

The RAL Tier-1 Network

James Adams
RAL, STFC, UKRI
GridPP49, Abingdon

2023-03-30

Overview

- Background
- What's Changed?
- What's Next?

Background

- Building out a new network for the Tier-1 alongside the “legacy” network
 - Fully-routed eBGP ECMP architecture
 - Mellanox switches running Cumulus Linux
 - Joined to legacy network by SCD SuperSpine
- Started work July 2021
- Connected to SCD SuperSpine October 2021
- Connected to RAL site November 2021
- First worker nodes live by December 2021

What's Changed?

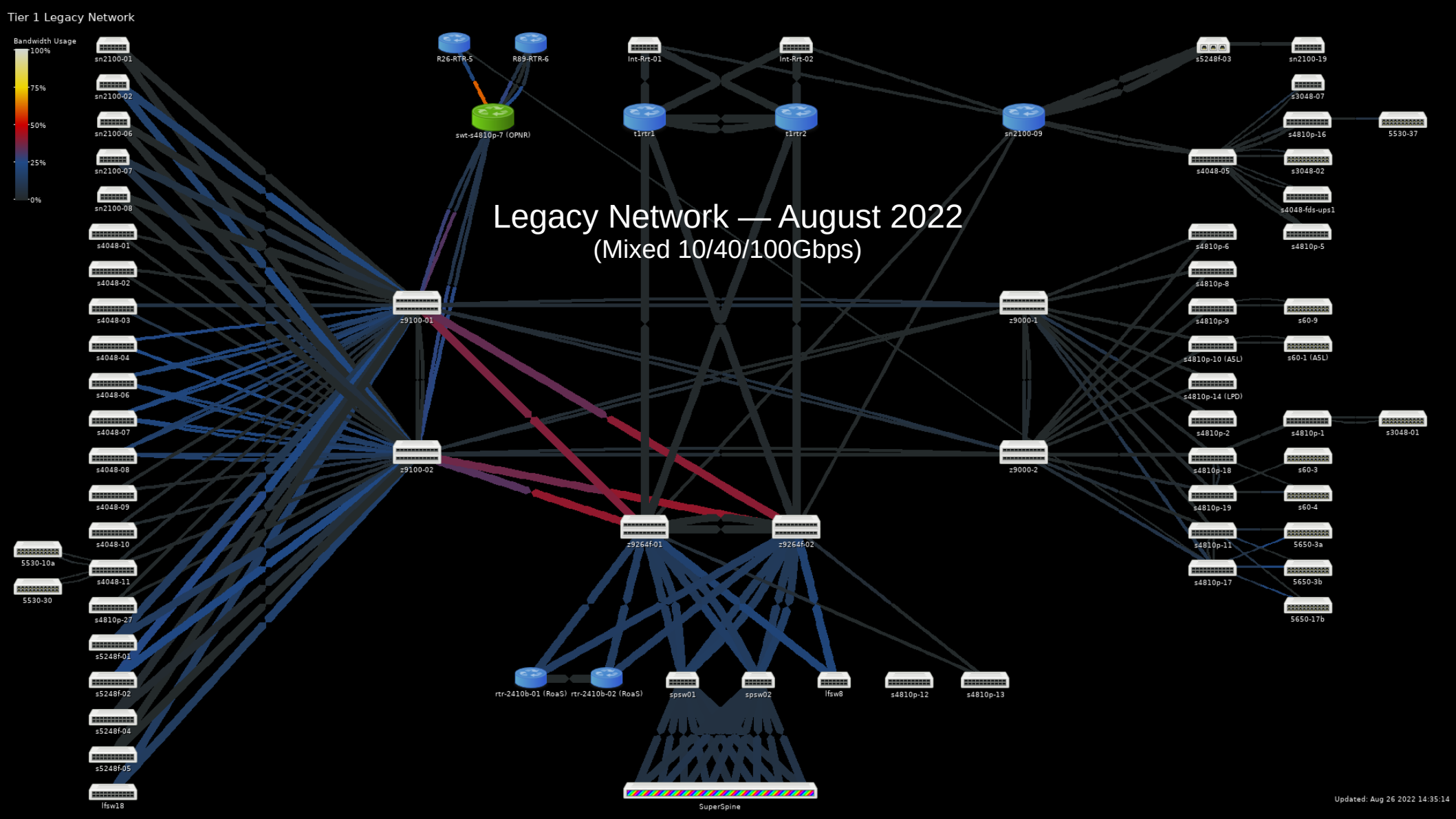
- Legacy network melt-down
- More hardware on new network
- Peering with LHCONE
- IPv6

