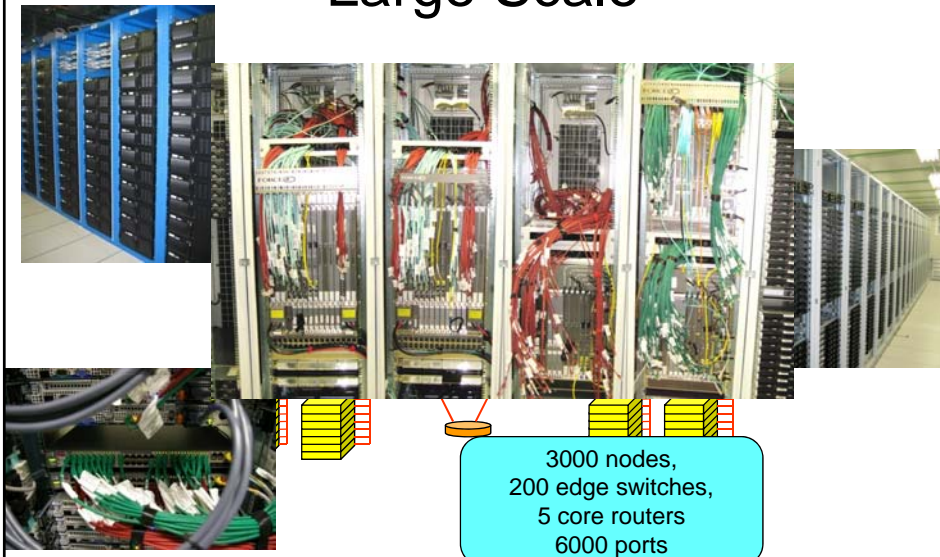


Monitoring the ATLAS TDAQ Network at CERN

Lucian LEAHU
Brasov, 15/01/2009

Large Scale



Plus physicists!

- Network dimensioned to meet 'requirements'
- Maximum average link occupancy = 60%
- **Should** mean peace of mind for Network Support
- **Actually** seen as a challenge by physicists
 - 40% for free! Turn up the wick until something breaks!
- Continuous running out of spec!
- Must distinguish between 'real' and 'self inflicted' problems

3

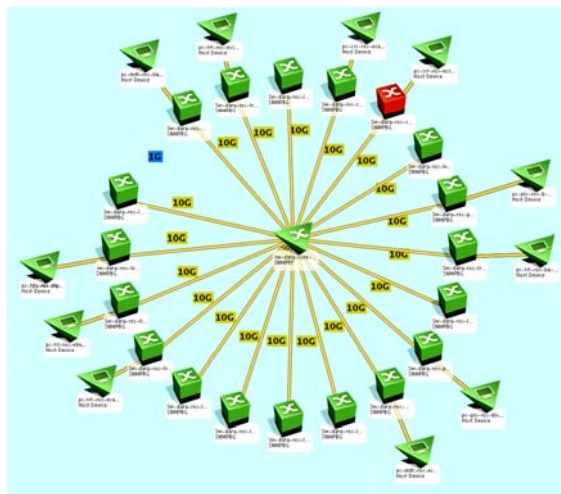
Commissioning

Basic question:
Does any of it work????
Monitor everything!

ICMP:
Internet Control Message
Protocol
SNMP:
Simple
Network Management
Protocol

SPECTRUM

Tells you it's alive
Tells you if it dies
Fetches status info..
- and hides it!



