# Monitoring and visualization of InfiniBand Fabrics

**InfiniBand Radar** – *InfiniBand Monitoring Tool*

Carsten Patzke

Spectrum Scale Strategy Days, Ehningen, March 2019

# About DESY

- Research institution

- ~2300 employees

- Over 3000 guest scientists yearly

- Research topics

  - Accelerator development

  - Photon science

  - Particle physics

  - Astroparticle physics

- 2 Sites
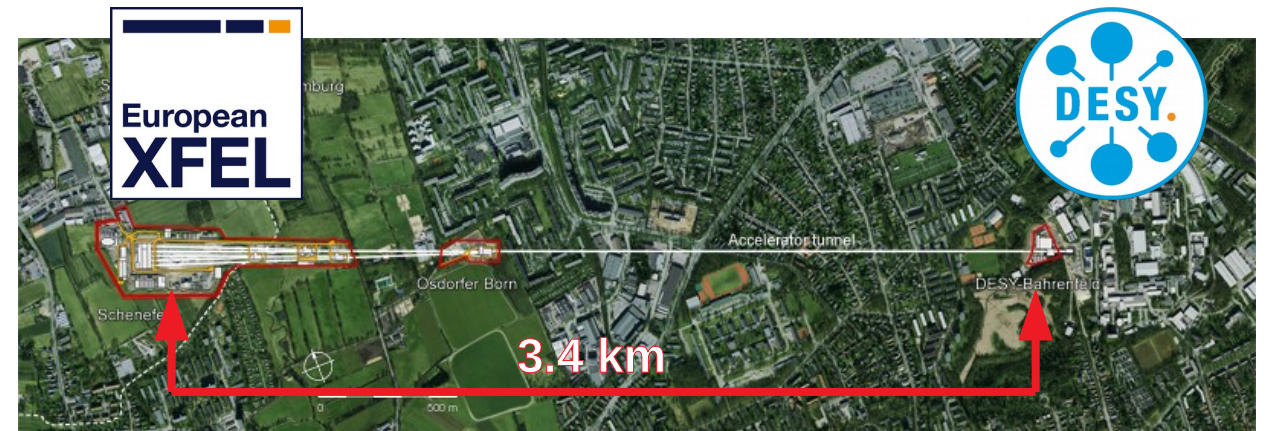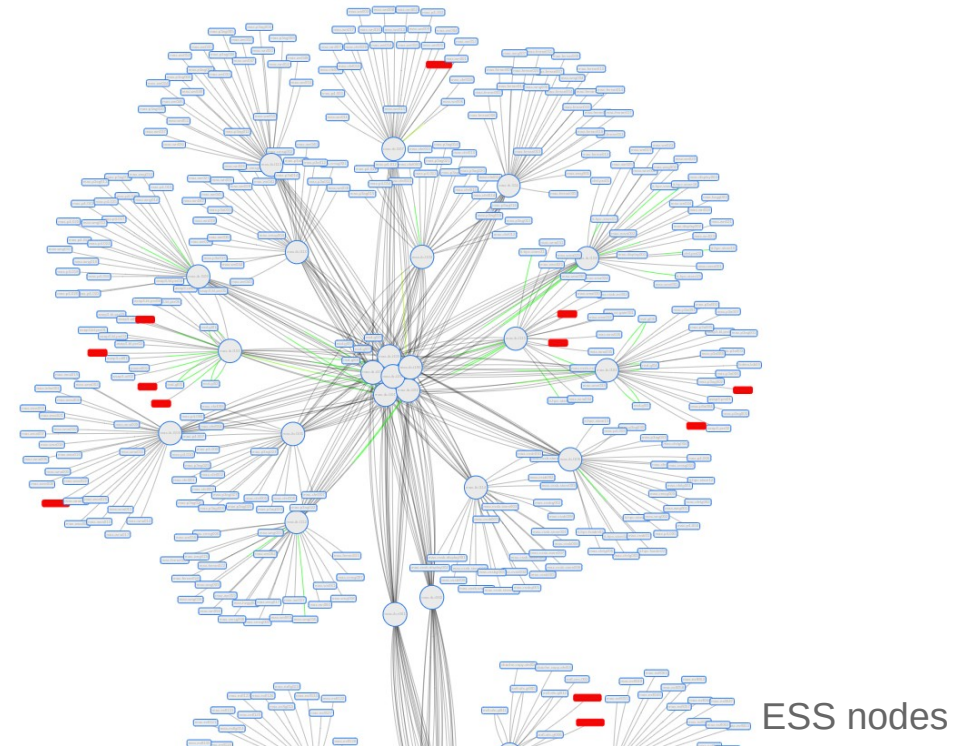
  - Hamburg

  - Zeuthen


Hamburg


Zeuthen

# Our GPFS/InfiniBand environment

- Using the GPFS since 2014

- Over 22 PB total storage capacity

- Only ESS Building blocks used (32 for 8 clusters)

- Metadata stored on SSDs

- GPFS is only available through InfiniBand

- Connected to over 900 individual server

- Some servers have access to two fabrics at once

- long-haul link to XFEL (MetroX)

- Deal with other traffic (MPI, BeeGFS)