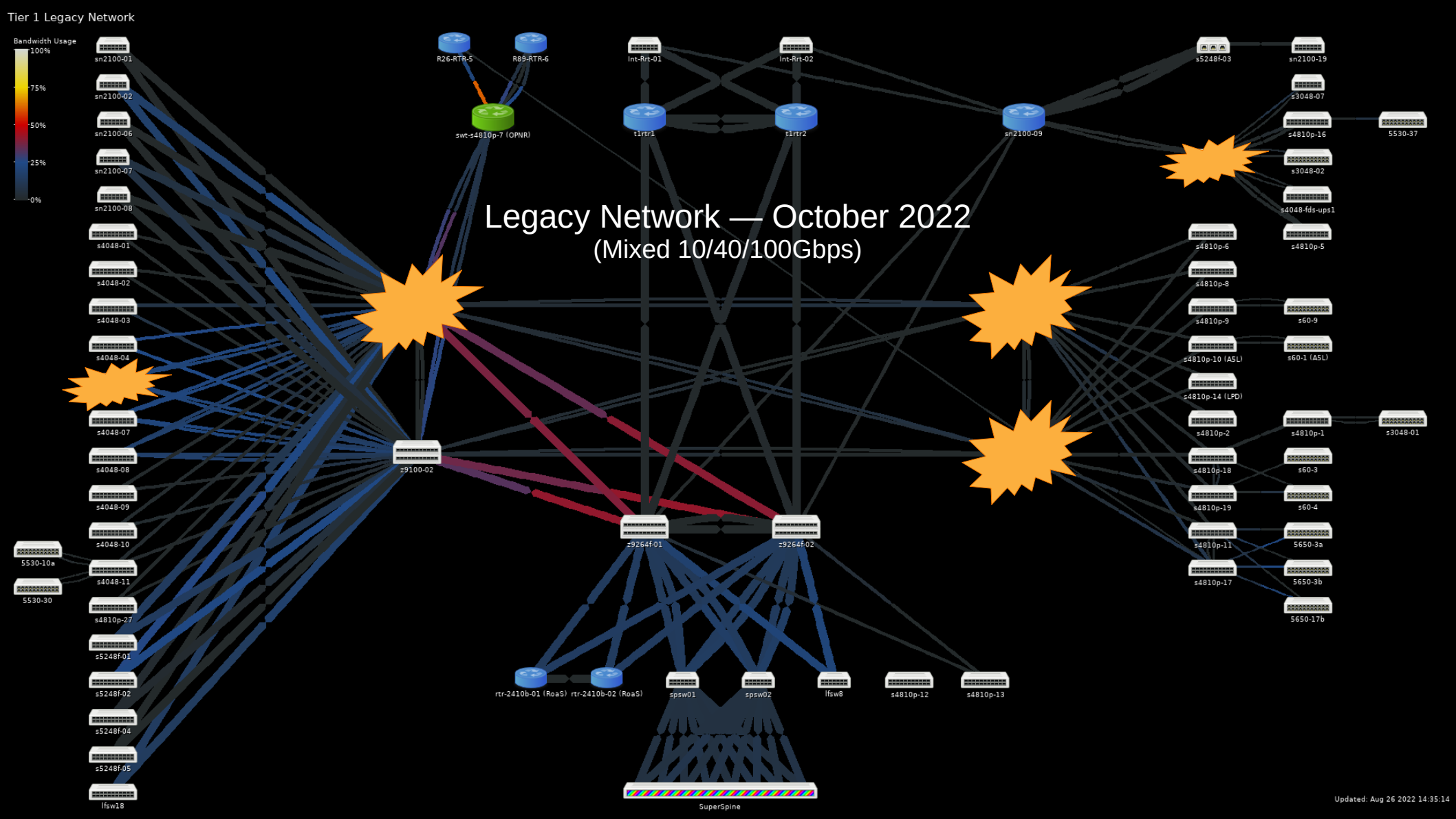
Legacy Network Meltdown

- Monday 17th October 2022 ~15:00 local time
 - One of the two core Z9100 failed
 - One of the Facilities database server switches suffered a hardware failure
 - Took out RT ticketing system, FTS and GOCDB
 - One of the echo storage switches crashed
 - Took out 25 storage nodes
 - Lots of spanning tree changes on management network
 - Rapid topology changes caused instability of MLAG & VRRP in core
 - Remaining core switches failed over to broadcast (i.e. hub) mode

