# Networks for ATLAS Trigger and Data Acquisition

S. Stancu*†, C. Meirosu†‡ , M. Ciobotaru*†, L. Leahu†‡, B. Martin†

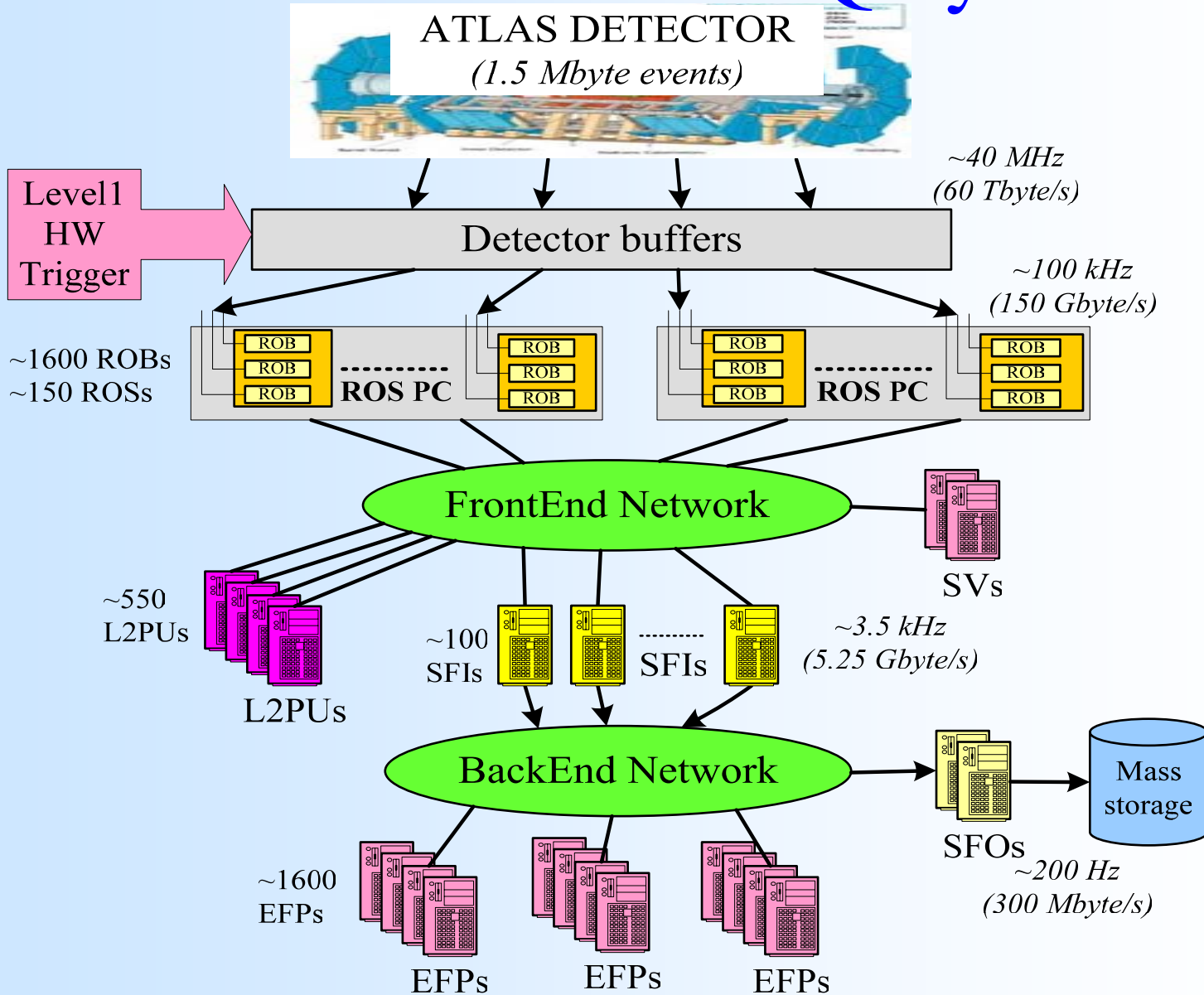\* University of California, Irvine
† CERN, Geneva
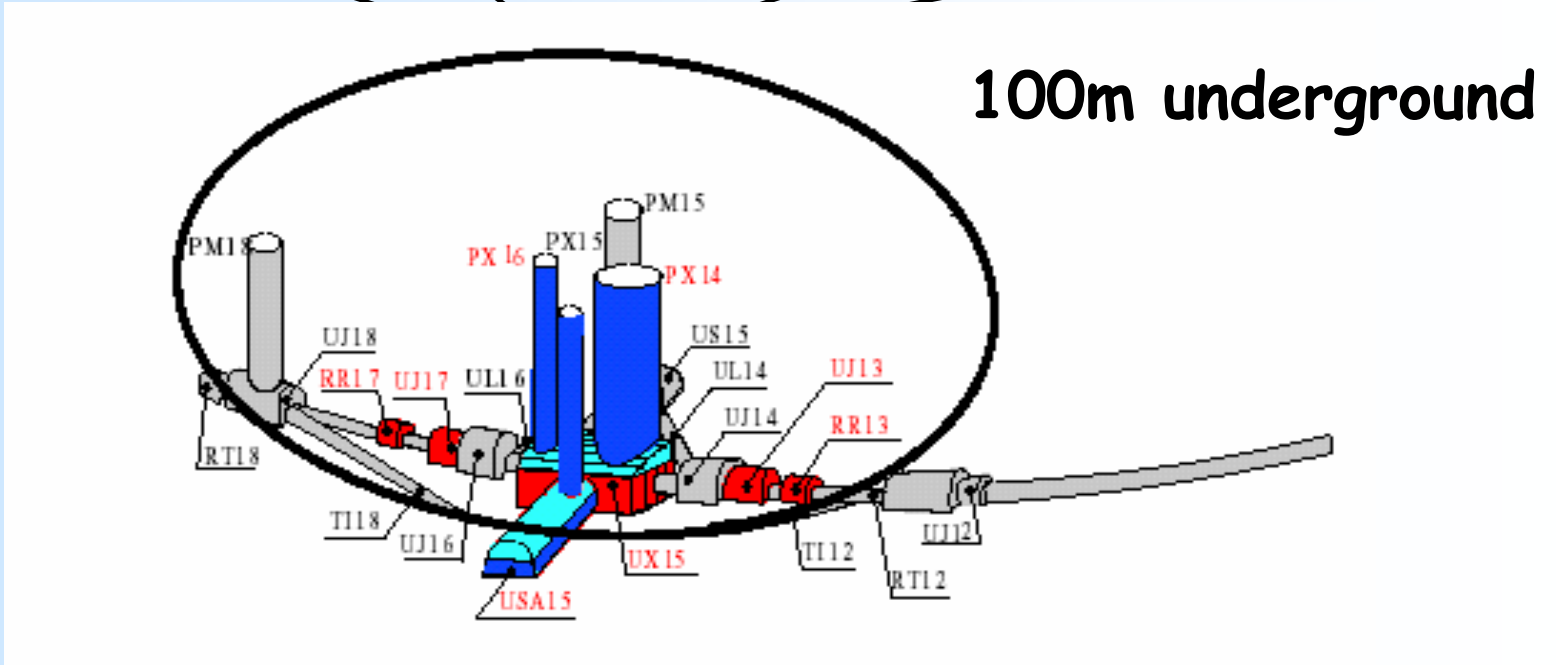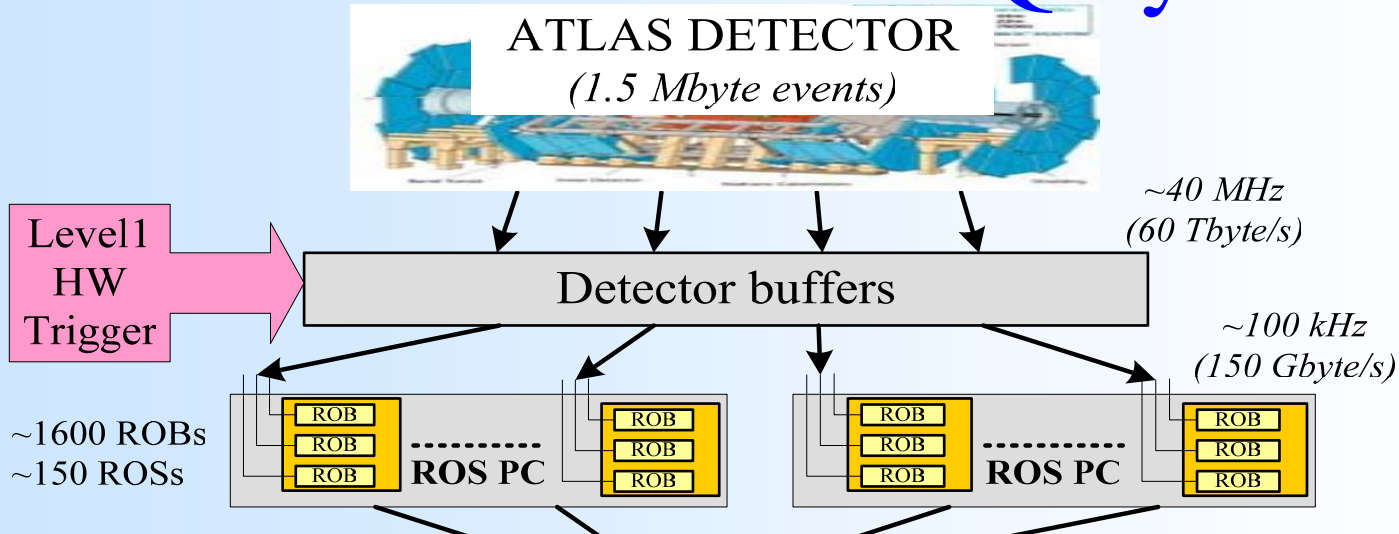‡ "Politehnica" University of Bucharest