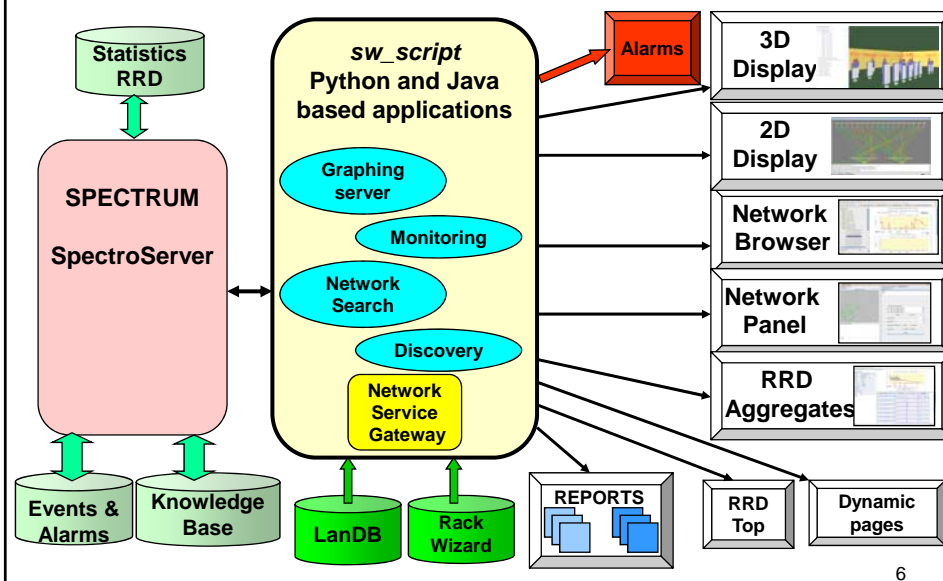
4

Commercial versus Proprietary

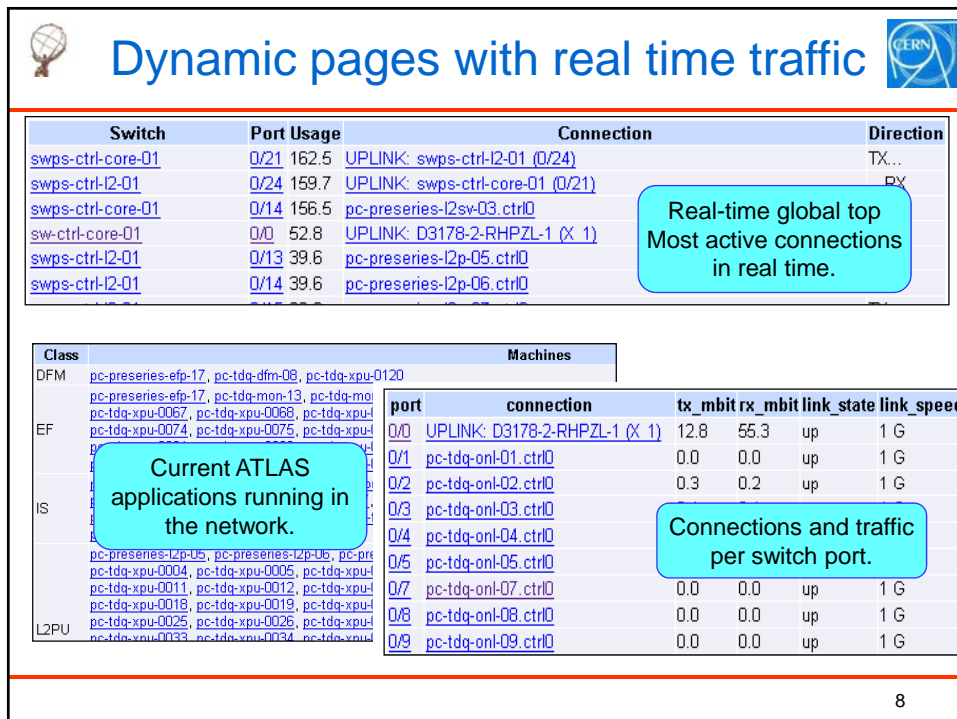
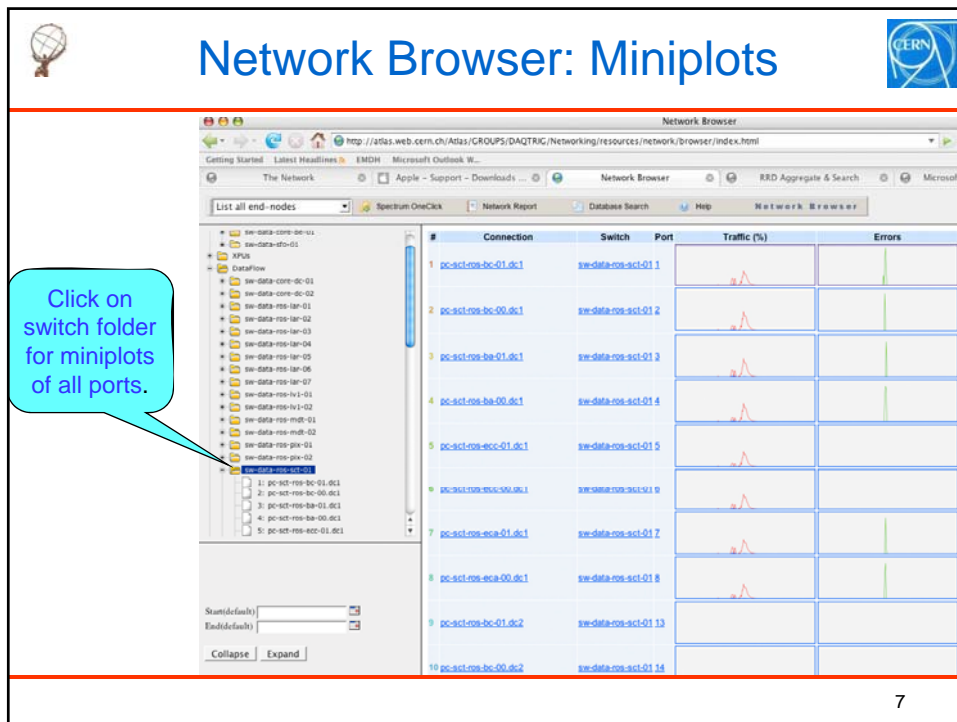
- Special needs:
 - Multiple networks per processor
 - Want to see whole picture for system analysis
 - Want to see all detail for component analysis
 - Want to see traffic volume visually
 - Want traffic/errors qualified
- Spectrum CORBA API clumsy
- Multiple requests hits the CPU hard

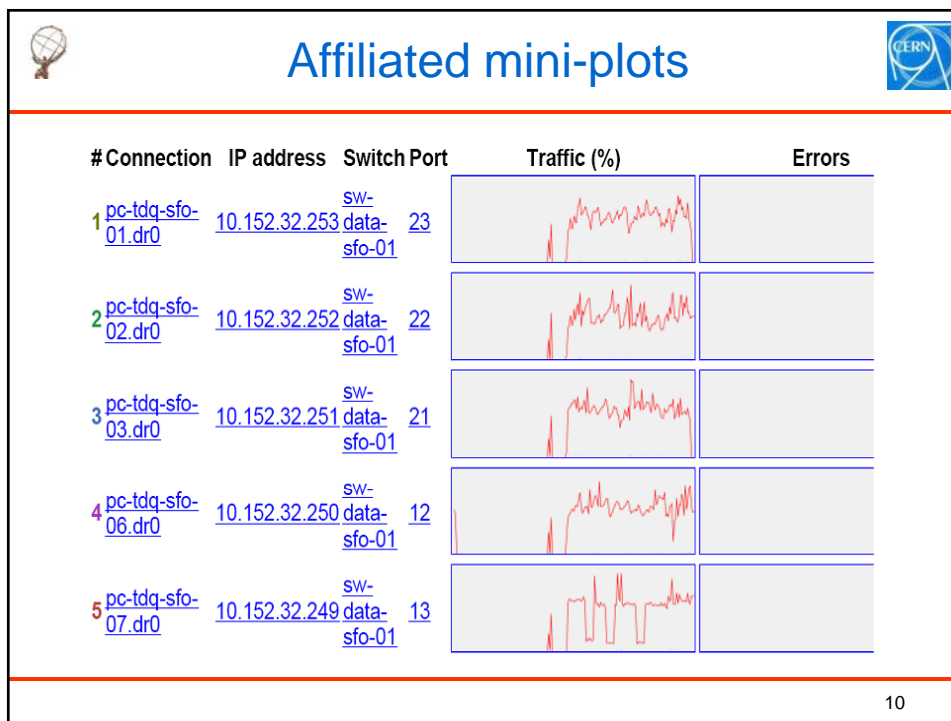
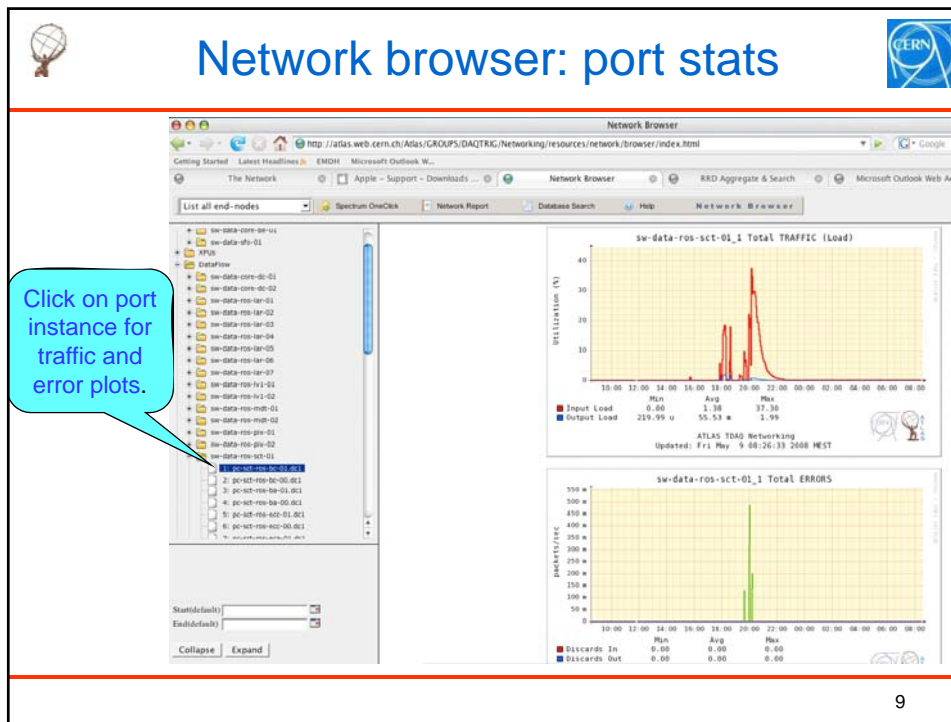
5

