# Swisscom Network Analytics

Visibility for a closed loop operated network

02.12.2021, Thomas Graf and Marco Tollini

« The customer knows before Swisscom that there is service interruption.

Unable to recognize impact and root cause when configurational or operational network changes occur.

Swisscom suffers reputation damage.
**We need to work together to mediate.** »

**Markus Reber**
Head of Networks at Swisscom

« **At IETF only 9.85% of the activities are related to network automation and monitoring.**
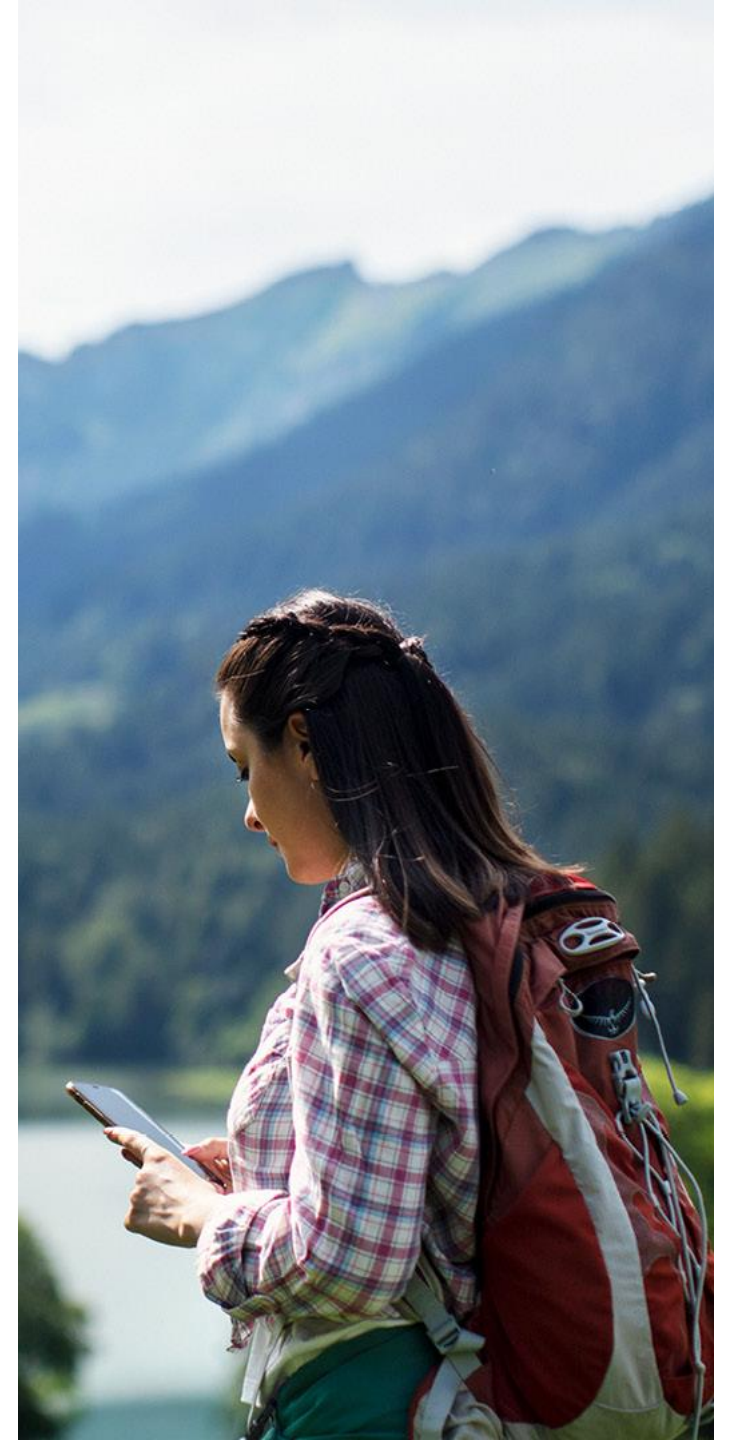
**We are still using protocols designed 40 years ago to manage networks.**

**IP network protocols are not made to expose metrics for analytics. <span style="color:red">IPFIX and BGP monitoring protocol are the rare exception.</span>** »
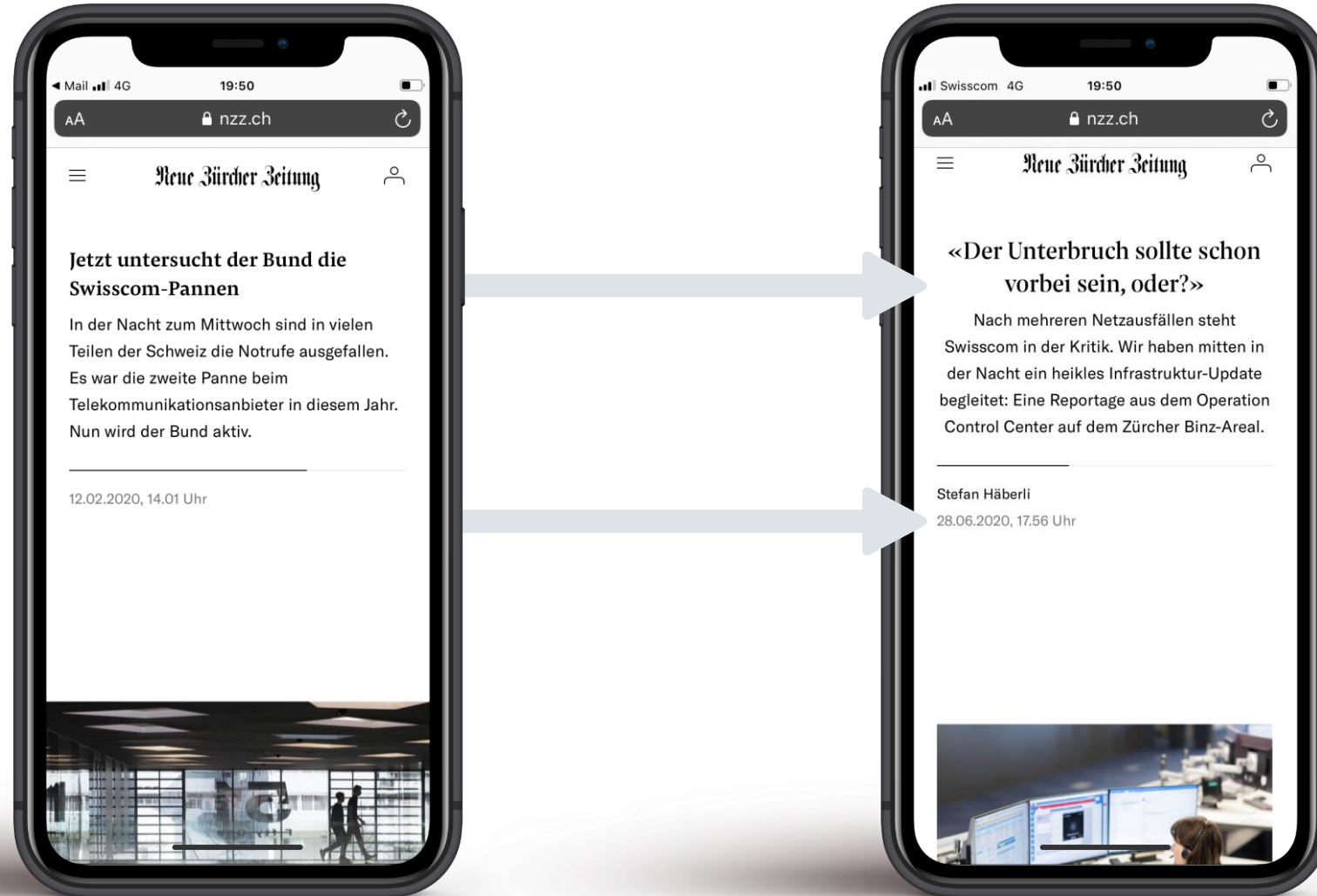
**Thomas Graf**
Distinguished Network Engineer
and Network Analytics Architect at Swisscom
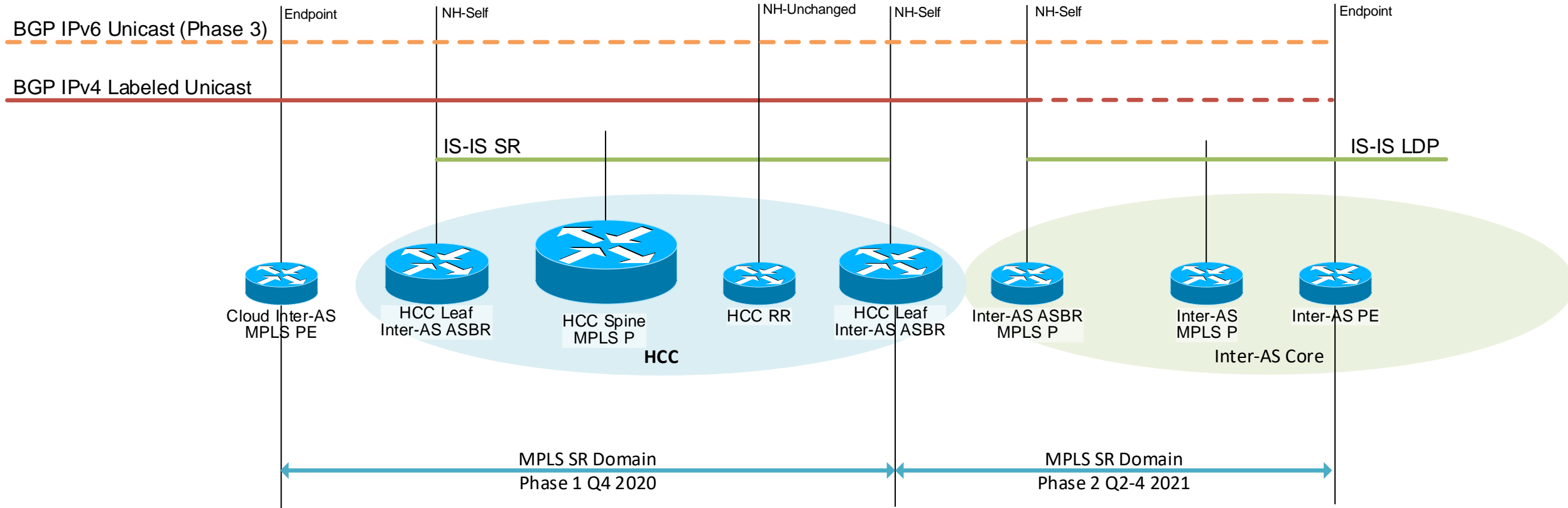
# Network Analytics Transformed Swisscom Media Reporting

Why networks and data lakes need to become one



**Jetzt untersucht der Bund die Swisscom-Pannen**

In der Nacht zum Mittwoch sind in vielen Teilen der Schweiz die Notrufe ausgefallen. Es war die zweite Panne beim Telekommunikationsanbieter in diesem Jahr. Nun wird der Bund aktiv.

