

The RAL Tier-1 Network

James Adams, RAL
GridPP48, Ambleside

2022-09-01

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?

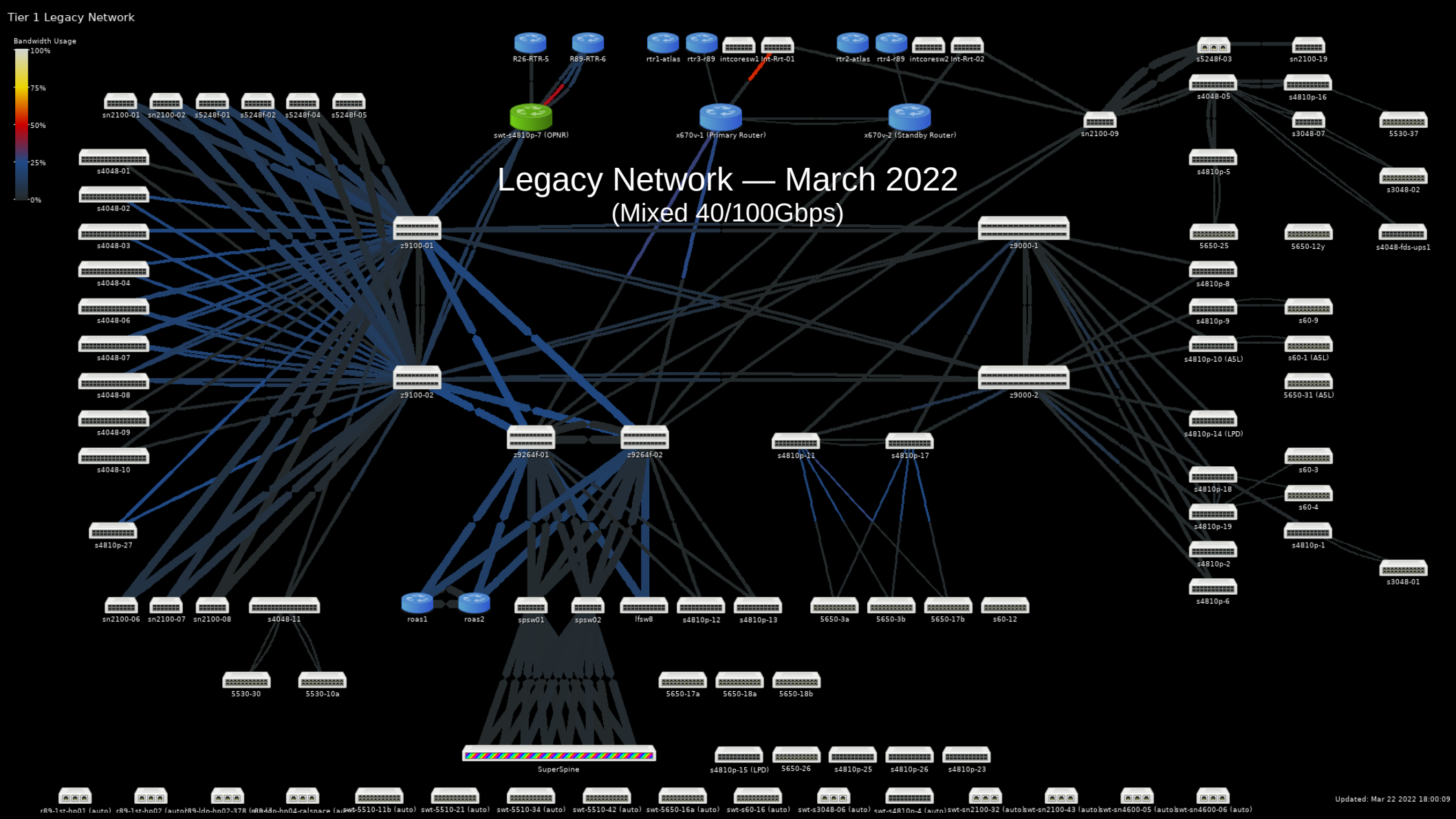
- RAL Macro-Segmentation
- Legacy network clean-up
- More hardware on new network
- More projects on SuperSpine
- Peering with LHCONe