# Outline

- Overview of the TDAQ system and networks
- Technology and equipment
- TDAQ networks:
  - ☆ Control network - no special bandwidth requirement
  - ☆ Dedicated data networks:
    - ✹ FrontEnd network – high bandwidth (~100Gbit/s cross-sectional bw.) and minimal loss
    - ✹ BackEnd  network – high bandwidth (~ 50Gbit/s cross-sectional bw.)
- Sample resiliency test
- Management/installation issues
  - ☆ Dedicated path for management and monitoring
  - ☆ Automatic topology/connectivity check
- Conclusions

# The ATLAS TDAQ System



ATLAS DETECTOR
*(1.5 Mbyte events)*

Level1 HW Trigger

Detector buffers

~40 MHz
*(60 Tbyte/s)*

~100 kHz
*(150 Gbyte/s)*

~1600 ROBs
~150 ROSs

ROB
ROB
ROB
**ROS PC**

ROB
ROB
ROB

ROB
ROB
ROB
**ROS PC**

ROB
ROB
ROB

FrontEnd Network

SVs

~550
L2PUs

L2PUs

~100
SFIs

SFIs

~3.5 kHz
*(5.25 Gbyte/s)*

BackEnd Network

Mass storage

SFOs

~1600
EFPs

EFPs

EFPs

EFPs

~200 Hz
*(300 Mbyte/s)*

# The ATLAS TDAQ System

ATLAS DETECTOR
*(1.5 Mbyte events)*

~40 MHz
*(60 Tbyte/s)*

Level1
HW
Trigger

Detector buffers

~100 kHz
*(150 Gbyte/s)*

~1600 ROBs
~150 ROSs

| ROB |
| ROB |
| ROB |

ROS PC

| ROB |
| ROB |
| ROB |

| ROB |
| ROB |
| ROB |

ROS PC

| ROB |
| ROB |
| ROB |

**100m underground**

PM15
PM18    PX16    PX15    PX14
UJ18    US15
RR17  UJ17  UL16    UL14    UJ13
RT18    UJ14    RR13
TI18  UJ16    UJ12
UX15    TI12    RT12
USA15

# The ATLAS TDAQ System

**Surface buildings**

FrontEnd Network

~550
L2PUs

L2PUs

~100
SFIs

SFIs

SVs

~3.5 kHz
(5.25 Gbyte/s)

BackEnd Network

~1600
EFPs

EFPs    EFPs    EFPs

SFOs

Mass
storage

~200 Hz
(300 Mbyte/s)

# The ATLAS TDAQ System

ATLAS DETECTOR
*(1.5 Mbyte events)*

~40 MHz
*(60 Tbyte/s)*

Level1 HW Trigger

Detector buffers

~100 kHz
*(150 Gbyte/s)*

~1600 ROBs
~150 ROSs

ROB ROB ROB **ROS PC**

ROB ROB ROB **ROS PC**

FrontEnd Network

SVs

~550 L2PUs

L2PUs

~100 SFIs

SFIs

~3.5 kHz
*(5.25 Gbyte/s)*

BackEnd Network

Mass storage

SFOs

~200 Hz
*(300 Mbyte/s)*

~1600 EFPs

EFPs          EFPs          EFPs

# The ATLAS TDAQ System

**ATLAS DETECTOR**
*(1.5 Mbyte events)*

**Level1 HW Trigger**

**Detector buffers**

~40 MHz
(60 Tbyte/s)