ESS nodes

3.4 km

# Existing monitoring tools
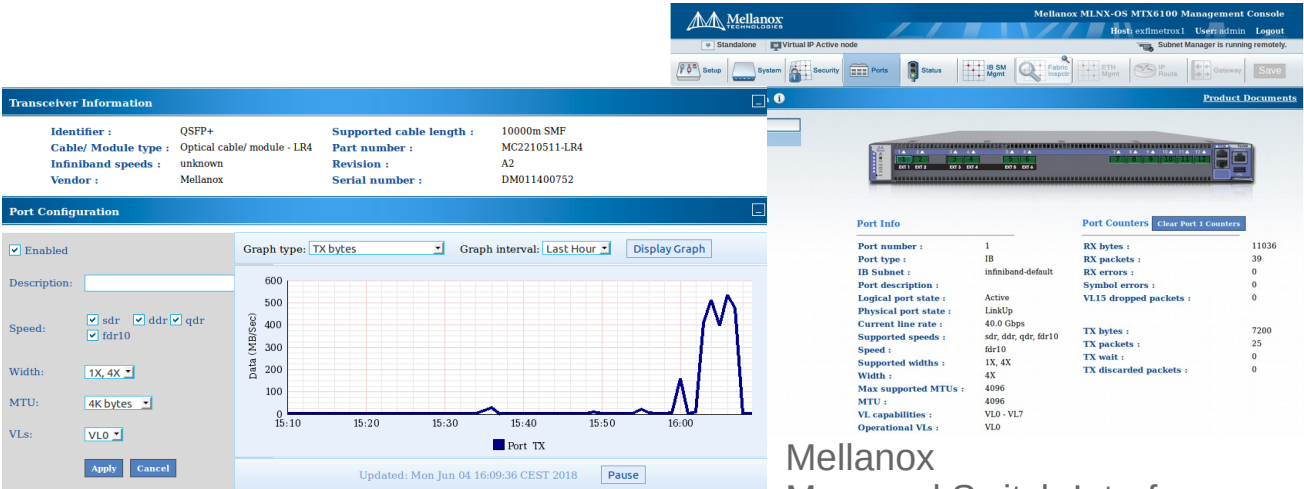
**Mellanox: Managed Switch Interface**

- Tracks a single system
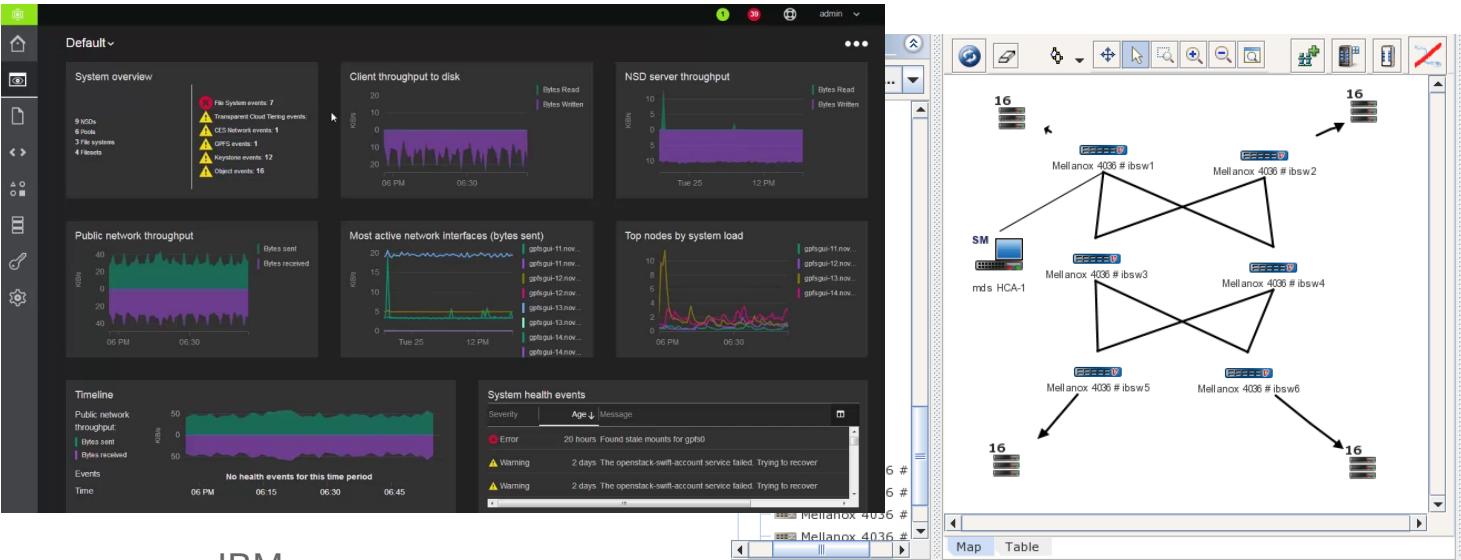
- For proprietary hardware

