

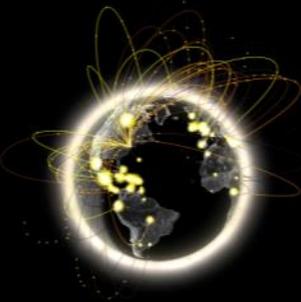
# Network Status Update: Where We Are Today

Shawn McKee

with input from Marian Babik and Edoardo Martelli

ATLAS ADC TIM

9 November 2016



# Why Networking?

- For an XRootD workshop, why talk about networking?
  - Of course we know we need, and heavily utilize, the network but what are the concerns?
  - To-date the main need has been to better support diagnosing, localizing and repairing network problems
- I will review our status and recent activities
- Then I will cover how networking is evolving and what may be changing in mid-to-long-term

# HEP Network Summary

- HEP (and especially LHC) networking is:
  - **Global**
  - **Foundational** to our computing models and infrastructure
  - Continuing an **exponential increase in bandwidth use**
  - **Functioning well** but facing some current and future challenges
- The HEP community has significantly benefited from the world-wide Research & Education (R&E) networking community
- There are a number of (relatively) small efforts in HEP engaged in network-related areas which I will try to cover
- While our wide-area networking needs are significant and have historically been the poster child for globally distributed e-Science, this may be changing over the coming years.

# First: A Little History

- There is a long history of work by physicists (led in large part by Harvey Newman) to enable HEP networking going back to 1986 (and actually starting in 1981 with a Caltech-CERN modem link)
- In the late 1990s the MONARC team developed a model of how LHC experiments might construct a suitable infrastructure accounting for compute, storage and networking
  - Model assumed the network was expensive, somewhat unreliable and not very performant.
  - The hierarchy of tiered computing centers was the output
- After the LHC turn-on the experiments found that the network was actually one of the most reliable and best performing components of our global infrastructure
  - And that excellent wide-area networking (WAN) was generally being provided without direct cost to the experiments
- Based upon the experience in Run-1 the LHC experiments evolved their computing models to **take better advantage of the network.**
  - The hierarchical model was replaced by more egalitarian access to data and sites
  - Direct access to data across the WAN became part of the toolkits (AAA, FAX, etc)

# Importance of Measuring Our Networks

- End-to-end network issues are difficult to spot and localize
  - Network problems are multi-domain, complicating the process
  - Standardizing on specific tools and methods allows groups to focus resources more effectively and better self-support
  - Performance issues involving the network are complicated by the number of components involved end-to-end.
- **perfSONAR** provides a number of standard metrics we can use
- Latency measurements provide one-way delays and packet loss metrics
  - Packet loss is almost always very bad for performance
- Bandwidth tests measure achievable throughput and track TCP retries (using Iperf3)
  - Provides a baseline to watch for changes; identify bottlenecks
- Traceroute/Tracepath track network topology
  - All measurements are only useful when we know the exact path they are taking through the network.
  - Tracepath additionally measures MTU but is frequently blocked

# Latency and packet loss matters

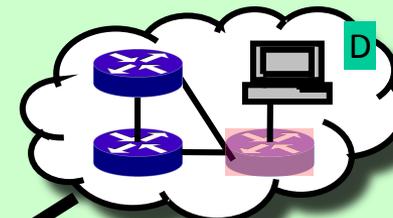
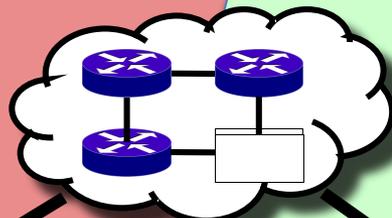
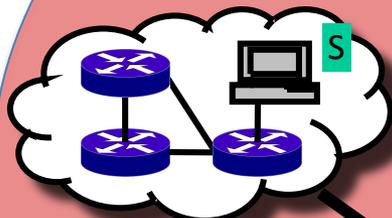
Performance is poor when RTT exceeds ~10 ms