~100 kHz
(150 Gbyte/s)

~1600 ROBs
~150 ROSs

ROB
ROB
ROB
**ROS PC**

ROB
ROB
ROB
**ROS PC**

ROB
ROB
ROB

**FrontEnd Net**

~550
L2PUs

L2PUs

~100
SFIs

SFI

**BackEnd Network**

~1600
EFPs

EFPs      EFPs      EFPs

SFOs

~200 Hz
(300 Mbyte/s)

Mass storage

**Control network**

- Infrastructure services
  - ☆ Network (DHCP, DNS)
  - ☆ Shared file systems
  - ☆ Databases, monitoring, information service, etc.
- Run Control

# The ATLAS TDAQ System

ATLAS DETECTOR
*(1.5 Mbyte events)*

Level1
HW
Trigger

**High availability**
24/7 when the accelerator is running

*Detector buffers*

*40 MHz*
*(60 Tbyte/s)*

*~100 kHz*
*(150 Gbyte/s)*

~1600 ROBs
~150 ROSs

ROB
ROB
ROB
**ROS PC**

ROB
ROB
ROB

ROB
ROB
**ROS PC**

ROB
ROB
ROB

FrontEnd Network

SVs

~550
L2PUs

L2PUs

~100
SFIs    SFIs    SFIs

*~3.5 kHz*
*(5.25 Gbyte/s)*

BackEnd Network

Mass
storage

SFOs
*~200 Hz*
*(300 Mbyte/s)*

~1600
EFPs

EFPs       EFPs       EFPs

# Technology and equipment

- *Ethernet* is the dominant technology for LANs
  - ☆ TDAQ's choice for networks (see [1])
    - ✸ multi-vendor, long term support, commodity (on-board GE adapters), etc.
  - ☆ Gigabit and TenGigabit Ethernet
    - ✸ Use GE for end-nodes
    - ✸ 10GE whenever the bandwidth requirements exceed 1Gbit/s
- Multi-vendor Ethernet switches/routers available on the market:
  - ☆ Chassis-based devices ( ~320 Gbit/s switching)
    - ✸ GE line-cards: typically ~40 ports (1000BaseT)
    - ✸ 10GE line-cards: typically 4 ports (10GBaseSR)
  - ☆ Pizza-box devices (~60 Gbit/s switching)
    - ✸ 24/48 GE ports (1000BaseT)
    - ✸ Optional 10GE module with 2 up-links (10GBaseSR)

# Resilient Ethernet networks

- What happens if a switch or link fails?
  - ☆ Phone call, but nothing critical should happen after a single failure.

- Networks are made resilient by introducing redundancy:
  - ☆ *Component-level redundancy*: deployment of devices with built-in redundancy (PSU, supervision modules, switching fabric)
  - ☆ *Network-level redundancy*: deployment of additional devices/links in order to provide alternate paths between communicating nodes.
    - ✷ Protocols are needed to *correctly* (and *efficiently*) deal with multiple paths in the network [2]:
      - ○ Layer 2 protocols: Link aggregation (trunking), spanning trees (STP, RSTP, MSTP)
      - ○ Layer 3 protocols: virtual router redundancy (VRRP) for static environments, dynamic routing protocols (e.g. RIP, OSPF).

# Control network

- ~3000 end nodes

- Design assumption: the instantaneous traffic does not exceed 1 Gbit/s on any segment, including up-link.

- One device suffices for the core layer, but better redundancy is achieved by deploying 2 devices.

- A rack level concentration switch can be deployed for all units except for critical services.

