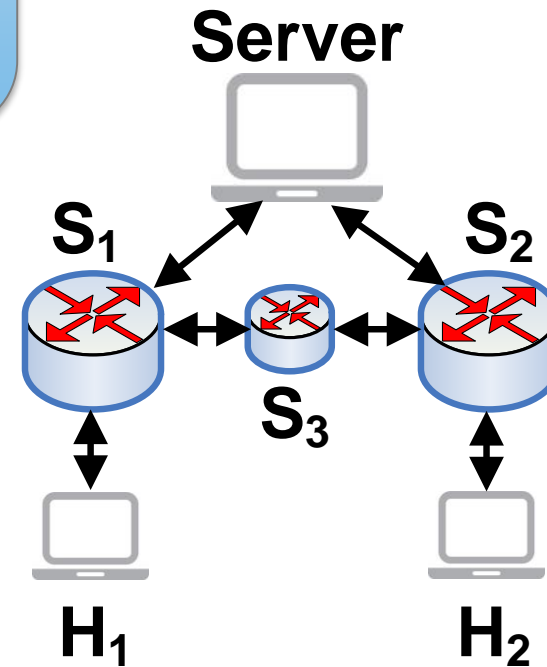
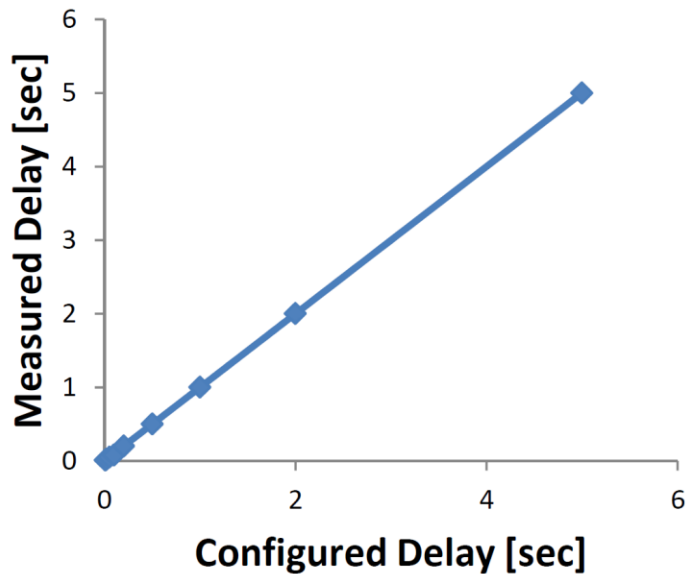
- TimeFlip is used to divide time into time
- The marking bit has a different interpretation in each time slot.
- Requires rough time synchronization, e.g., ~ 1 second.

# AM-PM Evaluation using Marvell Prestera Switches



# Software Implementation using P4

- Implemented in P4.
  - Time-of-day match field.
  - AM-PM in P4.
- Tested in Mininet.
- Open source code.