Current Infrastructure



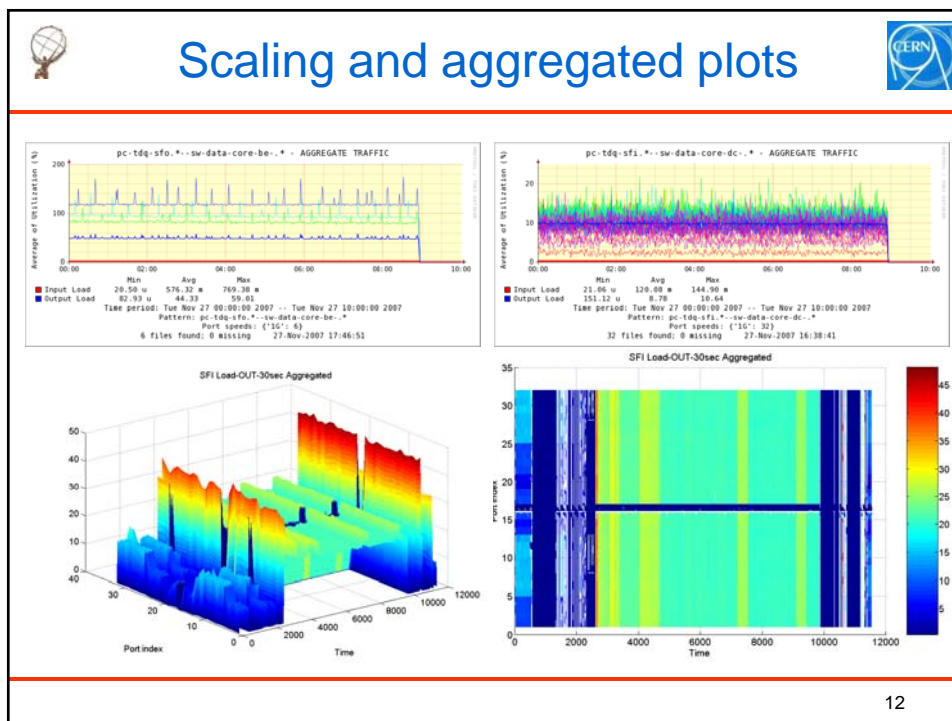
6







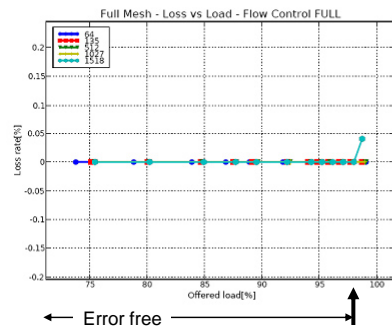
11



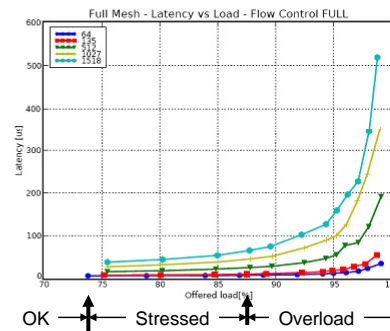
12



System Diagnostics



(a) Loss rate vs Load



(b) Latency vs Load

We have prior detailed measurements of losses and latencies against loads.
System dimensioned for a nominal operating maximum of 60% load on any port.

Naïve world view:

Flag Loads with high thresholds, above, say 85% throughput

Flag Errors with low thresholds, above, say, 0.05% or about 10pkts/sec at 1Gbps

Unfortunately it's not so simple

13

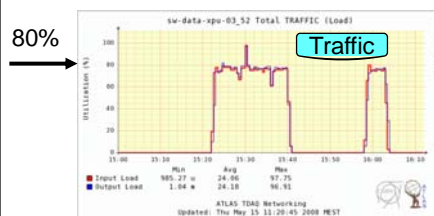


Where is the real error?