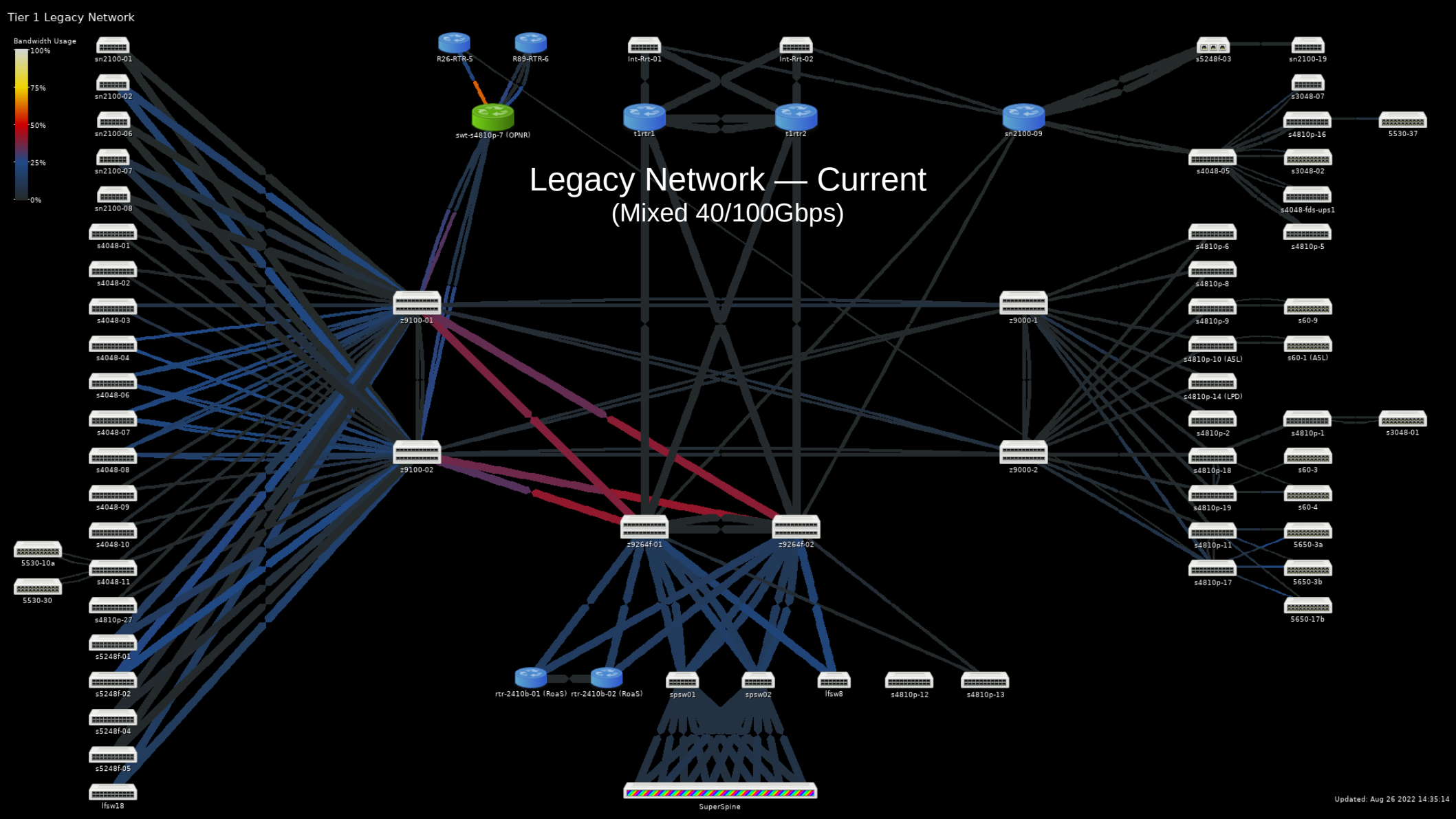
RAL Macro-Segmentation

- Essentially puts a firewall into the RAL Core
 - Designed to limit blast-radius of incidents
 - Breaks decades of connectivity assumptions
- Still discovering side effects
 - Active-Active router pairs no longer fail-out correctly
 - Long running SSH sessions dropped
- Somewhat mitigated by SuperSpine

Legacy Network

- Removed lots of old equipment
- Legacy path to LHCOPN increased to 80Gbps
 - Bottleneck to CERN improved
- Replaced Tier-1 Legacy Routers
 - Routes to RAL Campus and Janet
 - Removed 40Gbps bottleneck (now 200Gbps)
 - Will become the routers for facilities network later