- Layer 3 routed network
  - ☆ One sub-net per concentrator switch
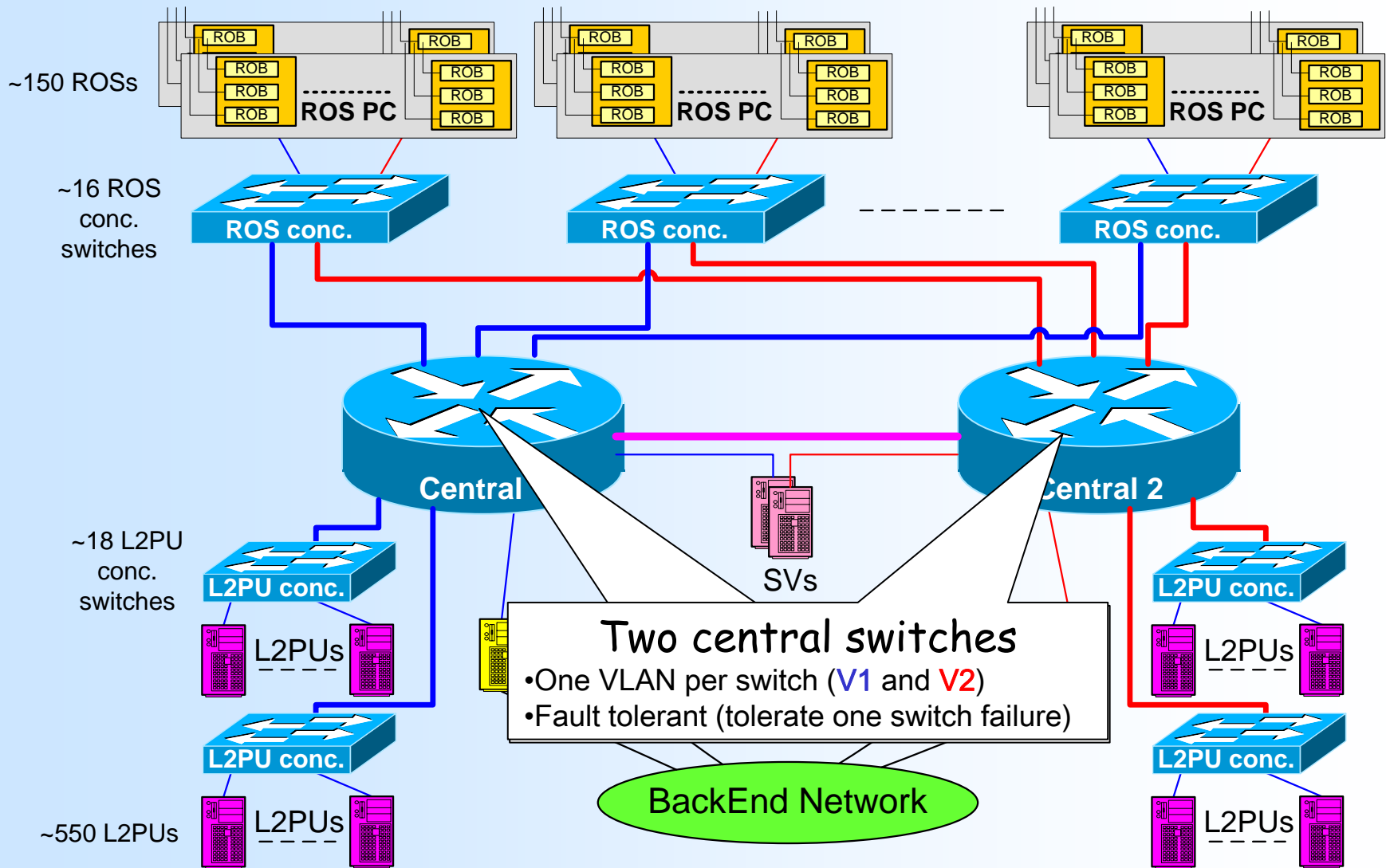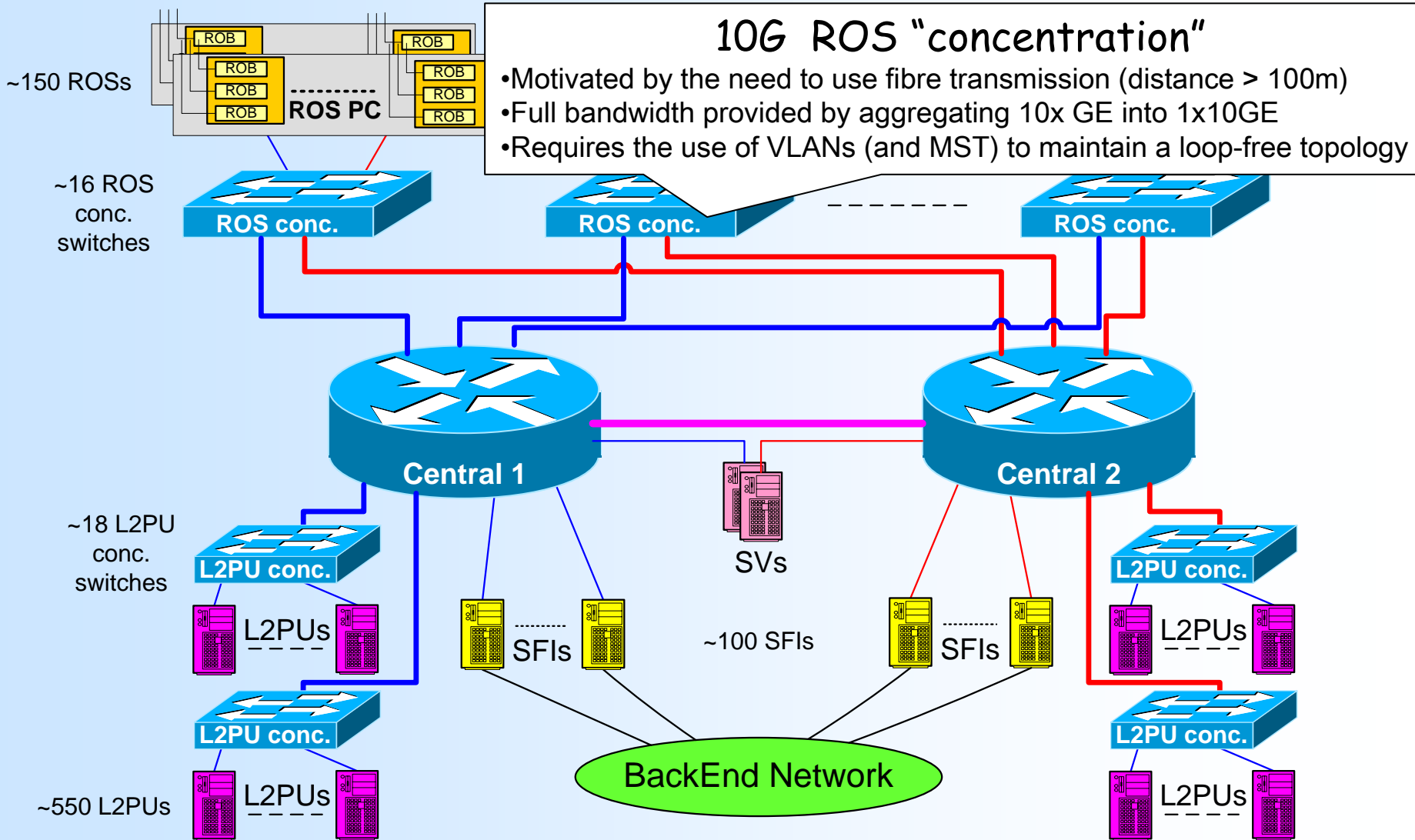  - ☆ Small broadcast domains → potential layer 2 problems remain local.
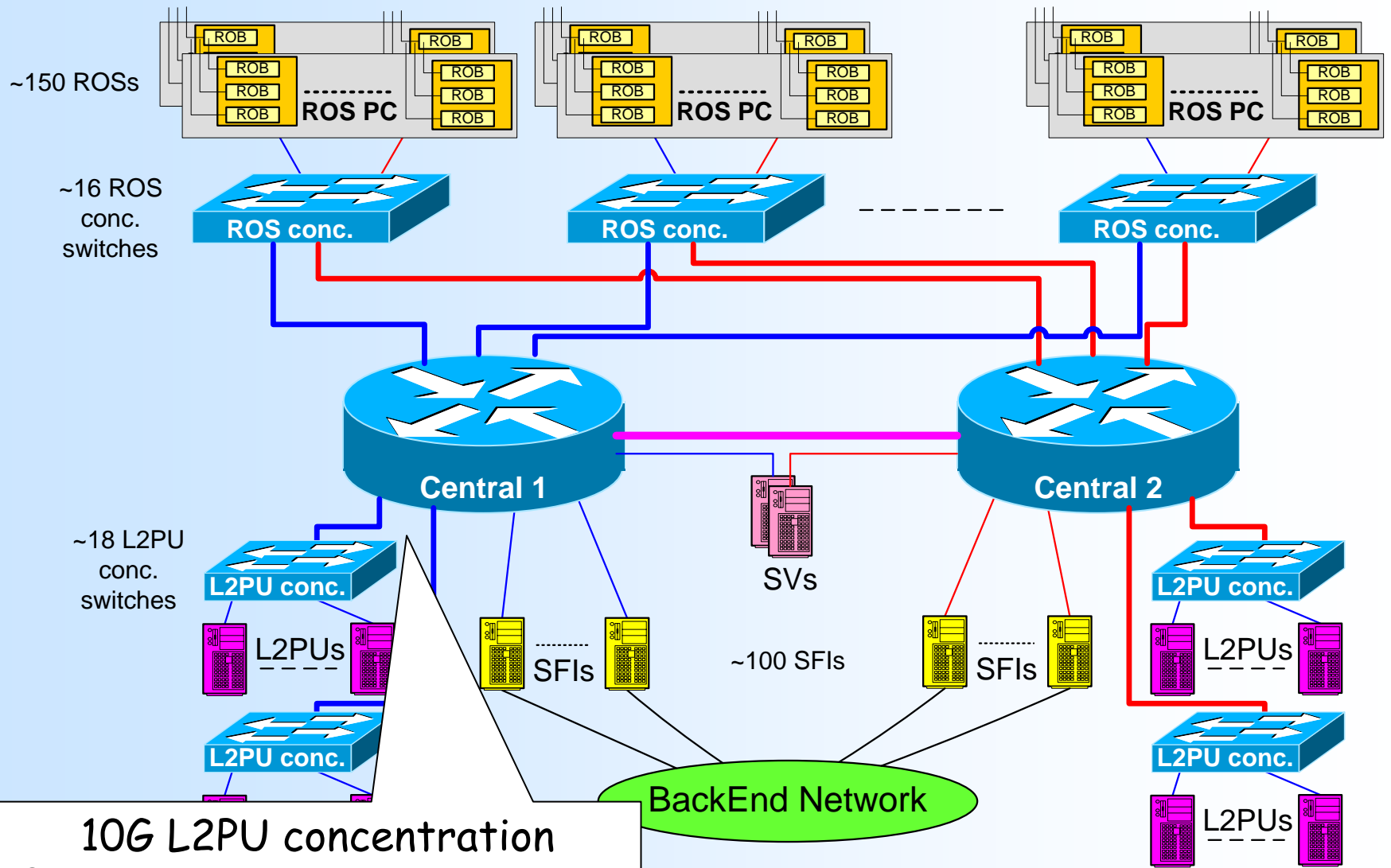