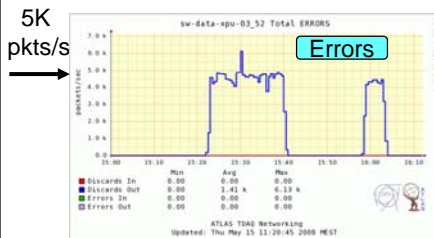


TCP

80%

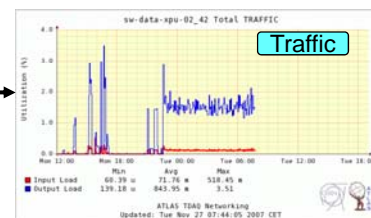


5K
pkts/s

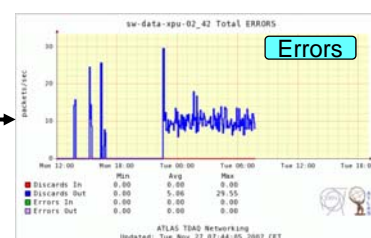


UDP

<2%



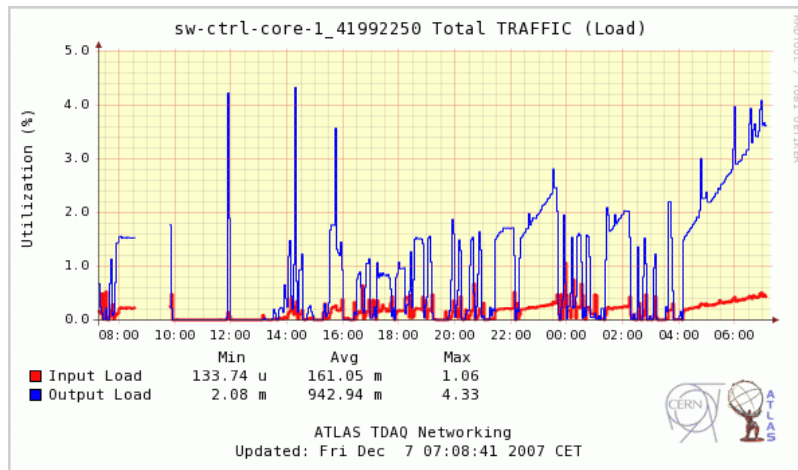
10
pkts/s



14



Anomalous traffic



...and from here to rules based expert systems

15

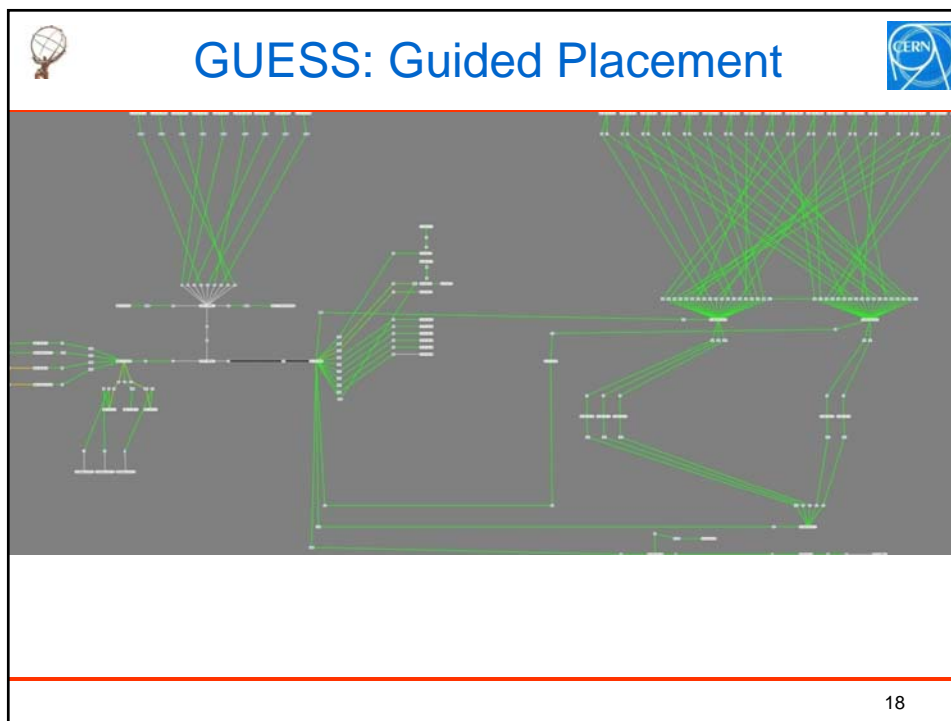
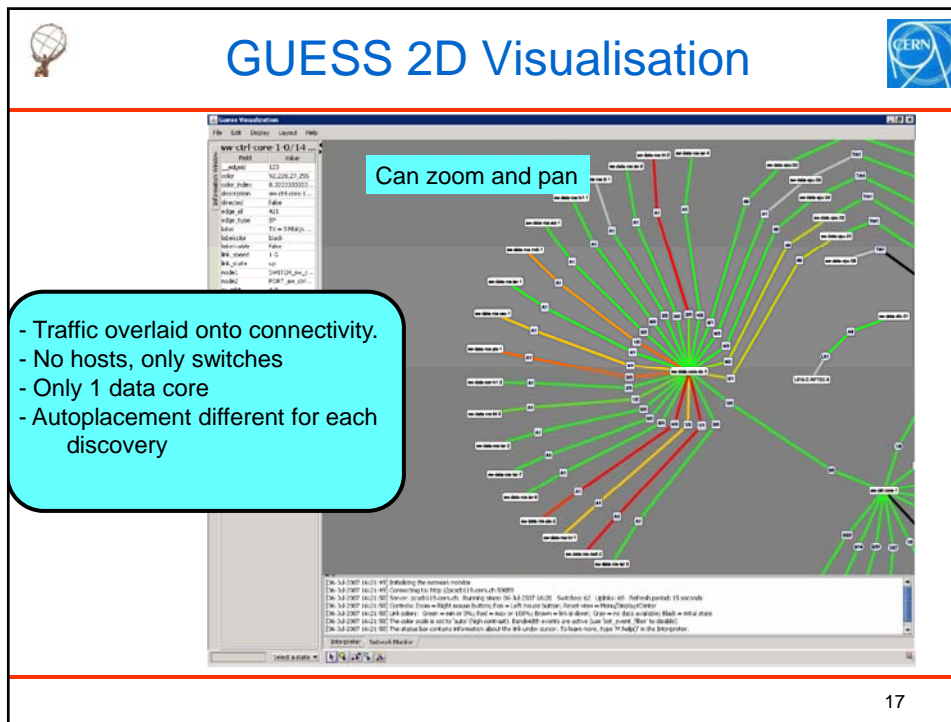