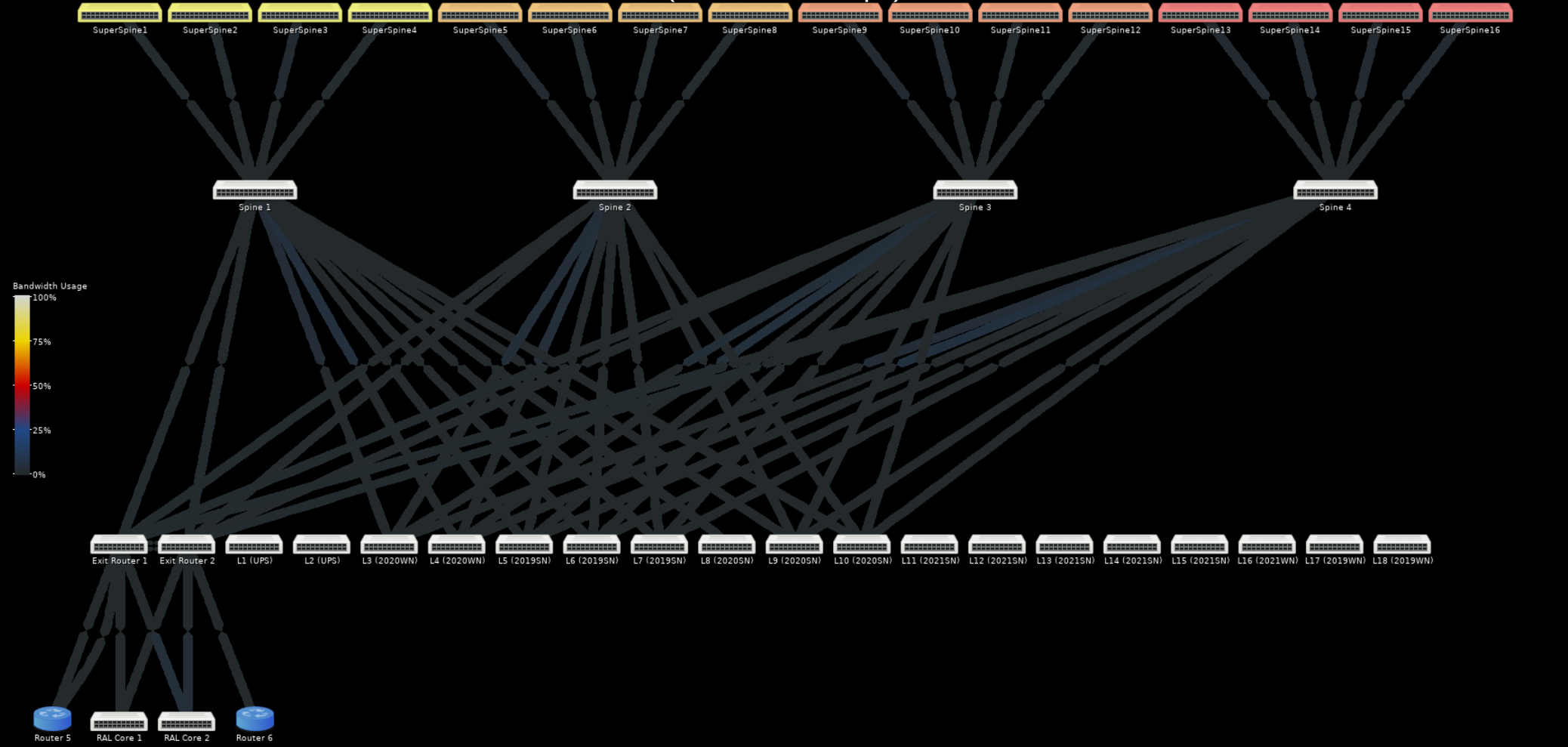
Updated: Aug 26 2022 14:35:14

New Network

- Three leaves of workers in production (+1)
 - 85% of pledge
- Five leaves of storage in production (+5)
 - 33% of Echo capacity
 - Data being migrated carefully
- Joined LHCONe
 - Initially with a very small prefix for PerfSonar

