



## New network infrastructure at CCIN2P3

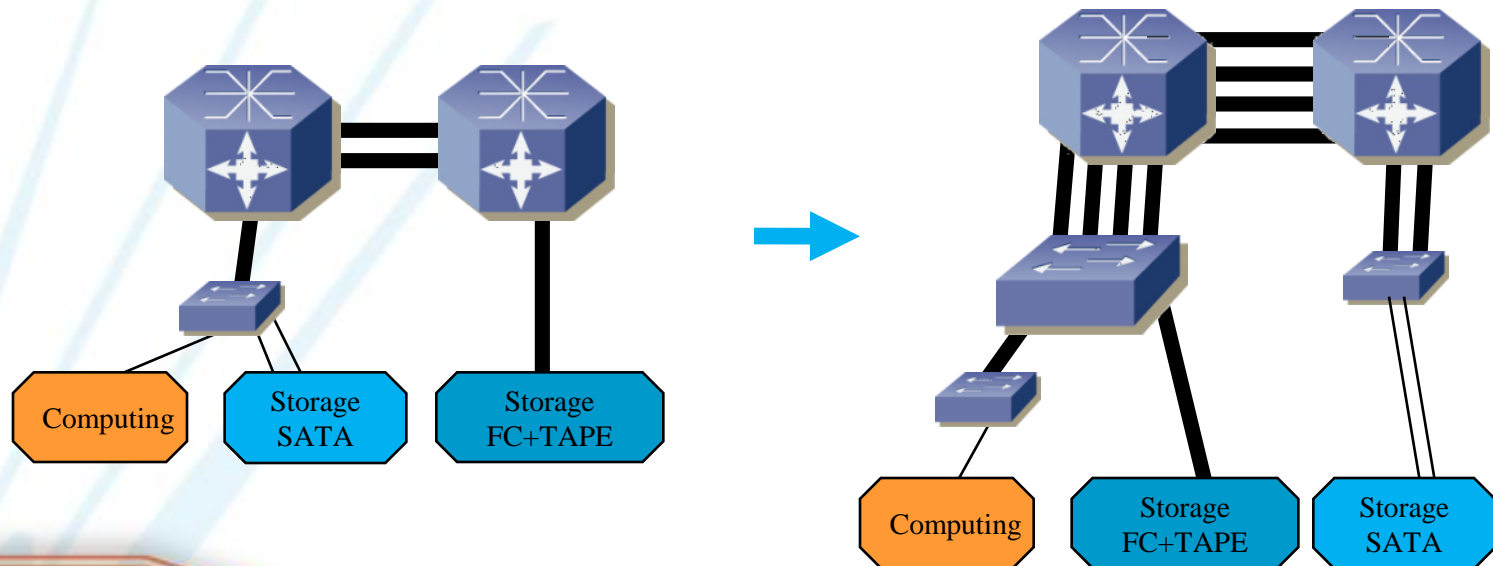
Guillaume Cessieux – CCIN2P3 network team  
Guillaume . Cessieux @ cc.in2p3.fr



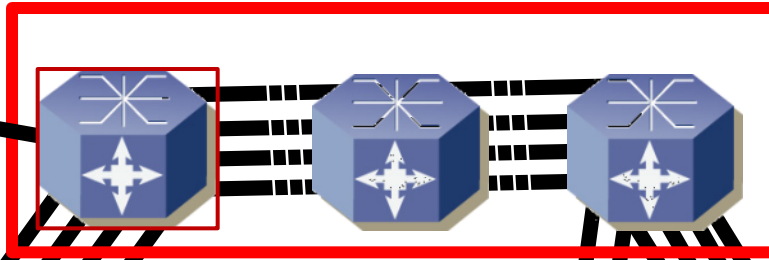
lrfu  
cea  
saclay

# Latest major network upgrade: Mid 2009

- Backbone 20G → 40G
  - No topology change, only additional bandwidth
- Linking abilities tripled
  - Distribution layer added