12.02.2020, 14.01 Uhr

**«Der Unterbruch sollte schon vorbei sein, oder?»**

Nach mehreren Netzausfällen steht Swisscom in der Kritik. Wir haben mitten in der Nacht ein heikles Infrastruktur-Update begleitet: Eine Reportage aus dem Operation Control Center auf dem Zürcher Binz-Areal.

Stefan Häberli
28.06.2020, 17.56 Uhr

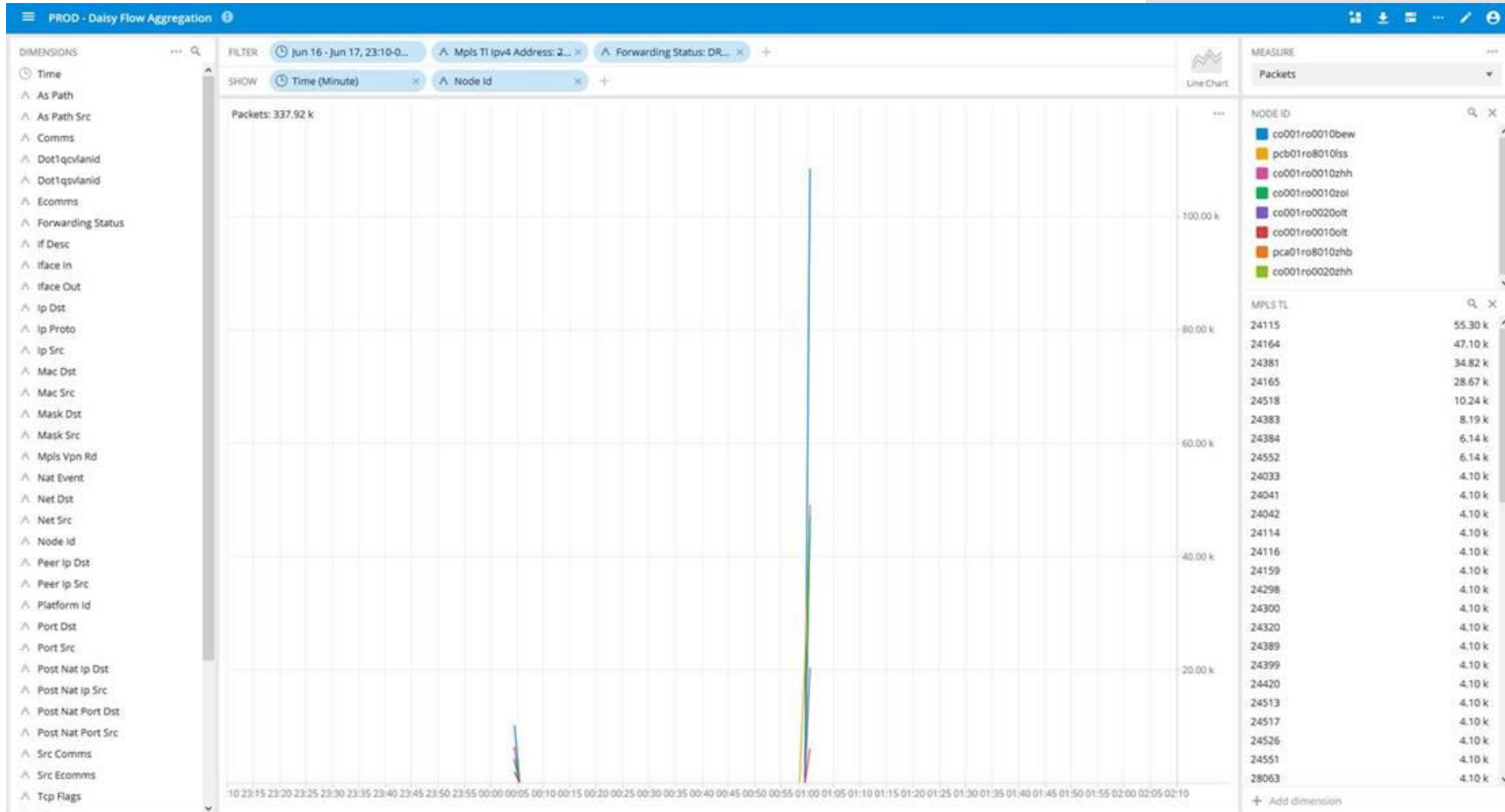# Transition to Segment Routing
## From MPLS over MPLS-SR to SRv6



Segment Routing **reduces the amount of routing protocols, simplifies forwarding-plane monitoring** while **enabling traffic engineering with closed loop** and increase scale.

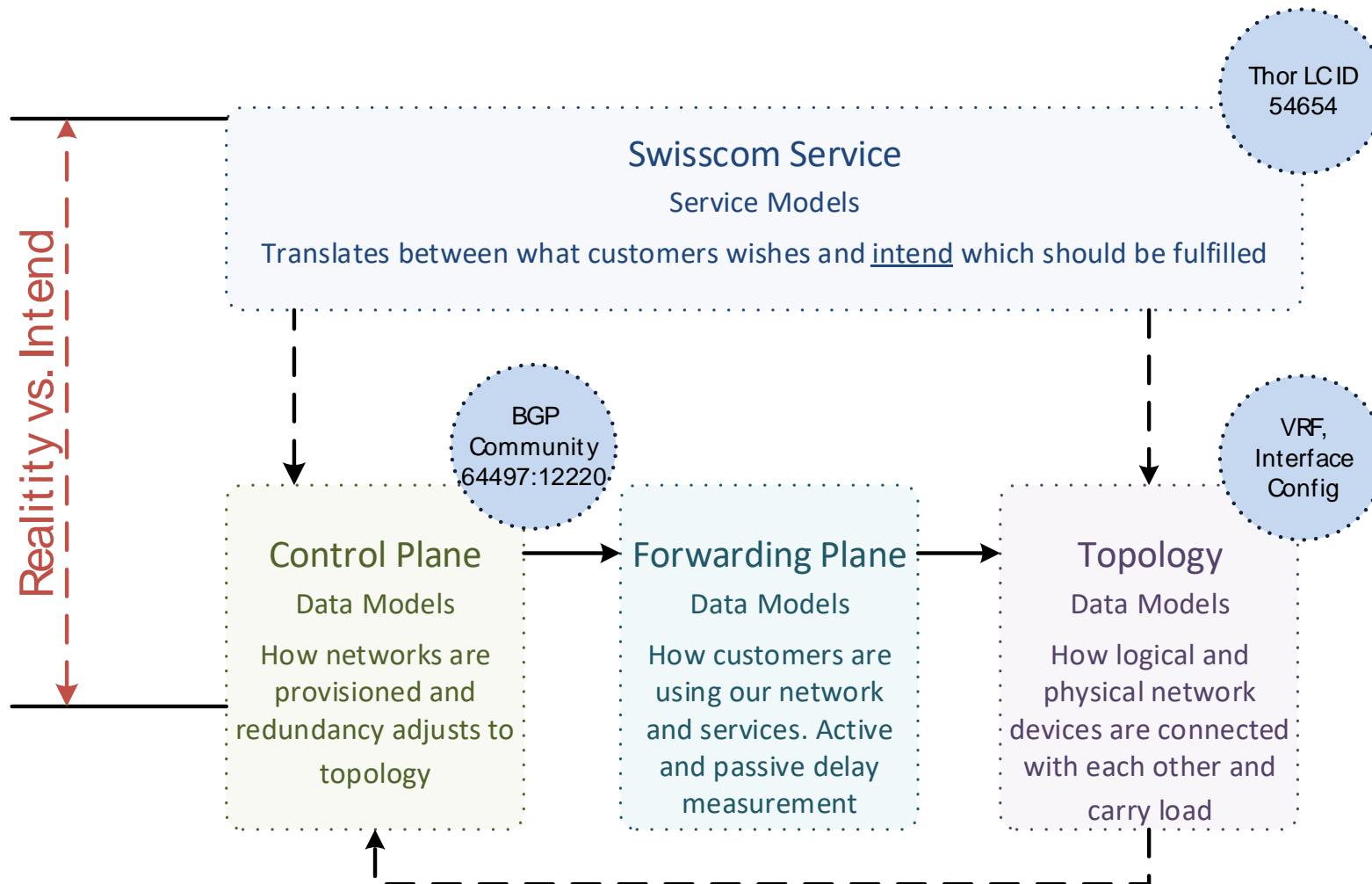# 337'920 Packets Dropped
Successfully migrated to a 3 label stack



6

# Data Collection with Network Telemetry

*Structured metrics enable informed decision-making*



**Reality vs. Intend**

**Swisscom Service**

Service Models

Translates between what customers wishes and <u>intend</u> which should be fulfilled

Thor LC ID
54654

BGP
Community
64497:12220

VRF,
Interface
Config

**Control Plane**

Data Models

How networks are provisioned and redundancy adjusts to topology

**Forwarding Plane**

Data Models

How customers are using our network and services. Active and passive delay measurement

**Topology**

Data Models

How logical and physical network devices are connected with each other and carry load

**Network Telemetry:**

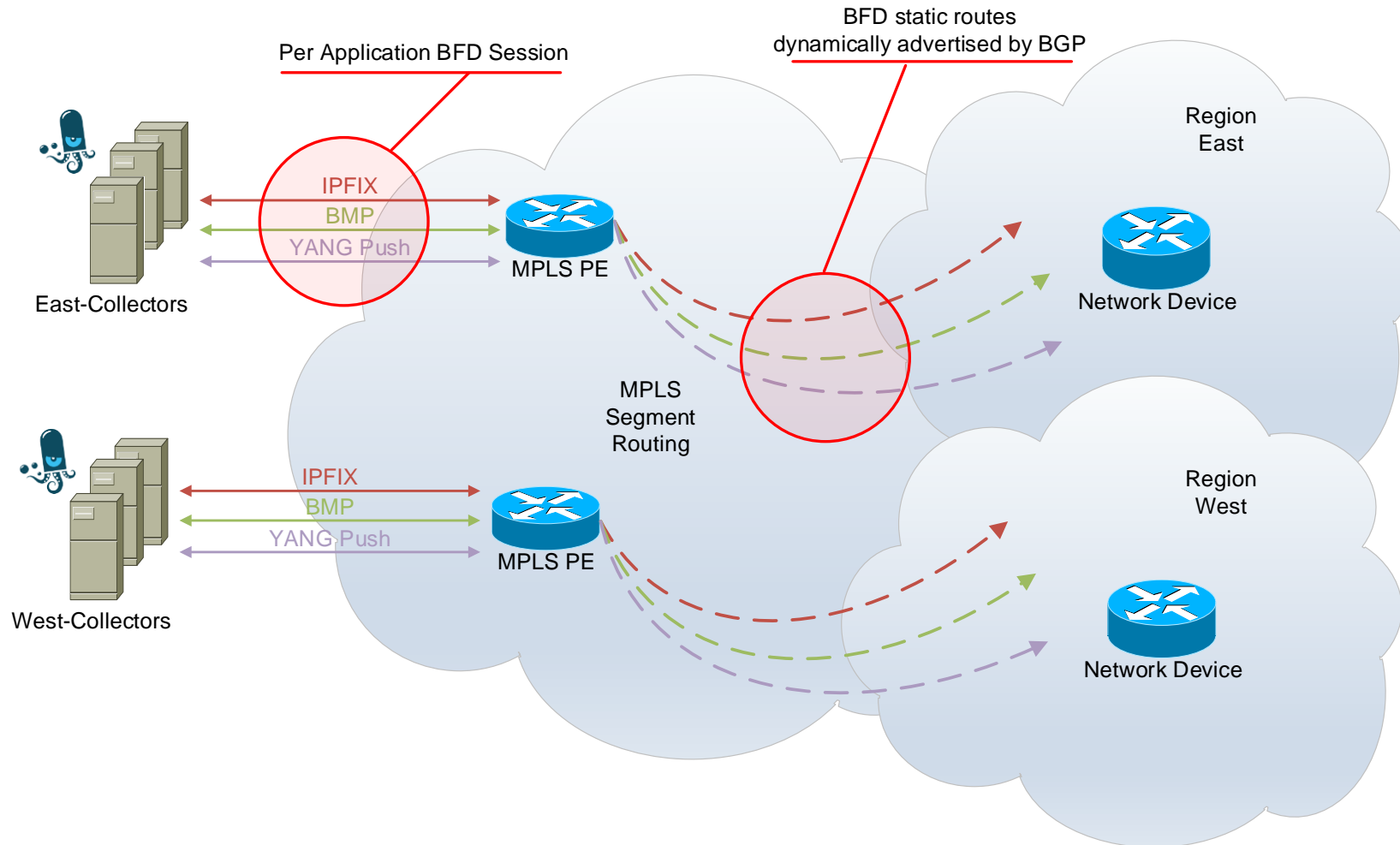> A data collection framework where the network device pushes its metrics to Big Data.