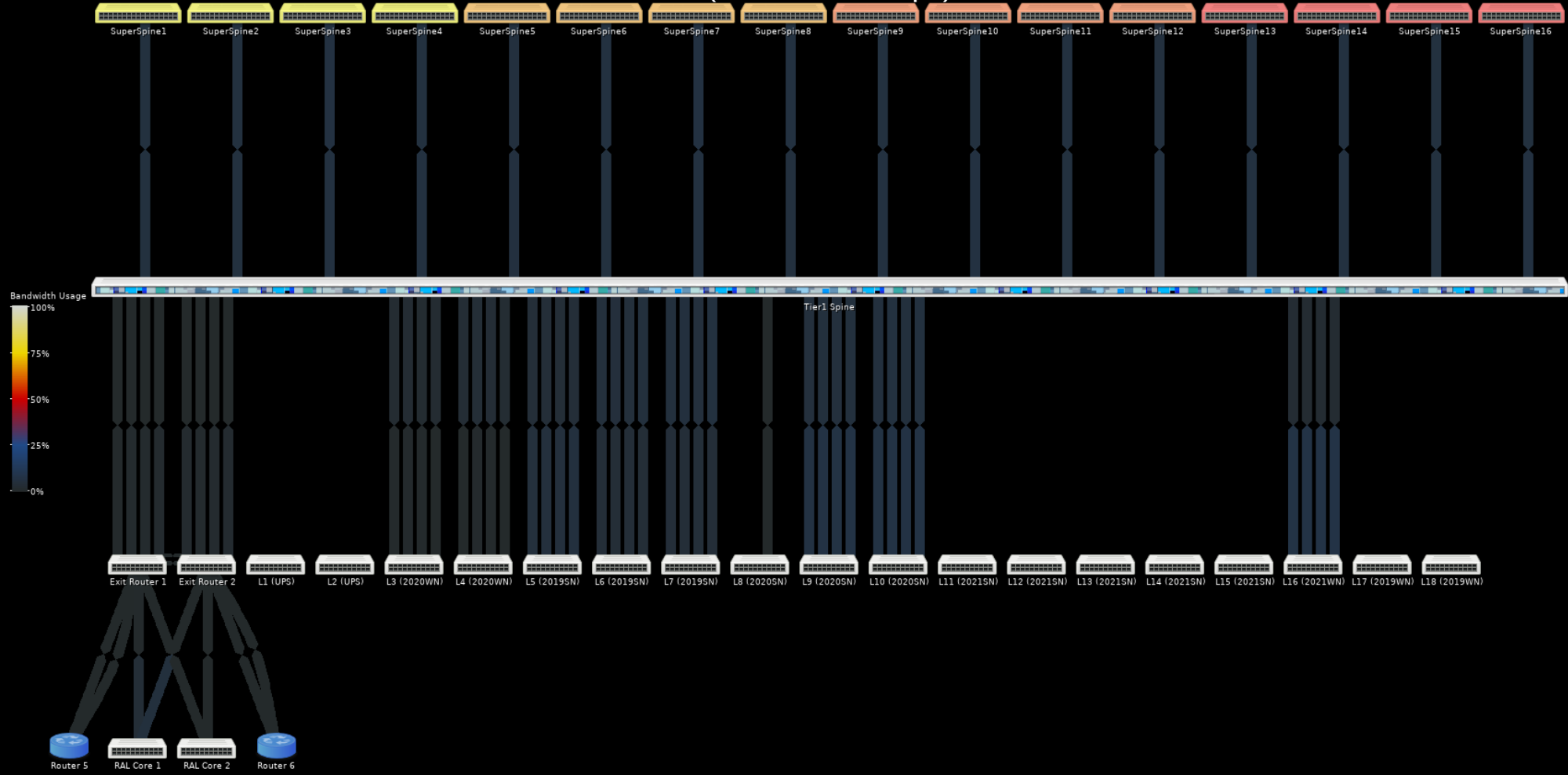
New Network — March 2022

(All links 100Gbps)



New Network — Current

(All links 100Gbps)