**Mellanox: UFM**

- Fabric wide monitoring

- Automatic fabric congestion detection

**IBM: Spectrum Scale GUI**

- Detailed information and management of Spectrum Scale clusters

- Only for Spectrum Scale
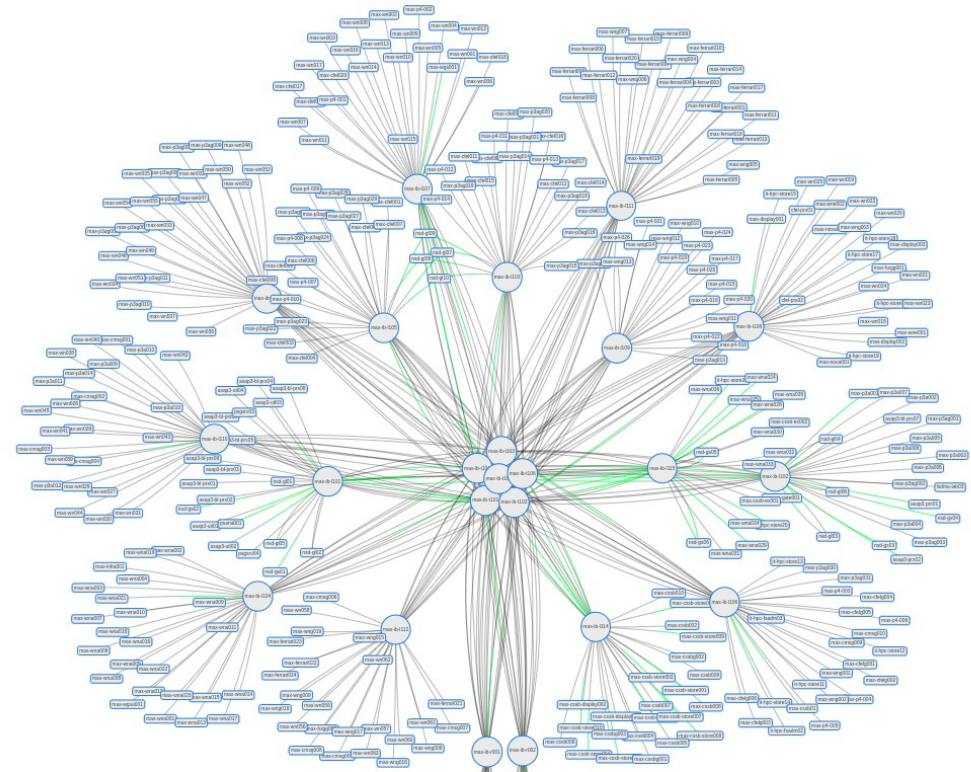


Mellanox
Managed Switch Interface



IBM
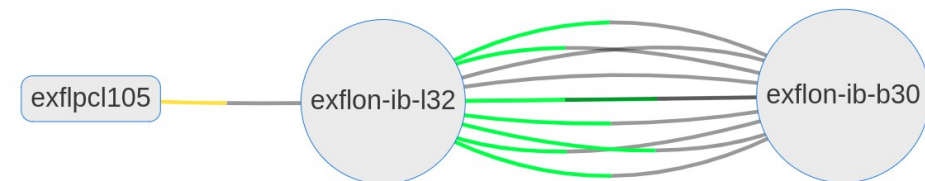Spectrum Scale GUI



Mellanox
Unified Fabric Manager

# InfiniBand Radar

## Features

- No proprietary hardware/software required

- Supports multiple fabrics at once

- Automatically detects topologies

- Web-based user interface

- Visualization via interactive map

- Traffic flow indicators

- Topology change detection



Fabric visualization



Traffic flow. Data send from left to right.
Green = low-, Yellow = medium-, Red = high-load

# InfiniBand Radar

## Features

- Diagrams of network utilization
  (Backed by a TSDB for history data)

- Detailed port information
  (Link speed, peer and CA-Name)

- Node search bar
  (Hostname, GUID or link speed)

- Search by tags
  (SM state or empty ports)