Performance is good when RTT is < ~10 ms

Source Campus

R&E Backbone

Destination Campus



0.0046% loss (1 out of 22k packets) on 10G link

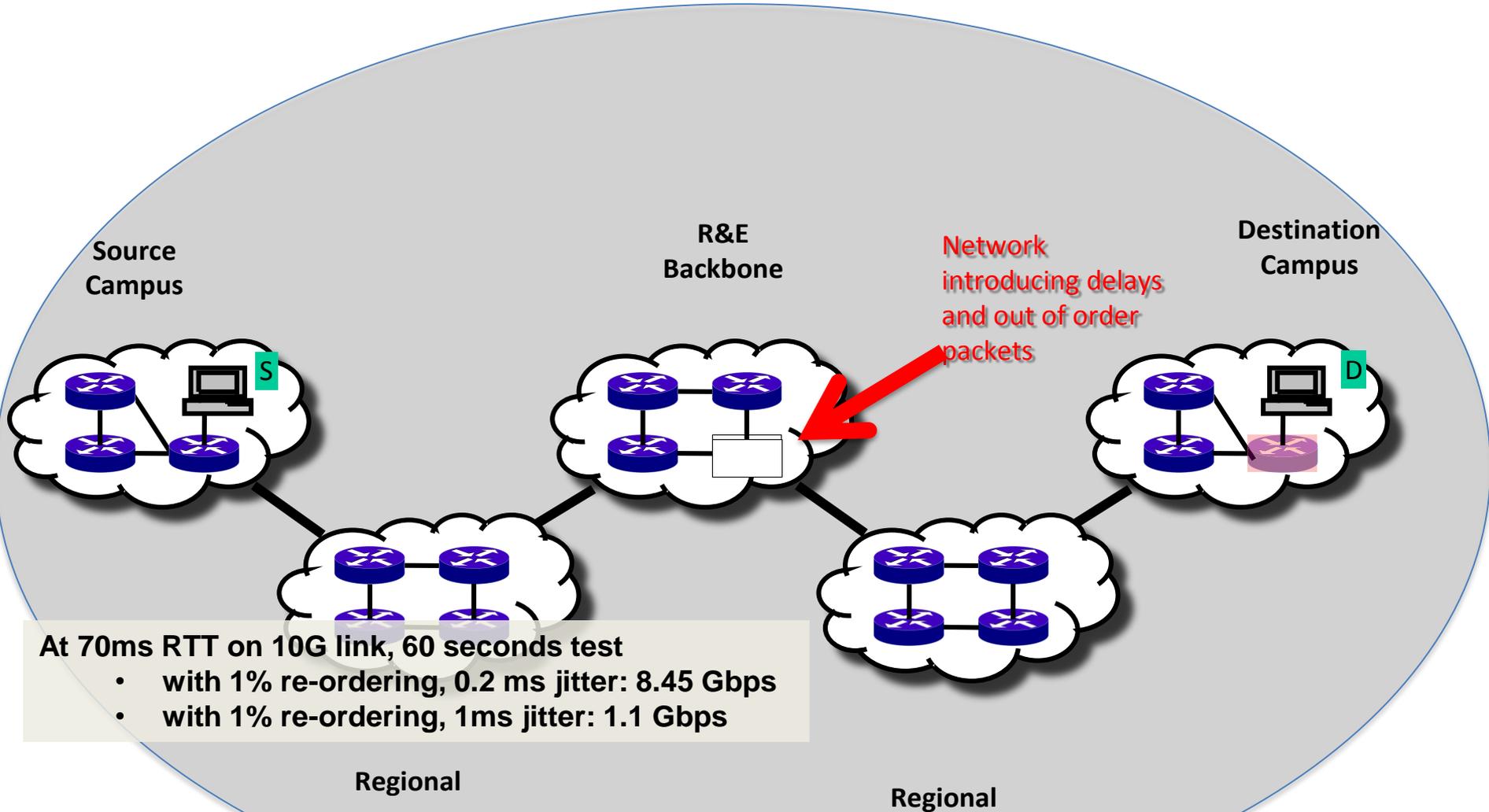
- with 1ms RTT: 7.3 Gbps
- with 51ms RTT: 122Mbps
- with 88ms RTT: 60 Mbps (factor 80)