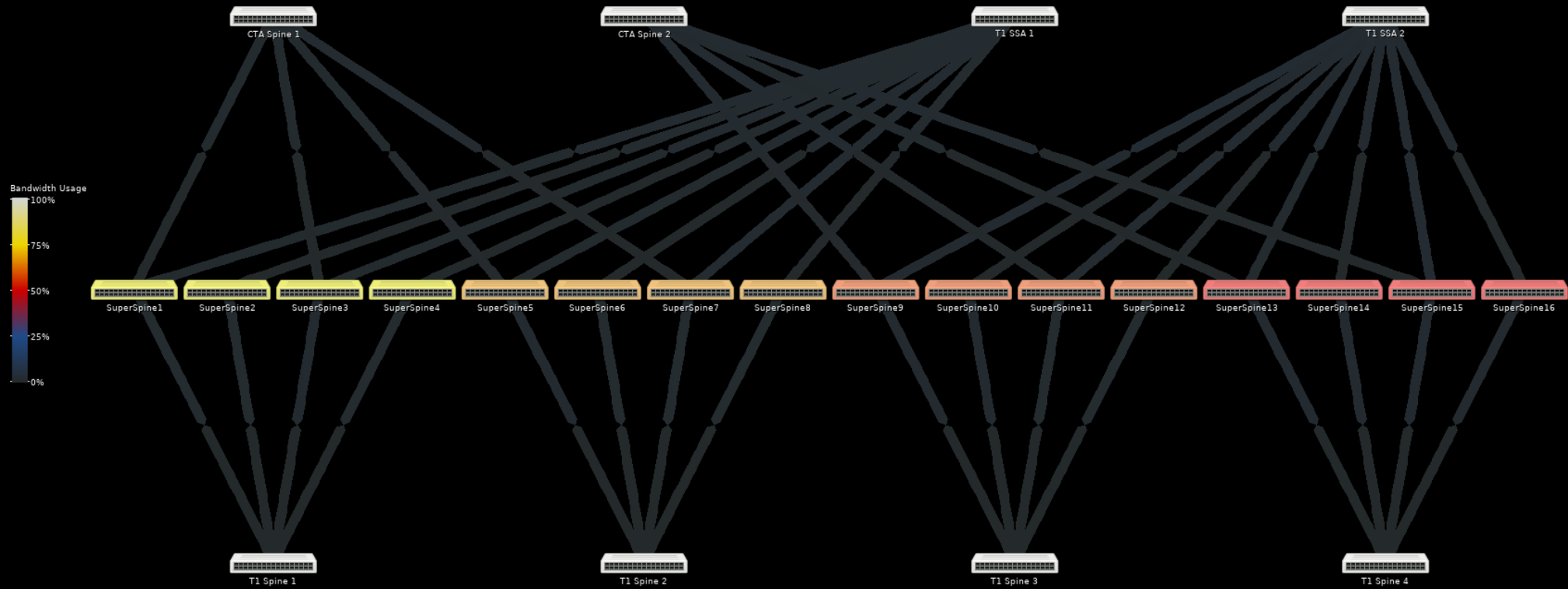


SuperSpine

- Joins SCD networks at up to 1.6Tbps (5-Stage Clos)
 - JASMIN
 - DAFNI
 - Tier1 legacy and new networks
 - Antares (CTA)
- SCD Cloud (R89 Pod) recently joined
 - R26 Pod joining soon