**Data Modelling:**

> Key for Big Data correlation to understand and react in the right context

> Are interface drops bad?

> How should we react?

# Network Distribution with BFD / Anycast
## Add as many servers as possible where you need them



- Each collection service is represented by an anycast service loopback address on the network.

- Depending if collection processes are running, service loopback installed in RIB with BFD static route and advertised by BGP at MPLS PE .

- Balances metrics from network devices to collectors depending on location **with 2 tuple hash (SRC/DST IP address).**
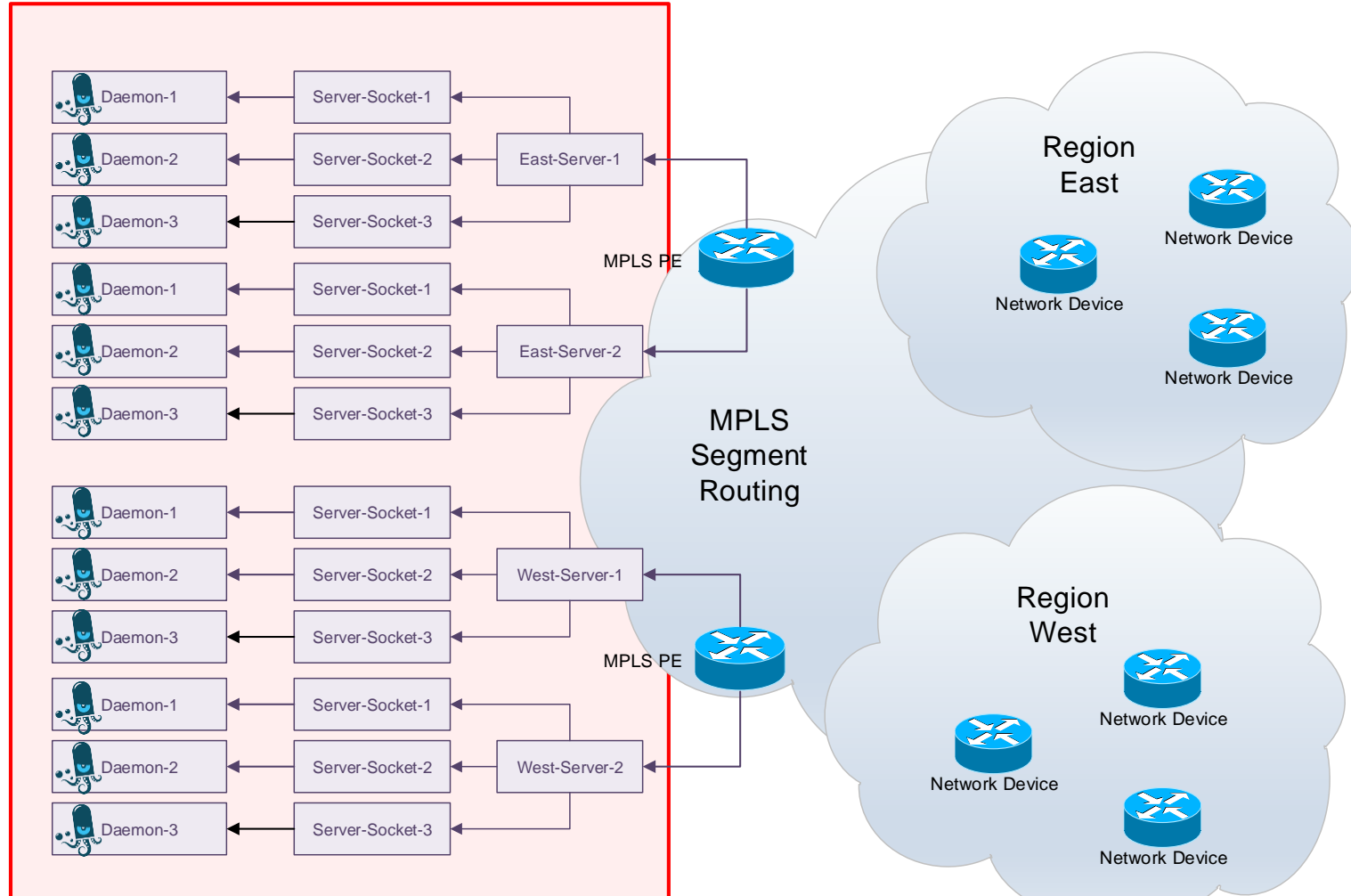
# Process Distribution with SO_REUSEPORT

Add as many daemons as possible where you need them
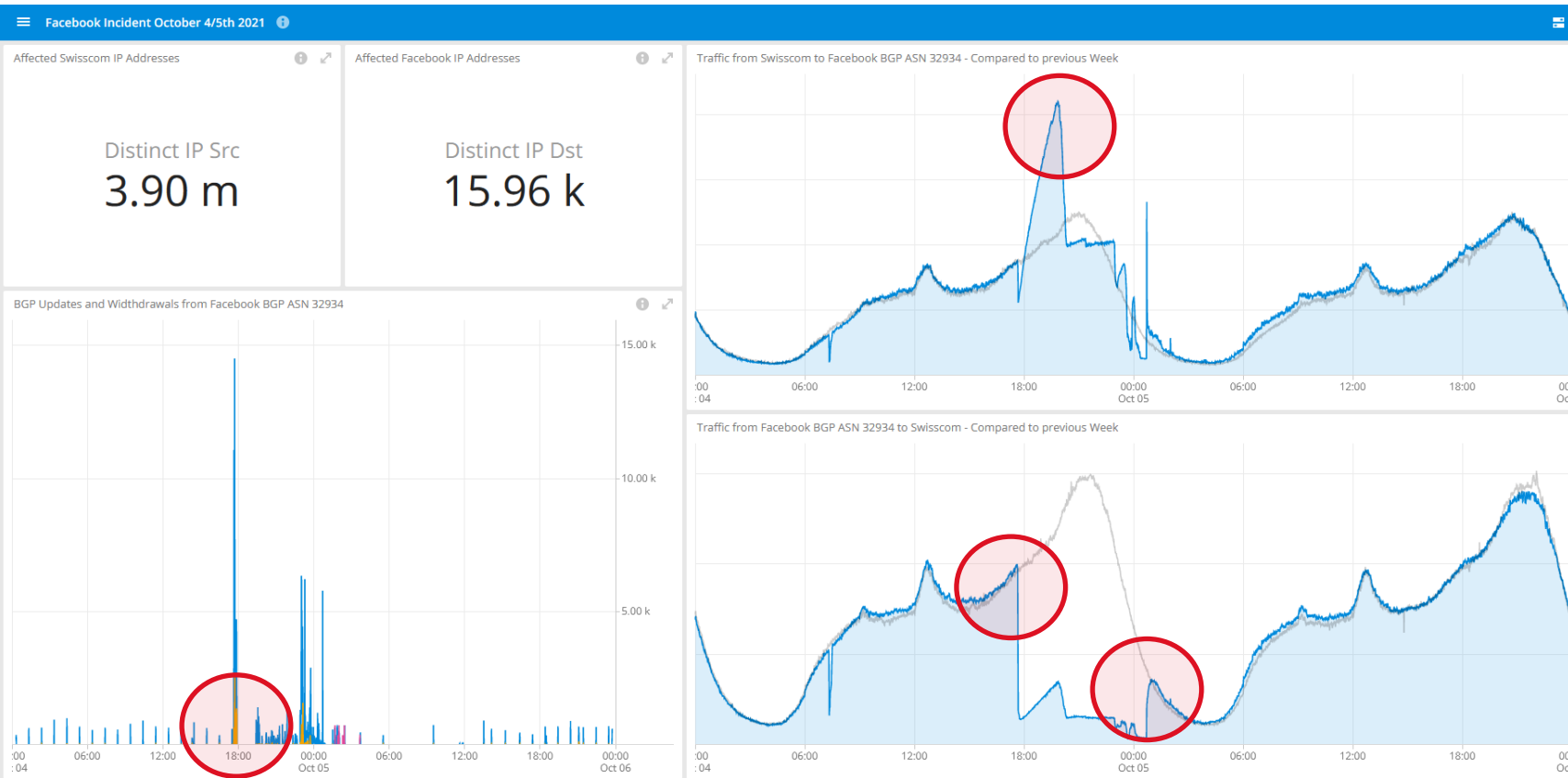
Linux Kernel SO_REUSEPORT Loadbalancing



- The Linux network socket is a bottleneck on servers with high connection/transaction rates.

- Linux kernel function SO_REUSEPORT allows infinite amount of network sockets per Layer4 port.

- Distributes metrics within server **with 2 tuple hash (SRC/DST IP address)** to daemons.

- Finite scale per server at lowest cost. CPU/memory resources per server is the only limitation factor.

# Facebook Incident October 4/5th
## The Swisscom perspective



**SOS** At 17:39 prefixes from Facebook BGP ASN 32934 where withdrawn. Outbound traffic steadily increased twofold until 20:20. Inbound traffic decreased by 85%.