# AM-PM: Where is it going?

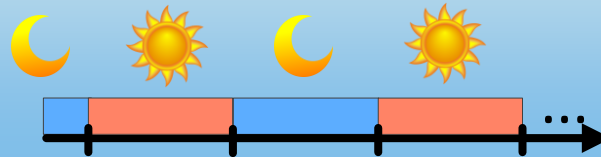
Network telemetry



Low overhead



AM-PM



Ongoing AM-PM work in the IETF:  
QUIC

MPLS

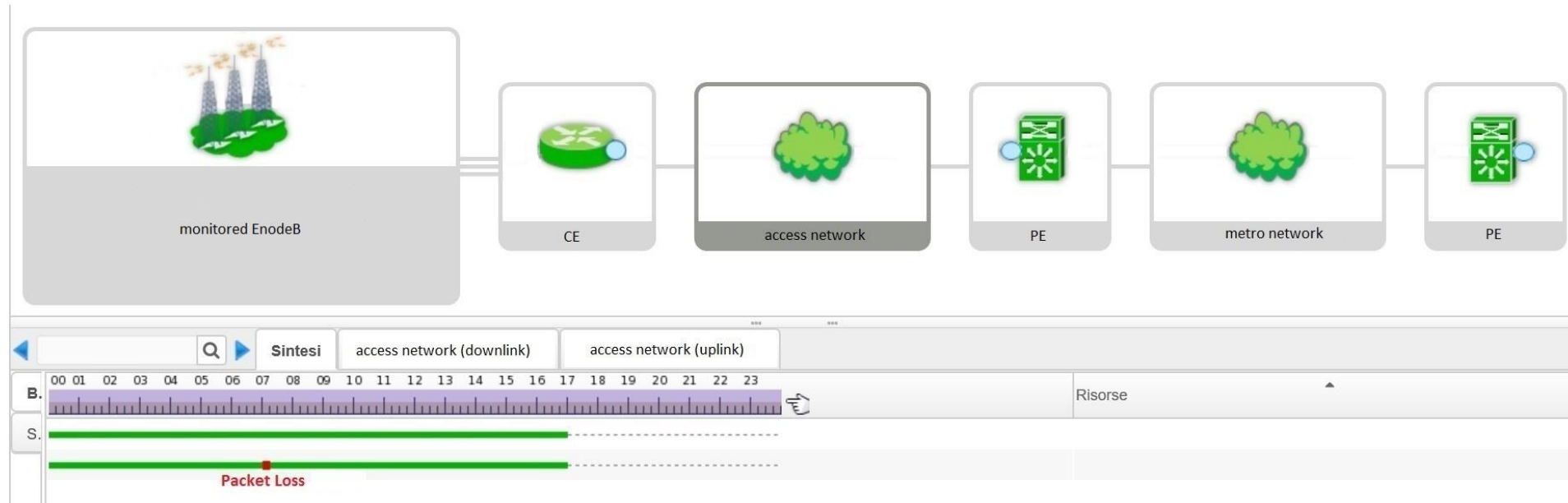
NSH

BIER

Geneve

AM-PM is under discussion in 6 working groups in the IETF...

# Large Scale Deployment in Telecom Italia



- **Mobile backhaul network ~ 1000 eNodeBs.**
- **AM-PM one bit (step-based) loss measurement.**
- **Uses unused bit in DSCP.**
- **Off-the-shelf network equipment.**

# Summary

Design and implementation of AM-PM

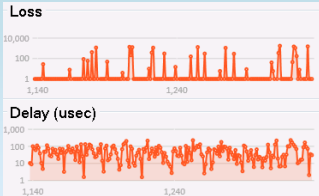
Hardware-based implementation using a Marvell switch.



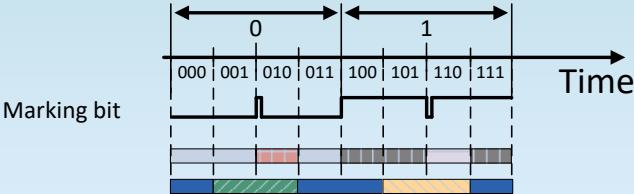
Software-based implementation in P4 – open source.

```
table look_for_flag {  
  reads {  
    intrinsic_metadata.time_of_day : ternary;  
    ipv4.flag_a : exact;  
  }  
  actions {  
    _look_for_flag;  
    _drop;  
  }  
  size: 256;  
}
```

Experimental results



Novel time-multiplexed parsing



12  
11 24  
10 23 1 13  
22 14 2  
9 21 **Thanks!** 15 3  
8 20 16 4  
19 17  
7 18 5  
6

# References

- [1] Fioccola, G., Capello, A., Cociglio, M., Castaldelli, L., Chen, M., Zheng, L., Mirsky, G., and T. Mizrahi, “Alternate Marking method for passive and hybrid performance monitoring”, [RFC 8321](#), 2018.
- [2] Mizrahi, T., Arad, C., Fioccola, G., Cociglio, M., Chen, M., Zheng, L., and G. Mirsky, “Compact Alternate Marking Methods for Passive and Hybrid Performance Monitoring”, [draft-mizrahi-ippm-compact-alternate-marking](#), work in progress, IETF, 2019.
- [3] Brockners, F., Bhandari, S., Pignataro, C., Gredler, H., Leddy, J., Youell, S., Mizrahi, T., Mozes, D., Lapukhov, P., Chang, R. and D. Bernier, J. Lemon, "Data Fields for In-situ OAM", [draft-ietf-ippm-ioam-data](#), work in progress, 2019.
- [4] C. Kim et al., “[In-band network telemetry \(INT\)](#)”, P4 consortium, 2015.
- [5] Mizrahi, T., Vovnoboy, V., Nisim, M., G. Navon, and A. Soffer, “[Network Telemetry Solutions for Data Center and Enterprise Networks](#)”, Marvell white paper, 2018.
- [6] Mizrahi, T., Rottenstreich, O. and Y. Moses, “TimeFlip: Scheduling Network Updates with Timestamp-based TCAM Ranges”, IEEE INFOCOM, 2015.
- [7] Mizrahi, T., Navon, G., Fioccola, G., Cociglio, M., Chen, M., and G. Mirsky, “AM-PM: Efficient Network Telemetry using Alternate Marking”, IEEE Network, 2019.
- [8] Riesenber, A., Kirzon, Y., Bunin, M., Galili, E., Navon, G., and T. Mizrahi, “Time-Multiplexed Parsing in Marking-based Network Telemetry”, ACM SYSTOR, 2019.
- [9] P4 AM-PM, <https://github.com/AlternateMarkingP4/FlaseClase>, 2018.