Time range picker



Network utilization



Search



Link selection

SM    = Subnet Manager
TSDB = Time Series Database

# Demo

# Fabric selection

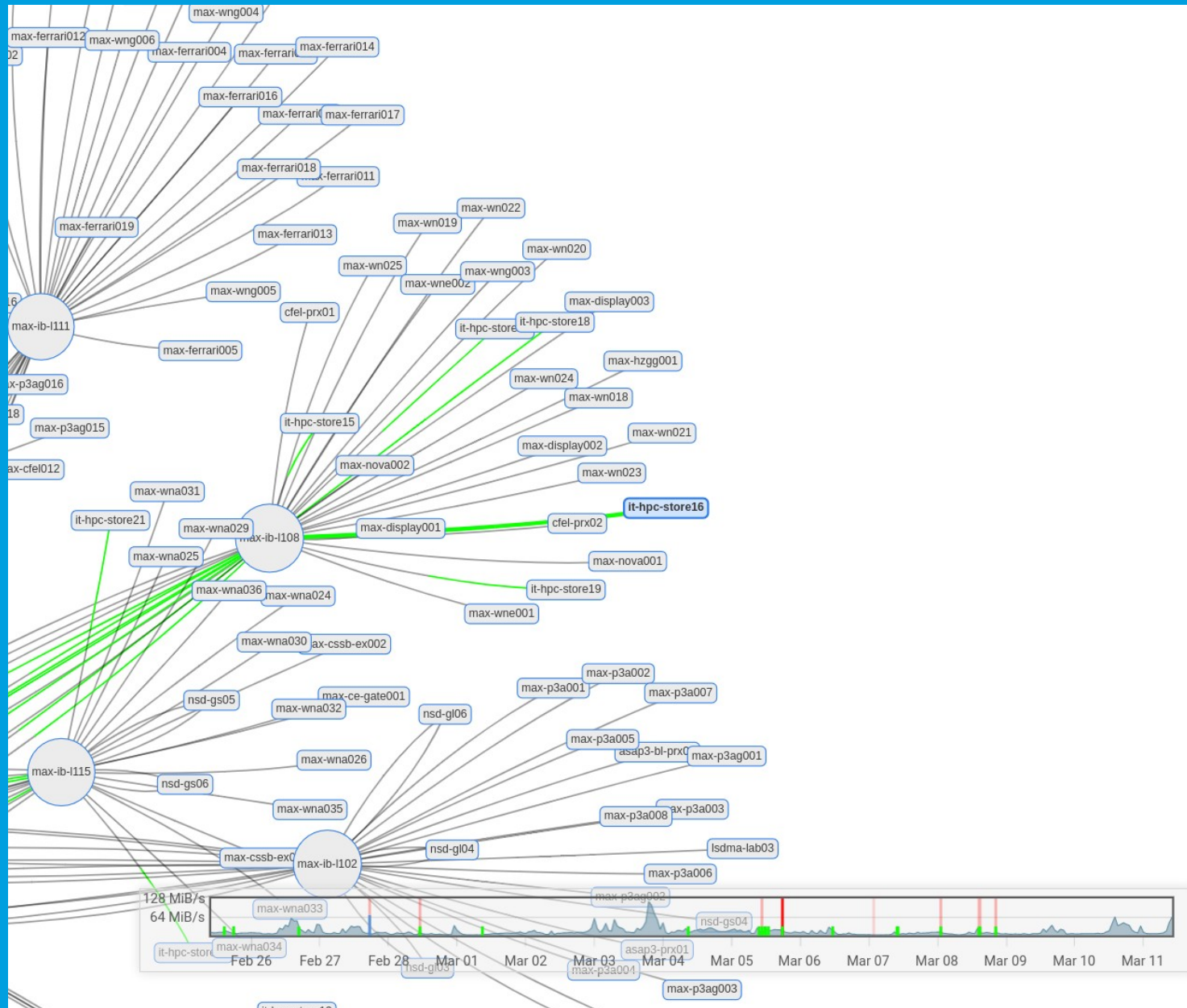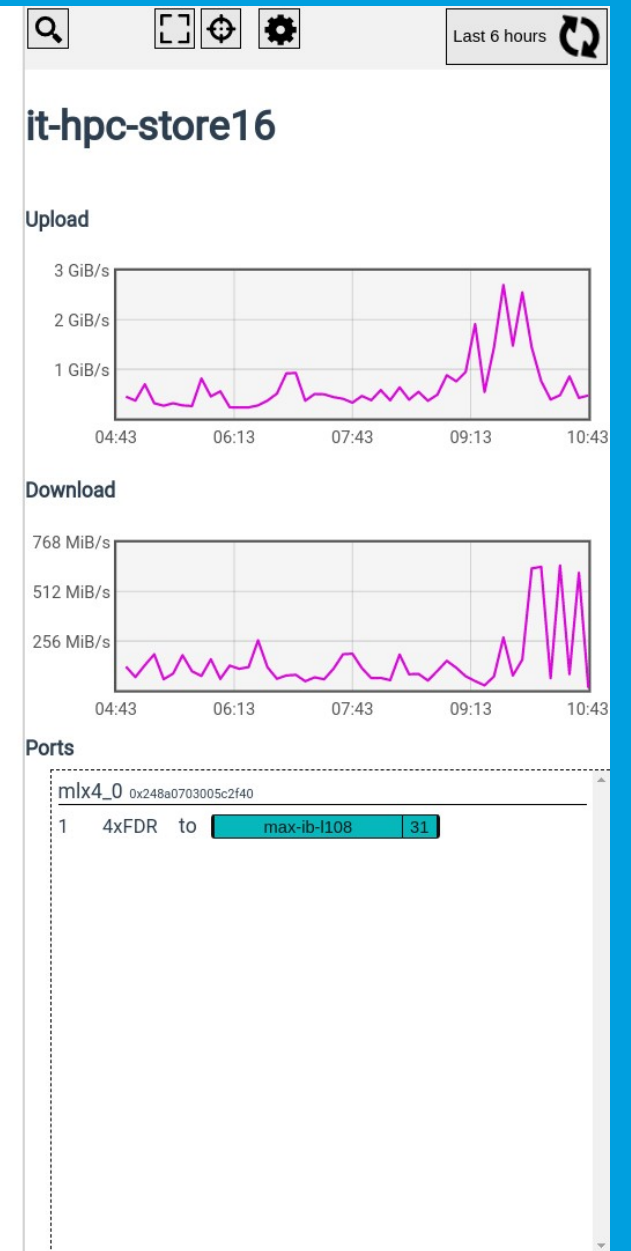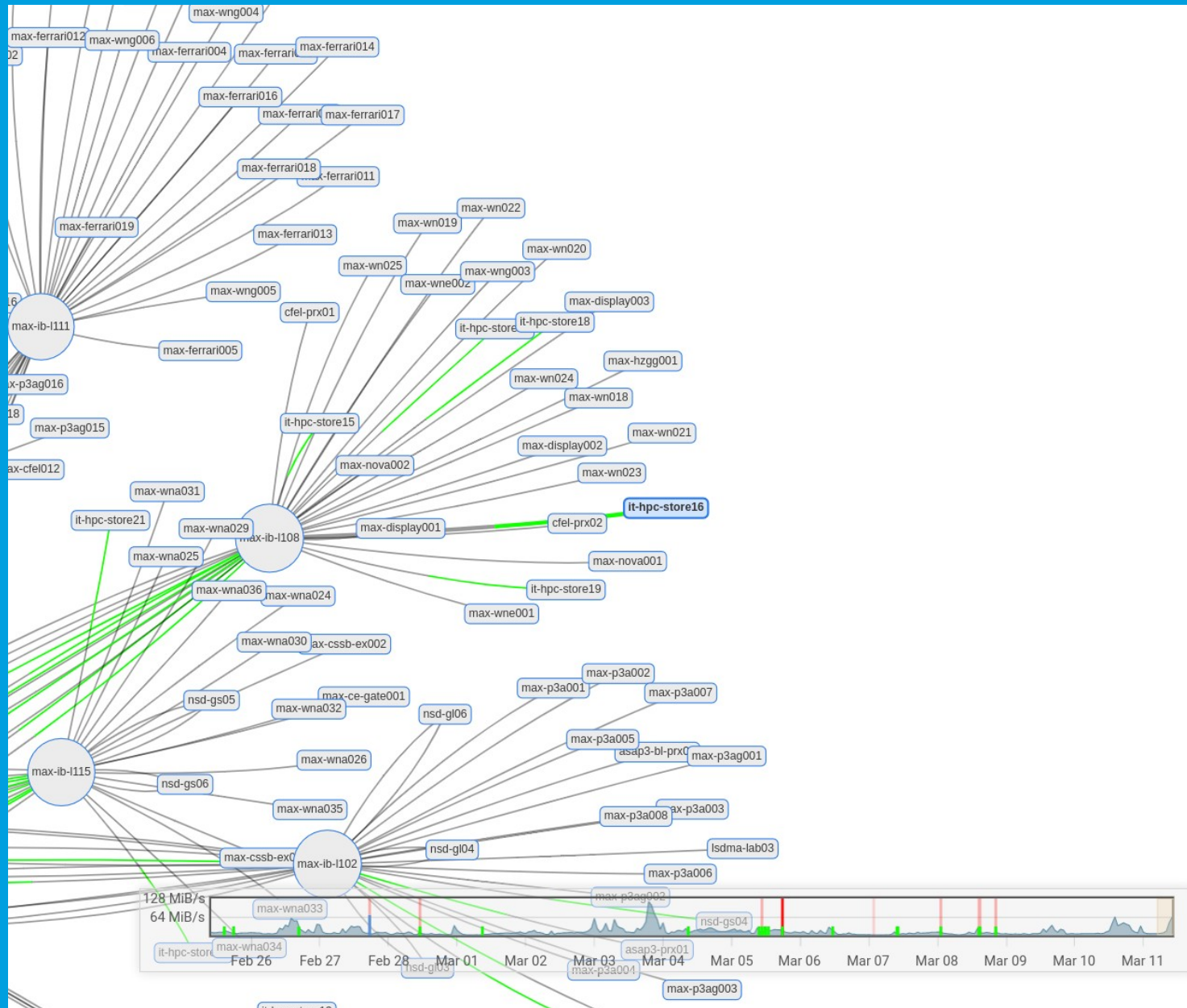| | |
|---|---|
|  | Maxwell |
|  | EU XFEL |
|  | PETRA III |
|  | DUST |

