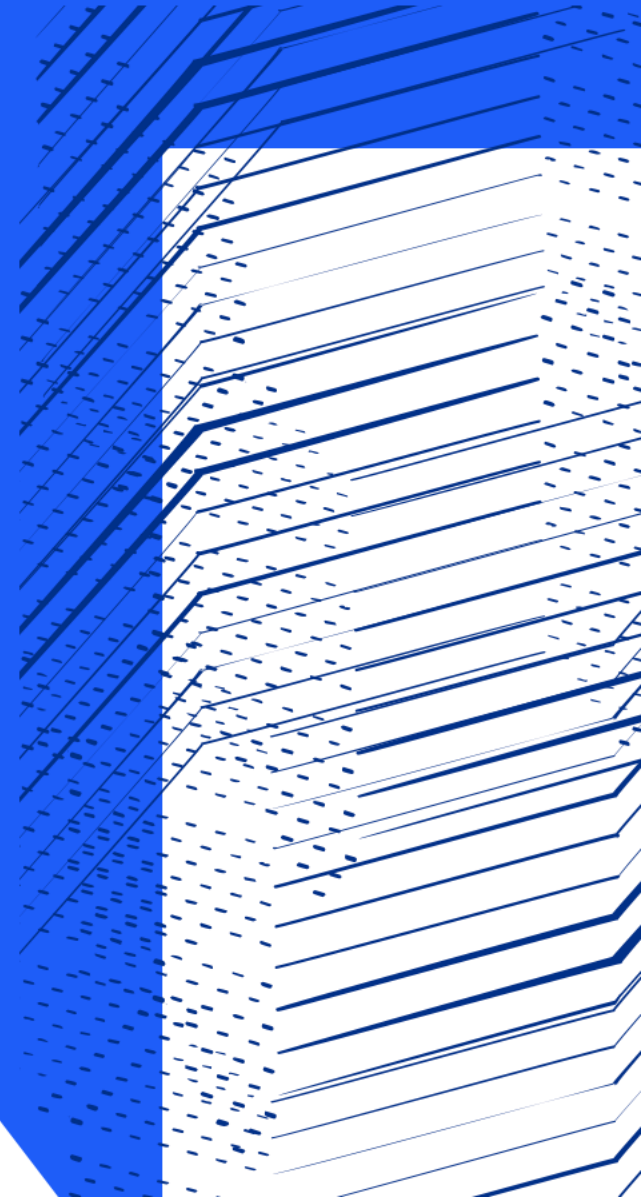




Science and
Technology
Facilities Council

RAL-LCG2 Network Outage October 2022

Alastair Dewhurst



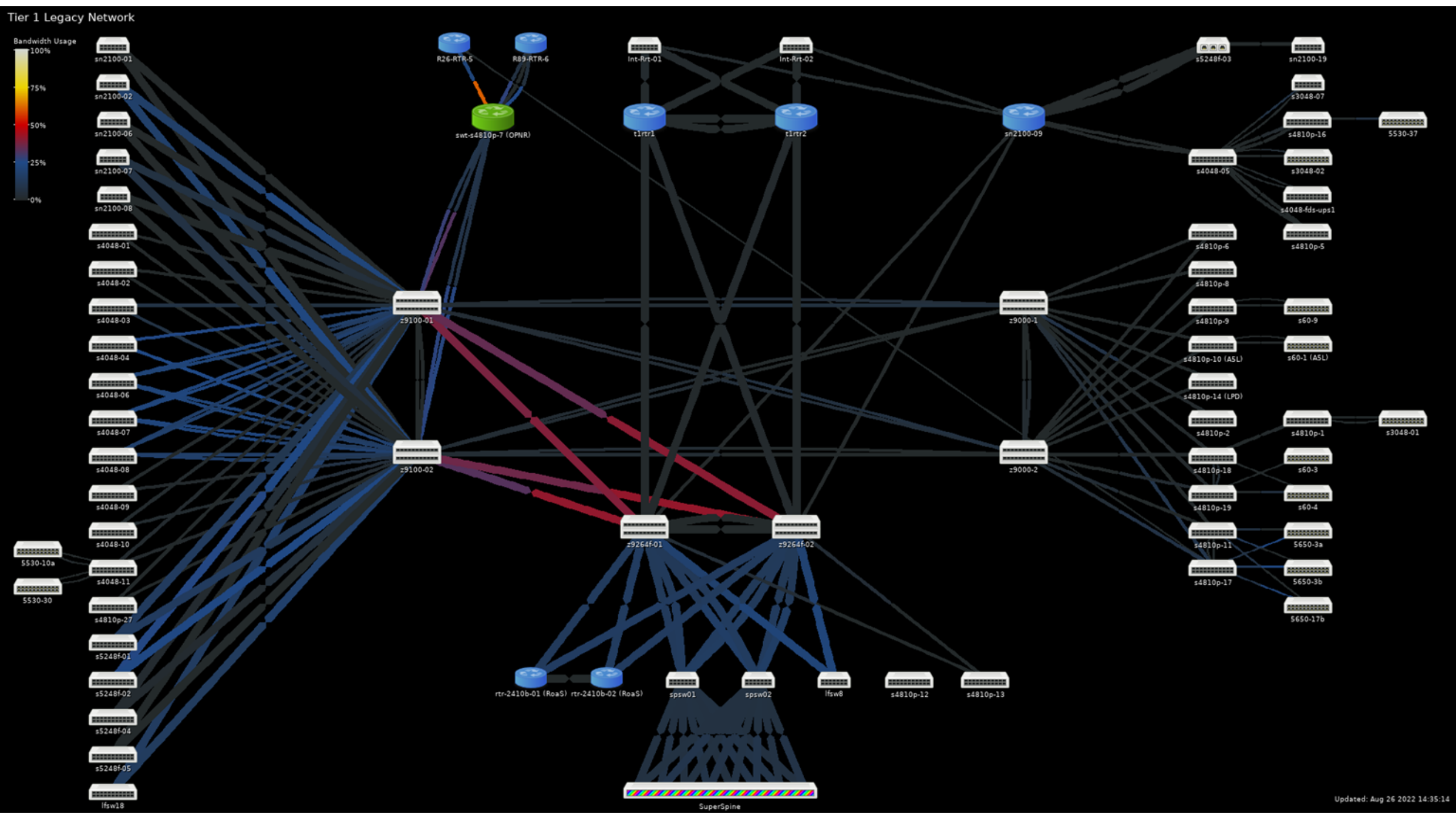
Introduction

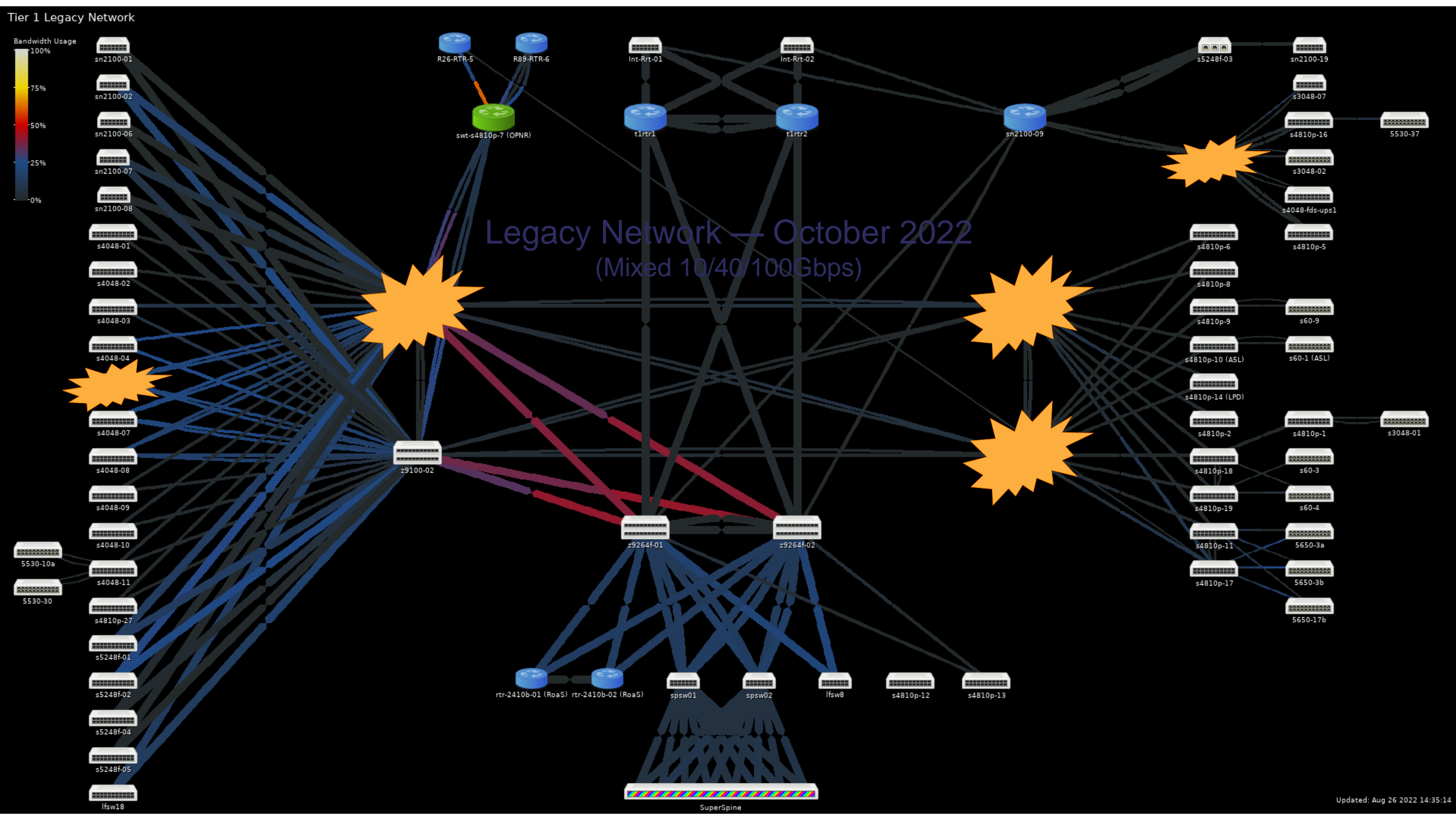
- This report summarises the RAL-LCG2 Network outage which involved two outages and caused significant degradation to services for ~10 days.
- Full SIR can be found [here](#).

Legacy Network Meltdown

On 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





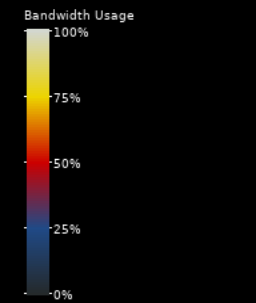
Initial Response

- 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
- We were able to get the network functional by the end of the day (17th) and many services returned to production the following day.
- However we had no resilience and batch and Echo needed to be limited.
- One good thing: Plans had been made for migrating off them