Regional

Regional

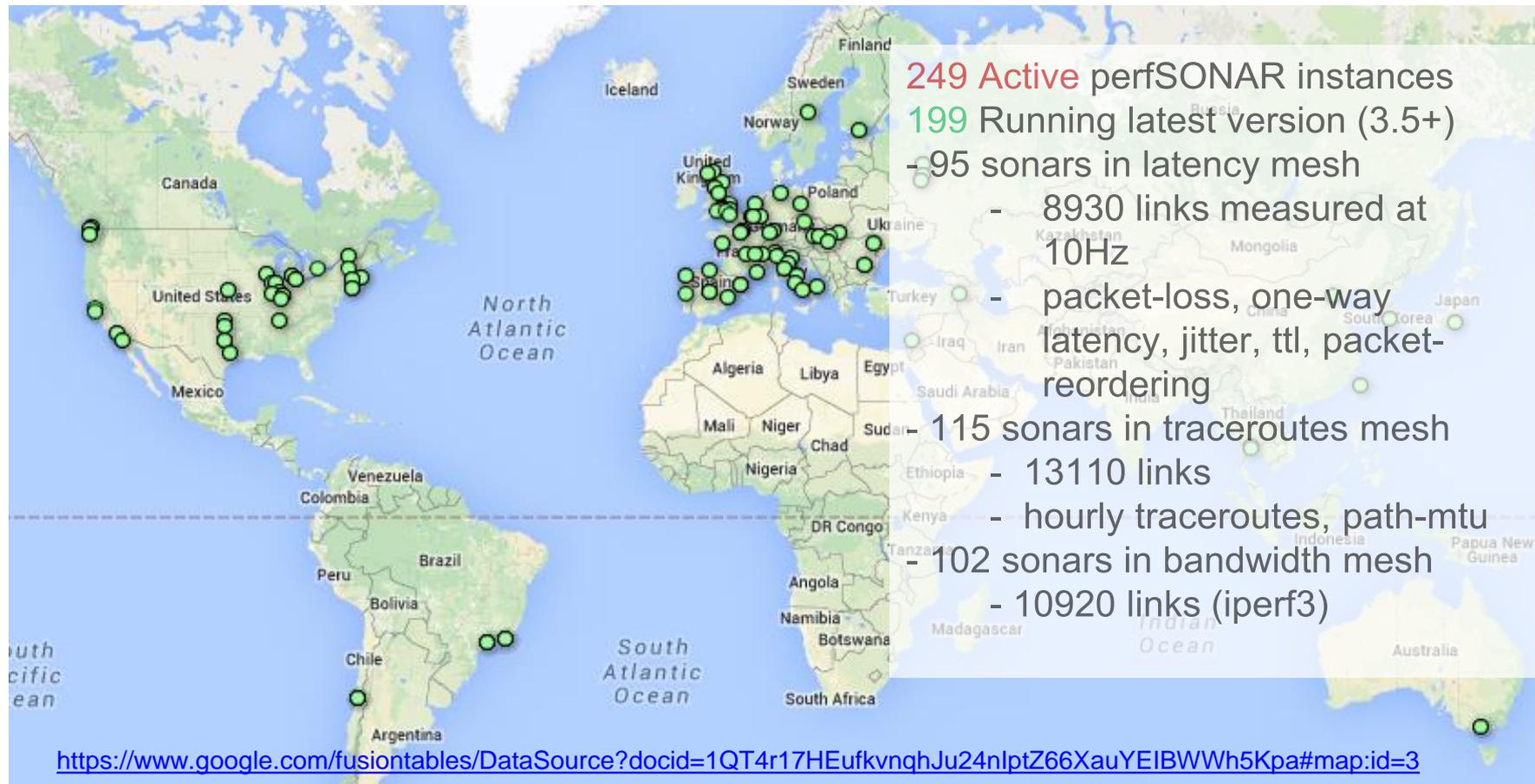
Switch with small buffers

# Packet ordering and jitter



# Current perfSONAR Deployment

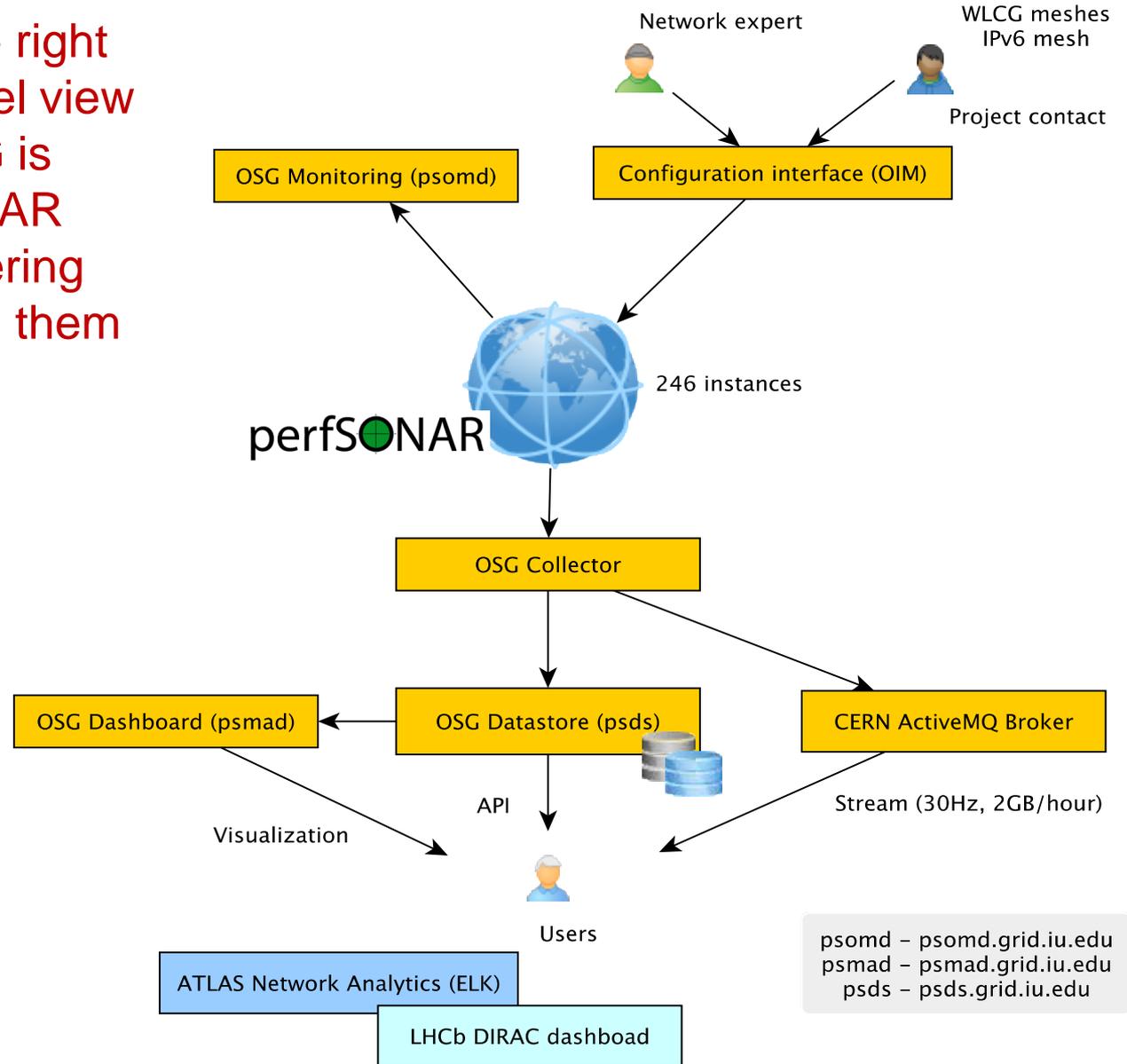
[http://grid-monitoring.cern.ch/perfsonar\\_report.txt](http://grid-monitoring.cern.ch/perfsonar_report.txt) for stats



- Initial deployment coordinated by WLCG perfSONAR TF
- Commissioning of the network followed by WLCG Network and Transfer Metrics WG

# Overview of perfSONAR Pipeline

The diagram on the right provides a high-level view of how WLCG/OSG is managing perfSONAR deployments, gathering metrics and making them available for use.



# HEP Networking

- Here is a quick snapshot of those involved in HEP Networking
  - The Open Science Grid
  - The WLCG Network and Transfer Metrics WG
  - Many institutions and communities supporting HEP networking
    - R&E backbone networks like ESnet, Internet2, GEANT,...
    - NRENs across the globe
    - Communities like LHCOPN/LHCONE, GLIF, perfSONAR Developers...
    - And all the many institutions around the world involved in network research relevant to HEP (**way too many to list!**)
    - **Our challenge is to incorporate this work into our infrastructure**



# OSG Networking Area Plans: Year 5

- Continue to do what we do now, and:
- Develop effective **Alarming and Alerting**
- Support higher-level network services
  - We have proto-typed a proximity service to find nearest SE given perfSONAR or to find the nearest perfSONAR give and SE
  - Create network cost prediction service to predict quality and capacity of source-destination paths for network decision support
- Improve the ability to manage and use network topology and network metrics: Analytics Platform?
- Prepare-for and integrate **Software Defined Networking**