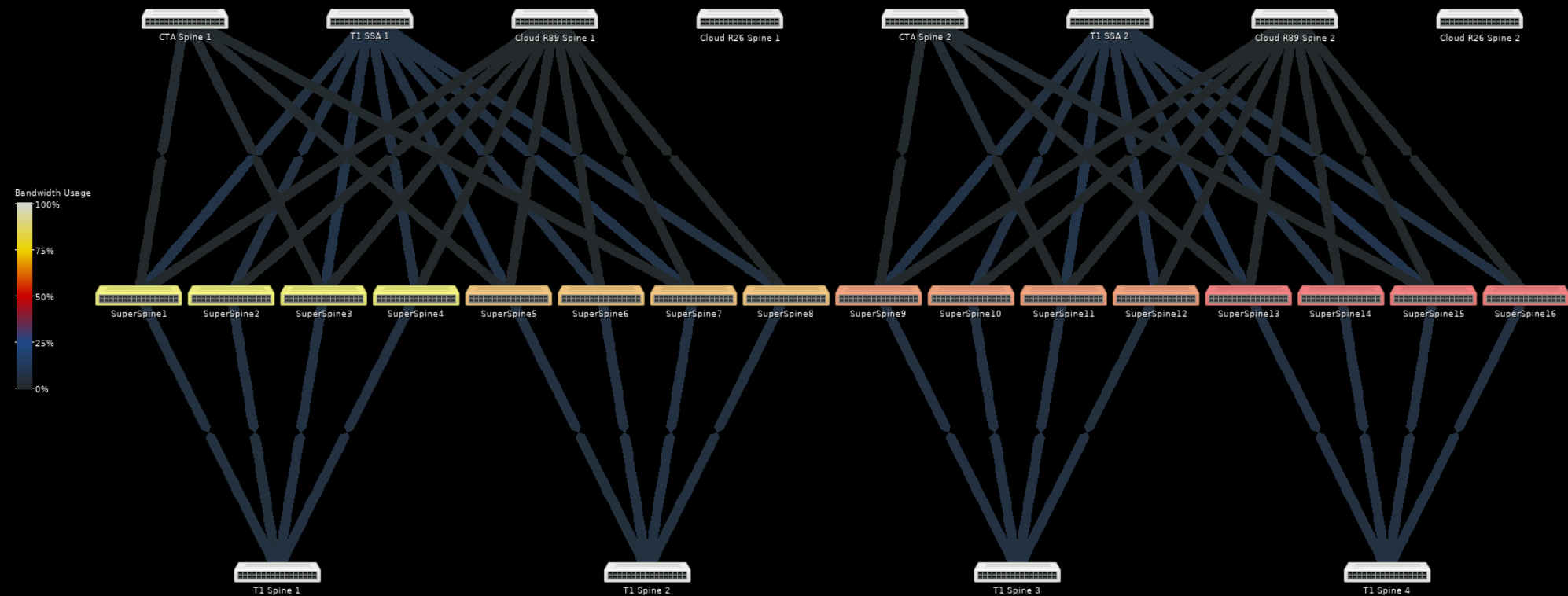
SuperSpine — March 2022

(All links 100Gbps, JASMIN & DAFNI not shown)



SuperSpine — Current

(All links 100Gbps, JASMIN & DAFNI not shown)

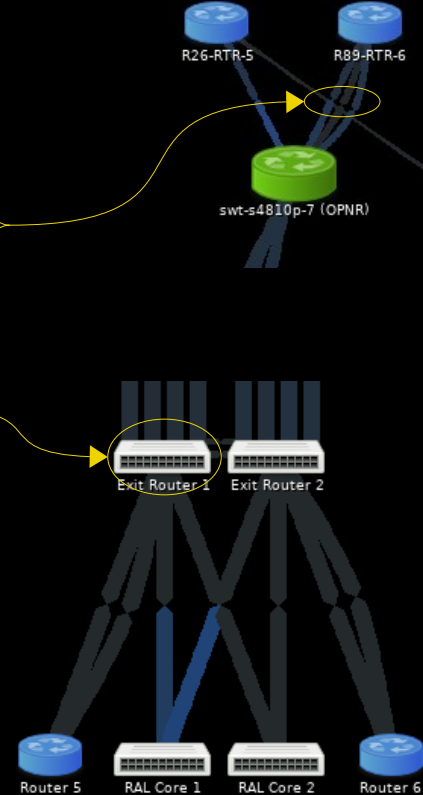


IPv6

- Needed more IPv6 space for new network
 - Legacy network has two /64
 - New network needs a /64 per leaf switch (currently twenty)
 - STFC addressing plan too restrictive (8-bits per department)
- A new /48 has been allocated by JISC for the Tier-1
 - Discussion happening internally to agree routing over SuperSpine
- New network will become “IPv6 first”

LHCOPN

- Currently peers with legacy OPNR via RAL Border Router 6.
- Move directly to Tier-1 Exit Router 1
 - Legacy network access over SuperSpine
 - Reverses existing traffic flow
- Pending IPv6 routing review
 - 60% OPN traffic is IPv6



LHCONE

- Currently peers with Tier-1 Exit Routers
- No services other than PerfSonar advertised yet
 - Reminder: LHCONE bypasses site firewalls
- Worker nodes will join after Condor 9 upgrade
 - Still need to roll out host firewalls
- Pending IPv6 routing review
 - Legacy network will access over SuperSpine

PerfSonar

- New nodes deployed on Tier-1 and Antares pods
- Problem with routes for the additional interfaces
 - Not being installed in the switching silicon
 - Being discussed with vendor
 - Alternative plan to move back to traditional method
 - Discouraged by PerfSonar developers
 - Remains the fall-back plan if not resolved by October

Questions?