# Fabric view

Last 5 minutes

Hostname, GUID, SM State

Fold all | Hosts: 607

+ asap3-bl-prx01
+ asap3-bl-prx02
+ asap3-bl-prx03
+ asap3-bl-prx04
+ asap3-bl-prx05
+ asap3-bl-prx06
+ asap3-bl-prx07
+ asap3-bl-prx08
+ asap3-bl-prx09
+ asap3-prx01
+ asap3-prx02
+ asap3-utl01
+ asap3-utl02
+ asap3-utl03
+ M asap3-utl04
+ cfel-prx01
+ cfel-prx02
+ dcache-copy-xfel01
+ dcache-copy-xfel02
+ dcache-copy-xfel03
+ dcache-copy-xfel04
+ S exfl-ces-001
+ exfl-ces-002
+ exfl-ofs-gl001
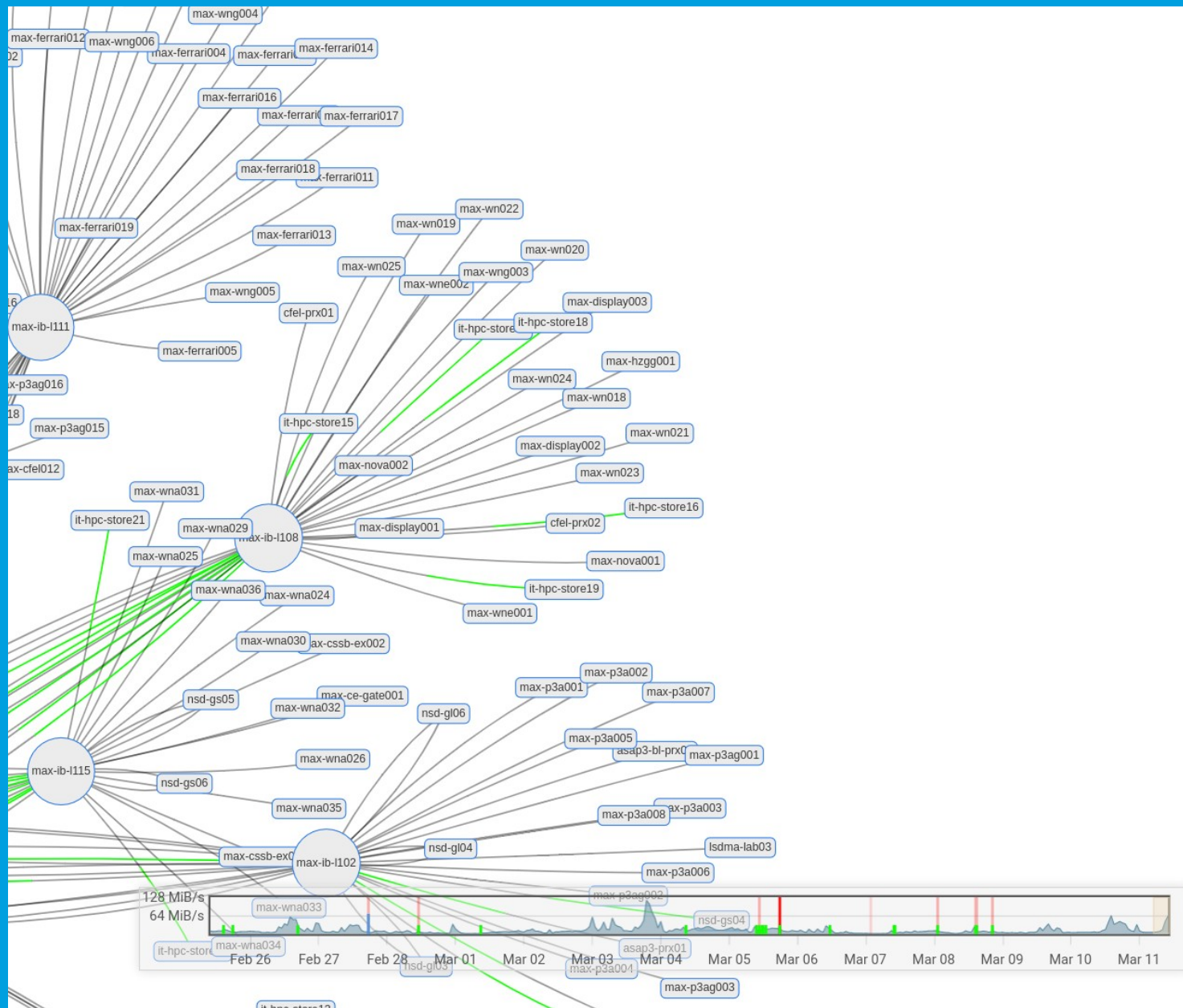+ exfl-ofs-gl002
+ exfl-ofs-gl003
+ exfl-ofs-gl004

128 MiB/s
64 MiB/s

Feb 26  Feb 27  Feb 28  Mar 01  Mar 02  Mar 03  Mar 04  Mar 05  Mar 06  Mar 07  Mar 08  Mar 09  Mar 10  Mar 11

