# WLCG perfSONAR Update

Marian Babik, CERN IT
on behalf of WLCG Network Throughput WG



WLCG
Worldwide LHC Computing Grid

Open Science Grid

# Outline

- WLCG perfSONAR infrastructure status

- 100Gbps Testing

- OSG/WLCG Network Monitoring Platform
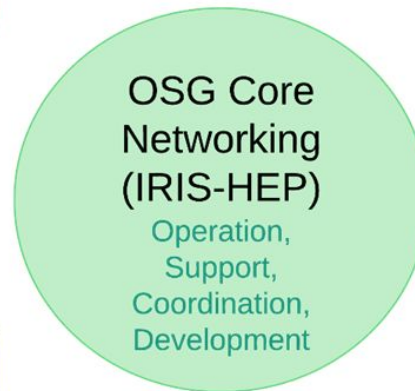
- New Analytics and Tools

- Summary

# OSG/WLCG networking projects

There have been 4 coupled projects around the core OSG Net Area

1. **SAND** (NSF) project for analytics (ended)

2. **HEPiX** NFV WG (finished work)

3. **perfSONAR** project

4. **WLCG Network** Throughput WG

**OSG Networking Components**

Ended July 2021

**SAND**
Analytics, VIsualization, Alerting/Alarming

**perfSONAR**
Framework, Metrics, Tools

**OSG Core Networking (IRIS-HEP)**
Operation, Support, Coordination, Development

**HEPiX Network Function Virtualization WG**
Technology exploration, Testing

**WLCG Throughput WG**
Configuration, Triage, Policy

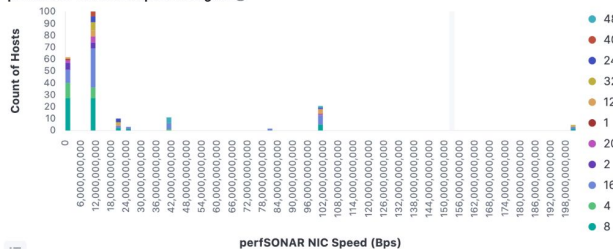WG Completed Work

# perfSONAR deployment



238 Active perfSONAR instances
- **207 production endpoints**
- T1/T2 coverage
- Dedicated latency and bandwidth nodes at each site
- Testing coordinated and managed from central place
- Continuously testing over 5000 links
- LHC experiments, DUNE, BelleII, LSST
- LHCOPN/LHCONE, ARCHIVER, StashCache, SLATE, CC*
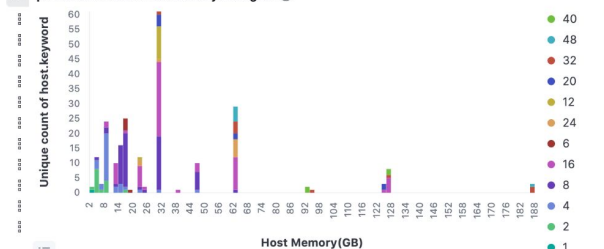- UK and FR meshes

# perfSONAR deployment

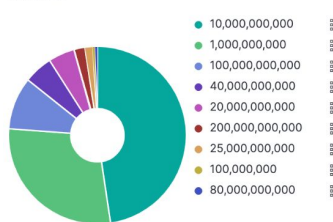238 Active perfSONAR instances - **207 production endpoints** - T1/T2 coverage
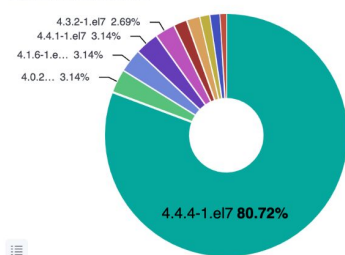


Our global toolkit deployment has a range of systems in terms of age and capability

Dashboard in ELK

Sites should remember to not only upgrade perfSONAR software but also the underlying **hardware,** as nodes become too old or are unable to test at the site storage speed.