Legacy Network Meltdown

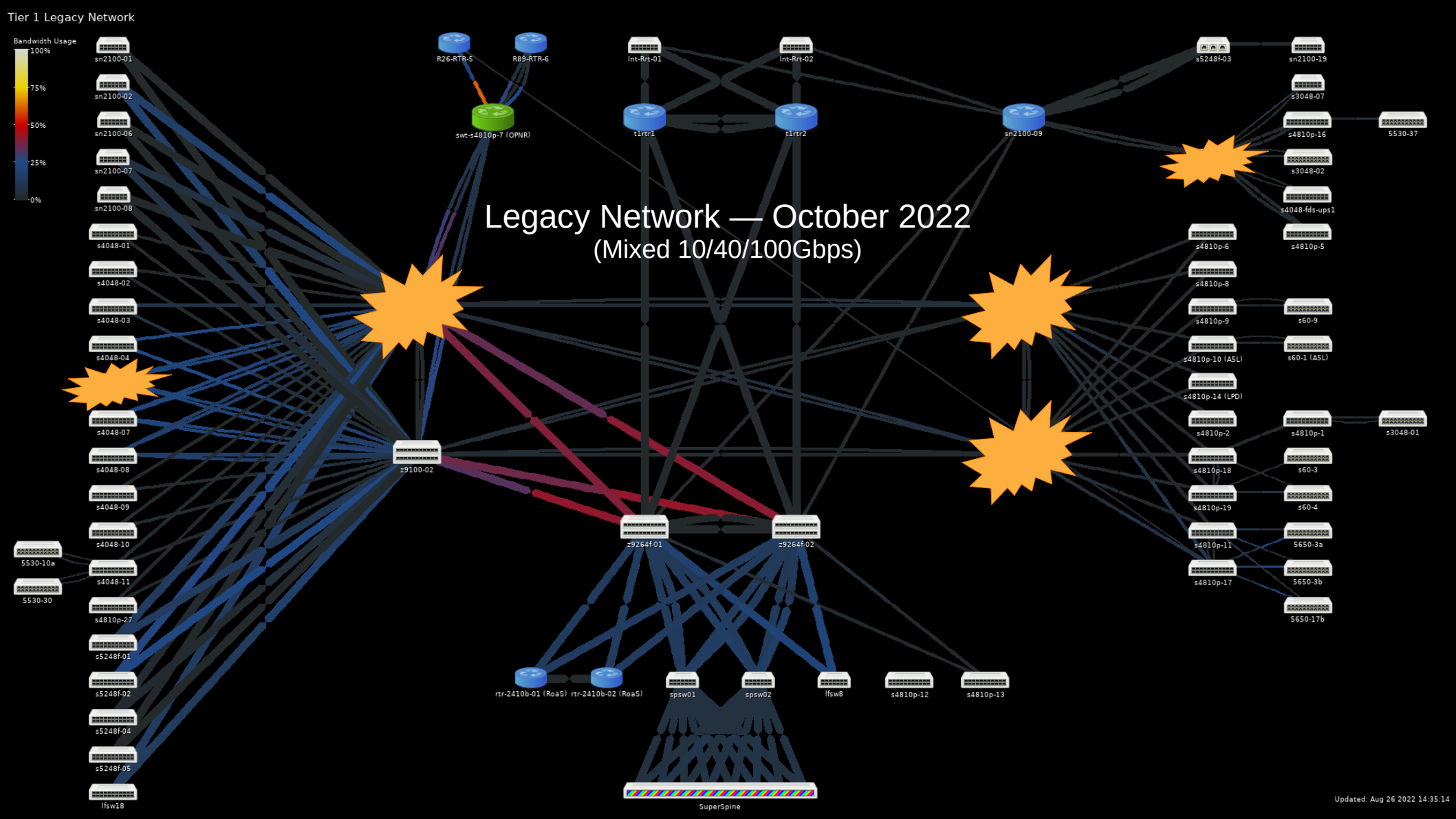
- Key staff away at HEPiX
- Additional faults found when trying to restart
 - Restoring VLT crashed switches immediately
 - Fan and power supply failed
- Z9000s & Z9100s out-of-maintenance
 - No spare hardware
 - Planning had already started for migrating off them