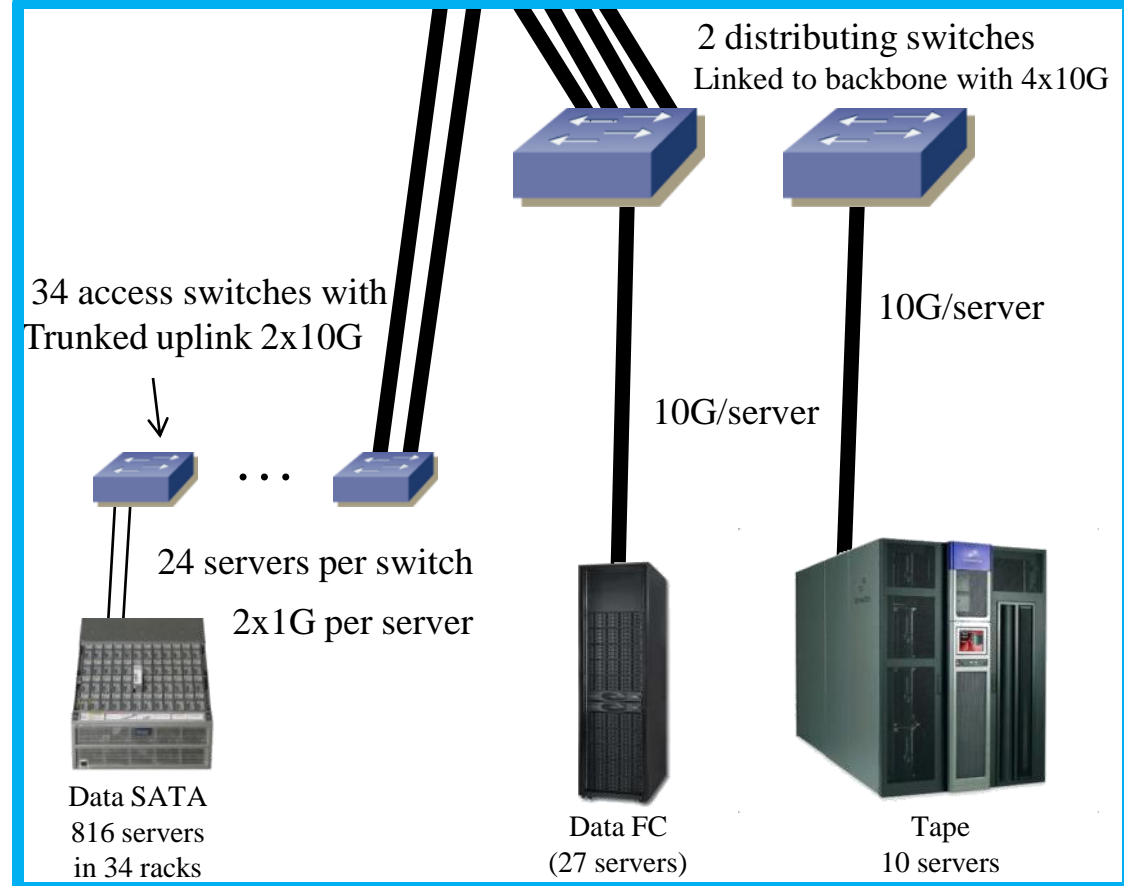
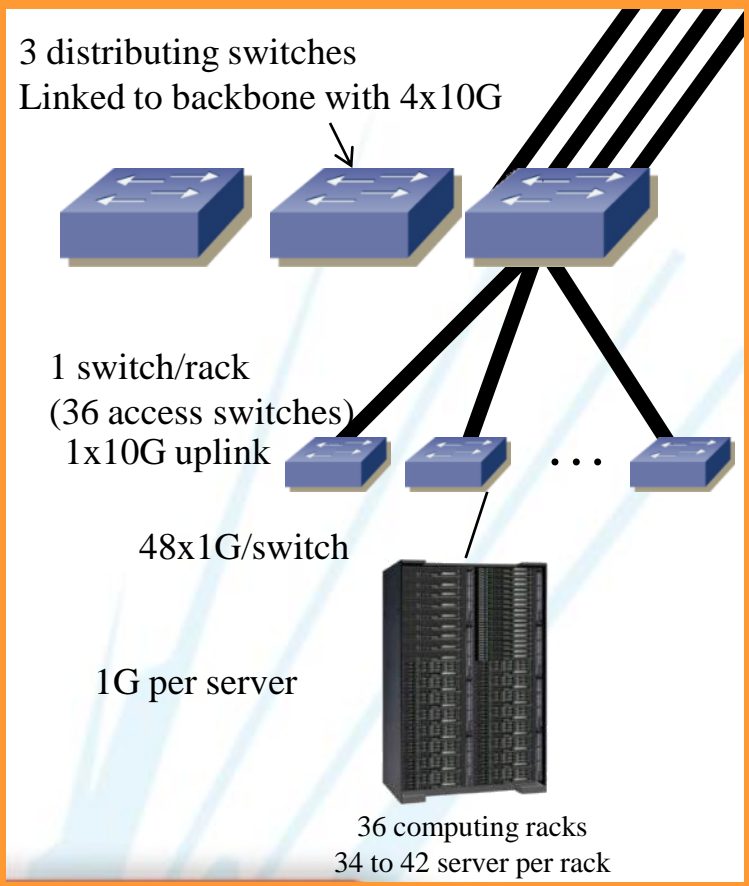
# Previous network architecture



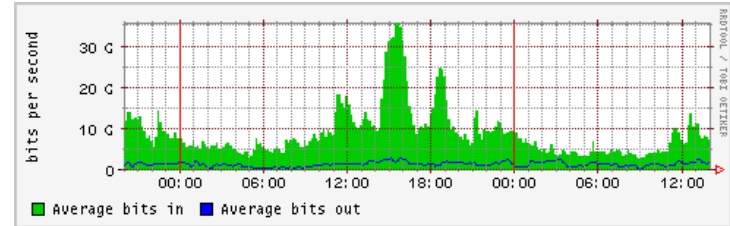
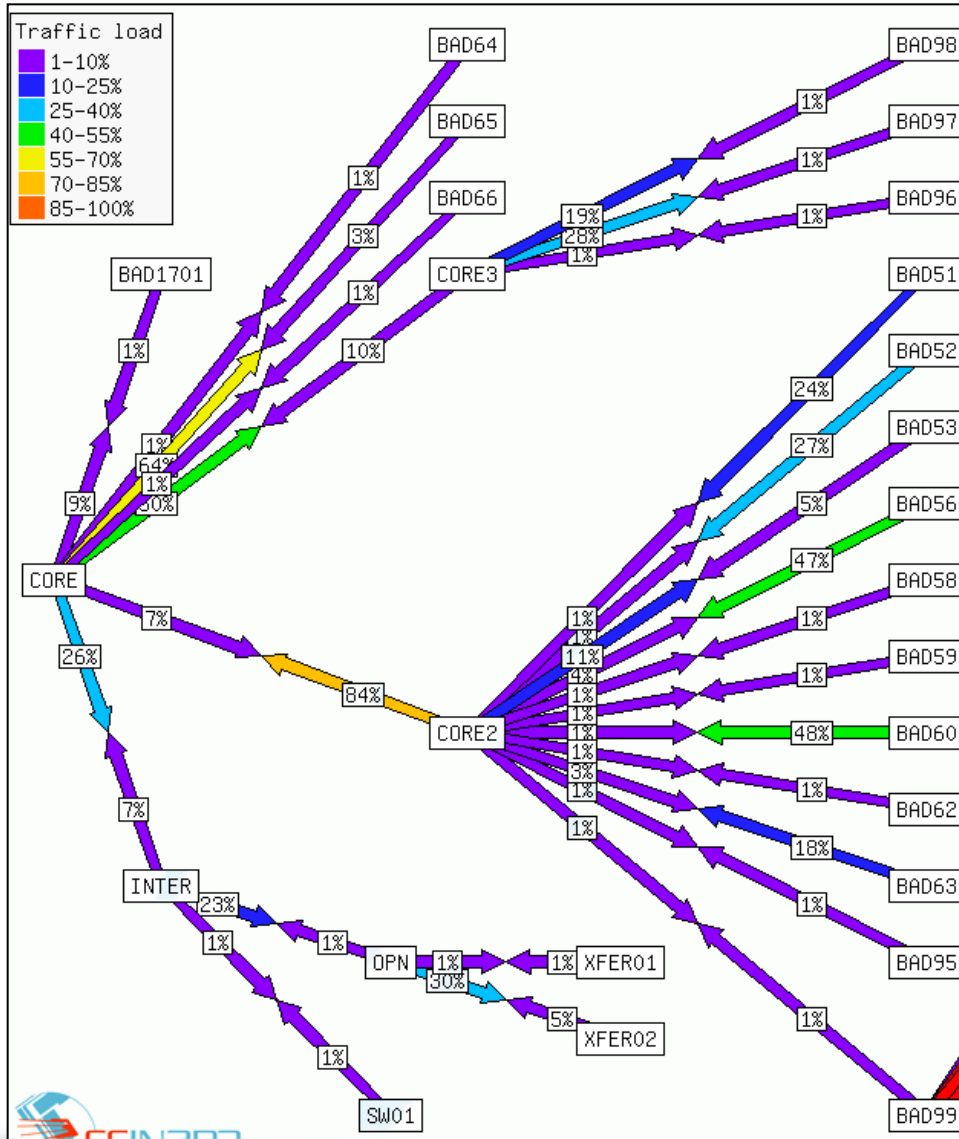
**Backbone 40G**