# Utilization graph

# Time range selection

# Utilization graph

# Search



Fabric List

store

Fold all | Hosts: 21

+ it-hpc-store11
+ it-hpc-store12
+ it-hpc-store13
+ it-hpc-store14
+ it-hpc-store15
+ it-hpc-store16
+ it-hpc-store17
+ it-hpc-store18
+ it-hpc-store19
+ it-hpc-store20
+ it-hpc-store21
+ max-cssb-store001
+ max-cssb-store002
+ max-cssb-store003
+ max-cssb-store004
+ max-cssb-store005
+ max-cssb-store006
+ max-cssb-store007
+ max-cssb-store008
+ max-cssb-store009
+ max-cssb-store010

Last 6 hours

**Fabric List**

Last 6 hours

store

Fold all | Hosts: 21

+ it-hpc-store11
+ it-hpc-store12
+ it-hpc-store13
+ it-hpc-store14
+ it-hpc-store15
− it-hpc-store16

mlx4_0 0x248a0703005c2f40

1    4xFDR    to    max-ib-l108    31

+ it-hpc-store17
+ it-hpc-store18
+ it-hpc-store19
+ it-hpc-store20
+ it-hpc-store21
+ max-cssb-store001
+ max-cssb-store002
+ max-cssb-store003
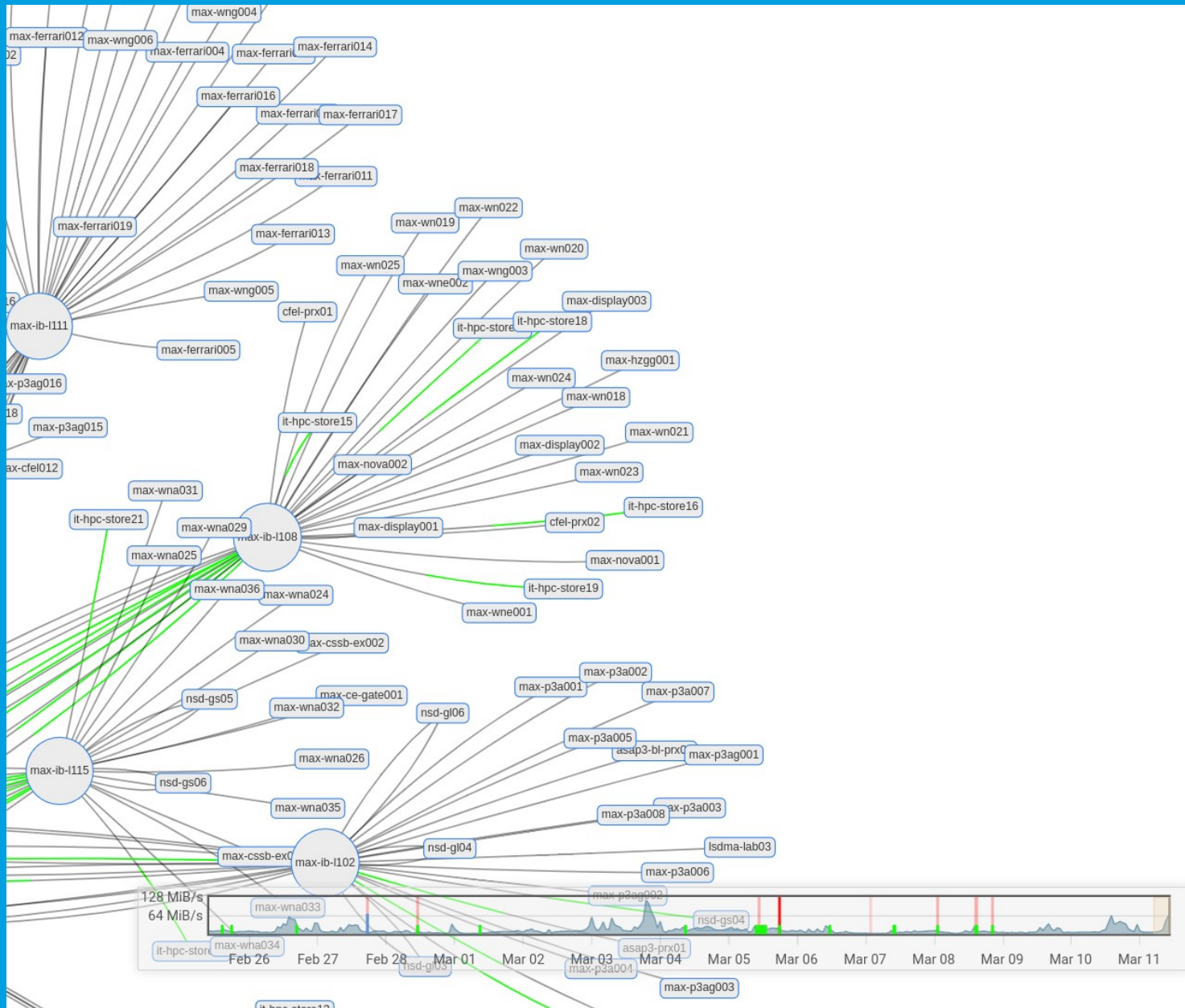+ max-cssb-store004
+ max-cssb-store005
+ max-cssb-store006
+ max-cssb-store007
+ max-cssb-store008
+ max-cssb-store009
+ max-cssb-store010