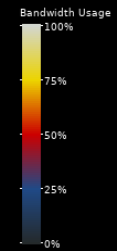




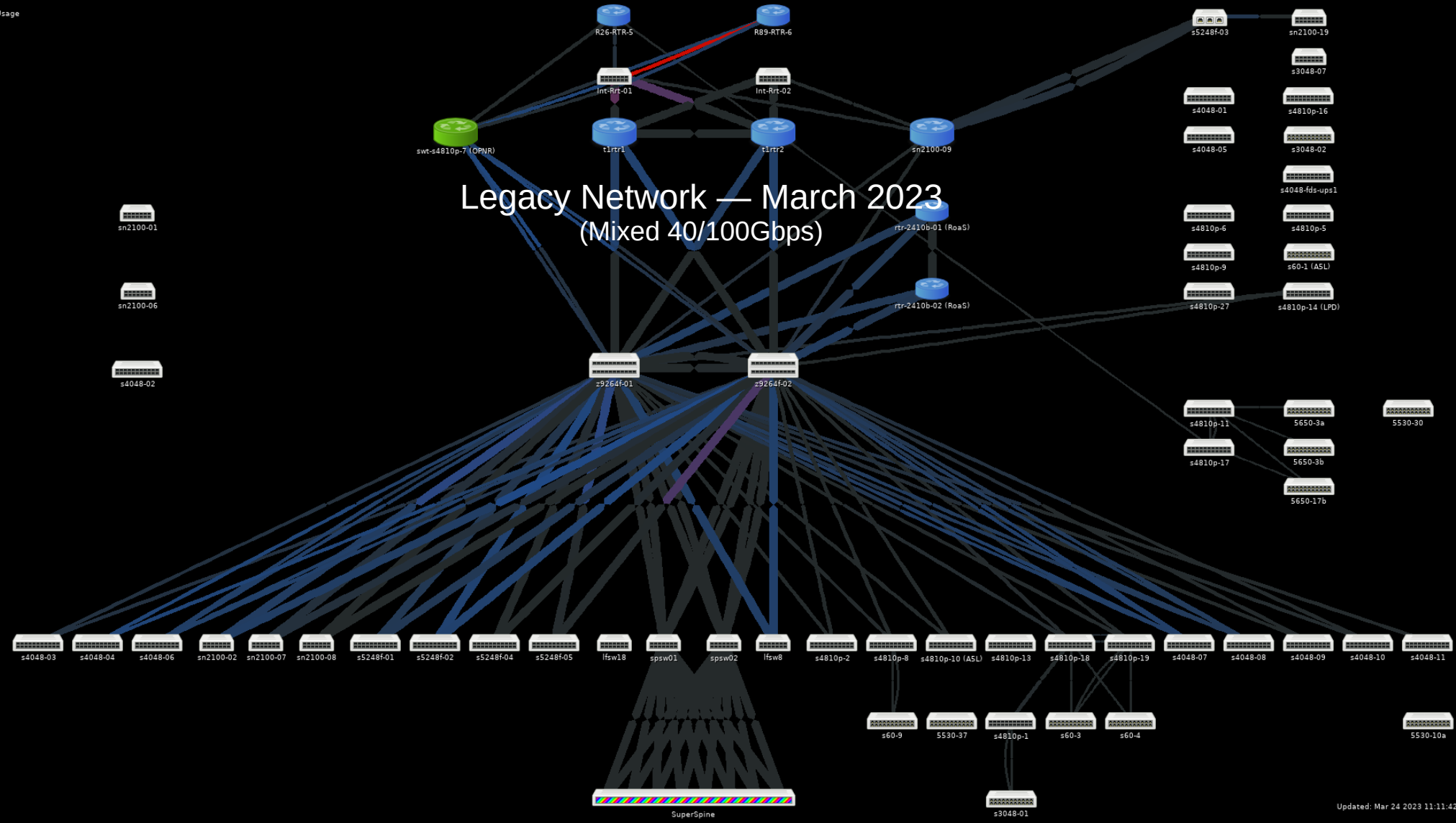
Legacy Network Meltdown

- Decided to make 12 months of planned changes in a single day
 - Moved 56 cables and config from four core switches
- Disabled management network entirely
 - Designed an interim replacement architecture
 - Replacement has been slow, but steady progress