## Computing

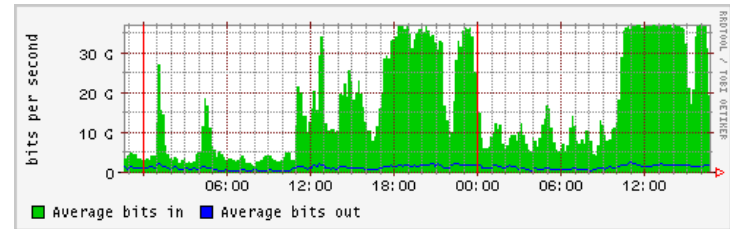
## Storage



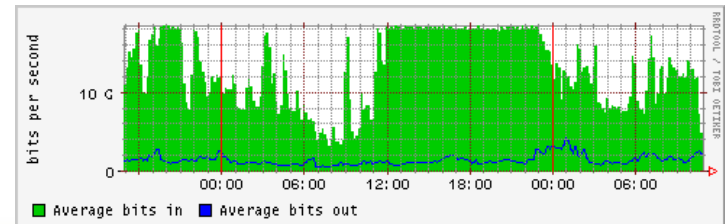
# Reaching limits (1/2)



Same 40G path 1 year later



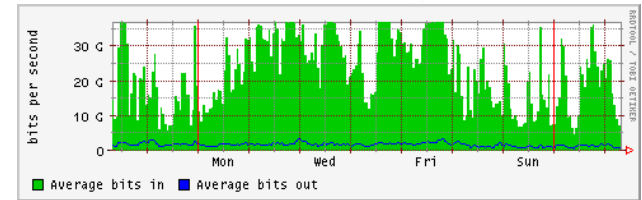
20G uplink of a distribution switch:



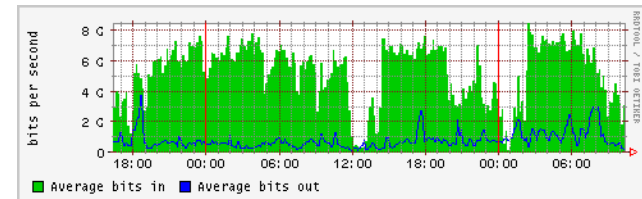
# Reaching limits (2/2)

- Clear traffic increase
  - More hosts exchanging more
  - More remote exchanges
- Limits appearing
  - Disturbing bottlenecks
  - Long path before being routed
- Upcoming challenges
  - New computing room, massive data transfers, virtualization, heavy Grid computation