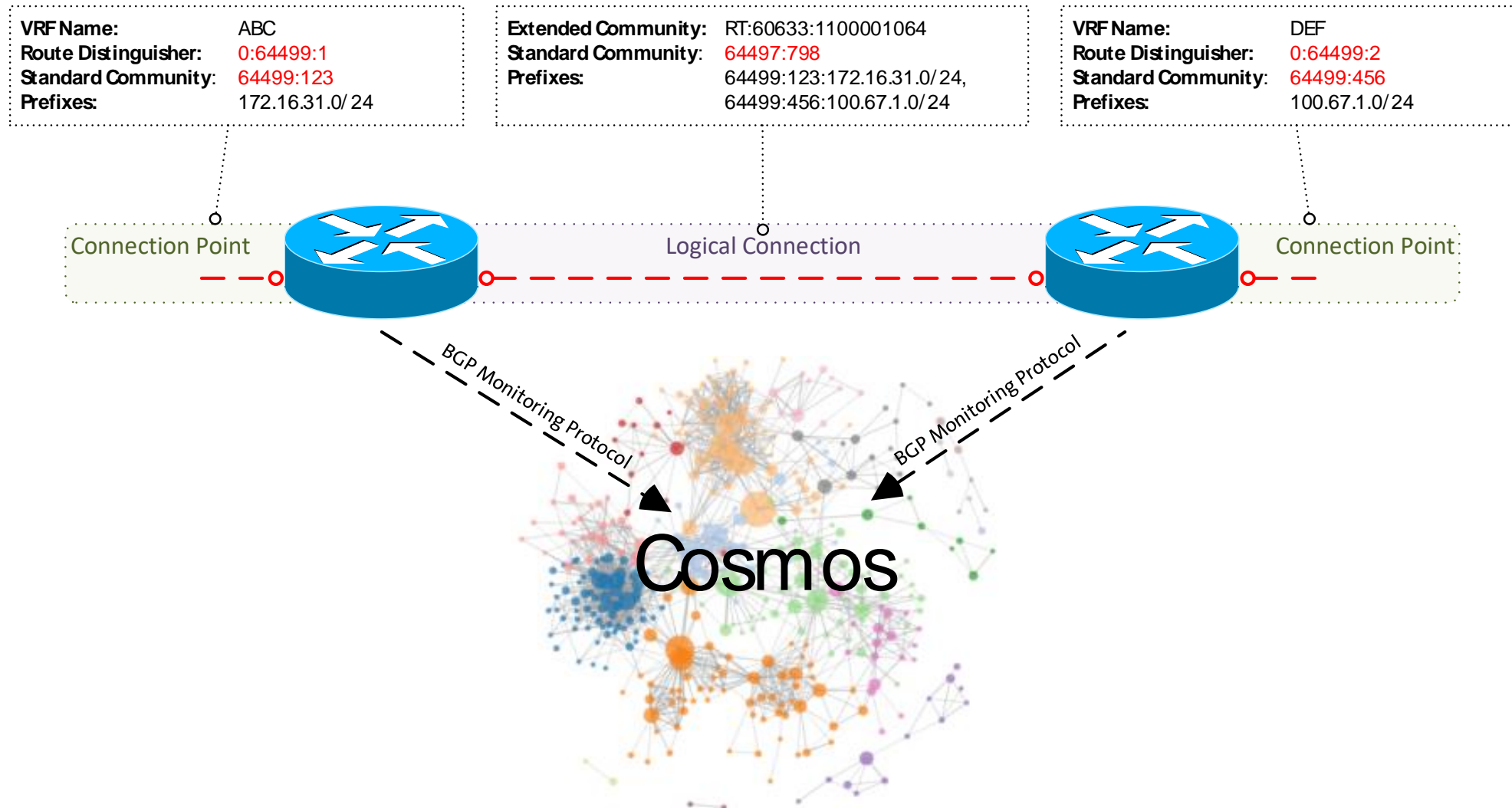Between 19:25 and 00:51, BGP updates and withdrawals where received.

At 00:41 traffic rate restored to normal.

# Visualizing Layer 3 VPN Topologies

Bringing Network Engineers visibility into topology changes

**VRF Name:** ABC
**Route Distinguisher:** 0:64499:1
**Standard Community:** 64499:123
**Prefixes:** 172.16.31.0/24

**Extended Community:** RT:60633:1100001064
**Standard Community:** 64497:798
**Prefixes:** 64499:123:172.16.31.0/24,
64499:456:100.67.1.0/24

**VRF Name:** DEF
**Route Distinguisher:** 0:64499:2
**Standard Community:** 64499:456
**Prefixes:** 100.67.1.0/24

Connection Point

Logical Connection

Connection Point

BGP Monitoring Protocol

BGP Monitoring Protocol

Cosmos

# The earth isn't flat, so are our networks
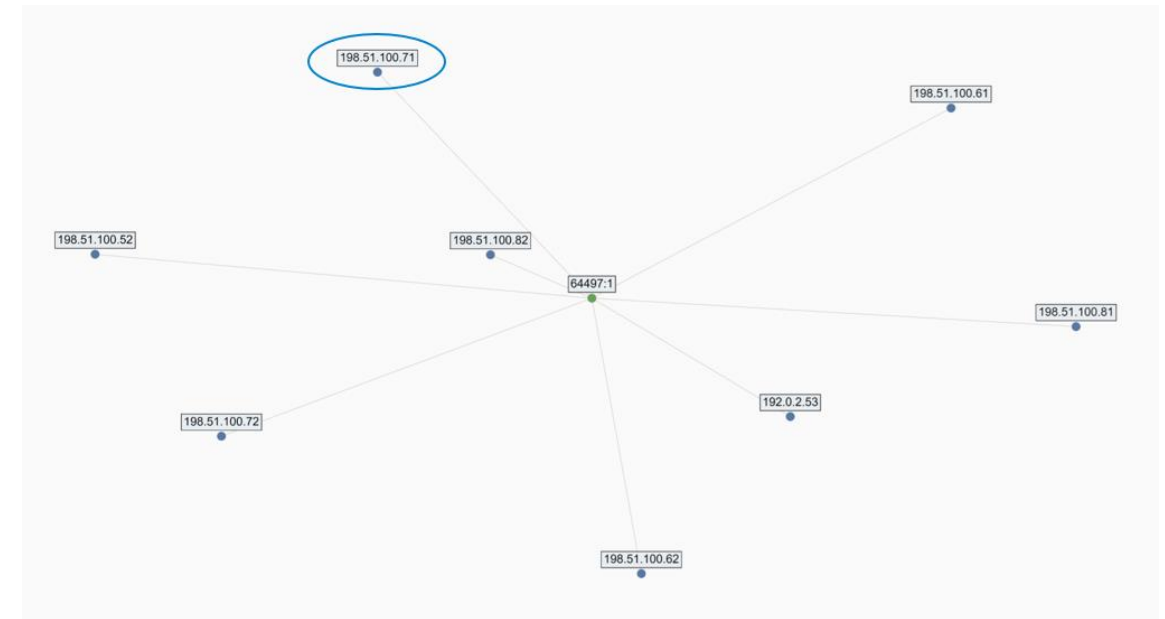BGP Communities are defining VPN's and Endpoints. Let's Visualize!



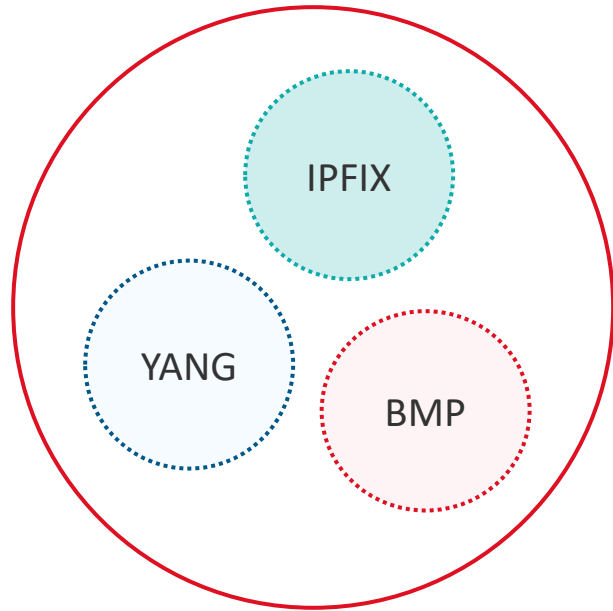**VS.**

**A table** showing
the **event** changes per node

**The Layer3 topology** visualization showing
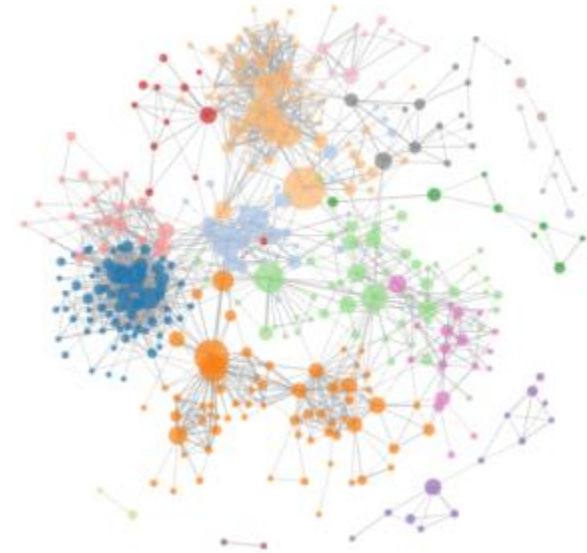the **state** changes per endpoint

# Objective of the project
Mapping the network, aiming for the stars



**Network Telemetry Protocols**

**Network Visualization**

# Event Based BGP Monitoring Protocol

BGP updates and withdrawals from the network

T1

T2

T3

T4 (current)

Timeline

**Advertise**
**1.1.1.1/24**

**Advertise**
**2.2.2.2/32**

**Withdraw**
**1.1.1.1/24**

# Challenge 1 – Obtaining the current BGP routing table
Keeping track of the changes

What is the current routing table state?

| T1 | T2 | T3 | T4 (current) |
|----|----|----|----|

Timeline

**Advertise 1.1.1.1/24**

**Advertise 2.2.2.2/32**

**Withdraw 1.1.1.1/24**

# Challenge 1 – Obtaining the current BGP routing table
Keeping track of the changes

Routing table:
- **2.2.2.2/32**

T1

T2

T3

T4 (current)

Timeline

**Advertise**
**1.1.1.1/24**

**Advertise**
**2.2.2.2/32**

**Withdraw**
**1.1.1.1/24**

# Challenge 2 – Retention time
## Do we have all the information?

We lost 2.2.2.2/32

T1
T2
T3
T4
T8 (current)

.....