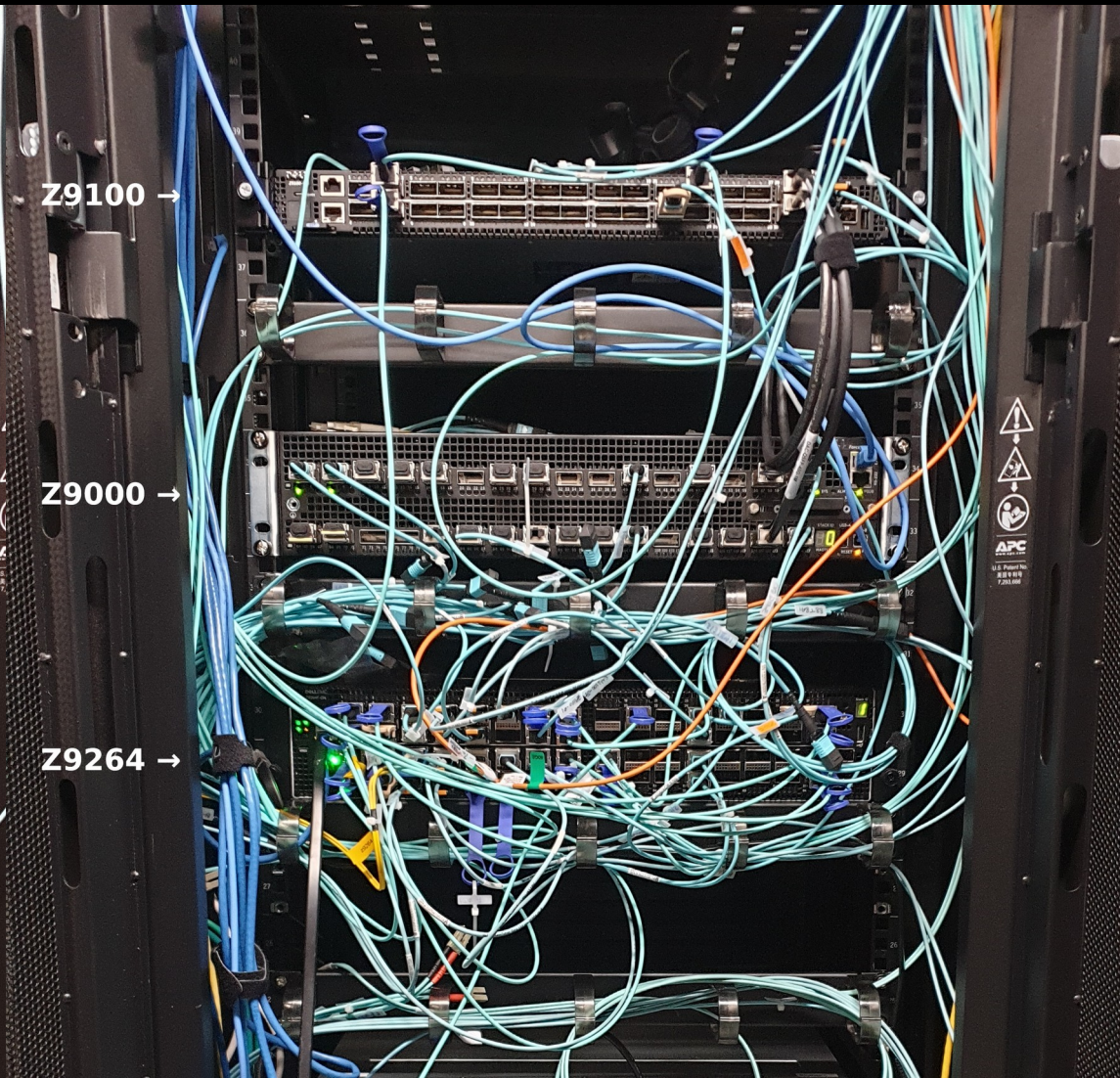
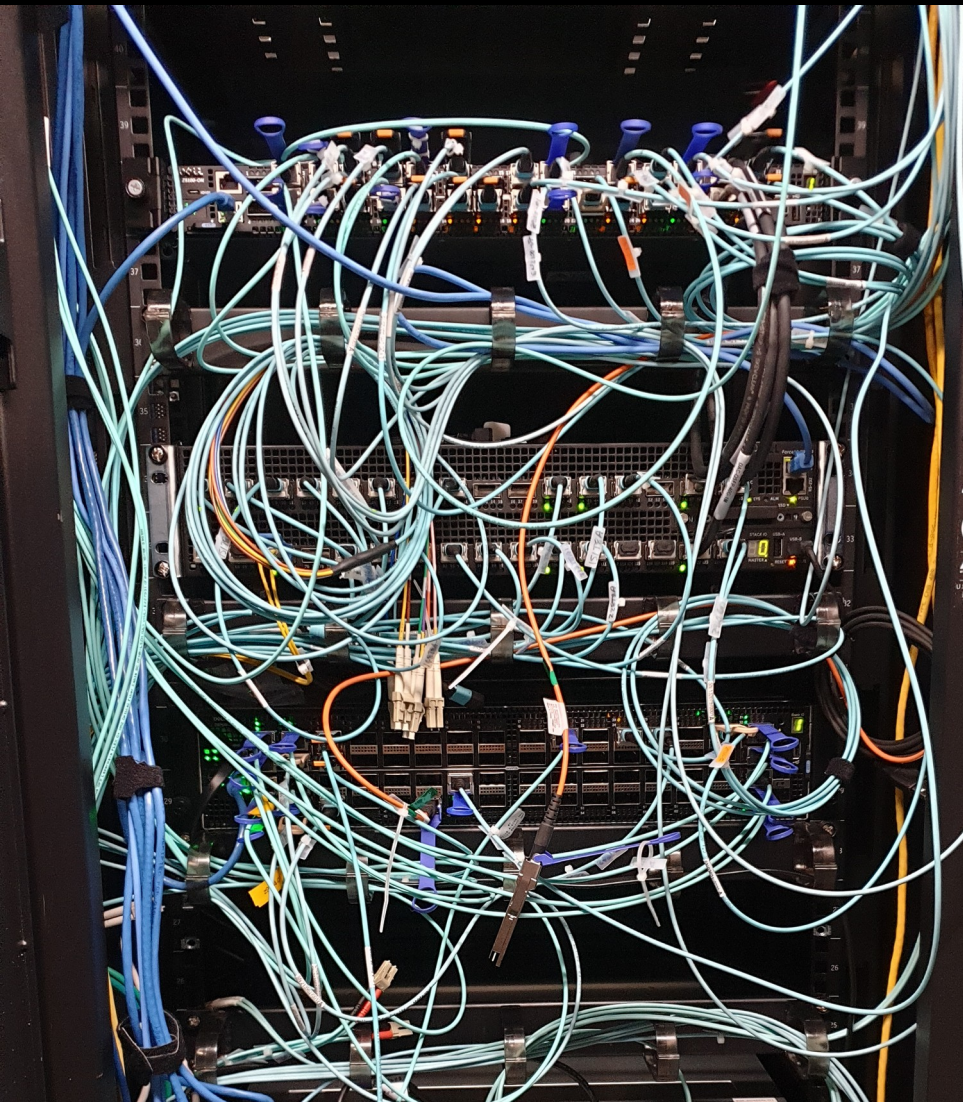
Legacy Network — March 2023 (Mixed 40/100Gbps)



Before

(all of this, twice)