# perfSONAR News

- perfSONAR 5 (beta out)
  - OpenSearch as local archive (replacing esmond/Cassandra) + Logstash
  - Grafana visualisations (dashboards)
  - Toolkit supports CC7, latest Debian 10, Ubuntu 18/20 and RHEL8 (Alma/Rocky)
    - CS8 will not be officially supported
    - Our recommendation is to wait for RHEL9 support
- [4.4.4 bug fix](#) released April 5th (WLCG baseline release)
  - Number of bug-fixes in pscheduler - please update if you haven't done already
- We're still seeing issues with some nodes hitting resource limits on very busy nodes (reboot resolves this, permanent fix is part of perfSONAR 5)
- HW updates on very old nodes might be needed
  - We now support configurations on a single node with two NICs
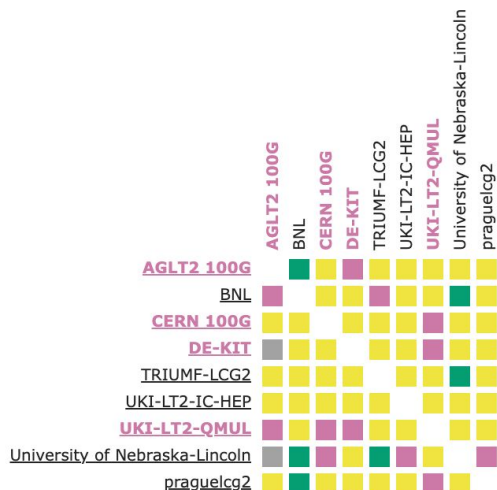  - Docker deployments (testpoint only) can also be considered

ESnet  GÉANT  INDIANA UNIVERSITY  INTERNET2  RNP ORGANIZAÇÃO SOCIAL DO MCTI  UNIVERSITY OF MICHIGAN

- WLCG 100Gbps mesh



WLCG 100G Mesh - WLCG 100G IPv6 Bandwidth - Throughput

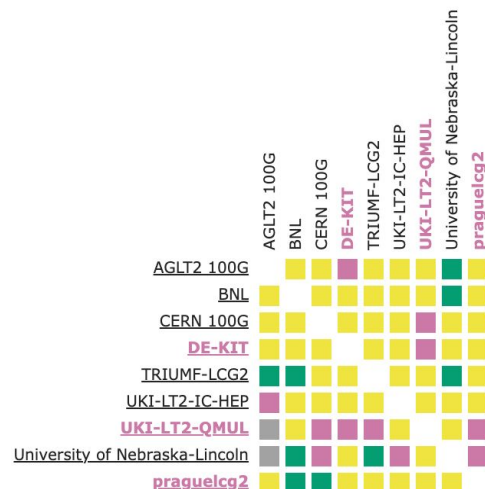Throughput >= 10Gbps    Throughput < 10Gbps    Throughput <= 1Gbps    Unable to fi

⚠ Found a total of 5 problems involving 4 hosts in the grid

WLCG 100G Mesh - WLCG 100G IPv4 Bandwidth - Throughput

Throughput >= 10Gbps    Throughput < 10Gbps    Throughput <= 1Gbps    Unable to fi