# The WLCG Network and Transfer Metrics Working Group

- Started in Fall 2014, it brings together network & transfer experts
  - Follows up on the WLCG perfSONAR Task Force goals
- **Mandate**
  - Ensure all relevant **network** and **transfer metrics** are identified, collected and published
  - Ensure sites and experiments can better understand and fix networking issues
  - Enable use of network-aware tools to improve transfer efficiency and optimize experiment workflows
- **Membership**
  - WLCG perSONAR support unit (regional experts), WLCG experiments, FTS, Panda, PhEDEx, FAX, Network experts (ESNet, LHCOPN, LHCONE)

<https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics>

# Coordinating Network Issue Response

- The working group has created a support unit to coordinate responses to potential network issues
  - Tickets opened in the support group can be triaged to the right destination
  - Many issues are potentially resolvable within the working group
  - Real network issues can be identified and directed to the appropriate network support centers
- Documented at [https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics#Network Performance Incidents](https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics#Network_Performance_Incidents)
- Example case CA<->EU [GGUS-118730](#)
  - resolved within hours of being reported
  - mainly due to our ability to narrow down using perfSONAR

# LHCOPN/LHCONE

- The LHCOPN working group was established by CERN, the WLCG Tier-1 sites and the various HEP related research and education networks to define, deploy and operate the LHC Optical Private Network interconnecting the Tier-1 and the Tier-0 at CERN
- The success of LHCOPN for the Tier-1s led to the creation of a similar network to support the Tier-2s and their interactions with the Tier-1s: The LHC Open Network Environment (LHCONE)
- The LHCOPN/LHCONE group meets jointly 2-3 times per year to discuss policy, operations and future evolution necessary to support the LHC (and now beyond) community.
  - This mostly volunteer effort has been very beneficial for LHC
  - There is a request to increase the participation from the experiments

# Network Analytics

- Ilija Vukotic/U Chicago has been leading an effort to get network metrics into an analytics platform (see HEPiX Talk <https://indico.cern.ch/event/531810/contributions/2321493/>)
- This analytics service indexes historical network related data while providing predictive capabilities for near term network throughput performance.
- Primary functions:
  - Aggregate, and index, network related data associated with WLCG “links”
  - Serve derived network analytics to ATLAS production, DDM & analysis clients
  - Provide a generalized network analytics platform for other communities in the OSG
- Part of ATLAS Analytics platform
  - <https://cds.cern.ch/record/2056257/files/ATL-SOFT-SLIDE-2015-752.pdf>

# Throughput predictions

- Throughput measurements are expensive so done at low frequency. Delays and packet loss rate are cheap.
- Idea is to use delays and packet loss rate to predict maximum possible throughput.
- Mathis formula is used to model impact of packet loss and latency on throughput
  - $\text{Rate} < (\text{MSS}/\text{RTT}) * (1 / \text{sqrt}(p))$ 
    - MSS – segment size
    - RTT – round trip time
    - p – packet loss
- Packet (re)ordering and jitter to be added as well