Timeline

**Advertise**
**1.1.1.1/24**

**Advertise**
**2.2.2.2/32**

**Withdraw**
**1.1.1.1/24**

Retention Time = 5T

# Challenge 3 – Computation Time
Are we able to merge all the information?

## Solution - Transforming events to state
How did a BGP RIB looked at a given time and how did it change over time?

Dump routing table:
- **2.2.2.2/32**

Routing table:
- **2.2.2.2/32**
- **3.3.3.3/16**

T1  T2  T3  T4  +  T5  +  ...  +  T8 (current)

Advertise  Advertise  Withdraw  **DUMP**  Advertise      Timeline
1.1.1.1/24  2.2.2.2/32  1.1.1.1/24  **Routing Table**  3.3.3.3/16

Retention Time = 5T

**Key:** Dump Interval <= Retention Time

# Network Visualization Data pipeline

Realtime, what else?



| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| Network Devices | Pmacct Data Collection | Apache Kafka Message Broker | Apache Kafka Connect JDBC Sink | TimescaleDB |

# Cosmos

2021-05-29 10:40:38

Select VPN

Select visualization

## Active Filters

Place your filters here

## All Filters

Search

## Playback

Speed    Resolution

No data loaded

# Cosmos

🕐 2021-05-29 10:40:38

| Now | 2021-05-29 | 10:40:38 |
| --- | --- | --- |

| 1 Minute | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| 5 Minutes | | | 2021 | May | | | |
| 15 Minutes | « | ‹ | | | › | » | |
| 1 Hour | Sun | Mon | Tue | Wed | Thu | Fri | Sat |
| 2 Hours | 25 | 26 | 27 | 28 | 29 | 30 | 1 |
| 6 Hours | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| 12 Hours | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| Yesterday | 16 | 17 | 18 | 19 | 20 | 21 | 22 |
| | 23 | 24 | 25 | 26 | 27 | 28 | **29** |
| | 30 | 31 | 1 | 2 | 3 | 4 | 5 |

Now    OK

Current

No data loaded

5
4
3
2
1
0

# Cosmos

🕐 2021-05-29 10:40:38

Select VPN ⌃

64497:1

64497:2

64497:3

64497:33

Search

## Playback

▶

🕐 Speed | 🖼 Resolution

5

4

Current

2

No data loaded

1

0

# Cosmos

🕐 2021-05-29 10:40:38

64497:1 ⌄

Select visualization ⌃

VPN Topology

VPN Routing Topology

Peering Topology

List

| afi | aigp | as_path |

| as_path_id | bgp_id |

| bgp_nexthop | bmp_msg_type |

| bmp_router | bmp_router_port |

| comms | ecomms | event_type |

| ip_prefix | is_filtered | is_in |

| is_loc | is_out | is_post | label |

| lcomms | local_ip | local_port |

| local_pref | log_type | med |

● All  ● Filtered  ⊕ ⊖ 🔍 ✋ 🏠 ☰

6

4

2

0

Current

05:00    06:00    07:00    08:00    09:00    10:00    11:00    12:00    13:00    14:00    15:00    16:00

# Cosmos

2021-05-29 10:40:38

64497:1

**VPN Topology**

Active Filters

Place your filters here

All Filters

Search

| afi | aigp | as_path |

| as_path_id | bgp_id |

| bgp_nexthop | bmp_msg_type |

| bmp_router | bmp_router_port |

| comms | ecomms | event_type |

| ip_prefix | is_filtered | is_in |

| is_loc | is_out | is_post | label |

| lcomms | local_ip | local_port |

| local_pref | log_type | med |

198.51.100.71

198.51.100.61

198.51.100.52    198.51.100.82

64497:1

198.51.100.81

198.51.100.72

192.0.2.53

198.51.100.62

● All  ● Filtered

6
4
2
0

Current

05:00  06:00  07:00  08:00  09:00  10:00  11:00  12:00  13:00  14:00  15:00  16:00

# VPN Topology view
*The macro view of a VPN*

| VPN | Loopback | RD | RD Origin |
|-----|----------|----|-----------|
| 1 | 10.10.0.1 | A | bmp |
| 2 | 10.10.0.2 | B | bmp |
| 1 | 10.10.0.3 | C | bmp |
| 1 | 10.10.0.3 | D | bmp |
| 1 | 10.10.0.4 | E | bgp |

**Steps:**

> Calculate Router Table of each router

> Filter by VPN and RD Origin

> Select needed data

# VPN Topology view
*The macro view of a VPN*

| VPN | Loopback | RD | RD Origin |
|-----|----------|-----|-----------|
| 1 | 10.10.0.1 | A | bmp |
| ~~2~~ | ~~10.10.0.2~~ | ~~B~~ | ~~bmp~~ |
| 1 | 10.10.0.3 | C | bmp |
| 1 | 10.10.0.3 | D | bmp |
| ~~1~~ | ~~10.10.0.4~~ | ~~E~~ | ~~bgp~~ |

**Steps:**

> Calculate Router Table of each router

> Filter by VPN and RD Origin

> Select needed data

# VPN Topology view
*The macro view of a VPN*

| VPN | Loopback | RD | RD Origin |
|-----|----------|-----|-----------|
| 1 | 10.10.0.1 | A | bmp |
| ~~2~~ | ~~10.10.0.2~~ | ~~B~~ | ~~bmp~~ |
| 1 | 10.10.0.3 | C | bmp |
| 1 | 10.10.0.3 | D | bmp |
| ~~1~~ | ~~10.10.0.4~~ | ~~E~~ | ~~bgp~~ |

**Steps:**

> Calculate Router Table of each router

> Filter by VPN and RD Origin

> Select needed data

# Cosmos

64497:1

VPN Routing Topology

**Active Filters**

Place your filters here

**All Filters**

Search

| afi | aigp | as_path |
|-----|------|---------|

| as_path_id | bgp_id |
|------------|--------|

| bgp_nexthop | bmp_msg_type |
|-------------|--------------|

| bmp_router | bmp_router_port |
|------------|-----------------|

| comms | ecomms | event_type |
|-------|--------|------------|

| ip_prefix | is_filtered | is_in |
|-----------|-------------|-------|

| is_loc | is_out | is_post | label |
|--------|--------|---------|-------|

| lcomms | local_ip | local_port |
|--------|----------|------------|

| local_pref | log_type | med |
|------------|----------|-----|

192.0.2.53-0:64499:23

192.0.2.53-0:64499:13

198.51.100.62-0:64499:42

198.51.100.61-0:64499:11

198.51.100.82-0:64499:32

198.51.100.71-0:64499:21

198.51.100.81-0:64499:31

198.51.100.72-0:64499:22

192.0.2.53-0:64499:43

198.51.100.62-0:64499:12

192.0.2.53-0:64499:33

● All   ● Filtered

6

4

2

0

Current

05:00  06:00  07:00  08:00  09:00  10:00  11:00  12:00  13:00  14:00  15:00  16:00

# VPN Routing Topology view
## *The routing view of a VPN*

| VPN | Loopback | RD | Next hop | Type |
|-----|----------|-----|----------|------|
| 1 | 10.10.0.1 | A | 10.10.0.1 (self) | out |
| 1 | 10.10.0.2 | B | 10.10.0.1 | local |
| 1 | 10.10.0.3 | C | 10.10.0.1 | local |
| 1 | 10.10.0.2 | B | 10.10.0.2 (self) | out |
| 1 | 10.10.0.3 | C | 10.10.0.2 | local |

**Steps:**

> Calculate Routing Table of each router

> Filter by VPN

> Find all routers which are advertising routes with next-hop self

> Find all routers which are importing locally the routes

> Join on different loopback but equal next-hop

# Cosmos

**Active Filters**

Place your filters here

**All Filters**

Search

| afi | aigp | as_path |
|---|---|---|
| as_path_id | bgp_id | |
| bgp_nexthop | bmp_msg_type | |
| bmp_router | bmp_router_port | |
| comms | ecomms | event_type |
| ip_prefix | is_filtered | is_in |
| is_loc | is_out | is_post | label |
| lcomms | local_ip | local_port |
| local_pref | log_type | med |

2001:db8:12::153-0:64499:23

2001:db8:11::153-0:64499:13

2001:db8:11::162-0:64499:12

2001:db8:12::171-0:64499:21

2001:db8:11::161-0:64499:11

2001:db8:12::172-0:64499:22

192.0.13.182-0:64499:32

2001:db8:21::153-0:64499:43

192.0.11.162-0:64499:12

192.0.13.153-0:64499:33

192.0.11.153-0:64499:13

2001:db8:21::161-0:64499:41

192.0.11.161-0:64499:11

2001:db8:21::162-0:64499:42

192.0.13.181-0:64499:31

192.0.12.172-0:64499:22

2001:db8:13::181-0:64499:31

192.0.21.162-0:64499:42

192.0.12.153-0:64499:23

192.0.21.161-0:64499:41

2001:db8:13::153-0:64499:33

192.0.12.171-0:64499:21

192.0.21.153-0:64499:43

2001:db8:13::182-0:64499:32

🔴 All 🔵 Filtered

⊕ ⊖ 🔍 ✋ 🏠 ☰

6

4

2

0

Current

05:00  06:00  07:00  08:00  09:00  10:00  11:00  12:00  13:00  14:00  15:00  16:00

# Peering Topology view
*The peering view of a VPN*

| Loopback | RD | Local IP | Peer IP |
|----------|----|----------|---------|
| 10.10.0.1 | A | 190.10.0.1 | 190.10.0.2 |
| 10.10.0.2 | B | 190.10.0.2 | 190.10.0.1 |
| 10.10.0.1 | A | 190.10.0.1 | 190.10.0.3 |
| 10.10.0.3 | C | 190.10.0.3 | 190.10.0.1 |
| 10.10.0.4 |   | 190.10.0.4 | 190.10.0.2 |
| 10.10.0.2 |   | 190.10.0.2 | 190.10.0.4 |
| 10.10.0.4 |   | 190.10.0.4 | 190.10.0.3 |
| 10.10.0.3 |   | 190.10.0.3 | 190.10.0.4 |

**Steps:**

> Calculate Routing Table of each router

> Filter by VPN

> Get all loopbacks participating the VPN

> Filter peer up events by loopbacks

> Join twice the generated peer up table to create connection on:
  a.LocalIP = b.PeerIP and
  a.PeerIP = b.LocalIP



10.10.0.2 (B)

10.10.0.1 (A)

10.10.0.4 P Router

10.10.0.3 (C)

# Cosmos

bmp_router ✕  rd ✕  ip_prefix ✕  bgp_nexthop ✕  comms ✕

🕐 2021-05-29 10:40:38

64497:1 ⌄

List ⌄

Active Filters

Place your filters here

All Filters

Search

| afi | aigp | as_path |

| as_path_id | bgp_id |

| bgp_nexthop | bmp_msg_type |

| bmp_router | bmp_router_port |

| comms | ecomms | event_type |

| ip_prefix | is_filtered | is_in |

| is_loc | is_out | is_post | label |

| lcomms | local_ip | local_port |

| local_pref | log_type | med |

| bmp_router ⇕ | rd ⇕ | ip_prefix ⇕ | bgp_nexthop ⇕ | comms ⇕ |
|---|---|---|---|---|
| 192.0.2.53 | 0:64499:13 | 2001:db8::10/128 | 2001:db8:11::153 | 64496:299, 64496:1001, 64497:1, 64499:10 |
| 192.0.2.53 | 0:64499:13 | 2001:db8::10/128 | ::1 | 64496:299, 64496:1001, 64497:1, 64499:10 |
| 192.0.2.53 | 0:64499:13 | 2001:db8::15/128 | 2001:db8:11::144 | 64496:299, 64496:1001, 64497:1, 64497:2, 64499:15 |
| 192.0.2.53 | 0:64499:13 | 2001:db8::20/128 | 2001:db8:11::144 | 64496:299, 64496:1001, 64496:1033, 64497:1, 64499:20 |
| 192.0.2.53 | 0:64499:13 | 2001:db8::30/128 | 2001:db8:11::144 | 64496:299, 64496:1001, 64496:1033, 64497:1, 64499:30 |
| 192.0.2.53 | 0:64499:13 | 2001:db8::40/128 | 2001:db8:11::144 | 60633:1033, 64496:299, 64496:1001, 64497:1, 64497:2, 64499:40 |
| 192.0.2.53 | 0:64499:13 | 2001:db8::40/128 | 2001:db8:11::161 | 60633:1033, 64496:299, 64496:1001, 64497:1, 64497:2, 64499:40 |
| 192.0.2.53 | 0:64499:13 | 2001:db8::40/128 | 2001:db8:11::162 | 60633:1033, 64496:299, 64496:1001, 64497:1, 64497:2, 64499:40 |

🔴 All  🔵 Filtered

⊕ ⊖ 🔍 ✋ 🏠 ☰

6

4

2

0

Current

05:00  06:00  07:00  08:00  09:00  10:00  11:00  12:00  13:00  14:00  15:00  16:00

# We need Network Analytics to meet the challenge

Maximize Uptime. Networks are using BGP to steer traffic and ensure redundancy.