⚠ Found a total of 4 problems involving 3 hosts in the grid

# 100Gbps Testing

- Monthly meetings since January
  - Aim to achieve 10% of avail. capacity (~10Gbps) on a regular basis
  - Discussing ways to tune the nodes and improve stability
  - wlcg-perfsonar-100g mailing list (join)
- Tunings
  - Used CheckMK monitoring along with ES/Kibana dashboards to check status
  - TCP buffers and MTU appear to have made the biggest difference
    - TCP buffers by default at ~ 200MB, need to be increased to 1GB
  - References:
    - https://fasterdata.es.net/host-tuning/linux/100g-tuning/
  - Tried FQ but that actually decreased the throughput in tests (not work-conserving)
  - NIC interrupts/core sync only possible via manual tests
- maddash shows by default avg. over 24 hours - extended to 4 days
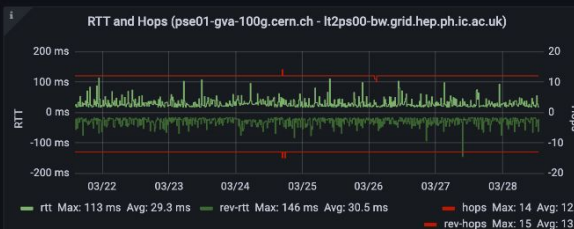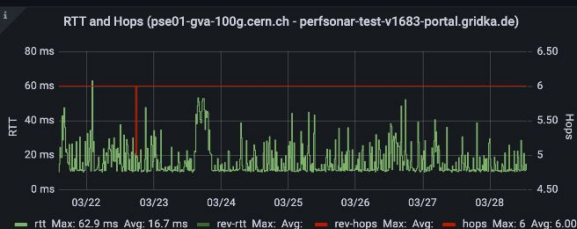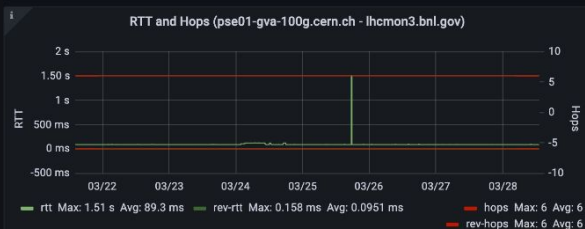- New host-based Grafana dashboard available

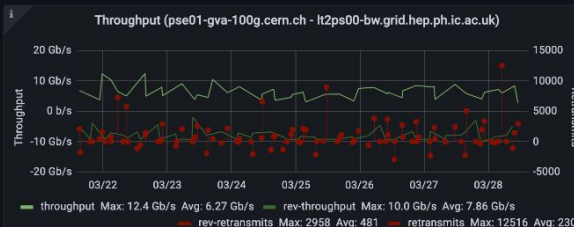# Grafana dashboard

# Grafana dashboard ([link](link))

# Grafana dashboard ([link](#))

- Collects, stores, configures and transports all network metrics
  - Distributed deployment - operated in collaboration
- All perfSONAR metrics are available via API, live stream or directly on the analytical platforms
  - Complementary network metrics such as ESNet, LHCOPN traffic also via same channels

# Network Measurement Platform Evolution

- **Collects, stores, configures and transports all network metrics**
  - Distributed deployment - operated in collaboration
- Planned evolution based on the perfSONAR 5
  - Directly publishing results from perfSONARs to ES@UC
  - High-level services provided to the experiments/users

# Tools and Applications for Network Data

- To organize access to all the various resources we have NEW homepage (https://toolkitinfo-nextjs.vercel.app/)
- We already have Kibana dashboards looking at
  - Bandwidth
  - Traceroute
  - Packetloss / Latency
  - Infrastructure
- With the completion of the SAND project, we have a few prototype tools that help us analyze and utilize our net data
  - We have a new perfSONAR focused dashboard: **ps-dash**
  - We have added a self-subscribe tool for network alarms call **AAAS**
  - ***Next two pages have the details on these two apps***

# pS (perfSONAR) Dash



https://ps-dash.uc.ssl-hep.org/

**Purpose**: provides a user dashboard to explore analyzed and summarized perfSONAR data.

Currently:
- Allows users to monitor their sites
- Provides tools for detecting basic problems

**Future plans:**
- Add today's Alarms
- Add traceroute data & plots
- Refine ranks
- Deduct possible cause for found issues

# ATLAS Alarms & Alerts Service



https://aaas.atlas-ml.org/

**Purpose**: provides user-subscribable alerting for specific types of network issues found by analyzing perfSONAR data

Currently available:
- Main packet loss issues
- Main throughput issues

**Future plans:**

- Add traceroute alarms:
  - Destination never reached
  - **Network path changes**
  - Node causes issues with multiple sites

# Bandwidth Alarms

Detecting changes in measured throughput wrt. 21-day average (ipv4, ipv6), e.g. see below a sample alarm