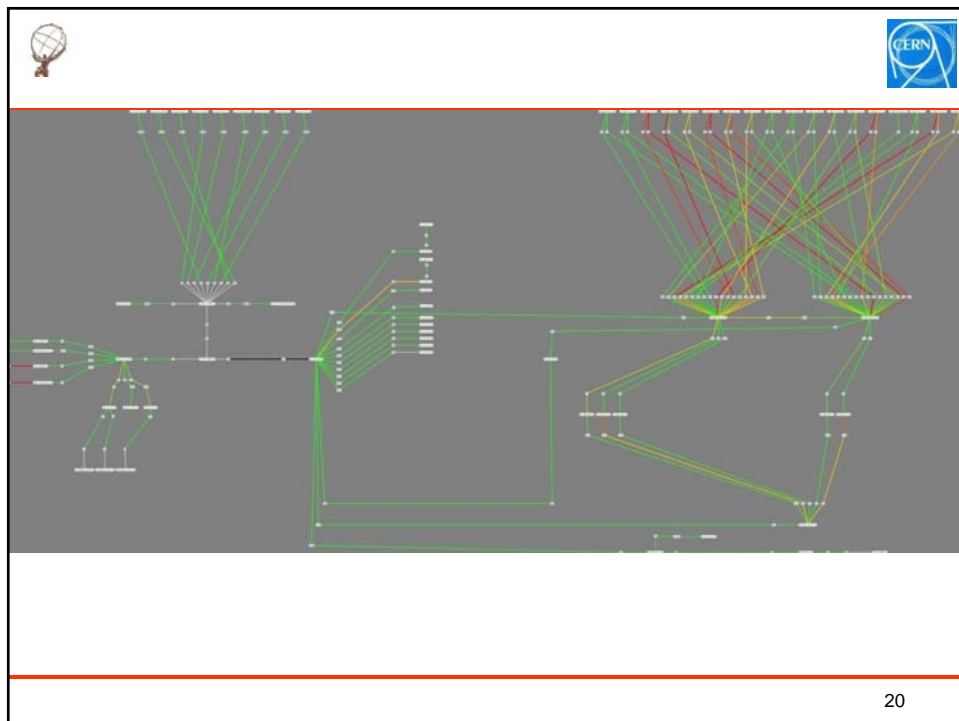
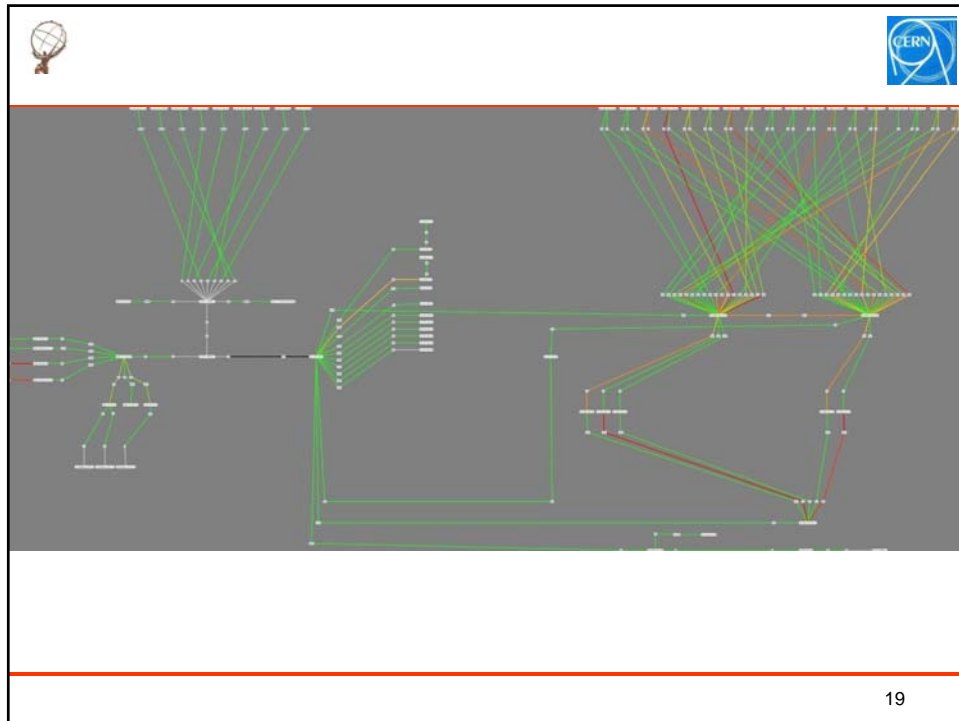
Transition to displays

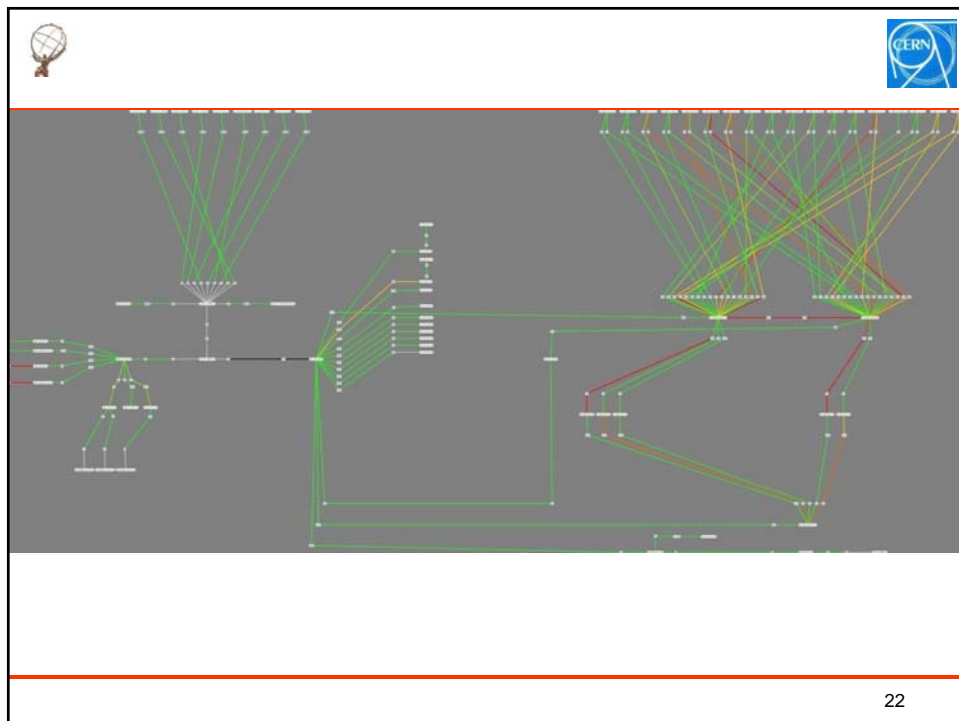
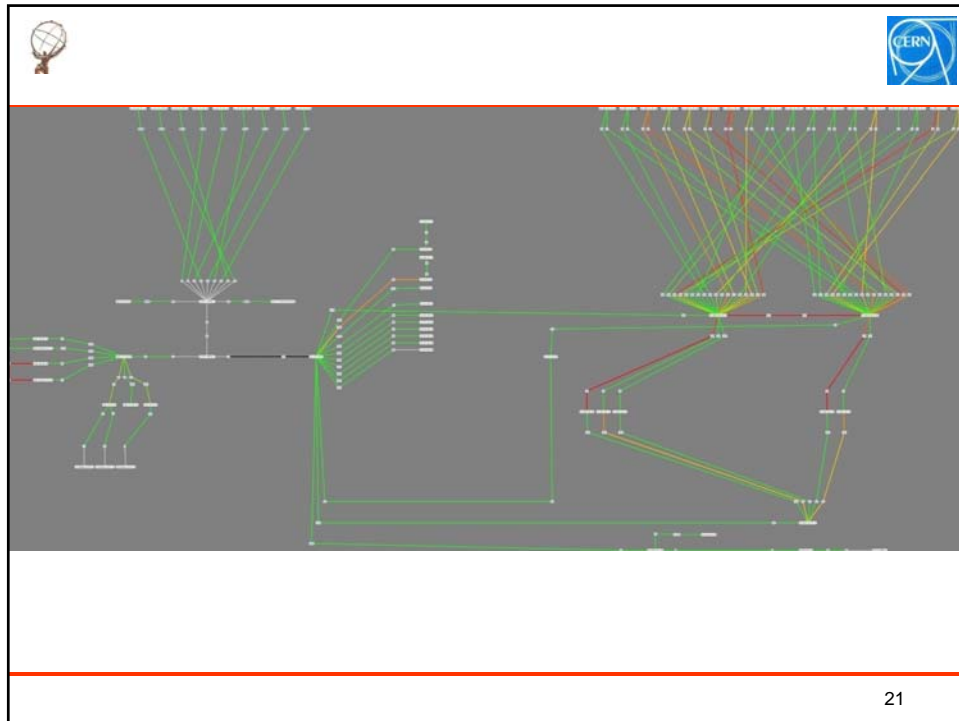