Infrastructure and TDAQ services

Core 1    Core 2

Conc.    Conc.    Conc.

~550 L2PUs    SVs    ~100 SFIs    SFOs    ~1600 EFPs

# FrontEnd network (see [3])



~150 ROSs

ROB · ROB · ROB · ROB · ROB · ROB · ROS PC

~16 ROS conc. switches

ROS conc.

Central 1

Central 2

SVs

~18 L2PU conc. switches

L2PU conc.

L2PUs

L2PU conc.

L2PUs

~550 L2PUs

SFIs

~100 SFIs

SFIs

BackEnd Network

# FrontEnd network



~150 ROSs

~16 ROS conc. switches

ROS PC

ROS conc.

~100 Gbit/s

Central 1

Central 2

SVs

~18 L2PU conc. switches

L2PU conc.

L2PUs

SFIs

~100 SFIs

SFIs

L2PU conc.

L2PUs

L2PU conc.

~550 L2PUs

L2PUs

BackEnd Network

L2PUs

# FrontEnd network



~150 ROSs

~16 ROS conc. switches

ROB · ROB · ROS PC

ROS conc.

~18 L2PU conc. switches

L2PU conc.

L2PUs

L2PU conc.

~550 L2PUs

L2PUs

Central

Central 2

SVs

**Two central switches**
- One VLAN per switch (V1 and V2)
- Fault tolerant (tolerate one switch failure)

L2PU conc.

L2PUs

L2PU conc.

L2PUs

BackEnd Network

# FrontEnd network



~150 ROSs

**ROB** ... **ROB**
**ROB** ... **ROB**
**ROB** ... **ROB**
**ROB** **ROS PC** **ROB**

### 10G ROS "concentration"
- Motivated by the need to use fibre transmission (distance > 100m)
- Full bandwidth provided by aggregating 10x GE into 1x10GE
- Requires the use of VLANs (and MST) to maintain a loop-free topology