# ATLAS Network Analytics

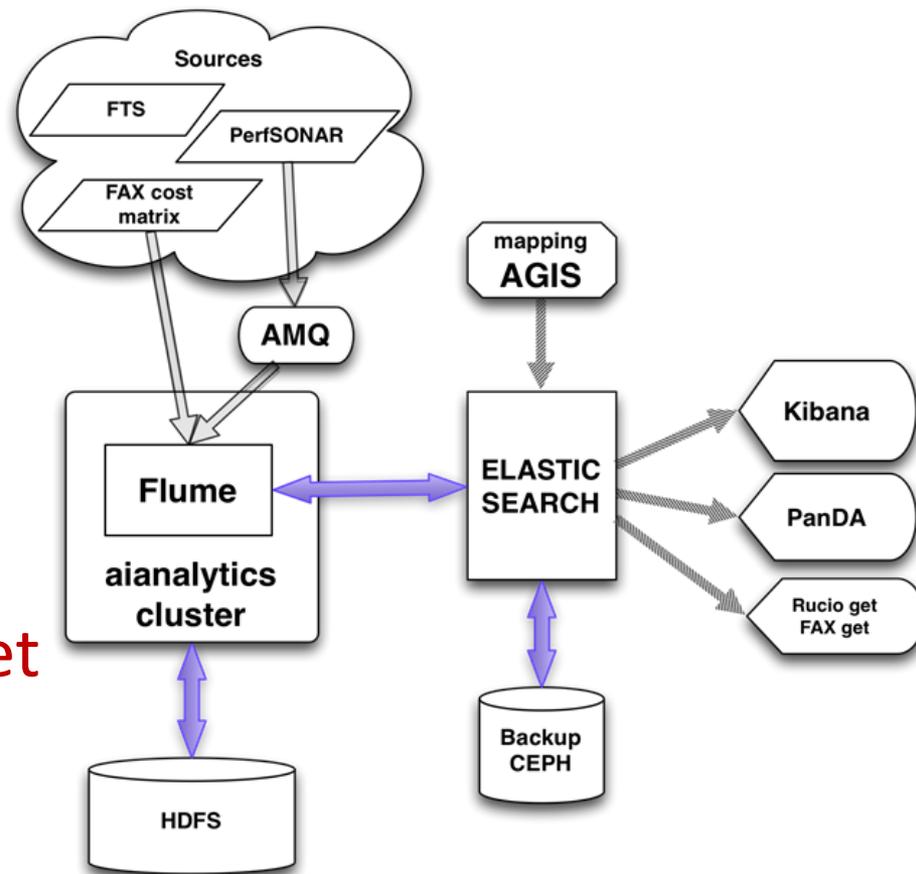
- Diagram shows the flow

- End-to-end+perfSONAR data both available to jointly analyze

- Kibana can be used to get customized views

<http://cl-analytics.mwt2.org:5601>

- More details at: <http://tinyurl.com/gt92zwb>



# Ongoing work in Network Analytics

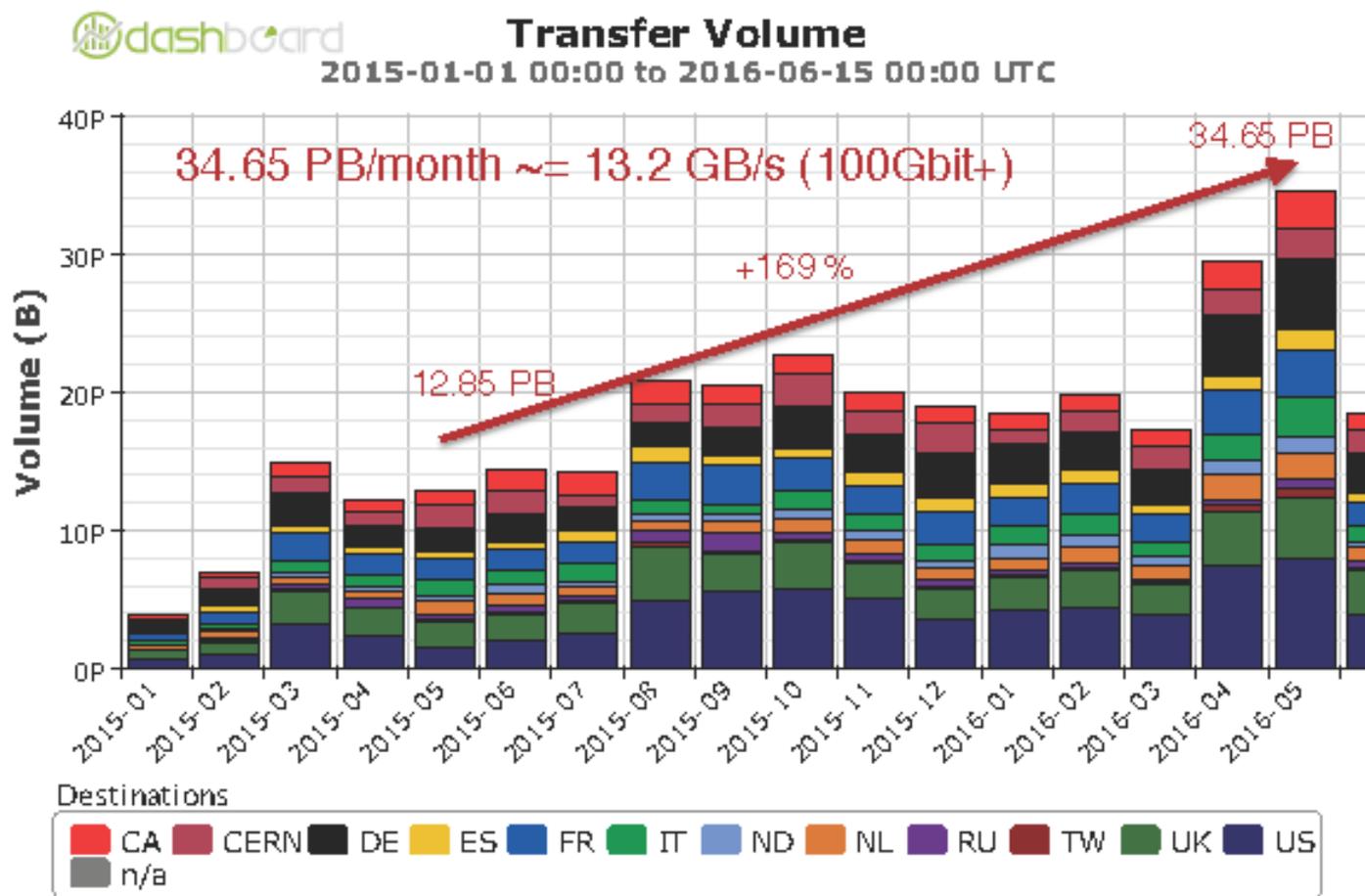
- The volume and complexity of network related data being collected by OSG and the experiments is challenging to use, but holds the promise of providing much deeper insights into our networks and hard to identify network problems.
  - To be most useful, the data requires cleaning, augmenting, transforming & correlating
- Ilija Vukotic (Univ. of Chicago) has developed ELK/jupyter stack for ATLAS Analytics and worked with Xinran Wang on [anomaly detection and advanced alerting/notifications](#) for network problems (See Track 5 talk Thursday afternoon)
  - Also looked at detection of the anomalies based on machine learning models
- Jerrod Dixon and Brian Bockelman (UNL) exploring network analytics in CMS
- Henryk Giemza (NCBJ), Federico Stagni integrating perfSONAR in DIRAC for LHCb
- Shawn McKee (Univ. of Michigan) working on real-time root cause analysis ([PuNDIT](#)) in collaboration with perfSONAR developers
- Hendrik Boras and Marian Babik (CERN) working on developing models for network cost-matrix - determine performance of network paths



# Network Evolution

- Historically the Wide-Area Network capacity has not always had a stable relationship compared to the data-center or end-node
  - In early days network links (on modems) significantly lagged the local speeds achievable within and between computers
  - The WAN technologies grew rapidly and for a while outpaced LAN and even local computing bus capacities
  - Today 100Gbps WAN links are the typical high-performance network speed but LANS are also in the same range.
    - Last Fall I bought a 32 port 100G switch, 4 dual-ported 100G NICs, 4 dual-ported 50G NICs, 4 dual-ported 25G NICs and all cables for \$18K
    - This summer I ordered a 100G NIC (Qlogic) for \$397 (for xrootd testing actually 😊)
- Today it is easy to oversubscribe our WAN links (in terms of \$ of local hardware at many sites)
- Will our R&E network providers be able to keep up with our needs?
  - So far, not a problem....
  - CERN currently testing 200Gbps waves
  - By 2020 800 Gbps waves will be available (assuming you buy the new hardware to support it)

# Network Use



ATLAS (and LHC in general) has been transferring an exponentially increasing amount of data since startup. This trend is likely to continue and is driven by increasing data volumes, more capable infrastructures and the excellent networks supporting our needs.