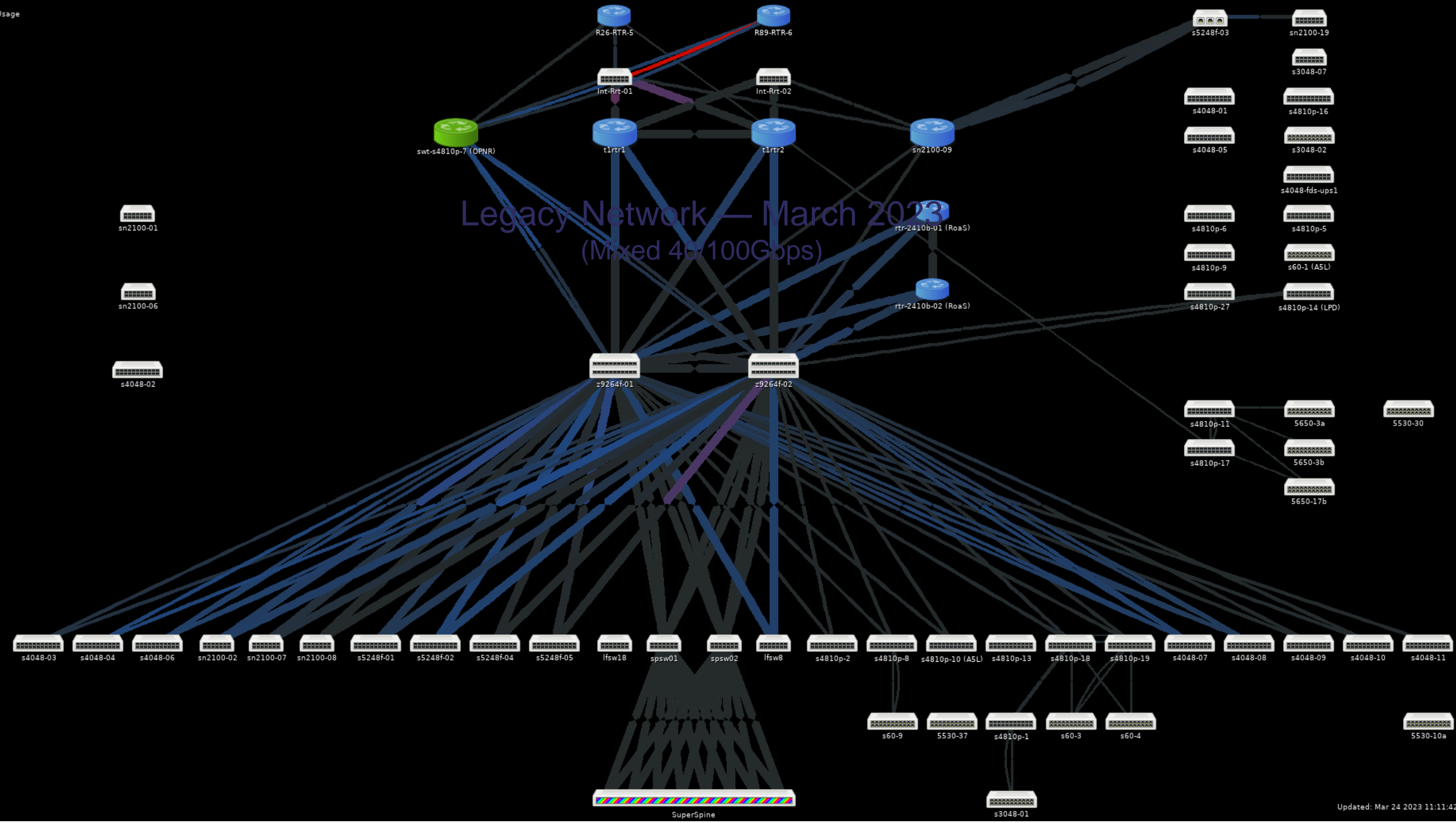
Follow up response

- Decided to make 12 months of planned changes in a single day
 - Scheduled an intervention for 26th October.
 - Moved 56 cables and config from four core switches
- Not everything came back correctly so downtime had to be extended.
 - It wasn't bad for the complexity of the change and the time pressure.
- Disabled management network entirely
 - Designed an interim replacement architecture
 - Replacement has been slow, but steady progress.