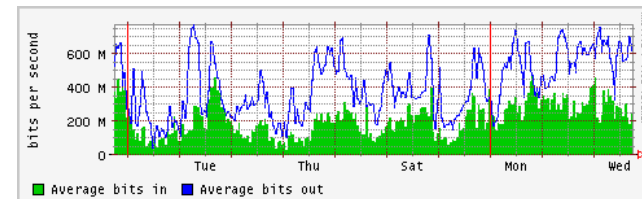
Usage of one 40G backbone etherchannel



10G direct link with CERN



Sample host average: 620M on 1G



# Complete network analysis performed

## Inventory

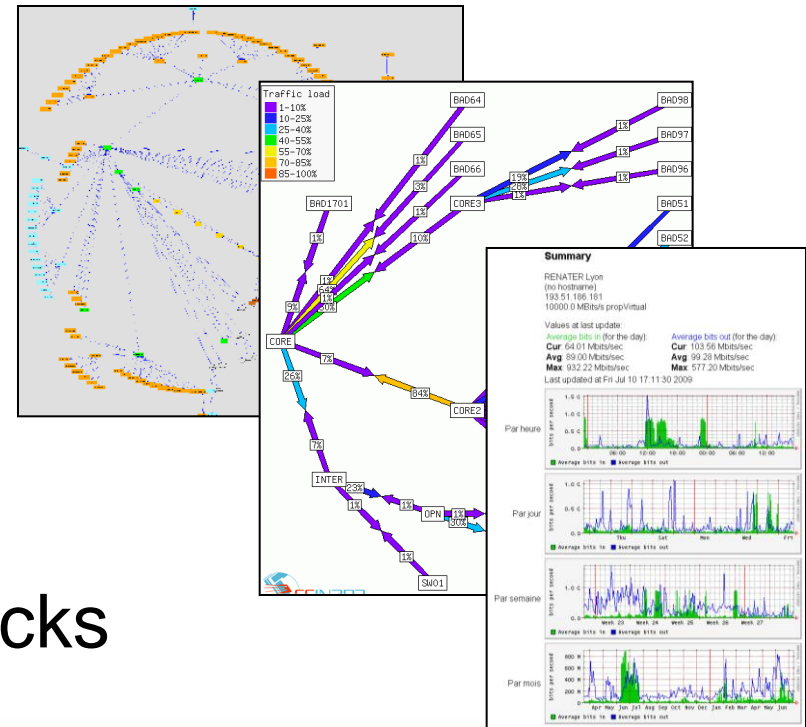
- 20 network devices found
  - Thanks discovery protocols...
- Software and features not harmonised

```
switch>show cdp neighbors
Capability Codes: R - Router, T - Trans Bridge, B - Source Route Bridge
S - Switch, H - Host, I - IGMP, r - Repeater, P - Phone,
D - Remote, C - CVTA, M - Two-port Mac Relay

Device ID    Local Intrfce  Holdtme  Capability Platform  Port ID
s1.in2p3.fr. Ten 2/1       146      R S I WS-XXXXX- Ten 1/50
s2.in2p3.fr. Ten 3/3       130      R S I WS-XXXXX- Ten 7/6
s3.in2p3.fr. Ten 3/4       150      R S I WS-XXXXX- Ten 6/6
```

## Topology

- A map worth anything



## Usage

- Traffic patterns, bottlenecks

# Requirements for new network architecture

- More bandwidth!
- Able to scale for next years
  - Allowing non disruptive network upgrade
  - Particularly with new computing room

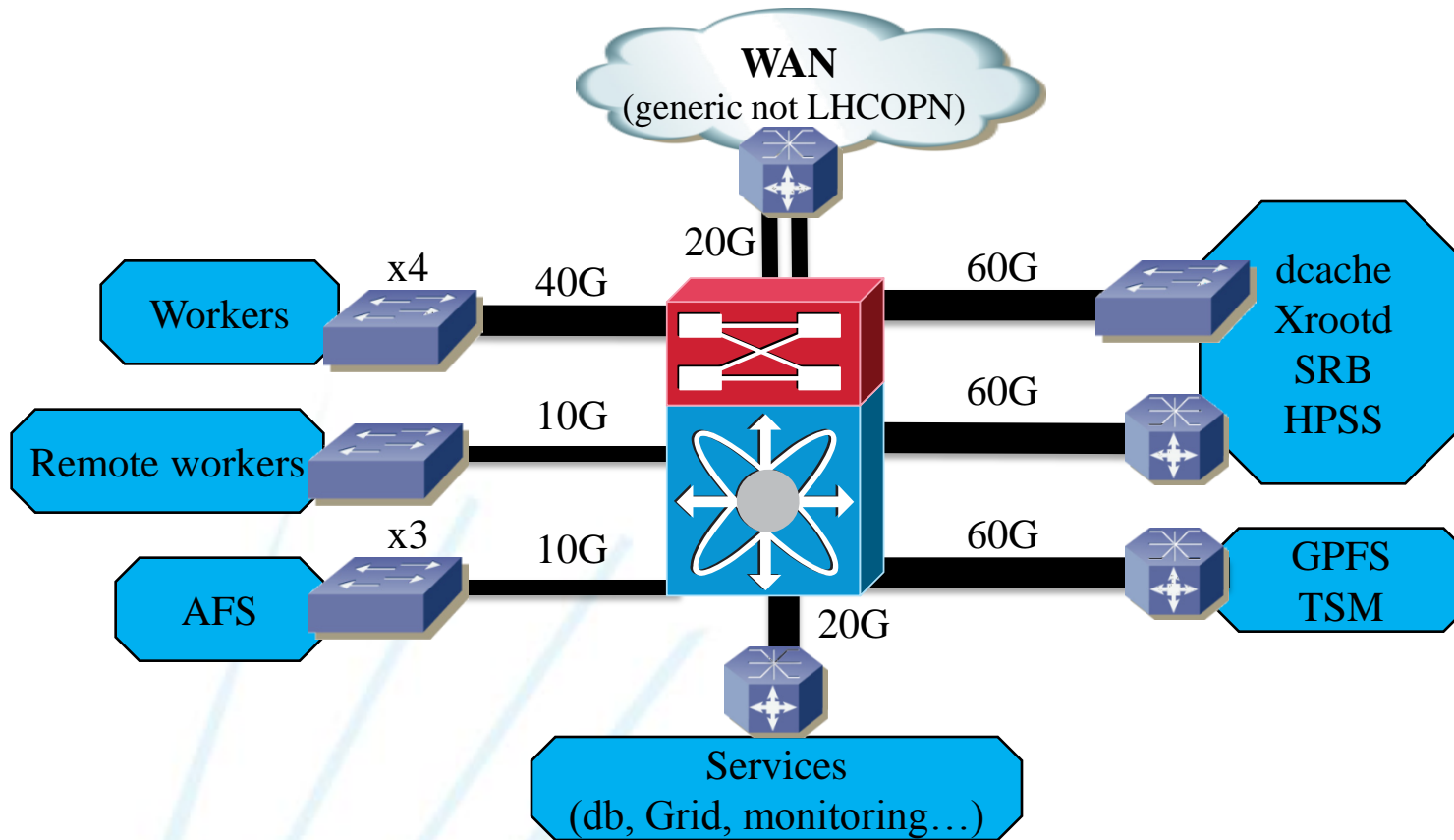


- Ease exchanges between major functional areas
- As usual: Good balance between risks, costs and requirements