With millions of routes in thousands of routing contexts and ten thousands of route-policies, to predict high availability, is for humans with CLI almost impossible.

> Which connection points are supposed to be highly available and are not?

> When a router or link is turned off for maintenance, which router or link will take over?

> Do all routers and links which are on standby have enough capacity to take over?

> When a route-policy is changed, how will the BGP attributes be affected in a logical connection and how will it affect the route propagation across the network?
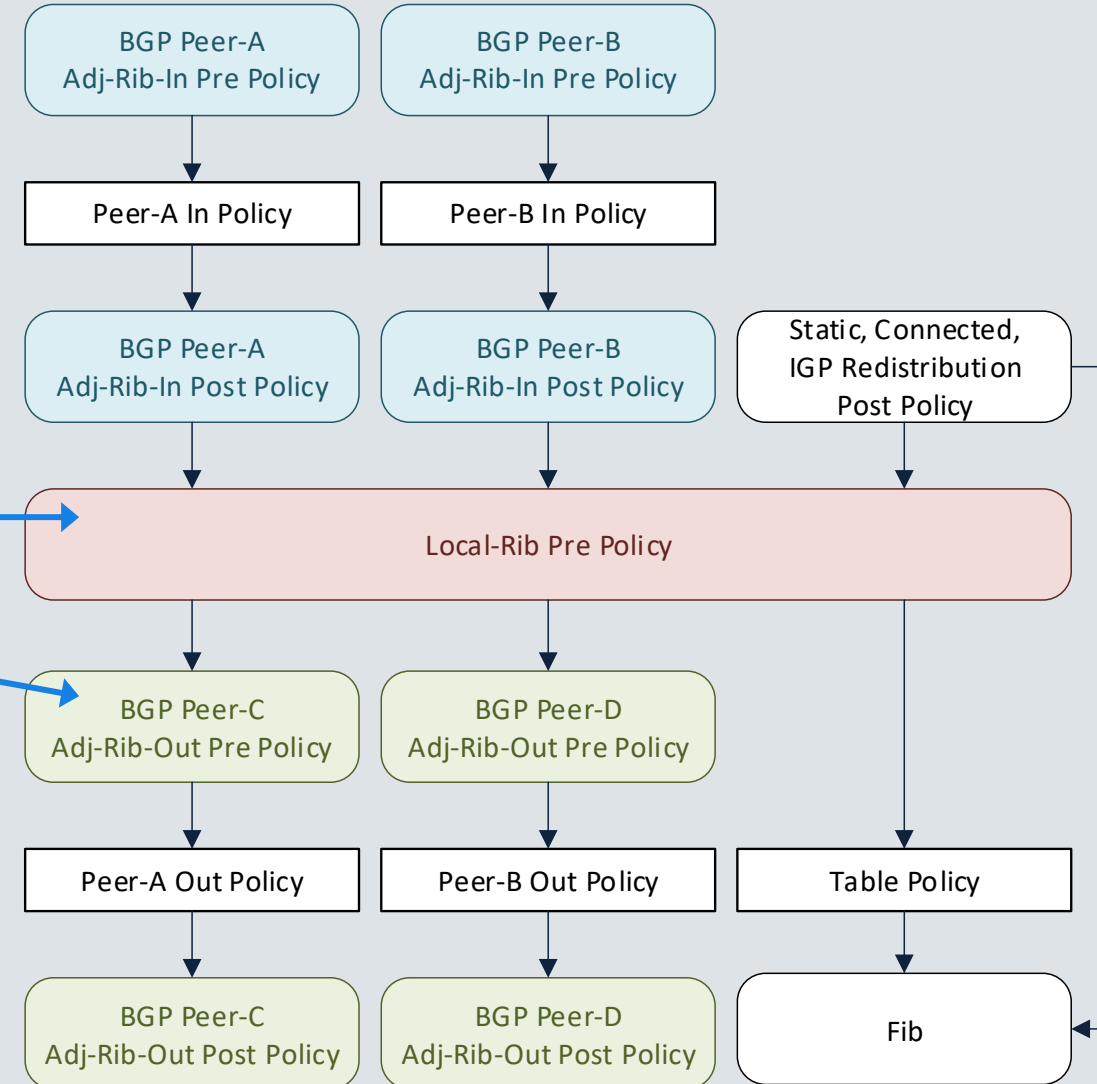
# BMP Covering all RIB's

Extends much needed RIB coverage

**BGP route exposure without BMP is a challenge of the first order:**

> Only best path is exposed (missing best-external and ECMP routes)

> Next-hop attribute not preserved all the time (allpaths)

> Filtering between RIB's not visible

- **Support for Local RIB in BGP Monitoring Protocol**
  https://tools.ietf.org/html/draft-ietf-grow-bmp-local-rib

- **Support for Adj-RIB-Out in BGP Monitoring Protocol**
  https://tools.ietf.org/html/rfc8671

Adj-RIB-Out an RFC since November 2019. Local RIB will follow soon. Juniper, Huawei and Nokia have public releases available supporting both. Cisco has test code available but haven't released yet.
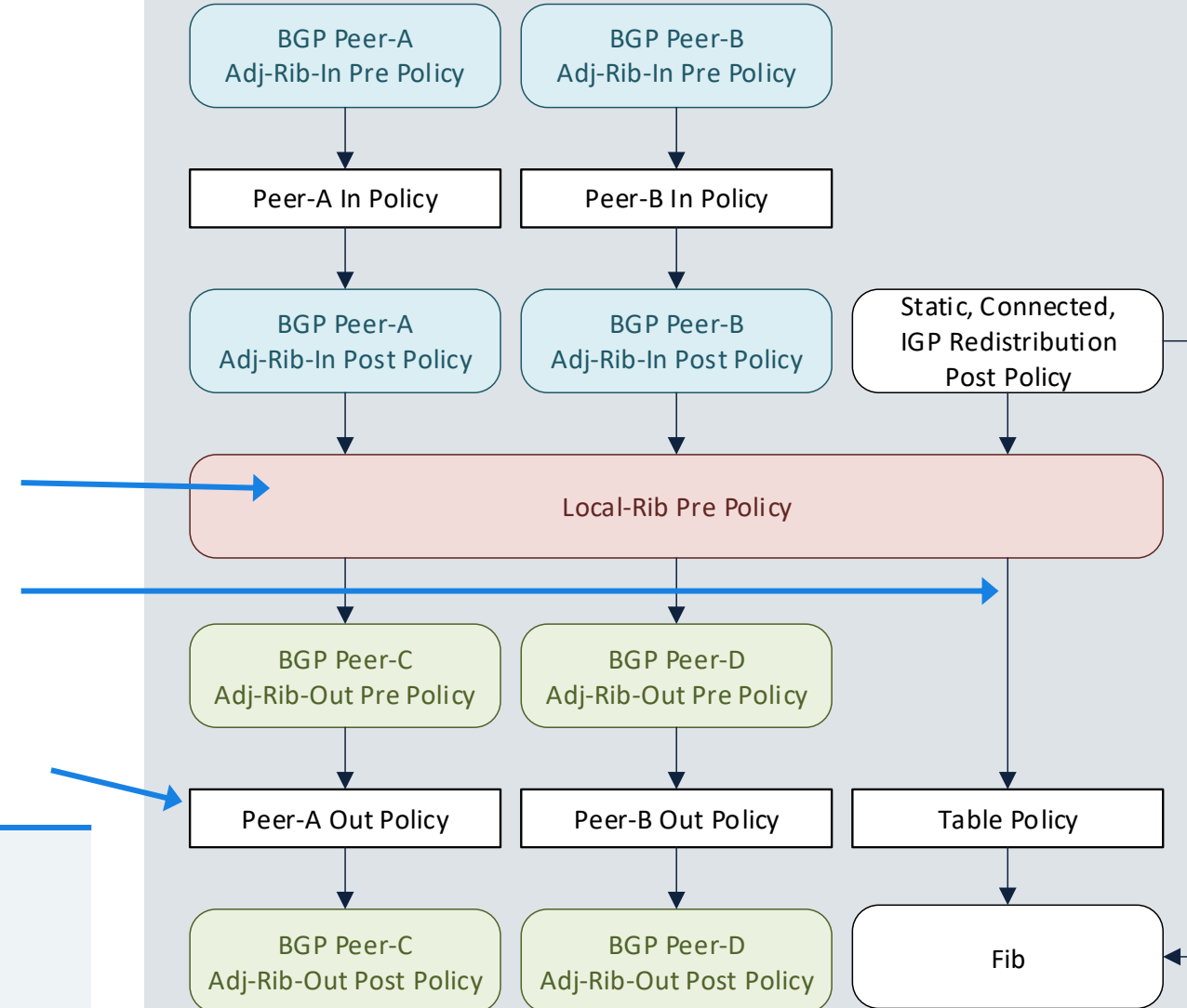
# BMP with extended TLV support
## Brings visibility into FIB's and route-policies

**Knowing all the routes in all the RIB's brings the new challenge**

> That we don't know how they are being used in the FIB/RIB (which one is best, best-external, ECMP, backup)

> That we don't know which route-policy permitted/denied/changed which prefix/attribute

• **TLV support for BMP Route Monitoring and Peer Down Messages**
https://tools.ietf.org/html/draft-ietf-grow-bmp-tlv

• **Support for Enterprise-specific TLVs in the BGP Monitoring Protocol**
https://tools.ietf.org/html/draft-lucente-grow-bmp-tlv-ebit

• **BMP Extension for Path Marking TLV**
https://tools.ietf.org/html/draft-cppy-grow-bmp-path-marking-tlv

• **BGP Route Policy and Attribute Trace Using BMP**
https://tools.ietf.org/html/draft-xu-grow-bmp-route-policy-attr-trace

For IETF 108 Hackathon, IETF lab network with Big Data integration has been further extended to collaborate development research with ETHZ, INSA, Imply, Huawei and pmacct (open source data-collection by Paulo Lucente).

# IPFIX Covering Segment Routing
## For MPLS-SR and SRv6

**Segment Routing coverage in IPFIX bring visibility for:**

> Which routing protocol provided the label in MPLS-SR.