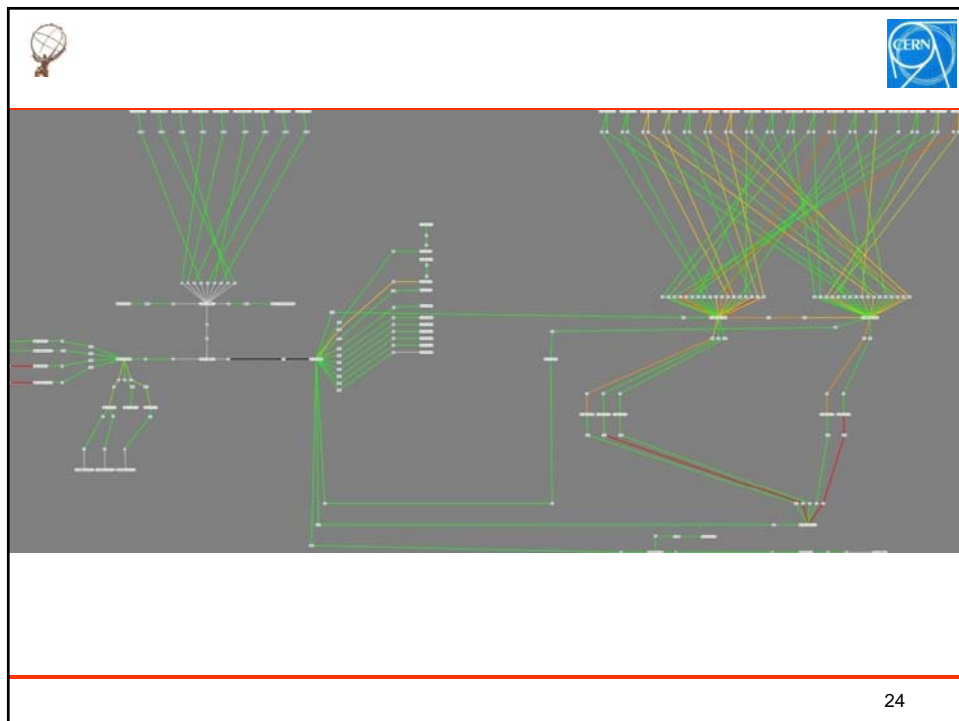
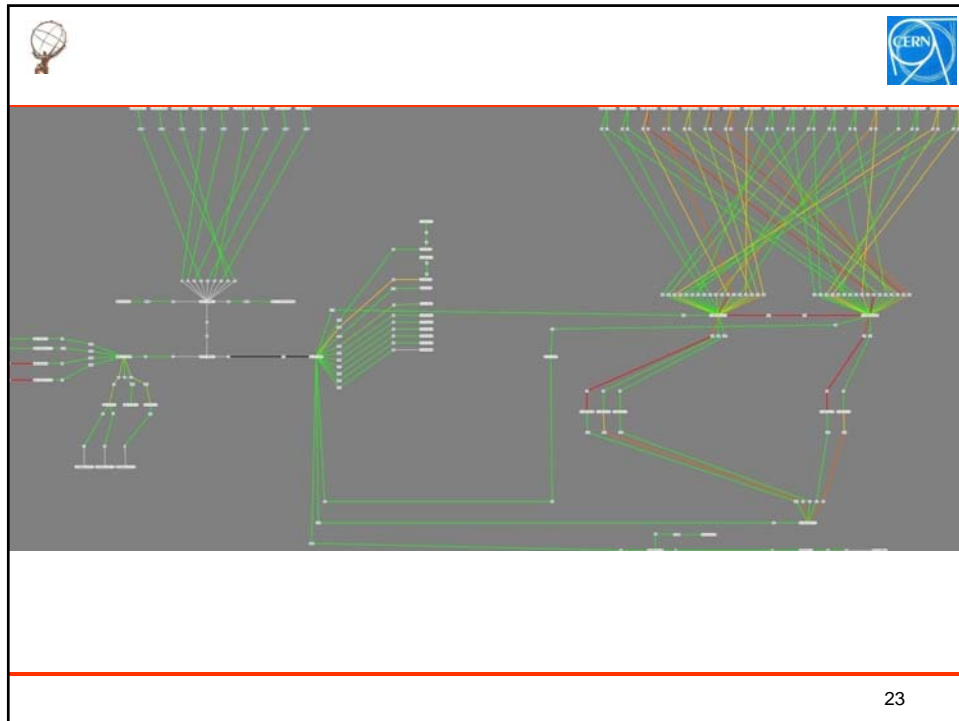
- So far have concentrated on low level port and switch monitoring and diagnostics.
- Identified scaling issues
- Want to have a display with a system view
- Want to retain architectural model
- Want to exploit real time stats