Currently working on creating high-level alarms (aggregating multiple alarms and running correlations with latencies and path alarms)

Sun, 05 Jun 2022 04:08:47    Networking/Perfsonar/Bandwidth decreased from/to multiple sites        Bandwidth decreased from/to multiple sites
tags: GRIF
Bandwidth decreased for ipv4 links between site GRIF to sites: ['INFN-T1', 'RO-14-ITIM', 'SAMPA'] change in percentages: [-50, -11, -27]; and from sites: ['IN2P3-CPPM', 'IN2P3-LAPP'], change in percentages: [-42, -44] with respect to the 21-day average.

Sun, 05 Jun 2022 04:08:47    Networking/Perfsonar/Bandwidth decreased from/to multiple sites        Bandwidth decreased from/to multiple sites
tags: IN2P3-CC
Bandwidth decreased for ipv4 links between site IN2P3-CC to sites: ['BNL-ATLAS', 'GLOW', 'INFN-T1', 'RRC-KI-T1', 'TOKYO-LCG2', 'TRIUMF-LCG2', 'UFlorida-HPC'] change in percentages: [-42, -28, -16, -10, -38, -30, -71]; and from sites: ['UKI-SOUTHGRID-OX-HEP'], change in percentages: [-51] with respect to the 21-day average.

Sun, 05 Jun 2022 04:08:47    Networking/Perfsonar/Bandwidth decreased from/to multiple sites        Bandwidth decreased from/to multiple sites
tags: IN2P3-CC
Bandwidth decreased for ipv6 links between site IN2P3-CC to sites: ['IN2P3-CPPM', 'NDGF-T1', 'Nebraska'] change in percentages: [-28, -81, -44]; and from sites: ['GLOW', 'IN2P3-CPPM', 'IN2P3-LAPP', 'TRIUMF-LCG2'], change in percentages: [-61, -42, -30, -21] with respect to the 21-day average.

Sun, 05 Jun 2022 04:08:47    Networking/Perfsonar/Bandwidth decreased from/to multiple sites        Bandwidth decreased from/to multiple sites
tags: IN2P3-CPPM
Bandwidth decreased for ipv4 links between site IN2P3-CPPM to sites: ['AGLT2', 'BU_ATLAS_Tier2', 'CA-VICTORIA-WESTGRID-T2', 'GRIF', 'IN2P3-LAPP', 'IN2P3-LPSC', 'INFN-T1', 'RO-03-UPB', 'UAM-LCG2', 'UKI-SOUTHGRID-OX-HEP'] change in percentages: [-99, -30, -99, -42, -65, -25, -97, -98, -89, -95]; and from sites: ['BEIJING-LCG2', 'IN2P3-LAPP', 'Taiwan-LCG2'], change in percentages: [-99, -48, -81] with respect to the 21-day average.

Sun, 05 Jun 2022 04:08:47    Networking/Perfsonar/Bandwidth decreased from/to multiple sites        Bandwidth decreased from/to multiple sites
tags: IN2P3-CPPM
Bandwidth decreased for ipv6 links between site IN2P3-CPPM to sites: ['BEgrid-ULB-VUB', 'BNL-ATLAS', 'CSCS-LCG2', 'IN2P3-CC', 'IN2P3-SUBATECH', 'MWT2', 'NDGF-T1', 'RO-16-UAIC', 'Taiwan-LCG2', 'praguelcg2'] change in percentages: [-97, -97, -96, -42, -53, -44, -95, -31, -84, -78]; and from sites: ['BEgrid-ULB-VUB', 'IN2P3-CC', 'IN2P3-LPSC', 'MWT2', 'NDGF-T1', 'RRC-KI-T1', 'SWT2_CPB', 'UKI-SOUTHGRID-OX-HEP'], change in percentages: [-99, -28, -13, -99, -54, -70, -99, -98] with respect to the 21-day average.