# Main directions

- Target non blocking mode
- No physical redundancy, meshes etc.
  - “A single slot failure in 10 years on big Cisco devices”
  - Too expensive and not worth for us
  - High availability handled at service level (DNS...)
  - Big devices preferred to meshed bunch of small
- Keep it simple
  - Ease configuration and troubleshooting
  - Avoid closed complex vendor solutions
    - e.g things branded “virtual”, “abstracted”, “dynamic”

# New network architecture

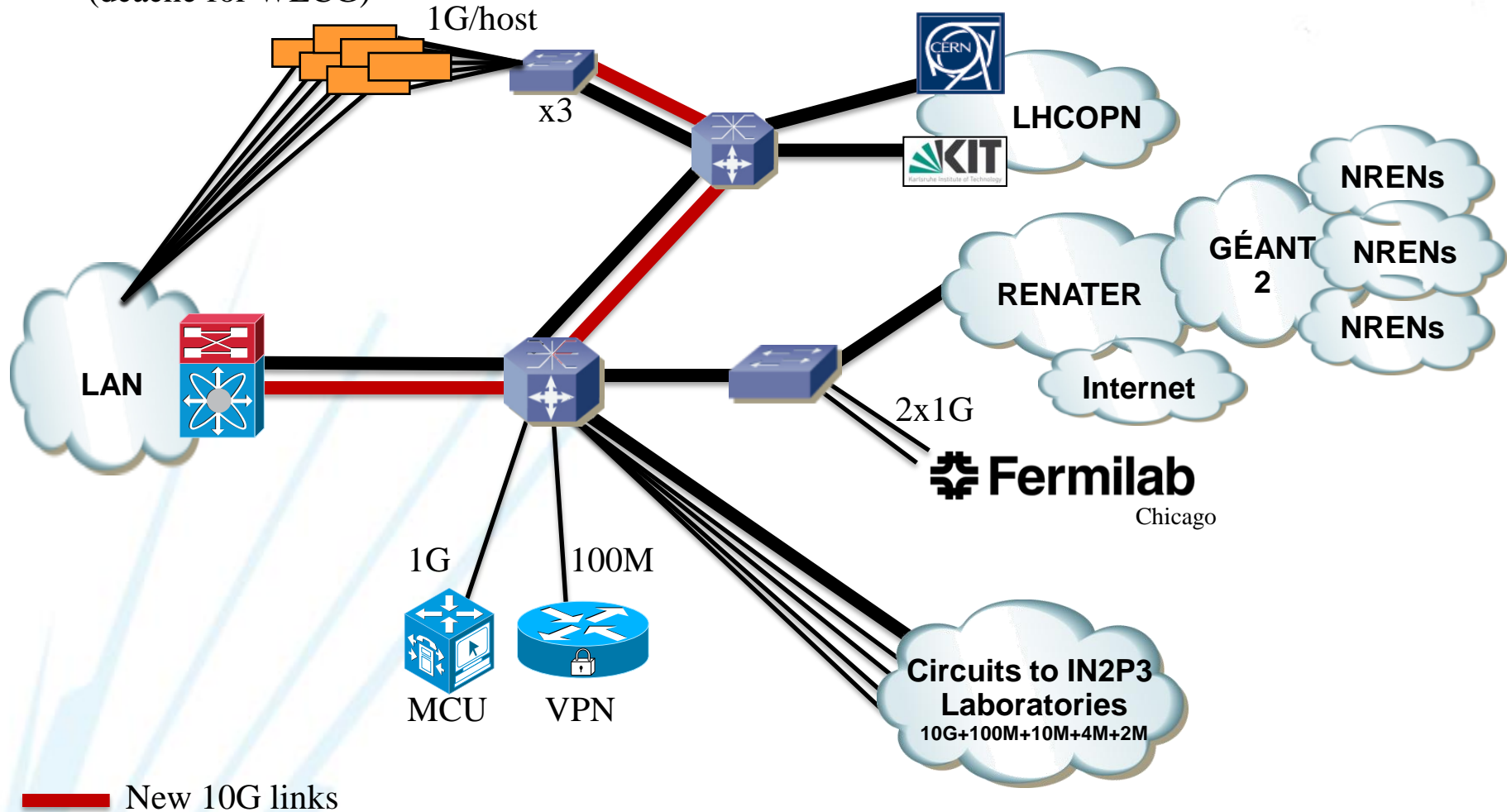


**From 160G often  
shared to 400G wirespeed**

Area	Old bandwidth	New bandwidth
AFS	30G <b>shared</b>	30G
GPFS TSM	40G <b>shared</b>	60G
Dcache Xrootd srb HPSS	40G <b>shared</b>	120G
Workers	40G <b>shared</b>	170G
WAN	10G	20G

# WAN upgrade

Dual homing of hosts doing massive data transfers  
(dcache for WLCG)

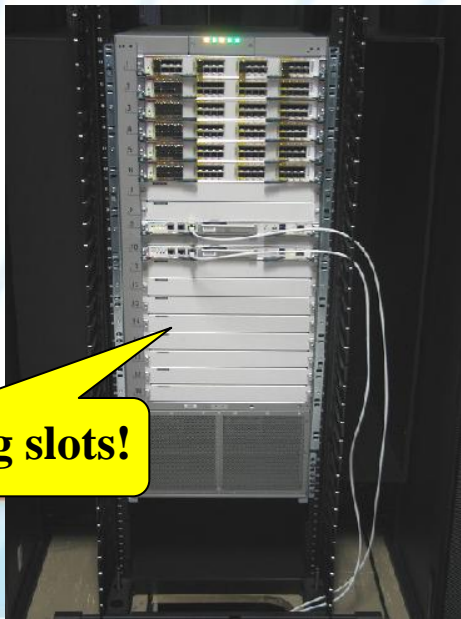


# Main network devices and configurations used



## ■ New core device: Nexus 7018

- High density device, really scalable
- Very modular: slots & switching engine cards
- 80G backplane per slot (8x10G non blocking)
  - Initial configuration: 6 slots, 3 switching engines (3x48G)
- This device is **vital**
  - 4 power supplies on 2 UPS, 2 managements slots



10 remaining slots!