After

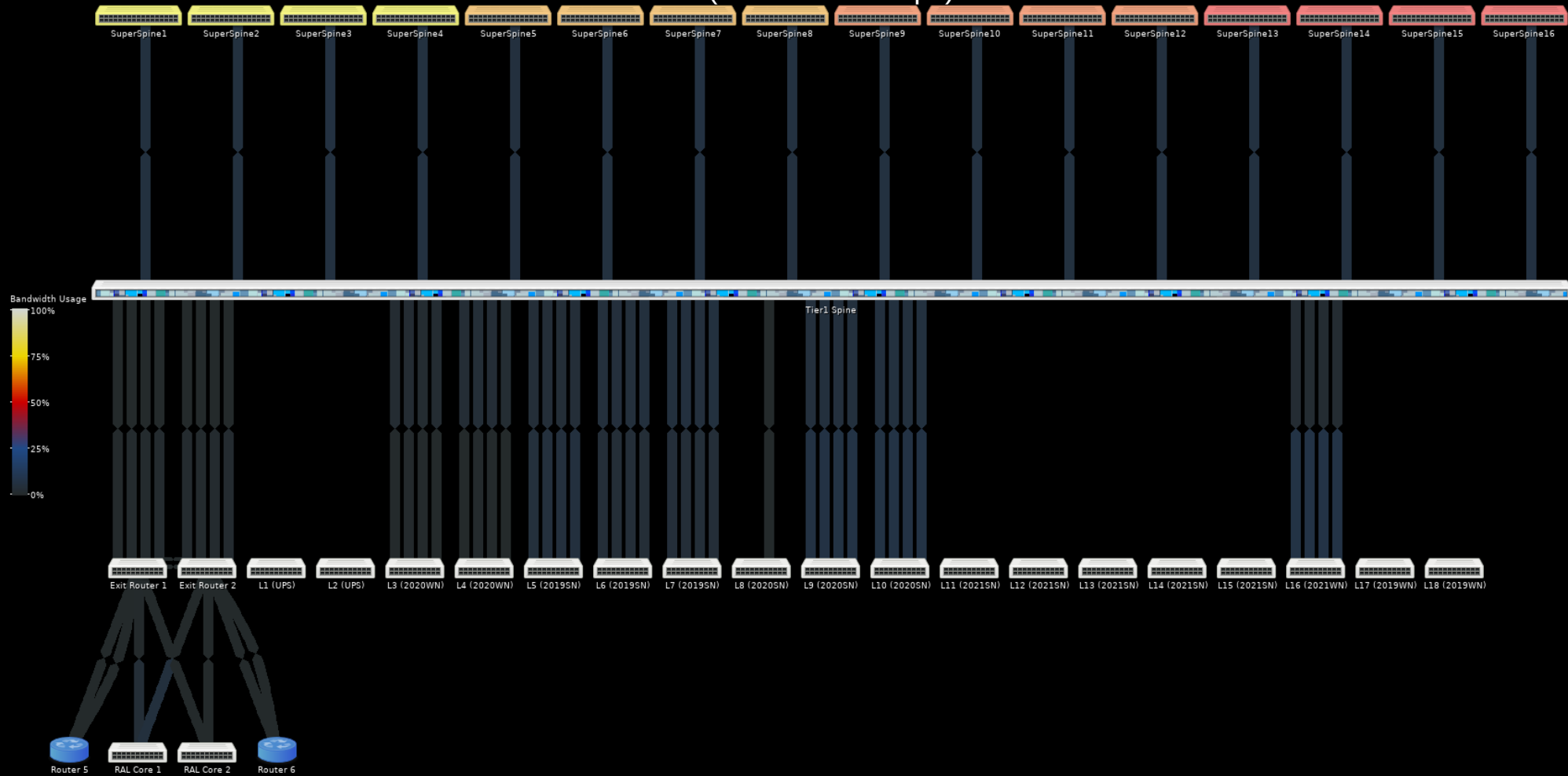


New Network

- Four leaves of workers (almost) **+1**
 - 104% of 2023 CPU pledge (146% soon)
- Eleven leaves of storage in production **+6**
 - 33% of Echo capacity
 - FY2022 nodes currently being deployed