Sun, 05 Jun 2022 04:08:47    Networking/Perfsonar/Bandwidth decreased from/to multiple sites        Bandwidth decreased from/to multiple sites
tags: IN2P3-LAPP
Bandwidth decreased for ipv4 links between site IN2P3-LAPP to sites: ['GRIF', 'IN2P3-CPPM', 'IN2P3-LPSC', 'RO-03-UPB'] change in percentages: [-44, -48, -27, -93]; and from sites: ['IN2P3-CPPM', 'TRIUMF-LCG2'], change in percentages: [-65, -19] with respect to the 21-day average.

Sun, 05 Jun 2022 04:08:47    Networking/Perfsonar/Bandwidth decreased from/to multiple sites        Bandwidth decreased from/to multiple sites
tags: IN2P3-LPSC
Bandwidth decreased for ipv4 links between site IN2P3-LPSC to sites: ['DESY-ZN', 'RRC-KI-T1', 'SWT2_CPB'] change in percentages: [-26, -70, -36]; and from sites: ['AGLT2', 'IN2P3-CPPM', 'IN2P3-LAPP'], change in percentages: [-84, -25, -27] with respect to the 21-day average.

# Network Path Anomalies Detection

## Detecting changes in ASNs sequences across all our traceroutes

Fri, 03 Jun 2022 17:26:32    Networking/Perfsonar/Path changed    Path changed

tags: CSCS-LCG2, SWT2_CPB, CA-VICTORIA-WESTGRID-T2, FMPhI-UNIBA, UKI-SCOTGRID-ECDF, RAL-LCG2, USC-LCG2, UKI-NORTHGRID-MAN-HEP, IFCA-LCG2, DESY-HH, GRIF, UKI-NORTHGRID-LANCS-HEP, WT2, KR-KISTI-GSDC-01, IN2P3-LPSC, TRIUMF-LCG2, pic, UKI-SCOTGRID-GLASGOW, TECHNION-HEP, IFIC-LCG2, CA-SFU-T2, FZK-LCG2, BNL-ATLAS, RO-16-UAIC, IN2P3-CPPM, RO-03-UPB, RO-14-ITIM, UKI-NORTHGRID-LIV-HEP, JP-KEK-CRC-02, DESY-ZN, UKI-SOUTHGRID-OX-HEP, NDGF-T1, CERN-PROD, UKI-LT2-QMUL, AGLT2, INFN-MILANO-ATLASC, BU_ATLAS_Tier2
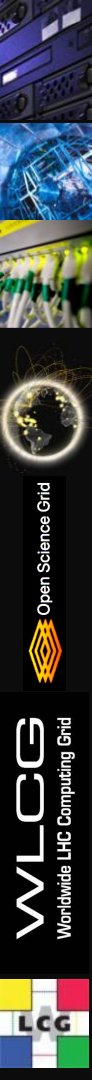
In the past 12 hours, path between 106 pairs diverged and went through ASN 6939 owned by HURRICANE, US. The change affected the following sites ['CSCS-LCG2', 'SWT2_CPB', 'CA-VICTORIA-WESTGRID-T2', 'FMPhI-UNIBA', 'UKI-SCOTGRID-ECDF', 'RAL-LCG2', 'USC-LCG2', 'UKI-NORTHGRID-MAN-HEP', 'IFCA-LCG2', 'DESY-HH', 'GRIF', 'UKI-NORTHGRID-LANCS-HEP', 'WT2', 'KR-KISTI-GSDC-01', 'IN2P3-LPSC', 'TRIUMF-LCG2', 'pic', 'UKI-SCOTGRID-GLASGOW', 'TECHNION-HEP', 'IFIC-LCG2', 'CA-SFU-T2', 'FZK-LCG2', 'BNL-ATLAS', 'RO-16-UAIC', 'IN2P3-CPPM', 'RO-03-UPB', 'RO-14-ITIM', 'UKI-NORTHGRID-LIV-HEP', 'JP-KEK-CRC-02', 'DESY-ZN', 'UKI-SOUTHGRID-OX-HEP', 'NDGF-T1', 'CERN-PROD', 'UKI-LT2-QMUL', 'AGLT2', 'INFN-MILANO-ATLASC', 'BU_ATLAS_Tier2']