2 extra switching engines possible

Compatibility check is done:			
Mod	boot	Impact	Install-type
1	yes	<b>non-disruptive</b>	rolling
2	yes	<b>non-disruptive</b>	rolling
3	yes	<b>non-disruptive</b>	rolling
4	yes	<b>non-disruptive</b>	rolling
5	yes	<b>non-disruptive</b>	rolling
6	yes	<b>non-disruptive</b>	rolling
9	yes	<b>non-disruptive</b>	reset
10	yes	<b>non-disruptive</b>	reset

# Main network devices and configurations used

6513



- 24x10G (12 blocking)  
+ 96x1G
- 48x10G (24 blocking)  
+ 96x1G

6509



- 64x10G (32 blocking)



Core & Edge

4900



16x10G



Distribution

4948



48x1G + 2x10G



Access

# A 3 months preparation

- Testing, preconfiguring, scripting, checklist
- Optical wiring: Fully done and tested before

[[telecom:checklist\_arret\_du\_2010-09-21]]

WIKI CC-IN2P3

Éditer cette page Antennes révisions

Vous êtes ici: start > serveur > idrftts\_cms/telecom-arret

1. Journée du 21 Septembre 2010:

- Arrêt des armoires électriques E, F, G, J, K, N pas en même temps et chacune pour environ
- armoires E et G: avertir DSI + CCSD + BV + LYRES (JB)
- armoire K: avertir REHATER (JB)
- armoires E et K: contrôler double alimentation des machines en racks Telecom, notamment

1. Soirée:

- Temps disponible: Arrêt telecom ~4h, coupure avec l'extérieur + service telecom: ~15 min
- Avant l'arrêt:
  - Déposer les nouveaux IOS (lent) et configurer les bonnes variables de boot [OK]
  - Mettre à jour le maximum d'équipements qui peuvent l'être (voir Xavier et al) pour tester
  - Précâbler ce qui peut l'être (coreA → core\*) [OK]
  - Mettre coreA dans les statistiques
  - Mettre à jour ou prévoir de le faire certains vieux équipements [OK]
    - ccnum, ccvialit, ccprn-bif, WS-C3524-XL à remplacer si possible par des 2960 : cccata-cccata-sg02
  - Avoir un modèle de configuration pour ccprn-coreA
  - Voir comment on se partage les tâches
- Pendant l'arrêt:
  - 19h-19h15:
    - Relead manuel des équipements en salle telecom: ccprn-inter, ccprn-opn
    - Changer cccata-sg01 (connecté à ccprn-inter) et sg02 (boite vers telecom & étage en d

2. 19h15-19h20:

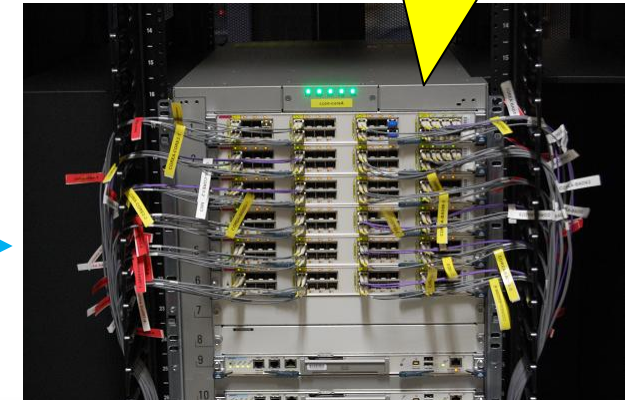
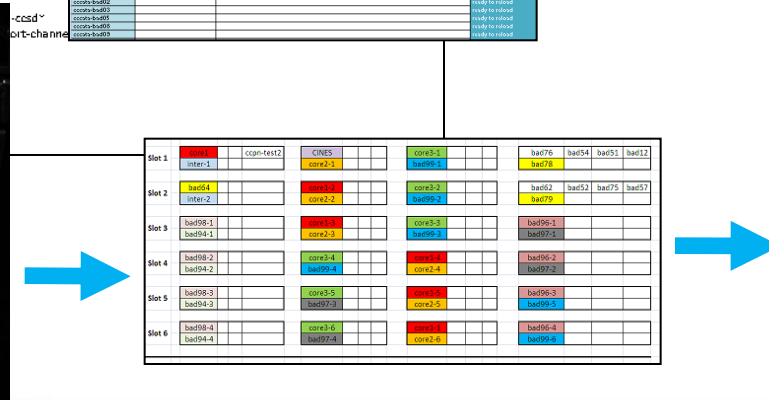
2. Soirée:

- Arrêt des armoires électriques E, F, G, J, K, N pas en même temps et chacune pour environ
- armoires E et G: avertir DSI + CCSD + BV + LYRES (JB)
- armoire K: avertir REHATER (JB)
- armoires E et K: contrôler double alimentation des machines en racks Telecom, notamment