# Making the Most of our Networks

- Much of our WLCG infrastructure is NOT tuned to take the best advantage of the networks we currently have
  - There are a wide range of **mis-configurations**, **non-optimal tunings** and **incorrect application** and **hardware** settings that lead to inefficient use of our networks
  - As mentioned, we have a wealth of data now available and ready for analysis to **identify bottlenecks** and **poor performance**.
- As we identify bottlenecks and poor performance we need to take the next step and work to improve our end-host's ability to effectively utilize the network we have
  - Doesn't require **SDN**, new hardware or new networks but can make a huge difference in network throughput for sites
  - Should we organize a near-term workshop to share best practices, tools and tuning information?

# Improving End-host Networking

- New operating systems and associated end-host improvements in hardware are making it easier to get high-performance on our wide-area networks
- TCP more stable in CC7, throughput ramp ups much quicker
  - Detailed [report](#) available from [Brian Tierney](#) / [ESNet](#)
- Fair Queueing Scheduler (FQ) available from kernel 3.11+
  - Even more stable, works better with small buffers
- Best single flow tests show TCP LAN at 79Gbps, WAN (RTT 92ms) at 49Gbps
  - IPv6 slightly faster on the WAN, slightly slower on the LAN
- New TCP congestion algorithm ([TCP BBR](#)) from [Google](#)
  - [Google](#) reports 2-4x performance improvement on path with 1% loss (100ms RTT)
  - Early testing from [ESNet](#) less conclusive, there is also question how tolerant BBR will be with other congestion algorithms on the same link.

# New Network Tools/Capabilities

- Some important and interesting possibilities for what we might provide in the future include the creation of tools and visualization systems which manage network topologies (which are time-dependent)
  - Combining topology and metrics is powerful for identifying and localizing network problems; **currently a very manual process.**
- Using these tools users can look for correlations with the metrics measured across those topologies.
  - This type of tool can be used to help localize problems.
- **Note it is only by using the complete set of OSG/WLCG network metrics that this becomes possible.**

# Network Tomography

Host A is getting poor performance to Host B and seeing **3% packet loss**  
Normally we would start to investigate partial paths to isolate the problem



However we also see Host D to Host C is having problems and **2% packet loss**:



And there is a third pair (Hosts E and F) having **1% packet loss**:



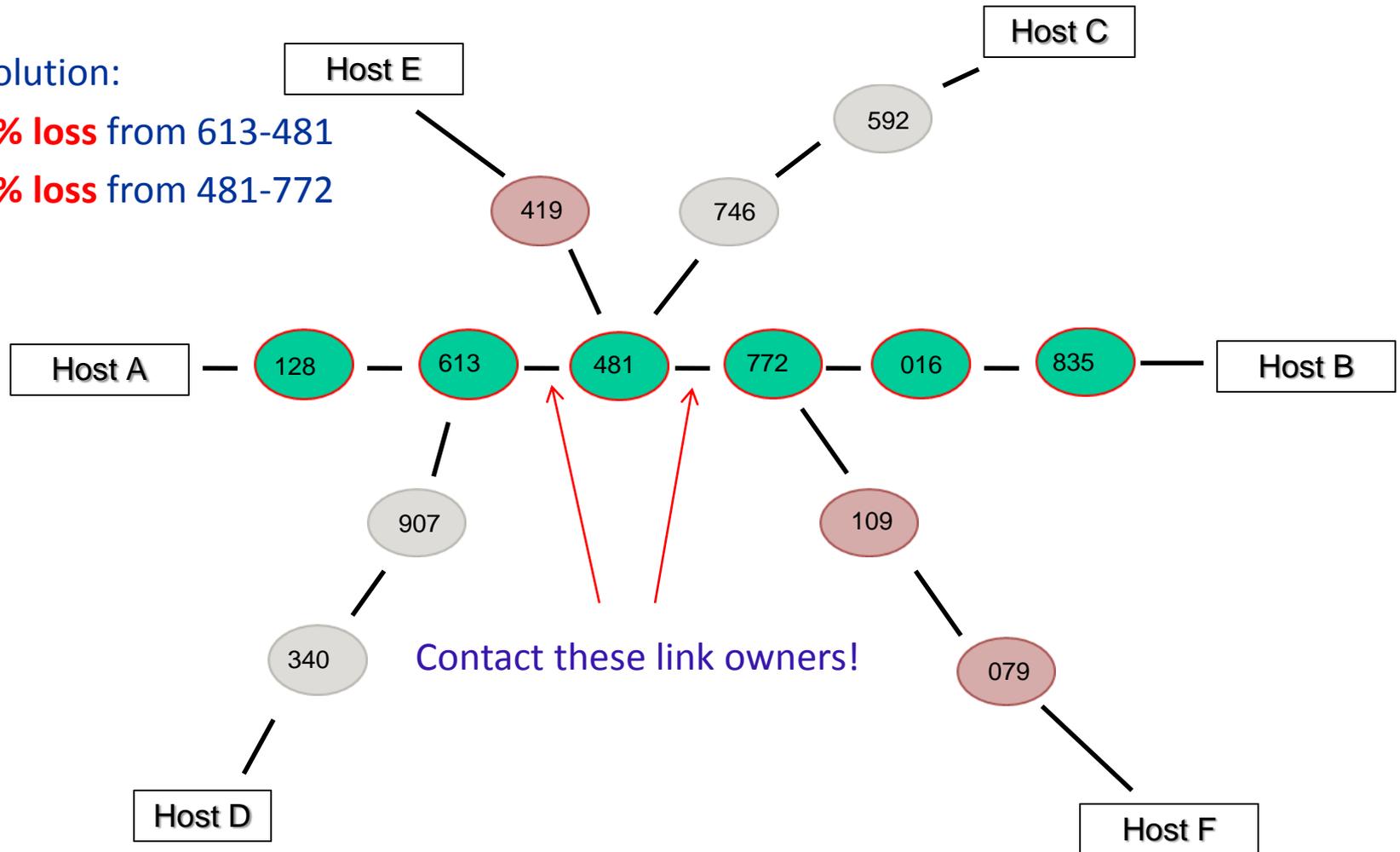
Let's correlate these paths

# Topology Problem Correlation

Solution:

2% loss from 613-481

1% loss from 481-772



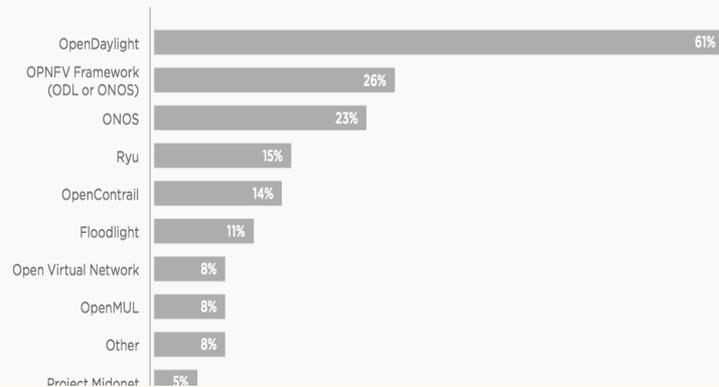
# R&E Networking

- High-Energy Physics (HEP) has significantly benefited from our strong relationship with Research and Education (R&E) network providers
  - To-date they have given us “infinite” capacity at relatively low (or no-direct) cost
    - They have been able to continually expanded their capacity to overprovision their networks relative to our needs and use.
- At the Terena network conference last spring SKA (Square Kilometer Array) noted **they will operate at data volumes 200xLHC scale** (<https://tnc16.geant.org/core/presentation/721>)
  - Besides Astronomy there are MANY science domains anticipating data scales beyond LHC: Health, Bioinformatics, Engineering...
- R&E network providers work closely with us in part because they view HEP as representative of future data-intensive science domains
  - HEP serves as the early prototype for such user communities
  - Network providers are concerned about what happens when there are **N** more HEP-scale science domains all wanting infinite capacity
    - ***Perhaps we should be too!***

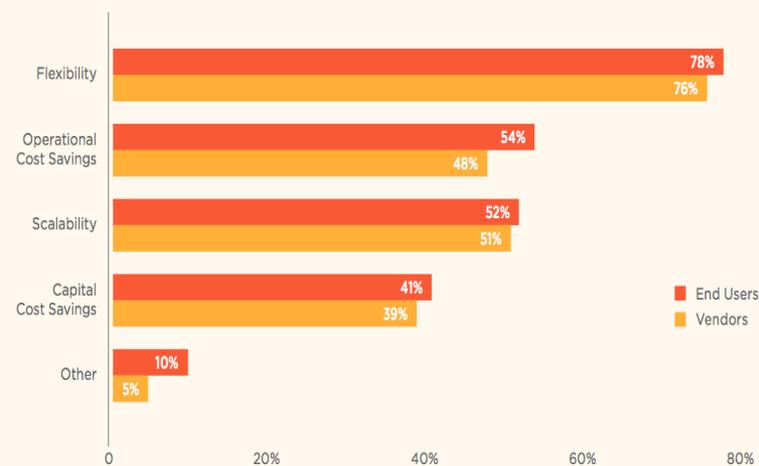
# Software Defined Networks (SDN)

- SDN is a set of technologies offering solutions for many of the future challenges
  - Current links might handle ~ 6x more traffic if we could avoid peaks and be more efficient
- Many different point-to-point efforts and successes reported within LHCOPN/LHCONE
  - The challenge remains getting this end-to-end
- While it's still unclear which technologies will become mainstream, it's already clear that software will play major role in networks in the mid-term (commercially driven)
  - Will experiments have effort to engage in the existing SDN testbeds to determine what impact it will have on their data management and operations ?

OPEN SOURCE NV/SDN SOLUTIONS DEPLOYED



BIGGEST BENEFITS OF NV



Respondents could choose multiple benefits of network virtualization. Results from 116 end-users and 85 technology vendors



sdxcentral.com

# Playing with SDN

- Future networks won't just have larger capacity but also programmability
- A group of people in the US from AGLT2, MWT2, SWT2 and NET2 are exploring SDN in ATLAS
  - Working with the LHCONe point-to-point effort as well
- We are deploying Open vSwitch on ATLAS production systems at these sites (<http://openvswitch.org/>)
  - IP addresses will be move to virtual interfaces
  - No other changes; verify no performance impact
  - Traffic can be shaped accurately with little CPU cost
- The **advantage** is the our data sources/sinks become **visible** and **controllable** by OpenFlow controllers like OpenDaylight
- Follow tests can be initiated to provide experience with controlling networks in the context of ATLAS operations.
- Interest from UVic, KIT and SurfSARA in participating
- Possible partnership with ESnet/CORSA in ~Dec timeframe
- *For more details talk to Rob Gardner or Shawn McKee*