16








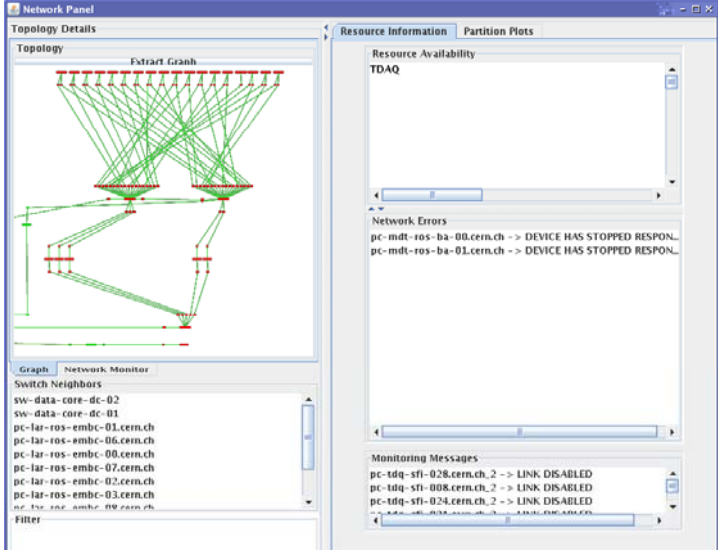





Atlas Users Network Panel



For operators we provide a summary of status per application set.



25



2D Display Limitations




2D Display is very good for switch to switch traffic visualisation

No 'level of detail' feature as you zoom in or out

Can't incorporate host details - - -would overwhelm the plot

BUT..



I want all the detail when something goes wrong

I want neighbouring views when examining individual behaviour

I want the system wide view for visual correlation

I want different levels of detail depending on my viewpoint

I want fly-through, and navigation and easily visible errors

I want ... I want ...

I want GOOGLE Earth for my network.

26



3D Display



Google earth as inspiration
Variable detail as a function of
viewing distance

Variable viewing angles
Intuitive navigation

Unfortunately Google doesn't
cope with our dynamic update
requirements

So we went looking for display
software that does.

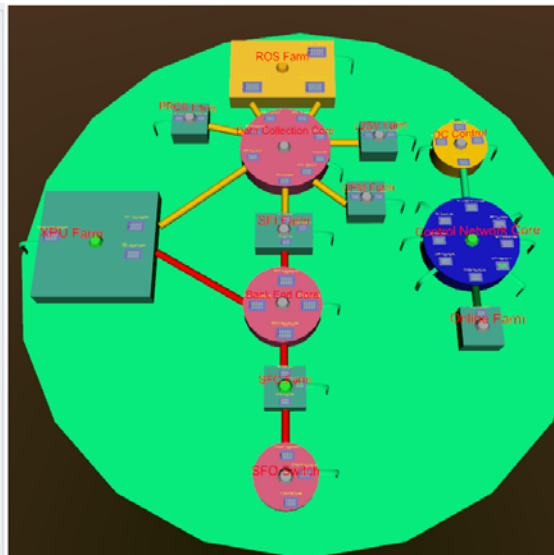
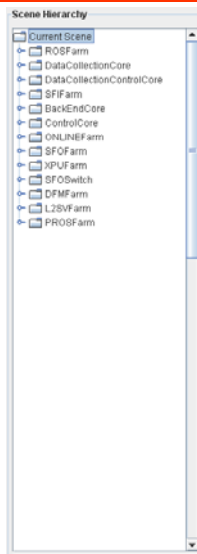
**X3D (enhanced VRML)
Octaga Player**