14 2001:48a8:68f7:8001:192:41:236:31-2001:630:441:905::b AGLT2-UKI-SOUTHGRID-OX-HEP Baseline: [20965, 231, 237, 11537, 786]

| | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | |
| 237 | 231 | 237 | 237 | 237 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | | | |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 |
| 237 | 231 | 237 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | |
| 237 | 231 | 237 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 |
| 237 | 231 | 237 | 237 | 237 | 237 | 11537 | 11537 | 20965 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | |
| 237 | 231 | 237 | 237 | 237 | 237 | 11537 | 11537 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | | |
| 237 | 231 | 237 | 237 | 237 | 237 | 11537 | 6939 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 | | |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 |
| 237 | 231 | 237 | 237 | 237 | 11537 | 11537 | 11537 | 11537 | 20965 | 20965 | 786 | 786 | 786 | 786 | 786 | 786 |

# Summary

- OSG in collaboration with WLCG operates a comprehensive network monitoring platform
  - Provides data and feedback to LHCOPN/LHCONE, HEPiX, WLCG and OSG communities
- The IRIS-HEP and SAND projects have produced some new tools for exploring and utilizing our network data
- Developing high-level services based on perfSONAR measurements that will help sites, experiments and R&Es receive targeted alarms/alerts on existing issues in the infrastructure
- We have to continue to watch our network monitoring infrastructure as it is a complex system with lots of areas for issues to develop.

# Acknowledgements

We would like to thank the **WLCG**, **HEPiX**, **perfSONAR** and **OSG** organizations for their work on the topics presented.

# Useful URLs

- OSG/WLCG Networking Documentation
  - https://opensciencegrid.github.io/networking/
- perfSONAR Infrastructure Dashboard
  - https://atlas-kibana.mwt2.org:5601/s/networking/goto/9911c54099b2be47ff9700772c3778b7
- perfSONAR Dashboard and Monitoring
  - http://maddash.opensciencegrid.org/maddash-webui
  - https://psetf.opensciencegrid.org/etf/check_mk
- perfSONAR Central Configuration
  - https://psconfig.opensciencegrid.org/
- Toolkit information page
  - https://toolkitinfo.opensciencegrid.org/
- Grafana dashboards
  - http://monit-grafana-open.cern.ch/
- ATLAS Alerting and Alarming Service: https://aaas.atlas-ml.org/
- The pS Dash application: https://ps-dash.uc.ssl-hep.org/
- ESnet WLCG DC Dashboard:
  https://public.stardust.es.net/d/IkFCB5Hnk/lhc-data-challenge-overview?orgId=1

# Backup Slides Follow

# WLCG Network Throughput Support Unit

Support channel where sites and experiments can report potential network performance incidents:

- Relevant sites, (N)RENs are notified and perfSONAR infrastructure is used to narrow down the problem to particular link(s) and segment. Also [tracking past incidents](#).
- Feedback to WLCG operations and LHCOPN/LHCONE community

**Most common issues**: MTU, MTU+Load Balancing, routing (mainly remote sites), site equipment/design, firewall, workloads causing high network usage

As there is no consensus on the MTU to be recommended on the segments connecting servers and clients, LHCOPN/LHCONE working group was established to investigate and produce a recommendation. (See coming [talk](#) :) )

# Importance of Measuring Our Networks

- **End-to-end network issues are difficult to spot and localize**
  - Network problems are multi-domain, complicating the process
  - Performance issues involving the network are complicated by the number of components involved end-to-end
  - Standardizing on specific tools and methods focuses resources more effectively and provides better self-support.
- **Network problems can severely impact experiments workflows and have taken weeks, months and even years to get addressed!**
- **perfSONAR provides a number of standard metrics we can use**
  - Latency, Bandwidth and Traceroute
  - These measurements are critical for network visibility
- **Without measuring our complex, global networks we wouldn't be able to reliably use those network to do science**