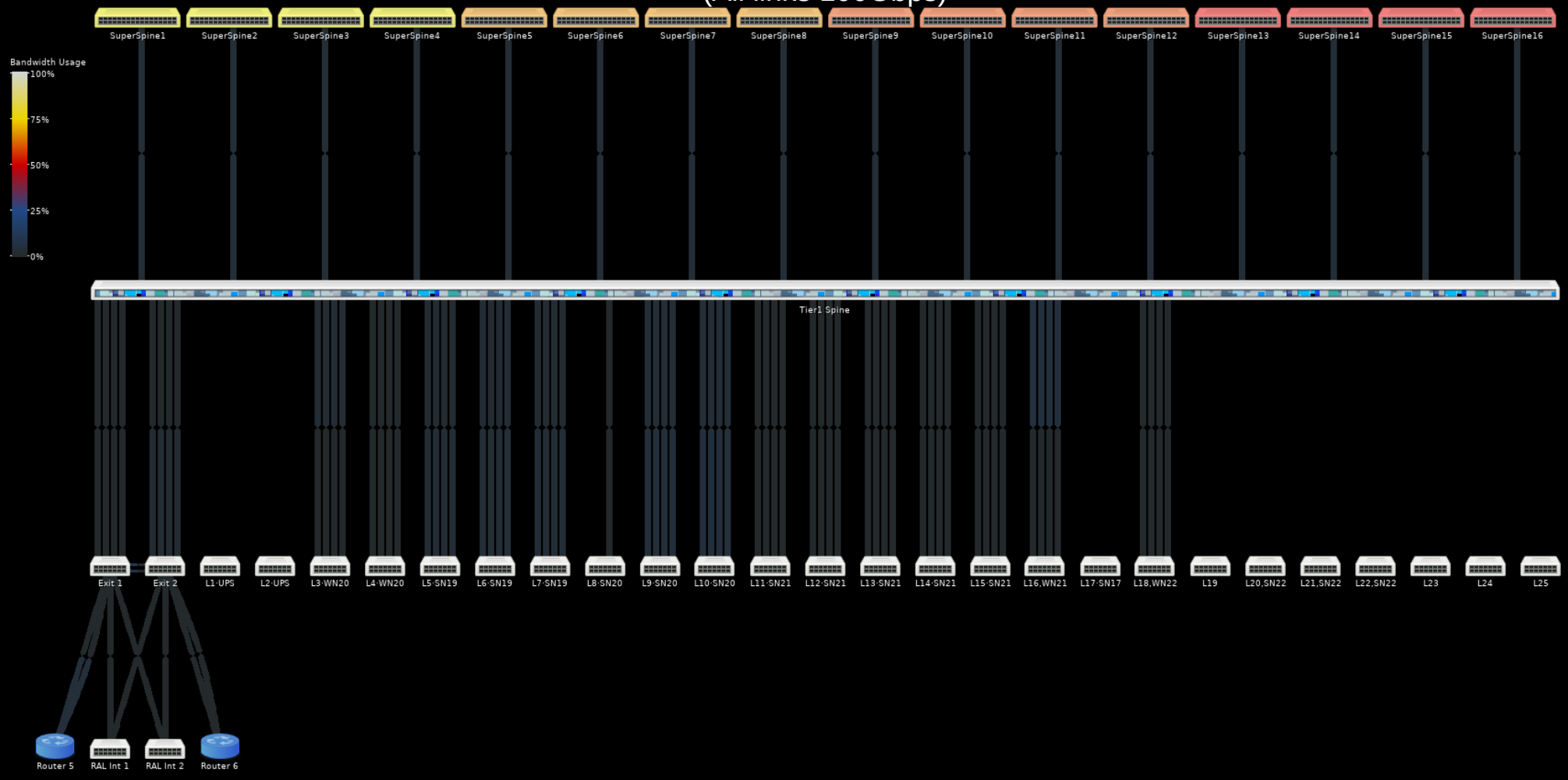
New Network — August 2022

(All links 100Gbps)

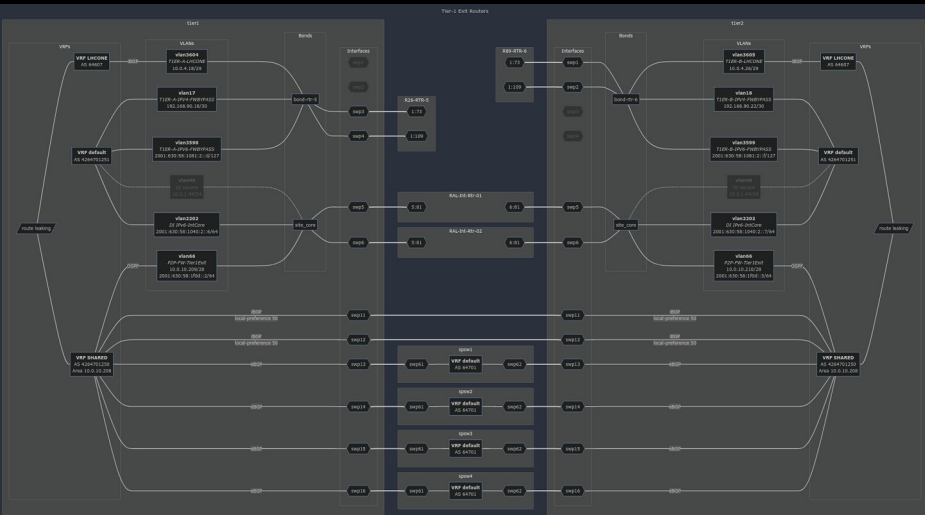


New Network — March 2023

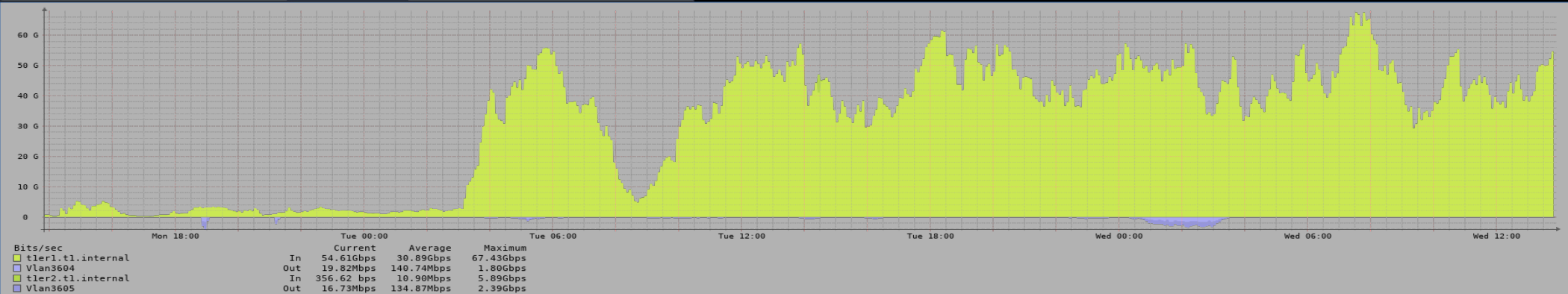
(All links 100Gbps)