> The IPv6 Segment where the packet is forwarded to in SRv6.

> The IPv6 Segments where the packet is going to be forwarded through in SRv6.

- **Export of MPLS Segment Routing Label Type Information in IPFIX**
  https://datatracker.ietf.org/doc/html/draft-ietf-opsawg-ipfix-mpls-sr-label-type

- **IPFIX export of Segment Routing IPv6 information**
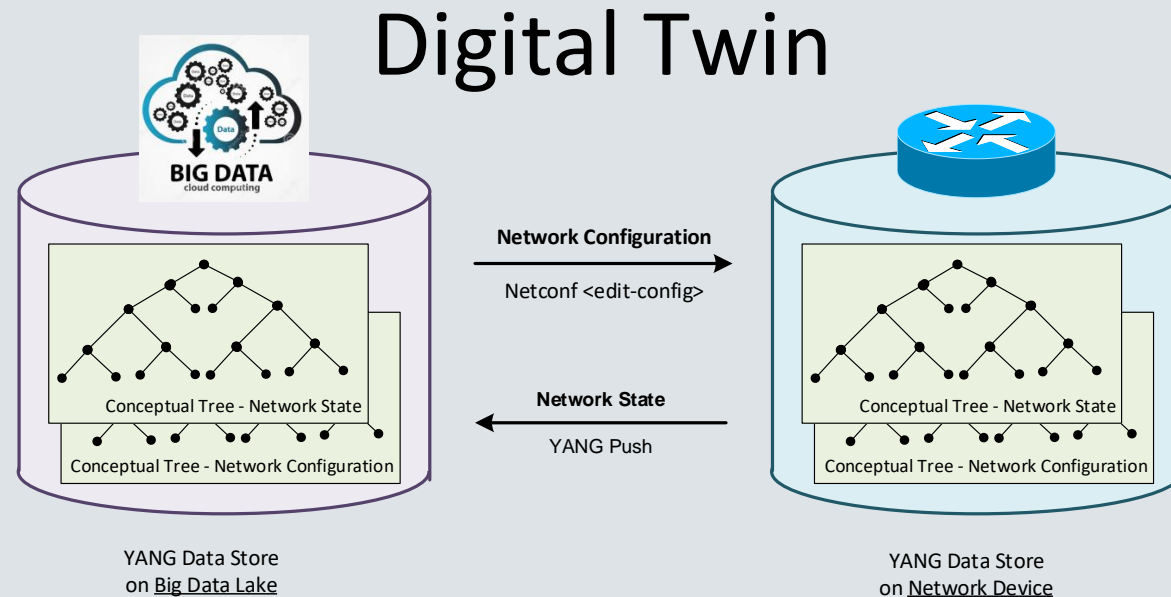  https://datatracker.ietf.org/doc/html/draft-patki-srv6-ipfix

draft-ietf-opsawg-ipfix-mpls-sr-label-type at final stage at IESG.
Driven by Swisscom. draft-patki-srv6-ipfix not being adopted yet.
Driven by Cisco. Swisscom is going to co-author.

```
> Frame 527: 182 bytes on wire (1456 bits), 182 bytes captured (1456 bits)
> Ethernet II, Src: Cisco_ea:ad:1c (00:32:17:ea:ad:1c), Dst: Vmware_21:95:d2 (00:0c:29:21:95:d2)
> Internet Protocol Version 4, Src: 138.187.57.63, Dst: 138.187.58.13
> User Datagram Protocol, Src Port: 44542, Dst Port: 9991
v Cisco NetFlow/IPFIX
    Version: 9
    Count: 1
    SysUptime: 516154.381000000 seconds
  > Timestamp: Feb 23, 2020 13:57:18.000000000 W. Europe Standard Time
    FlowSequence: 23685
    SourceId: 0
  v FlowSet 1 [id=313] (1 flows)
      FlowSet Id: (Data) (313)
      FlowSet Length: 120
      [Template Frame: 9]
    v Flow 1
      > MPLS-Label1: 17002 exp-bits: 0
      > MPLS-Label2: 24622 exp-bits: 0 bottom-of-stack
      > MPLS-Label3: 0 exp-bits: 0
      > MPLS-Label4: 0 exp-bits: 0
      > MPLS-Label5: 0 exp-bits: 0
      > MPLS-Label6: 0 exp-bits: 0
        InputInt: 87
        OutputInt: 111
        Octets: 216000
        Packets: 2000
      > [Duration: 5.753000000 seconds (switched)]
        TopLabelAddr: 138.187.57.13
        SrcAddr: ::
        DstAddr: ::
        ipv6FlowLabel: 0
        IPv6 Extension Headers: 0x00000000
        SrcAddr: 10.248.4.236
        DstAddr: 10.248.4.222
        SrcPort: 0
        DstPort: 2048
        MPLS Top Label Prefix Length: 32
        TopLabelType: LDP (5)
      > Forwarding Status
        Direction: Ingress (0)
        IP ToS: 0x00
        Protocol: ICMP (1)
      > TCP Flags: 0x00
        SamplerID: 1
        Ingress VRFID: 1610612736
        Egress VRFID: 1610612736
      Padding: 0000
```

# YANG Datastores enables Closed Loop Operation
## Automated data correlation – what else?



Digital Twin

Network Configuration
Netconf <edit-config>

Network State
YANG Push

Conceptual Tree - Network State

Conceptual Tree - Network Configuration

BIG DATA
cloud computing

YANG Data Store
on Big Data Lake

YANG Data Store
on Network Device

The IAB (Internet Architecture Board) at IETF took serious steps to bring automation into networks.
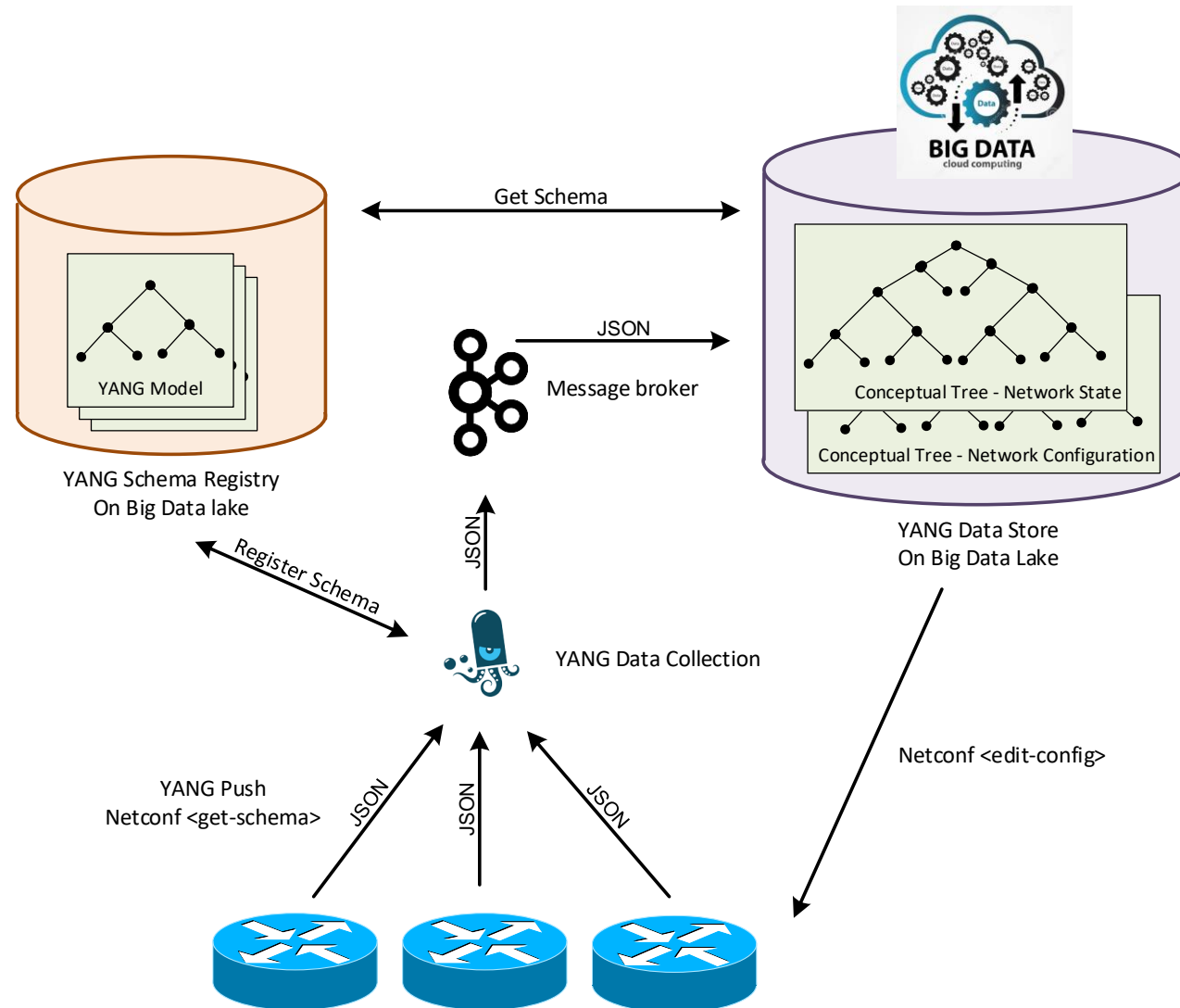
**At its core is YANG, a data modelling language which will not only transform how we managed our networks, it will transform also how we manage our services.**

**Automated networks can only run with  a common data model**. A digital twin YANG data store enables a comparison between intend and reality. Schema preservation enables closed loop operation. Closed Loop is like an autopilot on an airplane. We need to understand what the flight envelope is to keep the airplane within. Without, we crash.

# When Big Data and Network become one
A simple, scalable approach to YANG push



**Simplify** YANG push network data collection with high scale and low impact. **Suited for nowadays distributed forwarding systems.**

**Preserve YANG data model schema** definition throughout the data processing chain.

**Enable automated data correlation** among device, forwarding-plane and control-plane.

**UDP-based Transport for Configured Subscriptions**
https://datatracker.ietf.org/doc/html/draft-ietf-netconf-udp-notif

**Subscription to Distributed Notifications**
https://datatracker.ietf.org/doc/html/draft-ietf-netconf-distributed-notif

# Network Telemetry Overview
## Standards matter

### Why BMP?

Well established since June 2016 among major vendors and open-source community. Future proven thanks to encapsulating BGP PDU in BMP route-monitoring messages. Provides initial RIB state by providing initial state with subsequent updates for minimal performance penalty.