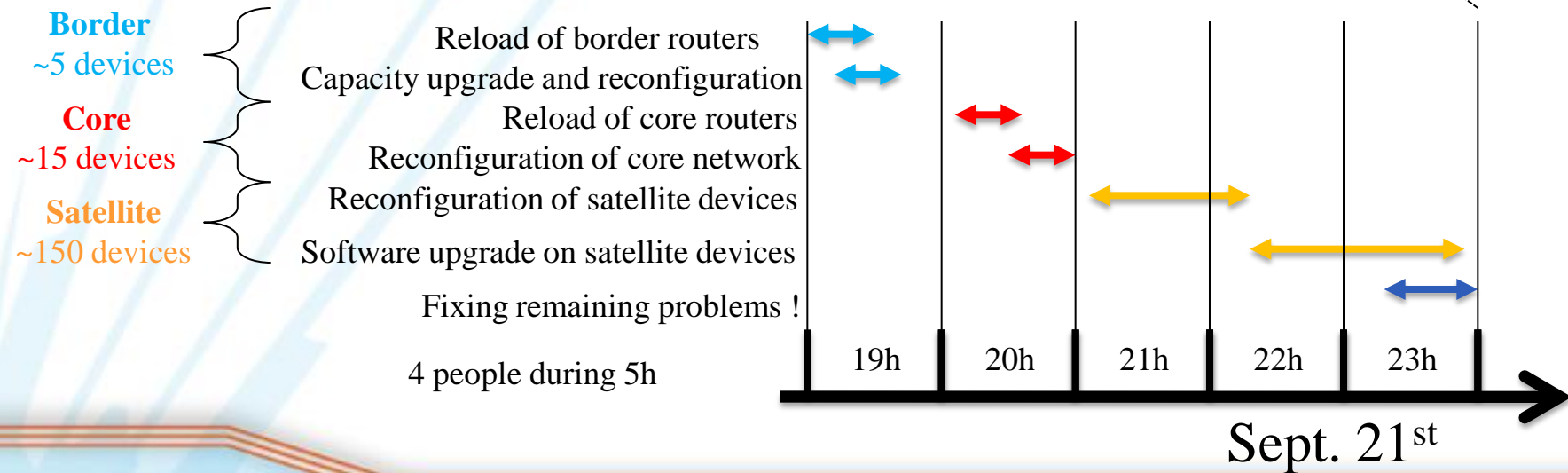
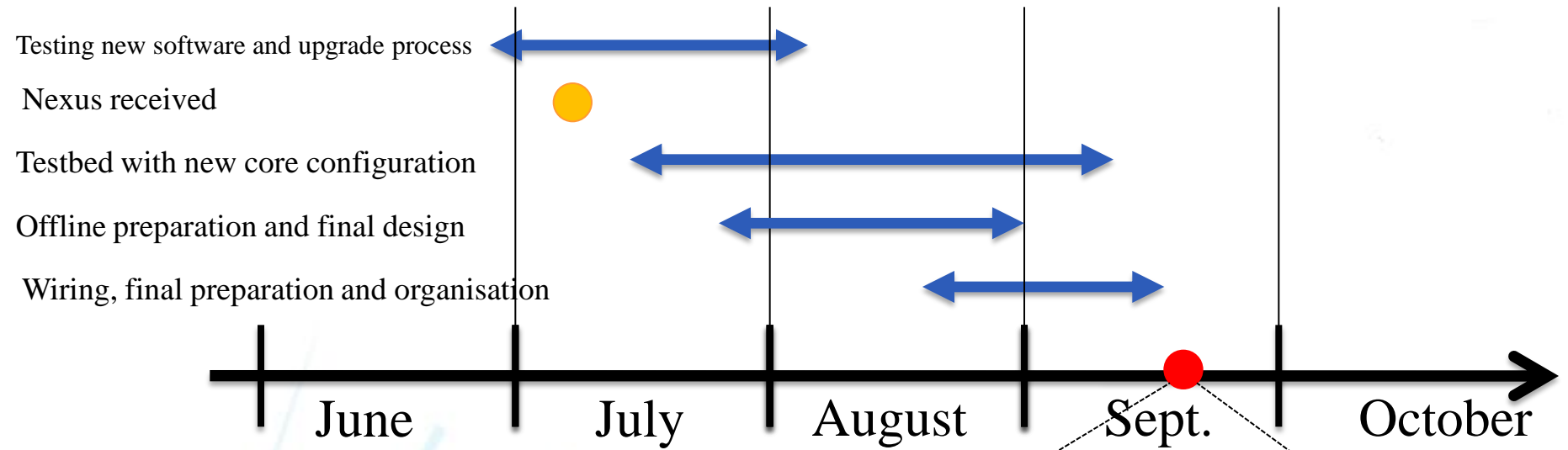
2. Soirée:

- Temps disponible: Arrêt telecom ~4h, coupure avec l'extérieur + service telecom: ~15 min
- Avant l'arrêt:
  - Déposer les nouveaux IOS (lent) et configurer les bonnes variables de boot [OK]
  - Mettre à jour le maximum d'équipements qui peuvent l'être (voir Xavier et al) pour tester
  - Précâbler ce qui peut l'être (coreA → core\*) [OK]
  - Mettre coreA dans les statistiques
  - Mettre à jour ou prévoir de le faire certains vieux équipements [OK]
    - ccnum, ccvialit, ccprn-bif, WS-C3524-XL à remplacer si possible par des 2960 : cccata-cccata-sg02
  - Avoir un modèle de configuration pour ccprn-coreA
  - Voir comment on se partage les tâches
- Pendant l'arrêt:
  - 19h-19h15:
    - Relead manuel des équipements en salle telecom: ccprn-inter, ccprn-opn
    - Changer cccata-sg01 (connecté à ccprn-inter) et sg02 (boite vers telecom & étage en d

~60 new fibres  
> 1km of fibre deployed!

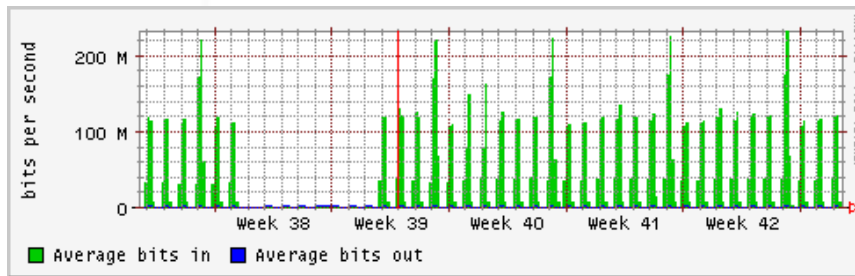


# Timeline



# Feedbacks

- We get no major surprise!
  - Heavy testing phase was fruitful
  - Main issue: Some routes not correctly announced
    - Not detected nor understood, but workaround found