# Search (Tags)

| Fabric List | | | | Last 6 hours |

**not connected**

Fold all  Hosts: 23

—**max-ib-l101**

(Switch) 0x7cfe90030095a7f0

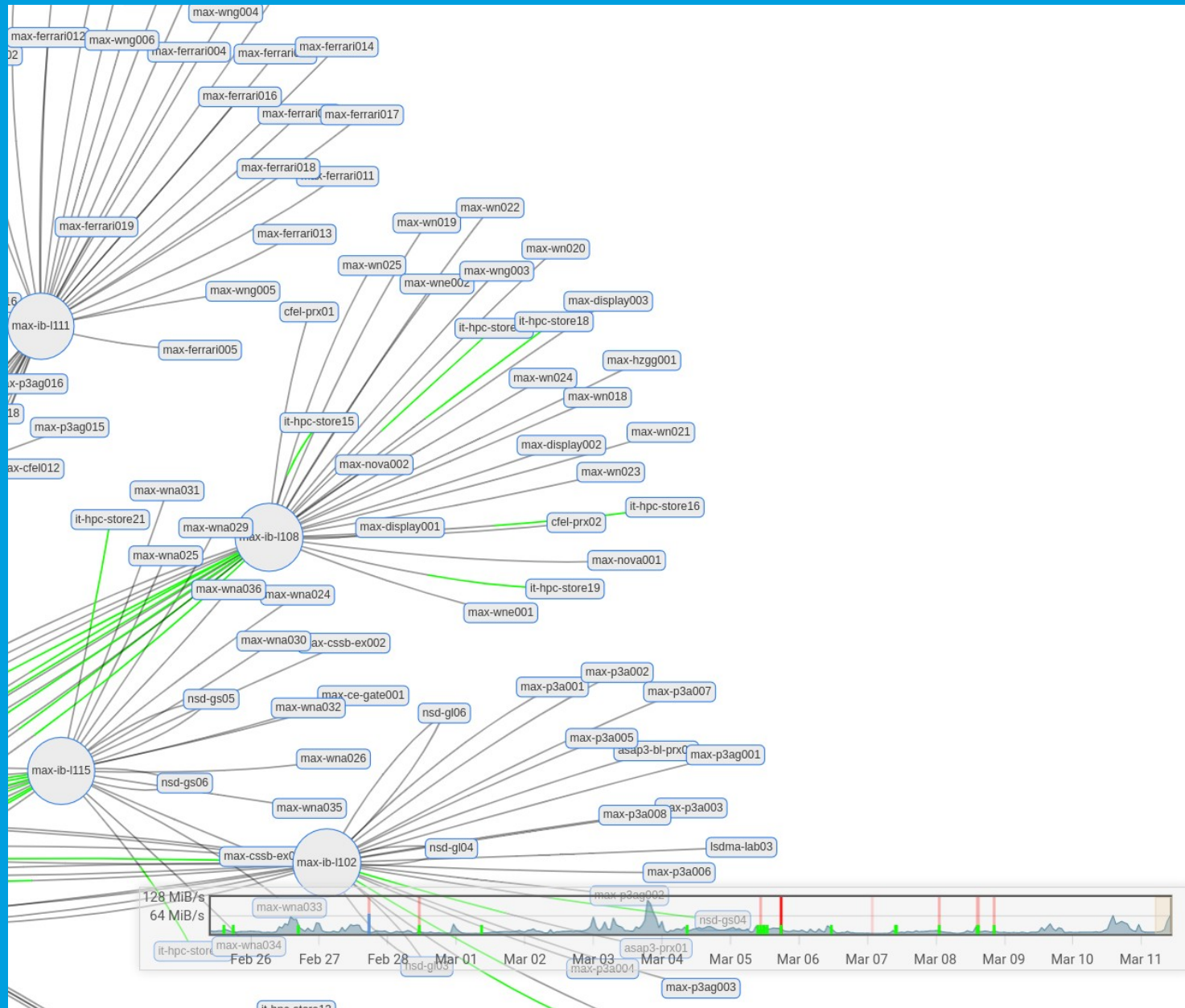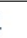| 1 | 4xFDR | to | nsd-gs02 | 1 |
| 2 | 4xFDR | to | nsd-gs01 | 1 |
| 3 | 4xFDR | to | asap3-bl-prx04 | 1 |
| 4 | 4xFDR | to | asap3-utl02 | 1 |
| 5 | 4xFDR | to | asap3-bl-prx05 | 1 |
| 6 | 4xFDR | to | psana001 | 1 |
| 7 | 4xFDR | to | asap3-utl01 | 1 |
| 8 | Not connected | | | |
| 9 | 4xFDR | to | asap3-bl-prx06 | 1 |
| 10 | 4xFDR | to | nsd-gl05 | HCA-1 | 1 |
| 11 | 4xFDR | to | nsd-gl05 | HCA-2 | 1 |
| 12 | 4xFDR | to | asap3-bl-prx09 | 1 |
| 13 | 4xFDR | to | asap3-bl-prx08 | 1 |
| 14 | 4xFDR | to | asap3-utl04 | 1 |
| 15 | 4xFDR | to | psgsrv04 | 1 |
| 16 | 4xFDR | to | asap3-utl03 | 1 |
| 17 | 4xFDR | to | psgsrv03 | 1 |
| 18 | 4xFDR | to | asap3-bl-prx03 | 1 |
| 19 | 4xFDR | to | asap3-bl-prx01 | 1 |
| 20 | 4xFDR | to | asap3-bl-prx02 | 1 |
| 21 | 4xFDR | to | nsd-gl02 | HCA-1 | 1 |
| 22 | 4xFDR | to | nsd-gl02 | HCA-2 | 1 |
| 23 | 4xFDR | to | nsd-gl01 | HCA-1 | 1 |
| 24 | 4xFDR | to | nsd-gl01 | HCA-2 | 1 |
| 25 | 4xFDR | to | max-ib-t103 | 18 |
| 26 | 4xFDR | to | max-ib-t103 | 22 |
| 27 | 4xFDR | to | max-ib-t104 | 17 |
| 28 | 4xFDR | to | max-ib-t104 | 14 |
| 29 | 4xFDR | to | max-ib-t101 | 29 |
| 30 | 4xFDR | to | max-ib-t101 | 30 |
| 31 | 4xFDR | to | max-ib-t102 | 29 |
| 32 | 4xFDR | to | max-ib-t102 | 30 |
| 33 | 4xFDR | to | max-ib-t105 | 13 |
| 34 | 4xFDR | to | max-ib-t105 | 22 |