BACKUP

STFC Addressing Scheme

Each project allocated one or more IPv6 /64

- 16 bits available to describe subnet

2001 : 0630 : 0058 : a b c d : 0000 : 0000 : 0000 : 0000

NETWORK HOST

JANET : RAL : a b c d : 0000 : 0000 : 0000 : 0000

a =STFC Address plan version (0-15)

b =Network Type

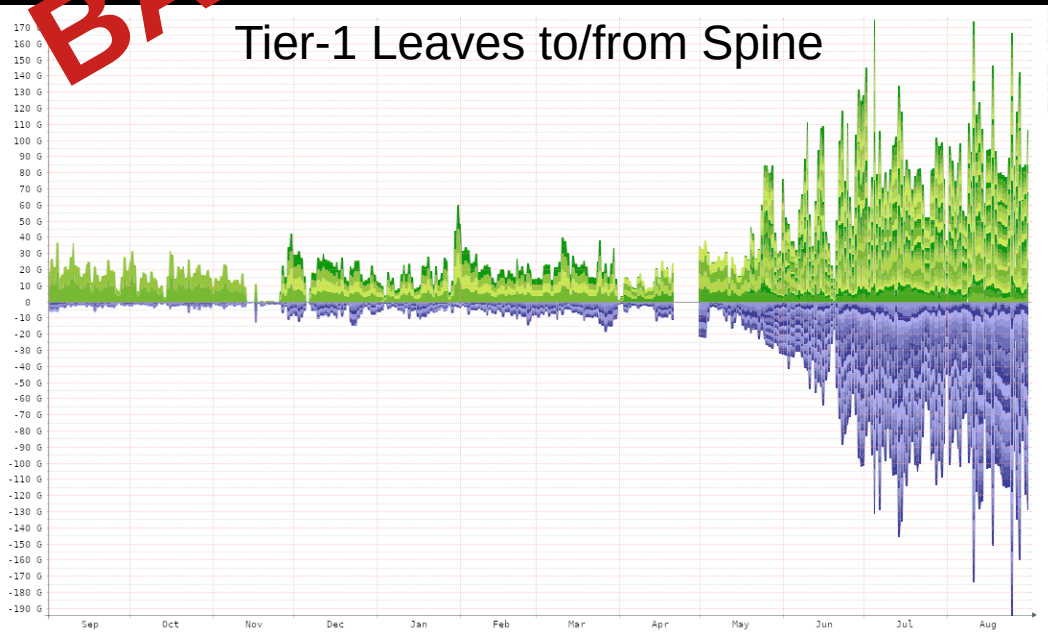
c =Network Subtype

d =Assigned by subnet owner (Tier 1 addressing scheme version)

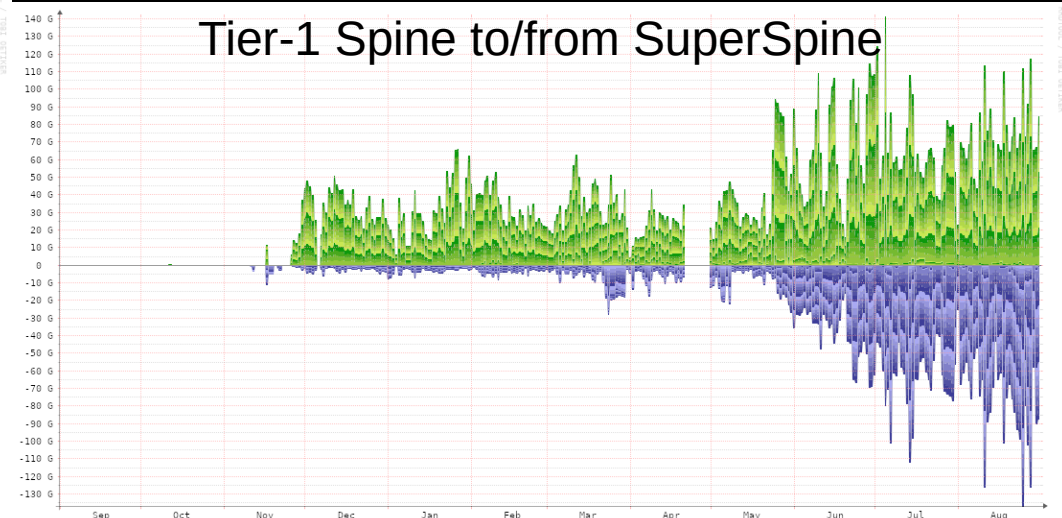
BACKUP

One Year of Traffic

Tier-1 Leaves to/from Spine



Tier-1 Spine to/from SuperSpine



BACKUP

One Year of Traffic