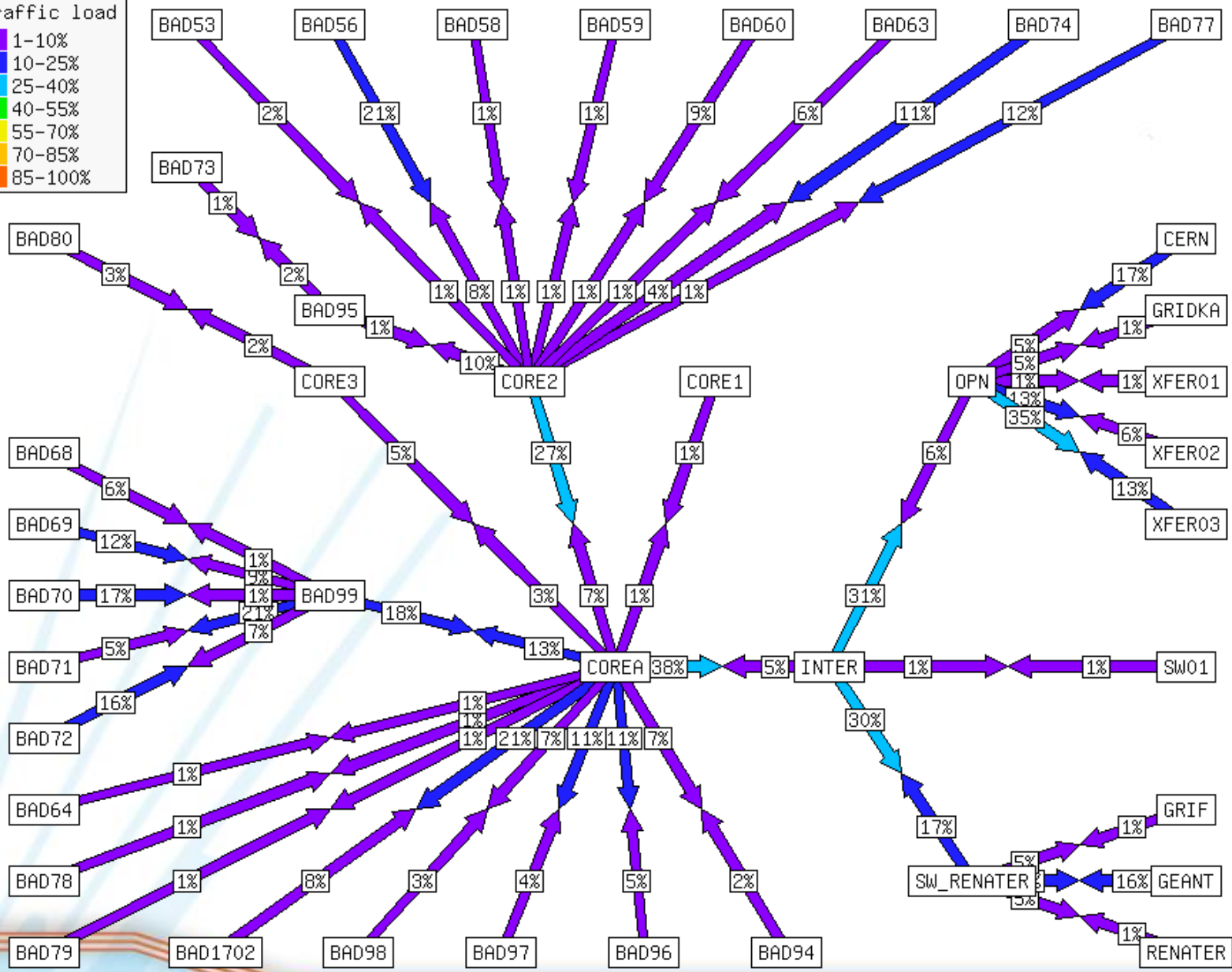
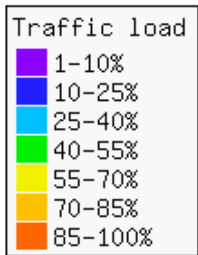
- Routing resilience hiding such problem
- Snapshot routing tables with a traceroute to each route!

- Keep monitoring, but deactivate alarms
  - Spare 800 SMS and 6k e-mails to each team members
- Do not parallelize actions too much
  - Hard to isolate faults or validate actions

## Key benefits

- Increased core capacity by 2.5
- Isolated areas, delivering wire speed, removed bottlenecks, shortened paths
- Seamless capacity upgrade now possible
- Harmonised softwares and features on 170 devices
  - From 37 different versions to 13

# Current status



# Upcoming

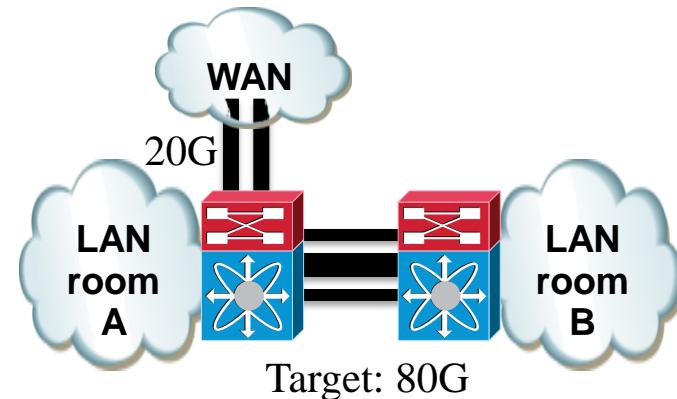
- **100G tests**

*This part was for eyes only*

- **Plug new computing room**

- Switch as far as you can

- 80G trunk to a similar architecture
- Not autonomous, this is an extension



# Conclusion

- Flat to starred network architecture
  - Closely matching our needs
- Average network usage down from 60% to ~15%
- Ready to face traffic increase for some more time
  - How long?