

# **LHCOPN/LHCONE perfSONAR Update**

**Shawn McKee/UM**

**LHCONE/LHCOPN Meeting**

**Brookhaven National Lab**

**April 4<sup>th</sup>, 2017**

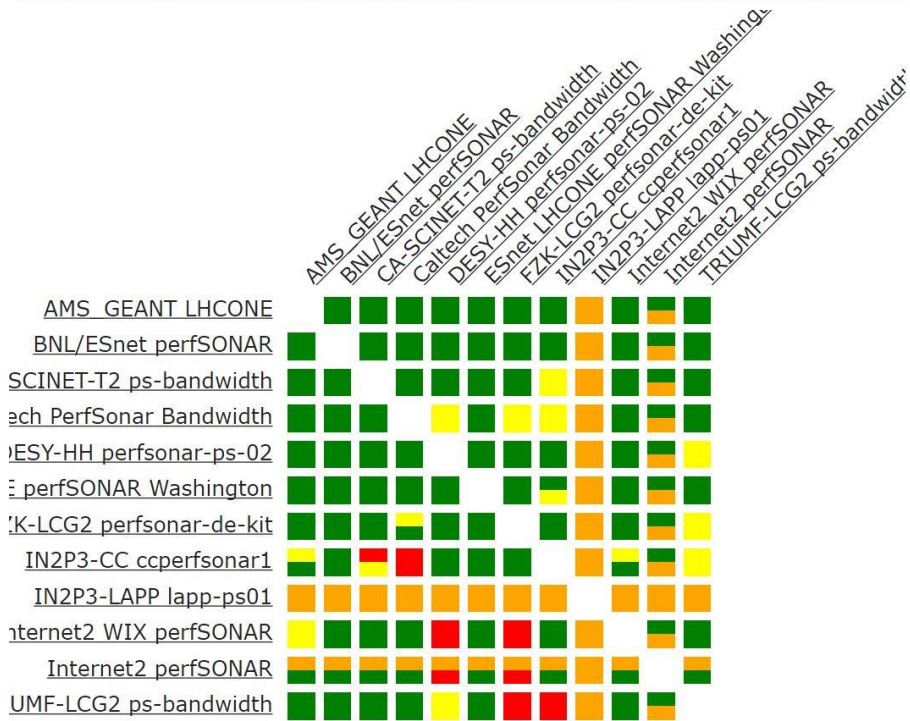
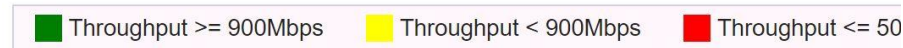
# Overview of Talk

- ❄ LHCONE and LHCOPN infrastructure overview
  - Status and changes in our meshes
- ❄ Tools Status and ongoing projects
  - New perfSONAR release, network analytics, mesh-config
- ❄ Example network debugging success story for LHCONE/LHCOPN
- ❄ Other Items

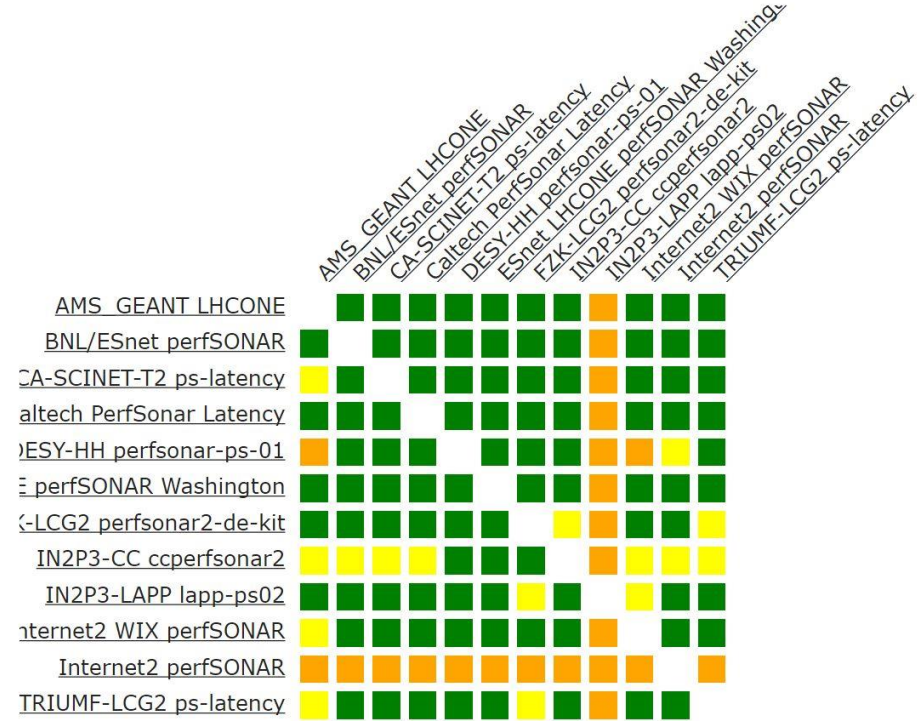
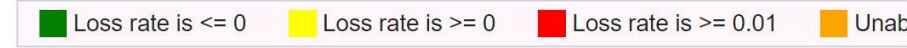
**Acknowledgements:** Marian Babik, Soichi Hayashi and Ilija Vukotic led the work on the new, interesting items I will cover.