~16 ROS conc. switches

**ROS conc.** **ROS conc.** **ROS conc.**

**Central 1** **Central 2**

SVs

~18 L2PU conc. switches

**L2PU conc.** **L2PU conc.**

L2PUs SFIs ... SFIs ~100 SFIs SFIs ... SFIs L2PUs

**L2PU conc.** **L2PU conc.**

~550 L2PUs L2PUs

BackEnd Network

L2PUs

# FrontEnd network



~150 ROSs

**ROB** **ROB** ... **ROS PC**

~16 ROS conc. switches

**ROS conc.** **ROS conc.** **ROS conc.**

**Central 1** **Central 2**

SVs

~18 L2PU conc. switches

**L2PU conc.**

**L2PUs**

**L2PU conc.**

**SFIs** ~100 SFIs **SFIs**

**L2PU conc.**

**L2PUs**

**L2PU conc.**

**L2PUs**

BackEnd Network

## 10G L2PU concentration
• One switch per rack

# BackEnd network

- ~2000 end-nodes
- One core device with built-in redundancy
- Rack level concentration with use of link aggregation for redundant up-links to the core.
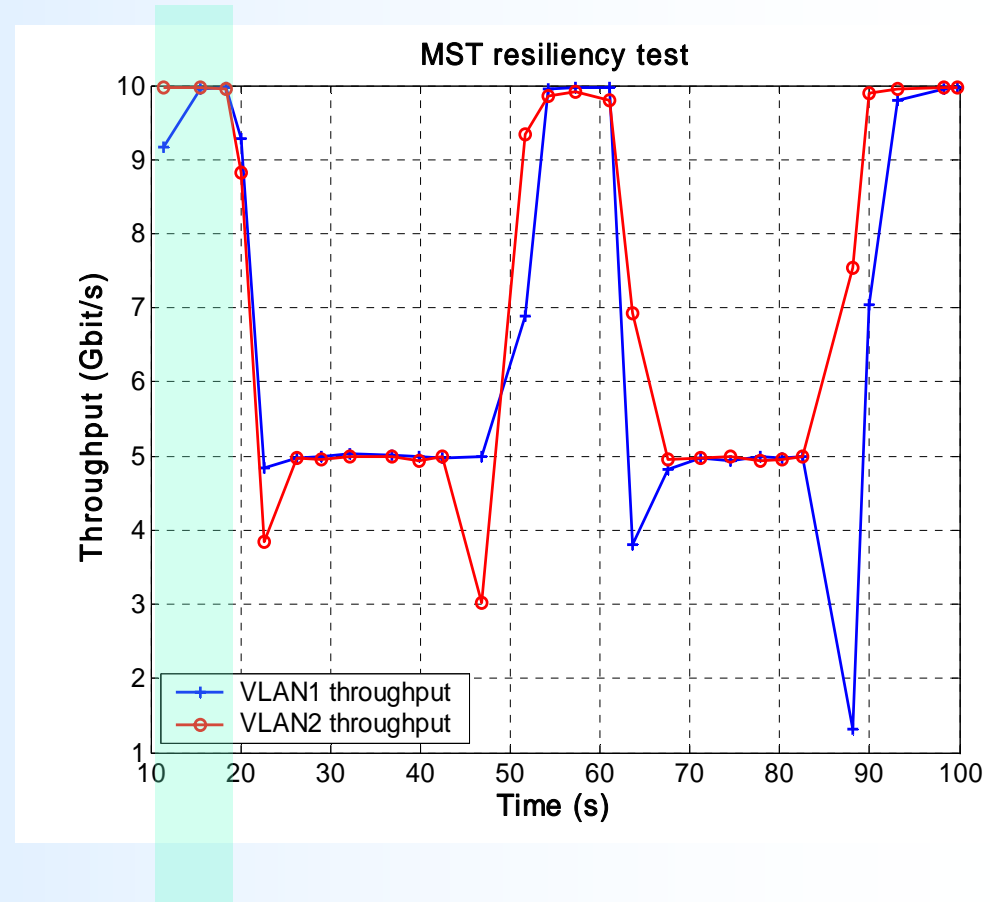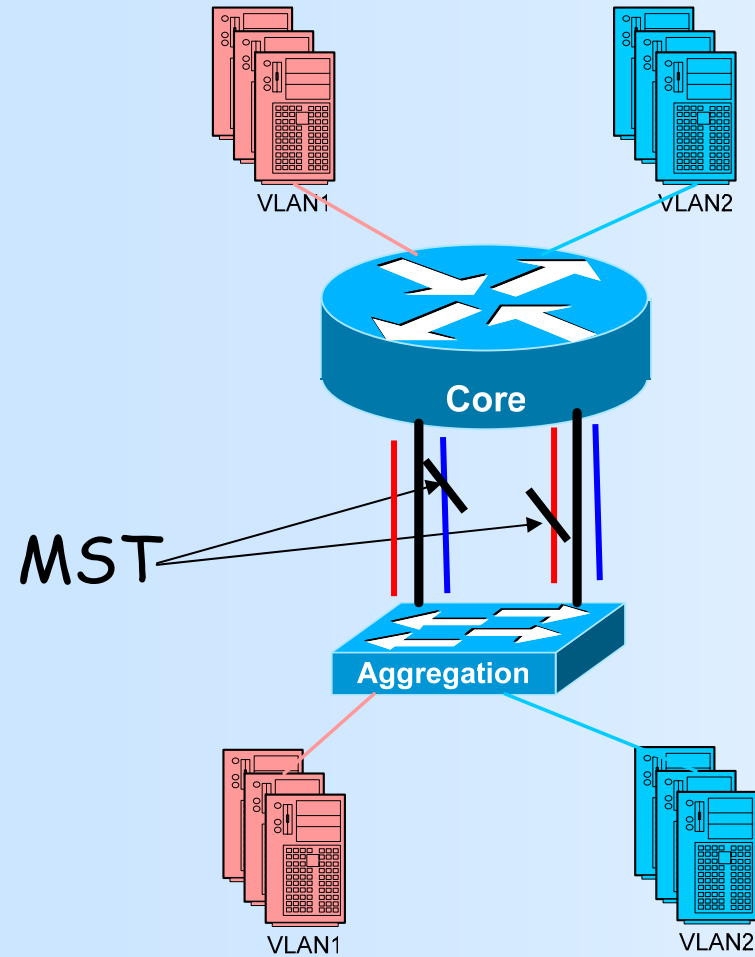- Layer 3 routed network to restrict broadcast domains size.