### Why IPFIX?

IPFIX succeeded well because of covering all three perspectives since day one: forwarding-plane, control-plane, device. In order to enable data correlation among different perspectives, key fields from other perspectives need to be present.
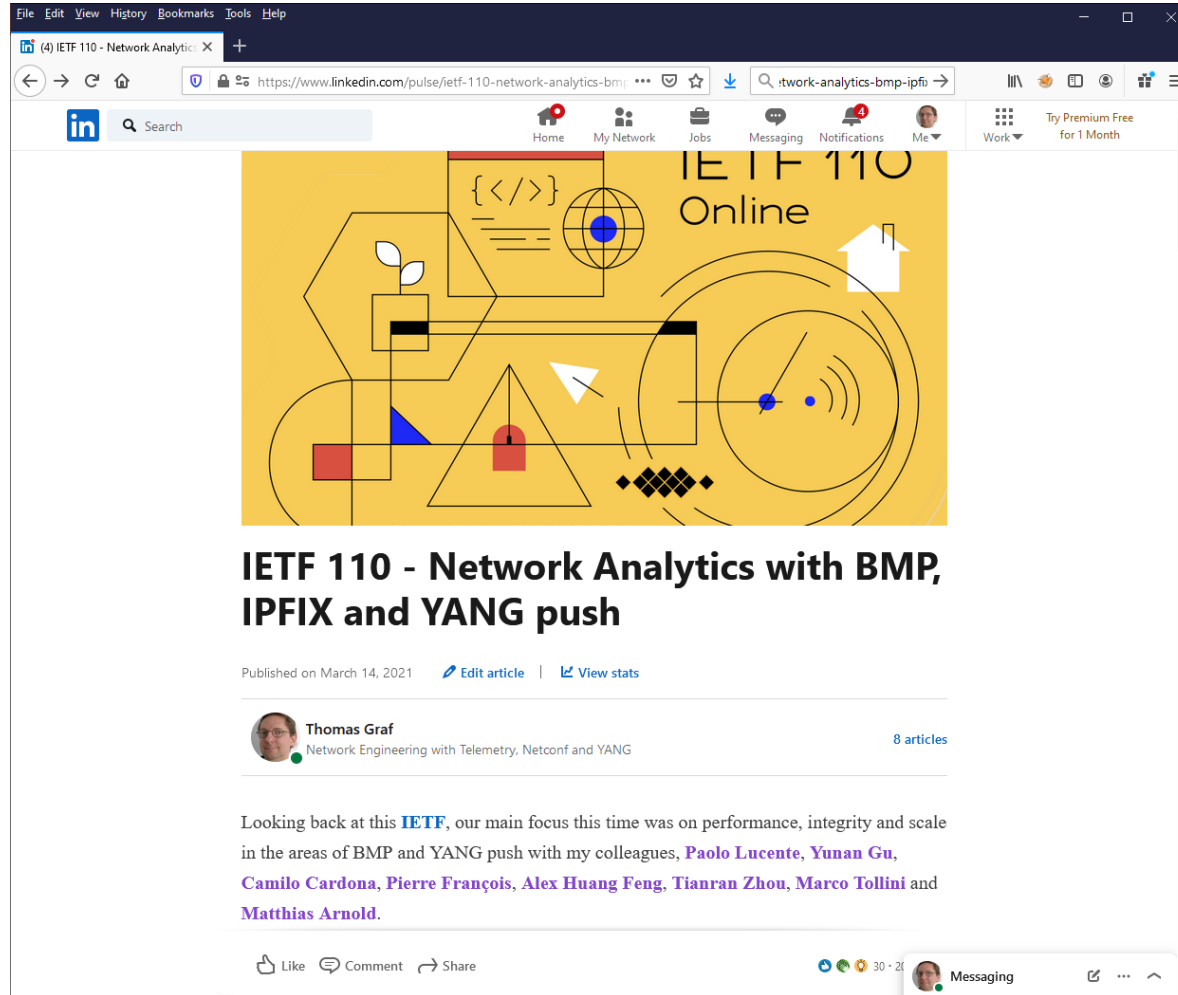
### Why YANG Push?

YANG is de-facto standard in network automatization since August 2016. With YANG push data-collection is going to be finally standardized and enabling closed loop operation frameworks.

# IETF 110 – Network Analytics
## With BMP, IPFIX and YANG Push

https://www.linkedin.com/pulse/ietf-110-network-analytics-bmp-ipfix-yang-push-thomas-graf/

👍 **6x BMP drafts** at GROW working group. Bringing RIB and route-policy dimensions into BMP and increase scale.

👍 **2x YANG push drafts** at NETCONF working group.

👍 **1x IPFIX MPLS Segment Routing draft** at OPSAWG working group.

✓ **Running code** being tested in IETF interoperability lab at 110 hackathon.

# ETH Zürich, Master Thesis Proposal – March 2022

High Availability with BGP monitoring Protocol Data Collection

## High Availability with BGP Monitoring Protocol Data Collection

Master thesis proposal with Swisscom

### Description

Swisscom collects millions of Network Telemetry [1] metrics every second with BMP [2], IPFIX [3] and YANG push [4] from thousands of network devices. In order to meet scaling demands of the data-collection, Swisscom uses a highly distributed load-balancing scheme across servers, Linux network sockets [5] and collector daemons. To further reduce the scale demands for the Big Data analysis, IPFIX and BMP metrics are highly aggregated [6] over a specified time bin during the data-collection.

This architecture imposes that the BGP [7] RIB state, which is collected through BMP route-monitoring messages, needs to be cached at the data-collection. The preservation of BGP RIB state caching across daemons is challenging, especially if faced with reload or migration events due to software upgrades or re-balancing decisions.

During this thesis you will first learn what metrics are collected with Network Telemetry, how they relate in terms of control-plane, forwarding-plane and device characteristics and how this allows us to distinguish between measurements and different dimensions. You will also understand why network schema needs to be preserved for a metric correlation which enables network-wide visibility. Finally, you will realize how Swisscom uses (i) Anycast [8] with ECMP [9] to distribute traffic across Layer 3 links and routers; and (ii) SO_REUSEPORT with an eBPF enhancement [10] to distribute incoming telemetry data to different collection processes on a server.

You will research and document how BMP-collected BGP RIBs (Routing Information Base) can be cached in a redundant fashion at the data collection layer, for the purpose of enriching Flow Aggregation [6], while saving persistently only the master copy at the database layer in order to avoid data duplication. Then you will implement your ideas in C and test them in a lab setup.

Experts from Swisscom, INSA [11] and Pmacct [12] will support you with a test environment and IETF level expertise in Network Telemetry data-collection, Linux network kernel and C development. You will be working in a well-supported group. Finally, you can present your thesis results at the IETF 115 GROW working group between November 5-11th 2022 to other network operators, vendors and universities.

https://nsg.ee.ethz.ch/fileadmin/user_upload/theses/2021-thesisproposal-bmp-high-availability.pdf

👍 **Research and document how BMP-collected BGP RIBs (Routing Information Base) can be cached in a redundant fashion at the data collection layer for the purpose of enriching Flow Aggregation.**
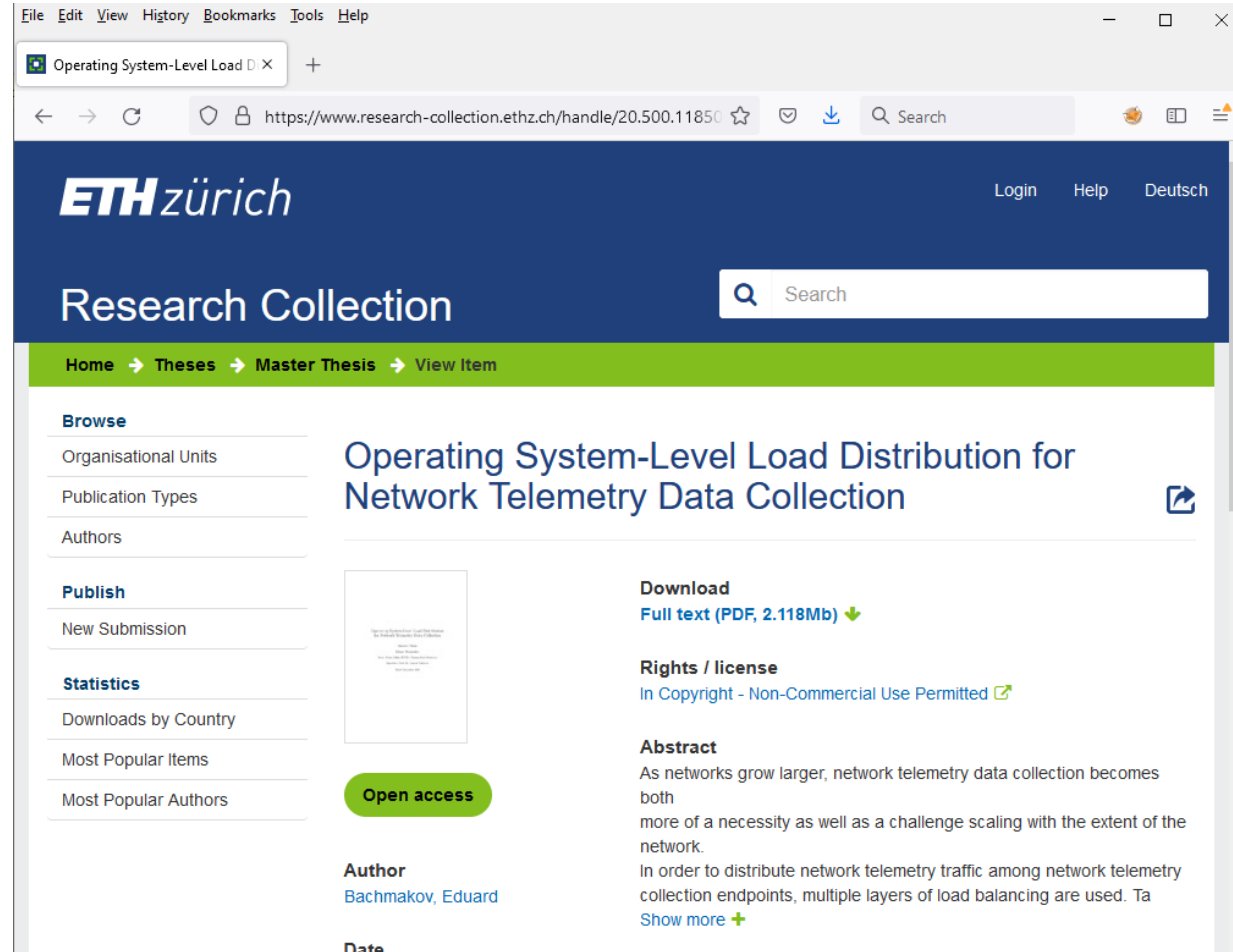
✓ **Develop and Test running code in C** and publish to the open-source and present to the IETF community.

# ETH Zürich, Eduard Bachmakov – Master Thesis

Operating System-Level Load Distribution for Network Telemetry Data Collection



👍 **From network data collection load distribution** with Anycast and ECMP on the network to SO_REUSEPORT with in the Linux network kernel.

👍 **Describes current load distribution challenges and extends** SO_REUSEPORT with cutome eBPF code.

✓ **Running code** on github at https://github.com/insa-unyte/ebpf-loadbalancer

https://www.research-collection.ethz.ch/handle/20.500.11850/507440

**Contact information**

Swisscom
Daisy Network Analytics


Thomas Graf
Binzring 17
8045 Zürich

Email thomas.graf@swisscom.com



Marco Tollini
Binzring 17
8045 Zürich

Email marco.tollini1@swisscom.com