# LHCONE MaDDash – 18 Sep 2016

LHCONE Mesh Config - TCP BWCTL Test Between



LHCONE Mesh Config - OWAMP Test Between LHCONE

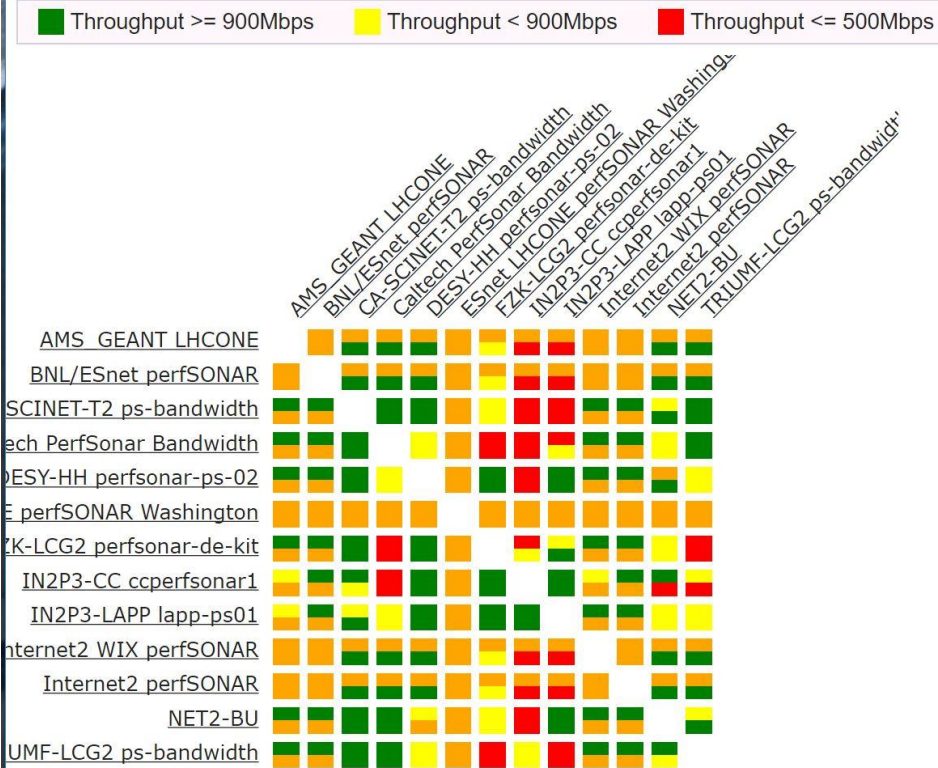


Compared to March 2016, things seem a bit better regarding the network quality.

# LHCONE MaDDash – 04 Apr 2017

LHCONE Mesh Config - TCP BWCTL Test Between LHCONE Mesh Config - OWAMP Test Between LHCONE