Tier 1 - Legacy Network



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





Science and
Technology
Facilities Council

Lessons learnt



Lesson learnt – legacy network

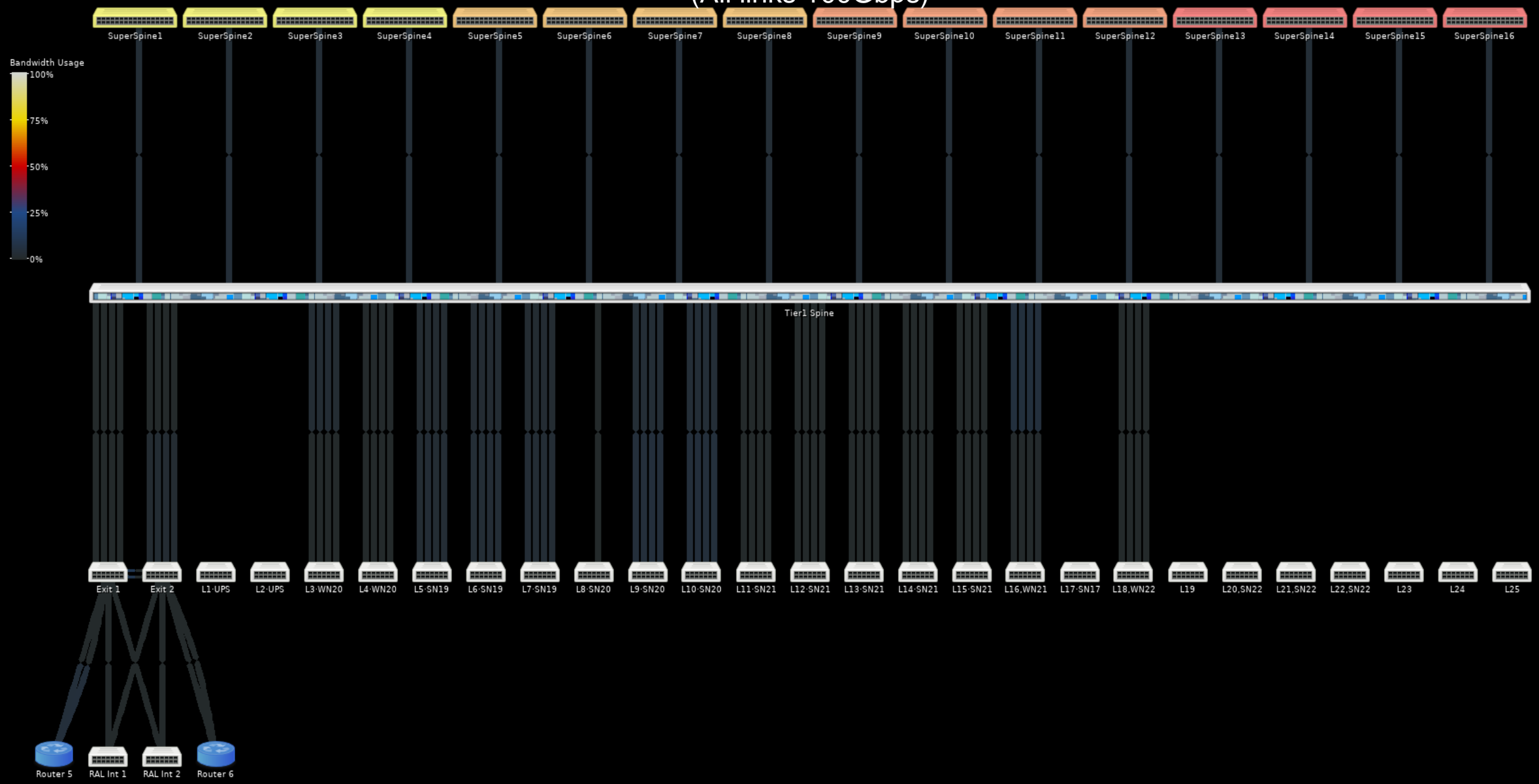
- We weren't surprised that the Legacy network failed.
 - We have been working to replace it since 2020.
- Between 2015 and 2018 there was no permanent RAL Tier-1 Manager
 - Strategic decisions about the network weren't made which delayed necessary upgrades and replacements.
- The legacy network has been rebuilt simpler and with less ancient hardware and will be completely phased out by end of 2024.
- The new network is starting to deliver significant benefits (e.g. LHCONe, multiple 100Gb/s OPN links)
- The criticality of the network has been recognized and effort to support it is a core part of the GridPP7 funding bid.

Lesson learnt – operational

- When RAL has a major outage it is likely to affect both the standard Tier-1 services and GOCDB.
 - Better documentation has been put in place to ensure that more people know how to announce a downtime when you can't use the GOCDB.
- One Antares server (CTA-Frontend) was setup on the legacy Tier-1 network.
 - This has now been moved to the correct location on the network.
- The criticality of the management network for disaster recovery had been underestimated.
 - We have invested £350k to purchase 110 new switches to modernize the hardware.
- We were already in the process of moving our documentation and ticket system to Jira and Confluence in the Cloud.
 - This should allow for more effective Disaster Recovery as it won't be affected by an outage.

New Network — March 2023

(All links 100Gbps)





Science and
Technology
Facilities Council

Questions?