# Search (Tags)

# Change detection

# Change detection (details)

max-wng004
max-ferrari012 max-wng006
max-ferrari004 max-ferrari
max-ferrari014
02
max-ferrari016

Fabric List

Last 6 hours

Hostname, GUID, SM State

Hosts: 607

max-ferrari019

max-ib-l111

Left 27.02.2019 17:07 ▼   Right 05.03.2019 16:34 ▼    Set *right* version as new *default*    X

—Hosts 0/0/0

Both versions are the same

—Connections 2/0/2

| Port A | Port B | Link |
|---|---|---|
| exfl-ofs-gl014/HCA-2/1 | max-ib-l209/(Switch)/16 | 4xFDR |
| exfl-ofs-gl014/HCA-4/1 | max-ib-l210/(Switch)/3 | 4xFDR |
| exfl-ofs-gl014/HCA-1/1 | max-ib-l209/(Switch)/16 | 4xFDR |
| exfl-ofs-gl014/HCA-3/1 | max-ib-l210/(Switch)/3 | 4xFDR |

max-p3ag016
18
max-p3ag015
ax-cfel012
max-wna
it-hpc-store21
max-wna

max-ib-l115
nsd

128 M
64 MiB/s

it-hpc-stor
max-wha034
Feb 26    Feb 27    Feb 28    Mar 01    Mar 02    Mar 03    Mar 04    Mar 05    Mar 06    Mar 07    Mar 08    Mar 09    Mar 10    Mar 11
nsd-gs03
asap3-prx01
max-p3a004
max-p3ag003

nsd-gs04

+ exfl-ofs-gl002
+ exfl-ofs-gl003
+ exfl-ofs-gl004

# InfiniBand Radar

## Architecture

Fabric

Web interface

One server per Fabric

Topology

Bandwidth statistics
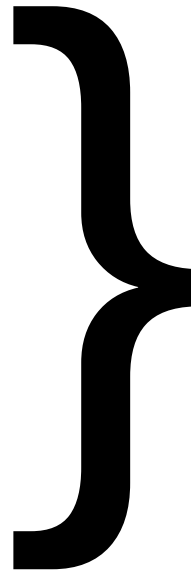
REST API

# InfiniBand Radar

**Analyze the fabric**

**ibnetdiscover**

- Collects fabric topology

**perfquery**

- Query performance counters

}

**Combined into a single executable**

- Sends data automatically
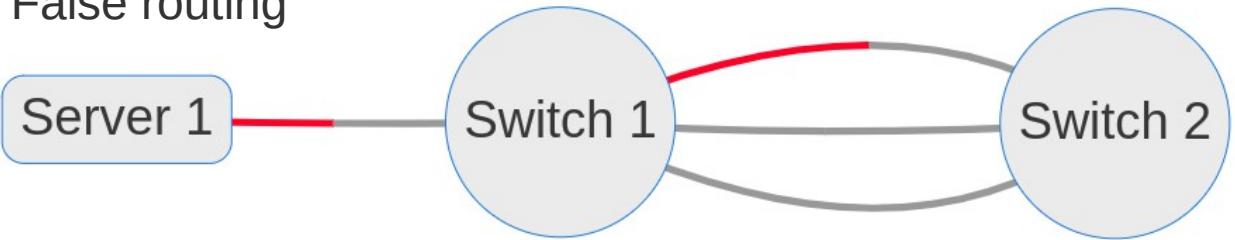
- In-house developed

- OFED libraries used

OFED = OpenFabrics Enterprise Distribution
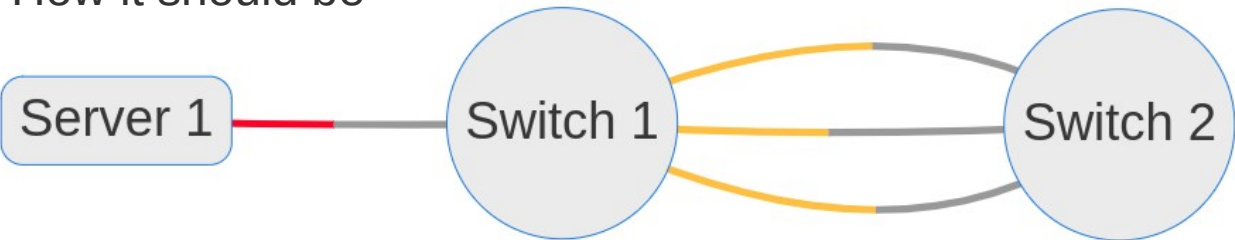
# Successful stories

## Balancing/Routing issue

- The traffic was not evenly distributed

- Caused by a miss configured SM

False routing



How it should be



## Invalid link state

- Lower link speed then expected
  (or the connection was not even established)

- Caused by broken cable / invalid handshake

−Connections 0/1/18

| Port A | Port B | Link |
|---|---|---|
| exflon-ib-l12/(Switch)/24 | exflfpcl01/mlx5_0/1 | 4xEDR=>4xFDR |
| exflon-ib-l15/(Switch)/21 | exfl-ons-gs106/HCA-3/1 | 4xEDR |
| exflon-ib-l15/(Switch)/23 | exfl-ons-gs106/HCA-1/1 | 4xEDR |

# InfiniBand Radar

## Easy installation

Complete source available:
https://github.com/infiniband-radar

### API Server

```
1 vim ./config/apiServer.json
2
3 docker-compose up -d
```

### Fabric Daemon

```
1 yum localinstall infiniband-radar-daemon.rpm
2
3 vim /etc/infiniband-radar/config.<FabricId>.json
4
5 systemctl enable infiniband-radar@<FabricId>
6 systemctl start infiniband-radar@<FabricId>
```

**Contact**

**DESY.** Deutsches
Elektronen-Synchrotron

www.desy.de

Carsten Patzke
carsten.patzke@desy.de
https://desy.de/~cpatzke