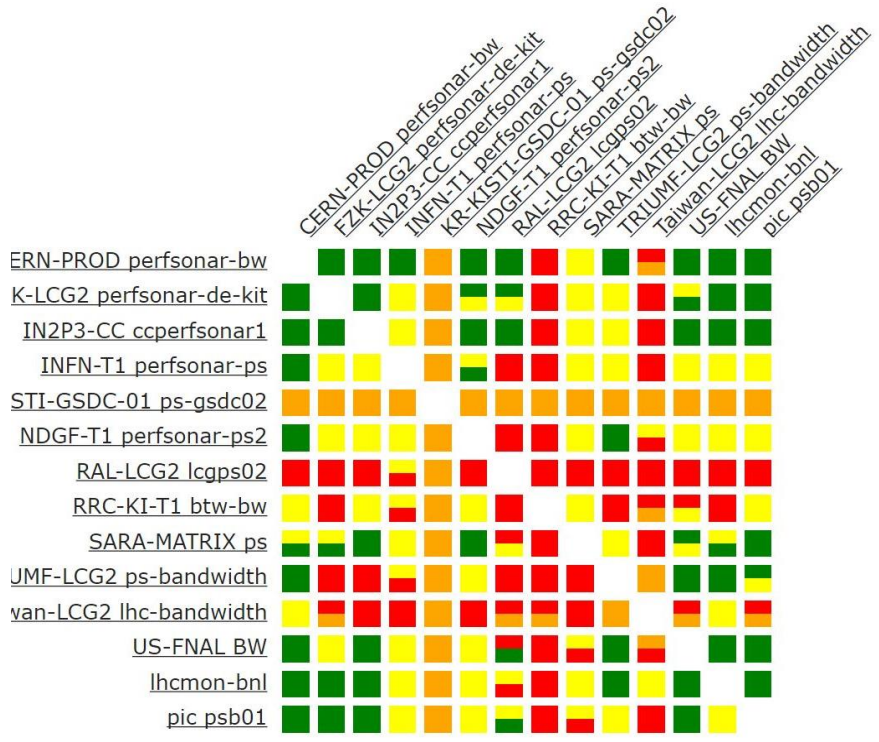
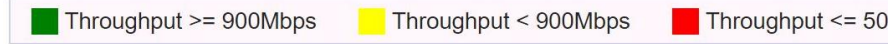
LHCONE Mesh Config - OWAMP Test Between LHCONE



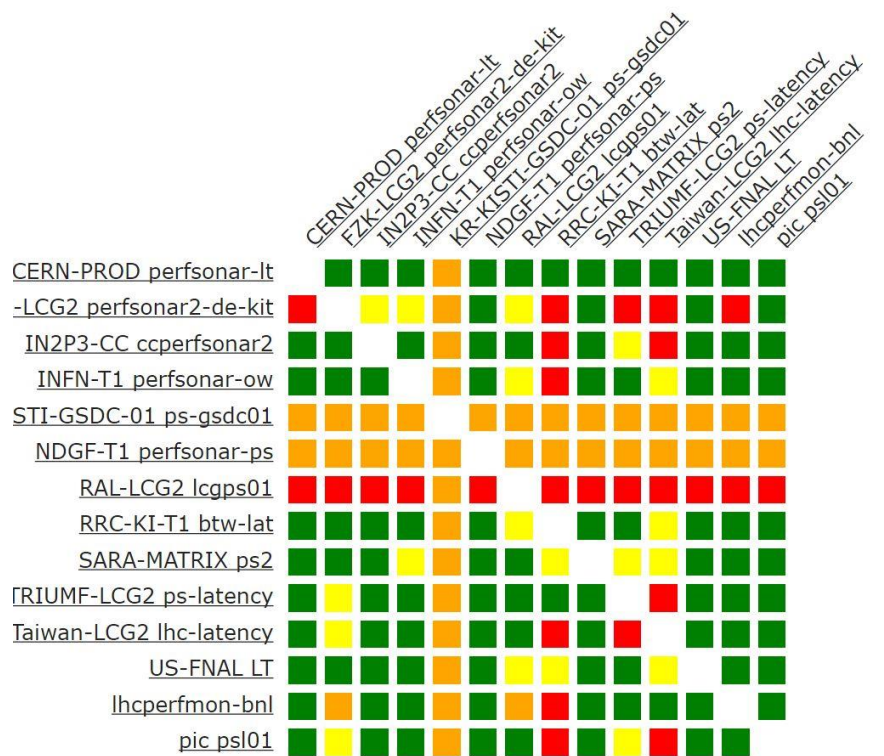
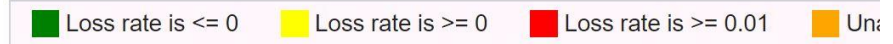
As you can see we have had lots of degradation in the SERVICE itself. Missing tests and results.