LHCONE



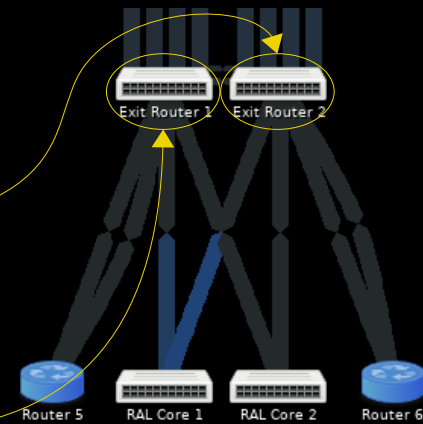
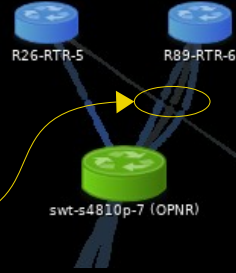
- Peering with both Exit Routers
- Everything on new network now advertised
 - 2020+ Worker nodes
 - 2019+ Storage
 - No echo gateways (yet)



LHCOPN



- Still peering with legacy OPNR via RAL Border Router 6.
- Second 100Gbps link now provisioned!
 - Will connect directly to Exit Router 2.
- Then move existing link to Exit Router 1.



IPv6

- In production on legacy network
 - Many services dual-stack for many years
- Rolled out to all leaves on new network
 - Started thinking about dual-stack worker nodes
 - Decisions to make around Docker integration
- Preparing leaves on CTA network (for Antares)
 - Required by ALICE

Questions?

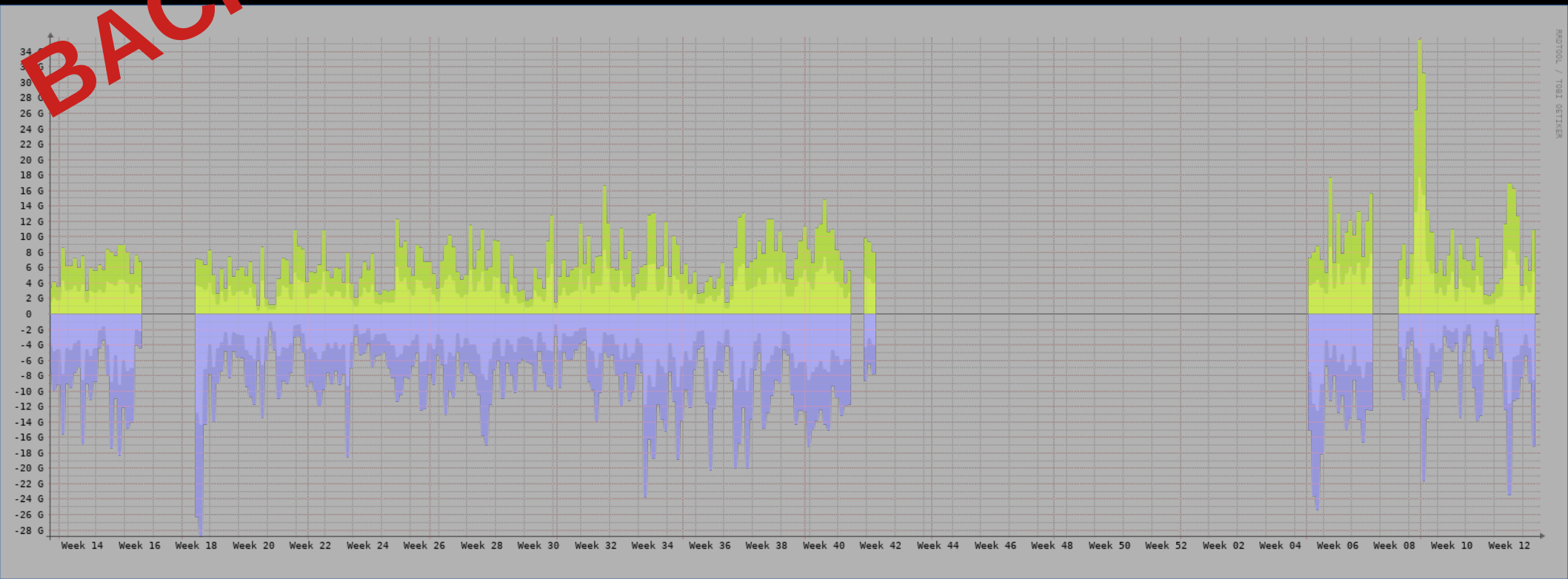
Last Year – Tier-1/SuperSpine

BACKUP



BACKUP

Last Year – LHCOPN



PROTON / TIBI GETIEN