# Current Measurement Efforts

- We have reorganized our bandwidth meshes
  - Previously one big mesh with 4 day cadence
    - Unable to finish tests due to nightly service restarts
  - New target: 3 (or more meshes) 24 hour cadence
    - **Proposal discussed this summer in European Throughput meeting**
    - **See [http://etf.cern.ch/perfsonar\\_meshes2.txt](http://etf.cern.ch/perfsonar_meshes2.txt)**
- We are also working on alerting when “obvious” problems are found
  - **Challenge:** getting appropriate contacts setup in check\_mk
  - Have initial version running in Analytics Platform now
  - Working with ETF to enable rule-based alerts (see <http://etf.cern.ch/docs/latest/user/overview.html#service> )

# Future Directions

- The WLCG efforts at CERN are being reorganized and this is an opportunity to chart future directions for the our networking efforts.
- We have a number of areas (projects; see next slide) we are considering and we need to understand where these efforts should be housed (Stay in WG, move to GDB, to LHCCONE)
  - It is important to note there is currently very little manpower for networking (much, much less than computing and storage)
  - To undertake all our plans will require identifying new effort
- We are planning a Pre-GDB meeting on January 10<sup>th</sup> 2017 focused on networking:  
<https://indico.cern.ch/event/571501/>
  - **Please REGISTER and ATTEND!**

# Possible Future Project Areas

- **Title:** LHCONE Traffic engineering
- **Areas:** LHCONE, routing, debugging, network orchestration
- **Title:** LHCONE L3VPN Looking Glass
- **Areas:** LHCONE, monitoring, debugging
- **Title:** Integration of network and transfer metrics to optimize experiments workflows
- **Areas:** FAX/Phedex, Rucio, perfSONAR, DIRAC
- **Title:** Advanced notifications/alerting for network incidents
- **Areas:** WAN, Advanced Notifications/Alerting, perfSONAR, Hadoop/Spark
- **Title:** Network performance of the commercial clouds
- **Areas:** Clouds, WAN connectivity, WAN performance (perfSONAR), establishing and testing network equipment at the cloud provider (VPN)
- **Title:** Software Defined Network Production Testbed
- **Areas:** WAN, SDN, LHCONE/LHCOPN, Storage/Data nodes

# Draft Perspective on Needed Effort

- **Short-term (1-2 years):** Focus on network monitoring, debugging and analytics. Find and fix network problems, improving our ability to utilize the networks we have.
- **Medium-term (3-7 years):** Plan for and evaluate the use of SDN for our infrastructures. Work on integration of those aspects deemed beneficial. Estimate the impact of other data-intensive science domains on our R&E networks and collaborate with them on their ramp-up to our scale.
- **Long-term(8-12 years):** Plan for and deal with the R&E network environment: **sharing, orchestration, automation** and the **implementation of smart networks**. Ensure our software can interact with smart network capabilities and agilely respond to dynamically changing infrastructure capacities and problems.

# Summary

- We have a working infrastructure in place to monitor and measure our networks
- perfSONAR provides lots of capabilities to understand and debug our networks
- Work on new applications is underway
  - Notifications/alerting
  - Predictive capabilities
  - Current utilization and capacity planning
  - Evaluating network performance of commercial clouds
- It is in HEP's best interest to stay aware of how the network is evolving and what the future landscape may look like
  - Important to start thinking of the network as something we will eventually be able to program/integrate into our architecture(s)

**Questions or Comments?**

# For Further Details

- WLCG network Use-cases document for experiments and middleware  
<https://docs.google.com/document/d/1ceiNITUJCwSuOuvbEHZnZp0XkWkwkPQTQic0VbH1mc/edit>
- Harvey's slides from Nordunet covering HEP networking history and ongoing work  
[https://www.dropbox.com/s/at2ky4rdc6szkmq/NGenIAGlobalNetworks\\_hbn091916.pptx?dl=0](https://www.dropbox.com/s/at2ky4rdc6szkmq/NGenIAGlobalNetworks_hbn091916.pptx?dl=0)
- OSG Network Documentation  
<https://www.opensciencegrid.org/bin/view/Documentation/NetworkingInOSG>
- WLCG Network and Transfer Metrics Working Group  
<https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics>
- perfSONAR deployment documentation for OSG and WLCG  
<https://twiki.opensciencegrid.org/bin/view/Documentation/DeployperfSONAR>
- WLCG workshop October 8, 2016 networking session presentations  
<https://indico.cern.ch/event/555063/sessions/203482/#20161008>

# References

- Network Documentation  
<https://www.opensciencegrid.org/bin/view/Documentation/NetworkingInOSG>
- Deployment documentation for OSG and WLCG hosted in OSG  
<https://twiki.opensciencegrid.org/bin/view/Documentation/DeployperfSONAR>
- Measurement Archive (MA) guide  
[http://software.es.net/esmond/perfsonar\\_client\\_rest.html](http://software.es.net/esmond/perfsonar_client_rest.html)
- Modular Dashboard and OMD *Prototypes*
  - <http://maddash.aglt2.org/maddash-webui>
  - [https://maddash.aglt2.org/WLCGperfSONAR/check\\_mk](https://maddash.aglt2.org/WLCGperfSONAR/check_mk)
- **OSG Production instances for OMD, MaDDash and Datastore**
  - <http://psmad.grid.iu.edu/maddash-webui/>
  - [https://psomd.grid.iu.edu/WLCGperfSONAR/check\\_mk/](https://psomd.grid.iu.edu/WLCGperfSONAR/check_mk/)
  - <http://psds.grid.iu.edu/esmond/perfsonar/archive/?format=json>
- Mesh-config in OSG <https://oim.grid.iu.edu/oim/meshconfig>
  - Being updated to a new standalone mesh-config application (ready for v4.0?)
- Use-cases document for experiments and middleware  
<https://docs.google.com/document/d/1ceiNITUJCwSuOuvbEHZnZp0XkWkwdkPQTQiC0VbH1mc/edit>