FrontEnd Network

~100 SFIs    SFIs
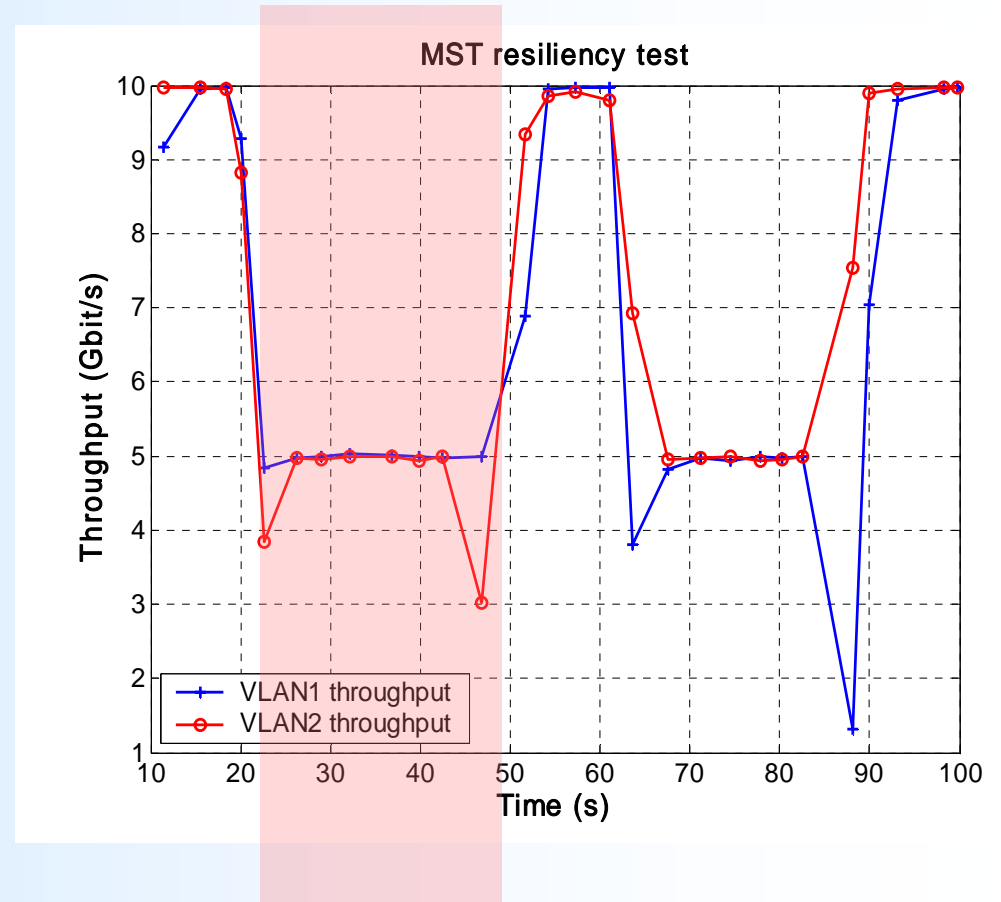
~30 SFOs

BE Central

SFOs

SFO conc.

Mass storage

~60 EF conc. switches

EF conc.    EF conc.    EF conc.

EFs    EFs    EFs

# BackEnd network



- ~2000 end-nodes
- One core device with built-in redundancy
- Rack level concentration with use of link aggregation for redundant up-links to the core.
- Layer 3 routed network to restrict broadcast domains size.

FrontEnd Network

~100 SFIs

SFIs

~30 SFOs

~50 Gbit/s

BE Central

SFOs

SFO conc.

Mass storage

~60 EF conc. switches

EF conc.

EF conc.

EF conc.

EFs

EFs

EFs

# BackEnd network

- ~2000 end-nodes
- One core device with built-in redundancy
- Rack level concentration with use of link aggregation for redundant up-links to the core.
- Layer 3 routed network to restrict broadcast domains size.



FrontEnd Network

~100 SFIs
SFIs

~30 SFOs

SFOs

~2.5 Gbit/s

BE Central

SFO conc.

Mass storage

~60 EF conc. switches

EF conc.

EF conc.

EF conc.

EFs

EFs

EFs

# BackEnd network



- ~2000 end-nodes
- One core device with built-in redundancy
- Rack level concentration with use of link aggregation for redundant up-links to the core.
- Layer 3 routed network to restrict broadcast domains size.

FrontEnd Network

~100 SFIs

SFIs

~30 SFOs

BE Central

SFOs

SFO conc.

Mass storage

~60 EF conc. switches

EF conc.

EF conc.

EF conc.

EFs

EFs

EFs

# Interchangeable processing power

- Standard processor rack with up-links to both FrontEnd and BackEnd networks.
- The processing power migration between L2 and EF is achieved by software enabling/disabling of the appropriate up-links.
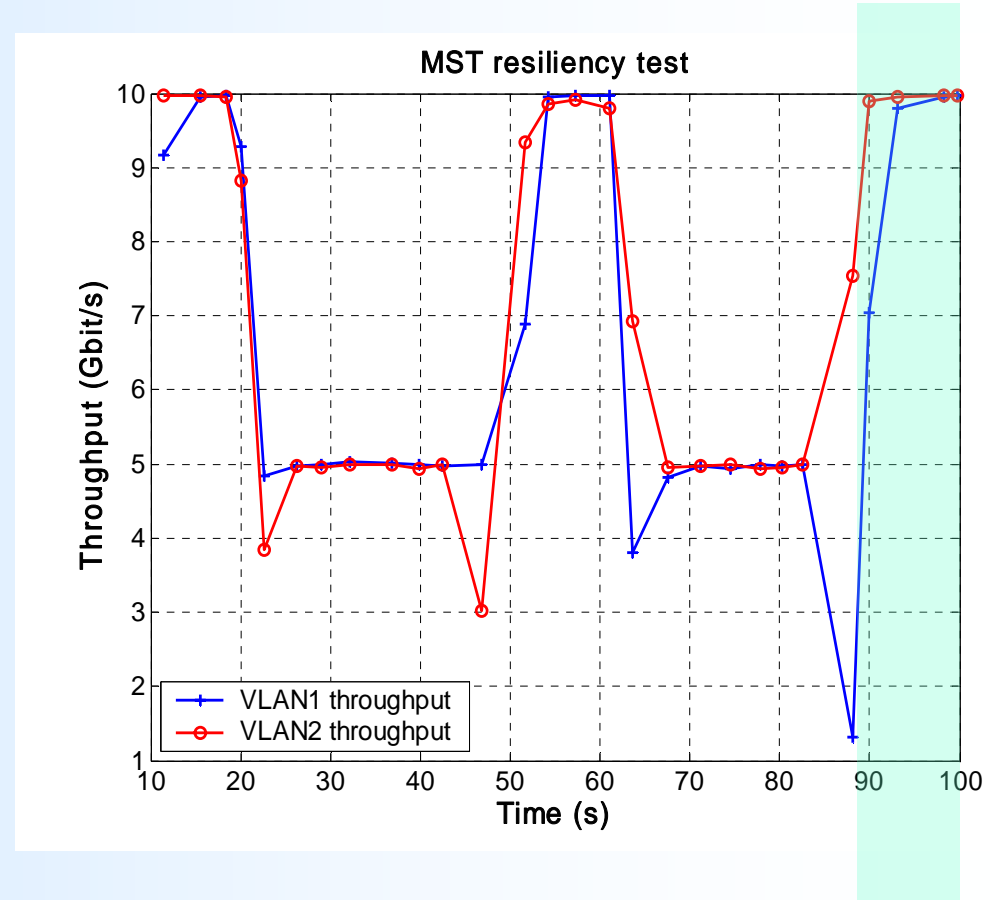
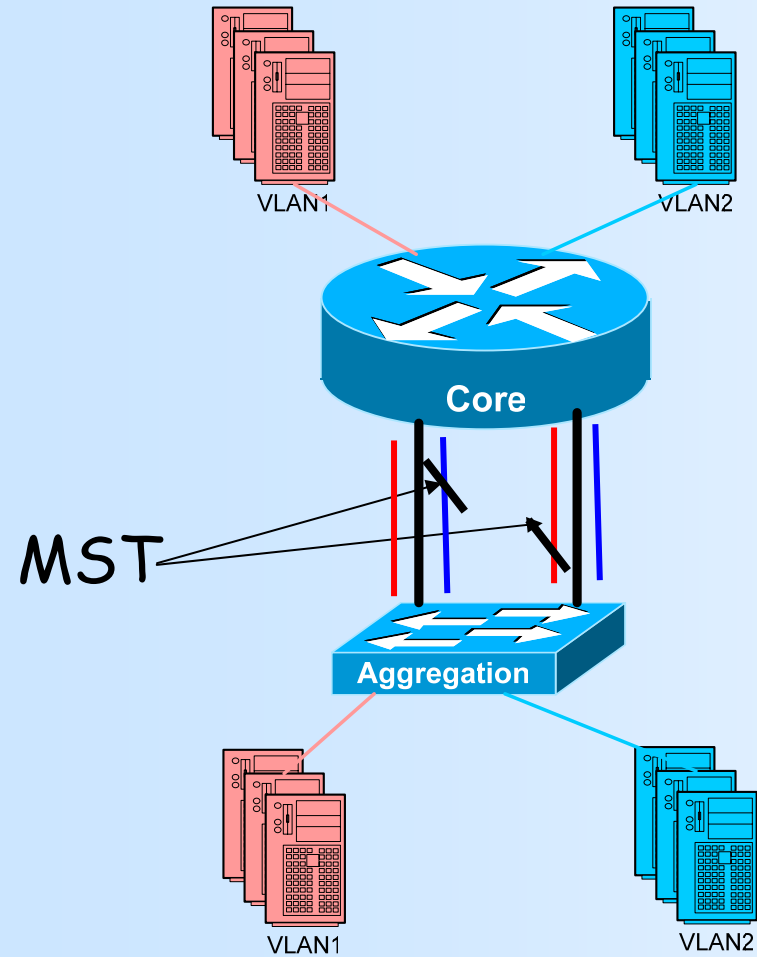# Sample resiliency test (see [4])