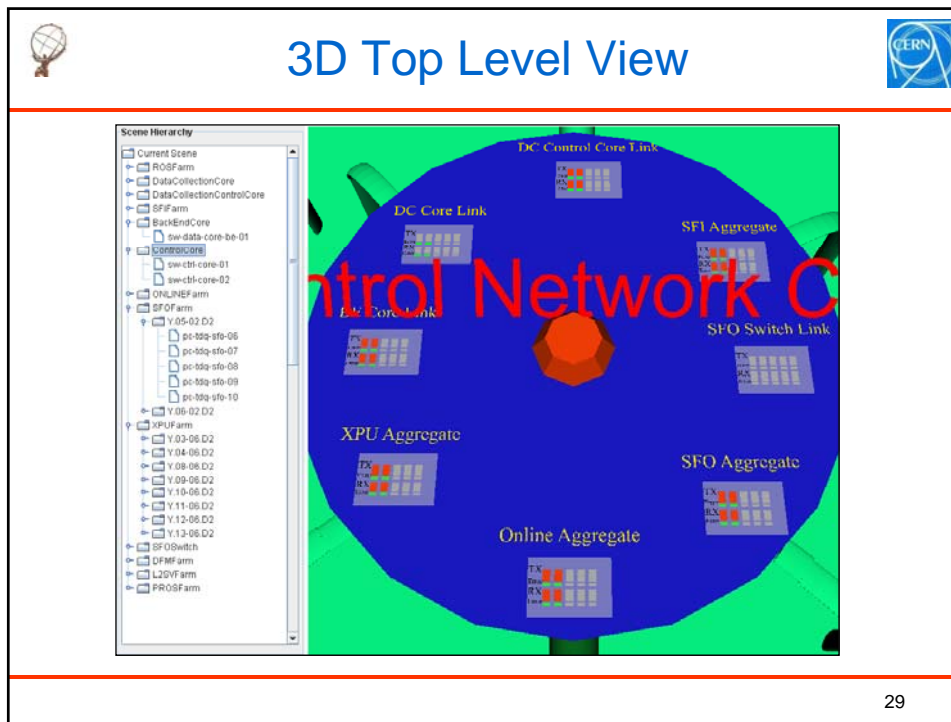
27



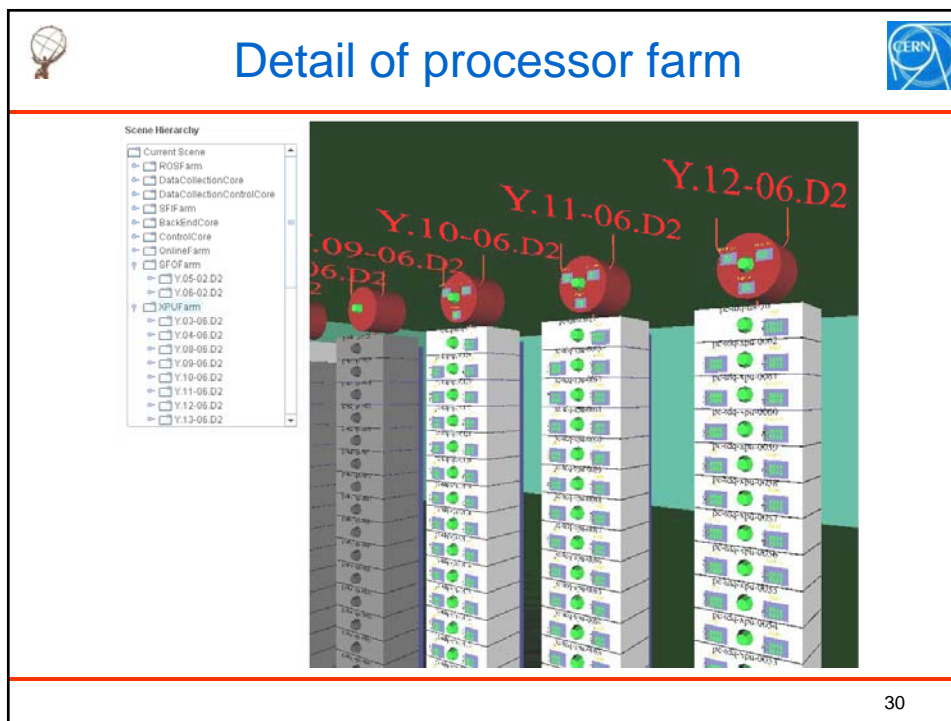
3D Top Level View



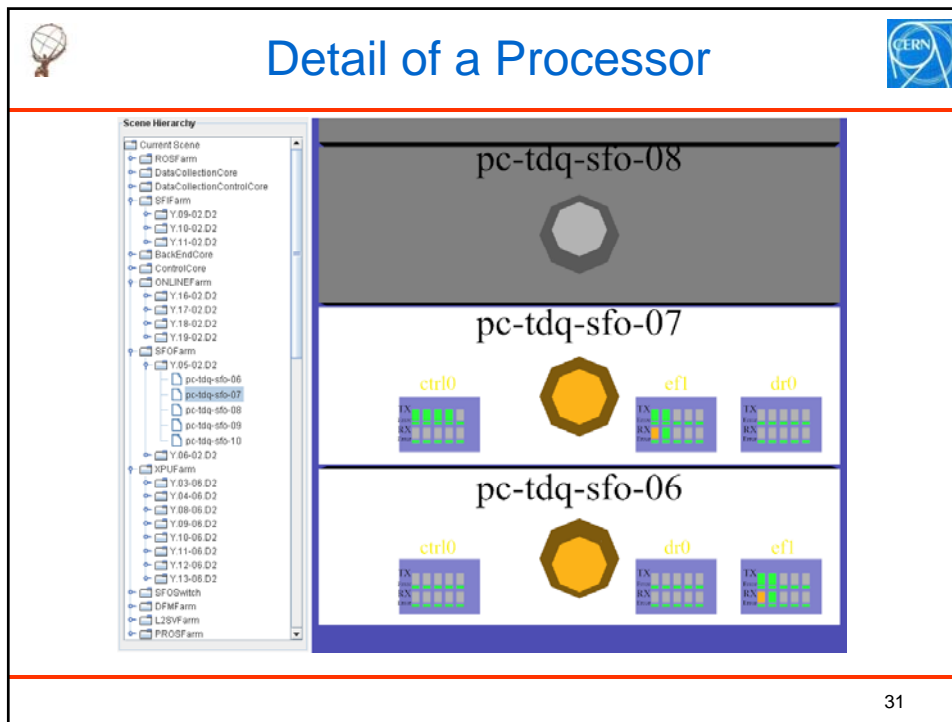
28



29



30



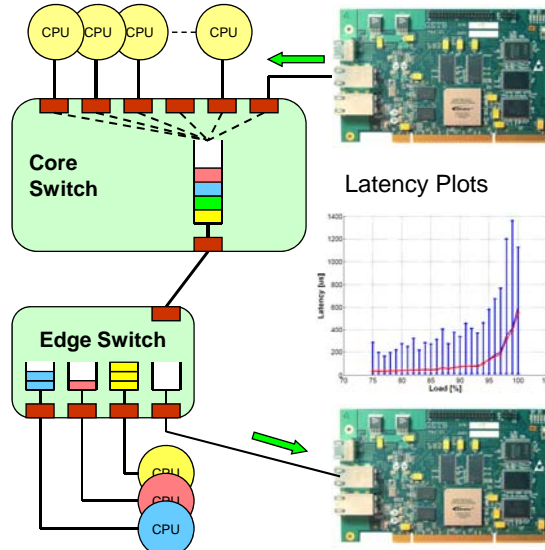
31

Diagnostic Tools: 1

- YATG (Yet Another Traffic Grapher)
 - High speed SNMP-based traffic monitoring (the switch is the limiting factor)
 - Fine time granularity statistics for selected device interfaces
- ATLAS-like traffic bandwidth measurements
 - Distributed applications replicate the transactional request-response transfer protocol
 - Demonstrate maximum achievable bandwidth

32

Diagnostic tools :2



Dynamic queue
growth
measurement

33

Diagnostic tools :3 SFLOW (1)

sflow

- sFlow is the standard for statistical packet sampling
 - Each network port → sampling system
 - All packet samples → central location (software)
 - Analysis → information about the *content* of the traffic
- By collecting packet samples, the packets can be classified into **flows**
 - A flow ~ network conversation between two applications
 - The bandwidth occupancy for each flow can be estimated
- We developed an sFlow analysis application in order to study the technology

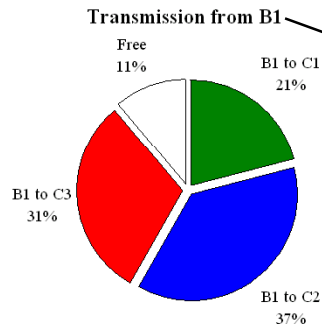
34

Diagnostic tools :3 SFLOW (2)

For one switch port

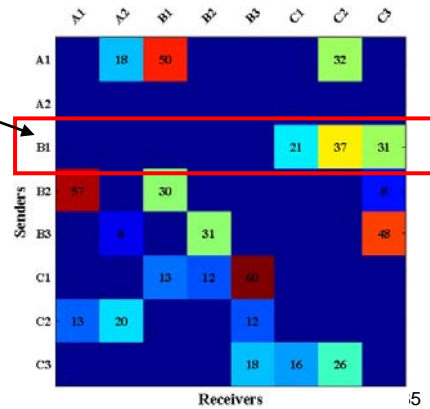
Using SNMP → "bandwidth usage on port B1 is 89%"

Using sFlow → Pie chart with traffic distribution



For the entire network

Traffic matrix with all senders and receivers



Summary

- large scale network
- monitor, diagnose
 - commercial tools
 - "in-house" built software
 - plots for 6000 ports
 - different levels of visual abstraction
- 24/7 operation
- errors not welcome