# Sample resiliency test

# Sample resiliency test
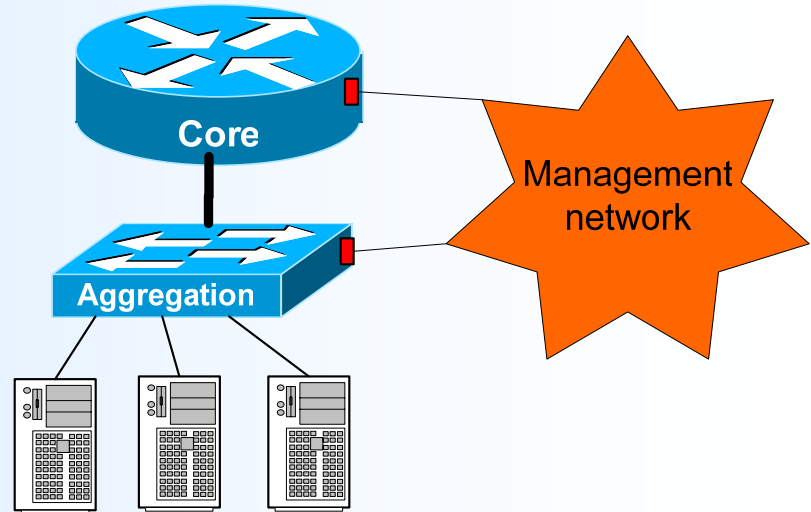
# Sample resiliency test

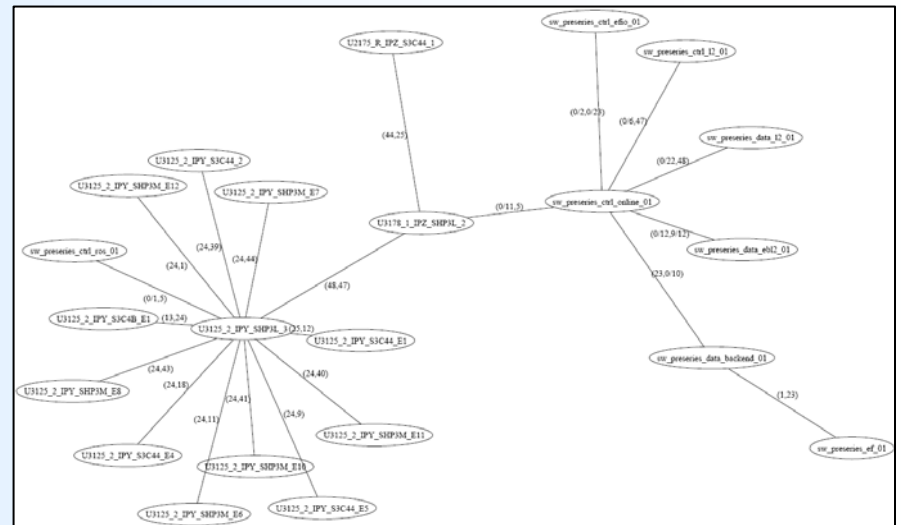# Sample resiliency test

# Installation/management issues



- Dedicated path for management
  - ☆ Each device will have an "out of band" interface dedicated for management.
  - ☆ A small layer 2 network will be used to connect to the "out of band" interfaces of devices

- Automatic topology discovery/check
  - ☆ Maintaining accurate active cabling information in the installation database is tedious
  - ☆ Developed a tool which constructs the network topology based on the MAC address table information
  - ☆ To do: interface the tool with the installation database.

# Conclusions

- The ATLAS TDAQ system (approx. 3000 end-nodes) relies on networks for both control and data acquisition purposes.
- Ethernet technology (+IP)
- Networks architecture maps on multi-vendor devices
- Modular network design
- Resilient network design (high availability)
- Separate management path
- Developing tools for automatic population/cross-checks of installation data-bases.
- Network operation → see Catalin Meirosu's talk [5].

# References

[1] S. Stancu, B. Dobinson, M. Ciobotaru, K. Korcyl, and E. Knezo, "*The use of Ethernet in the Dataflow of the ATLAS Trigger and DAQ*" in Proc. CHEP 06 Conference

[2] T. Sridhar, "*Redundancy: Choosing the Right Option for Net Designs*"

*http://www.commsdesign.com/showArticle.jhtml?articleID=25600515*

[3] S. Stancu, M. Ciobotaru, and K. Korcyl, "*ATLAS TDAQ DataFlow Network Architecture Analysis and Upgrade Proposal*" in Proc. IEEE Real Time Conference 2005

[4] Cisco whitepaper, "*Understanding Multiple Spanning Tree Protocol (802.1s)*" http://www.cisco.com/warp/public/473/147.pdf

[5] C. Meirosu, A. Topurov, A. Al-Shabibi, B. Martin, "*Planning for predictable network performance in the ATLAS TDAQ*", CHEP06 Mumbay, February 2006.