# LHCOPN MaDDash – 18 Sep 2016

OPN Config - TCP BWCTL Test Between OPN Band



OPN Config - OWAMP Test Between OPN Latency



Slight improvement since March 2016 in network. Specific issues noted last time.

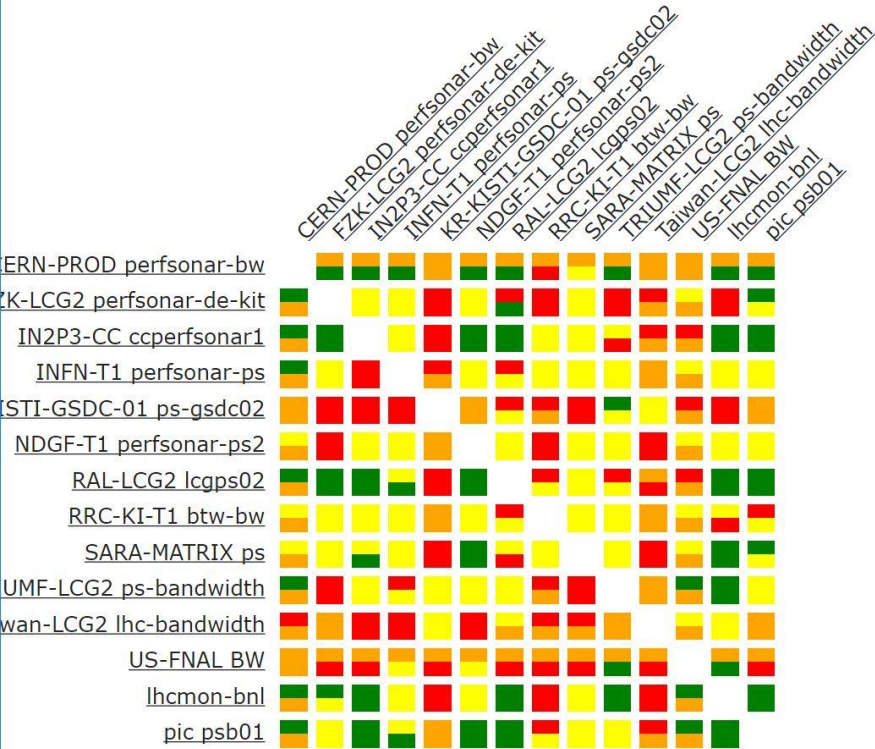
# LHCOPN MaDDash – 18 Sep 2016

OPN Config - TCP BWCTL Test Between OPN Bandwidth

■ Throughput >= 900Mbps 
 ■ Throughput < 900Mbps 
 ■ Throughput <= 50C

OPN Config - OWAMP Test Between OPN Latency

■ Loss rate is <= 0 
 ■ Loss rate is >= 0 
 ■ Loss rate is >= 0.01 
 ■ Un



Again degradation of the SERVICE. We should use the upcoming perfSONAR v4 release as an opportunity to get our installations fixed (and more resilient).

# perfSONAR Toolkit Status

- ❄ **The next major release of perfSONAR will be April 17**
  - ❑ RC3 came out a little more than 1 week ago
  - ❑ We will want to get ALL instances updated ASAP.
  - ❑ Default installations should auto-update.
  - ❑ Some LHCONe/LHCOPN instances may require manual intervention.
  - ❑ Need all pS hosts to provide JSON interface for monitoring
- ❄ **Along with perfSONAR updates, a number of other things are in process and I will cover in following slides**
  - ❑ New Mesh-config (MCA) is almost ready
  - ❑ New ETF monitoring for perfSONAR ready
  - ❑ New analytics capabilities using ELK + Jupyter

# perfSONAR v4.0 Update

## ❄ Focus is on “Control and Stability”

- ❑ New standalone mesh-config (MCA)
- ❑ New test scheduler (pScheduler):
  - ⌘ Shared by all tests and aware of the resources each uses
  - ⌘ Containing finer grained controls about who can run tests and what tests they are allowed to run.
  - ⌘ Increased visibility and control as to when tests will be run

## ❄ CentOS 7 and Debian 8 support

## ❄ New Endpoint selection capabilities

- ❑ Better metadata, topology information will be available

## ❄ New MaDDash w/MadAlert part of this release

## ❄ Minimum dual-core (2GHz+) and 4GB ram

## ❄ Toolkit now exposes info via JSON:

- ❑ <http://psum06.aglt2.org/toolkit/?format=json>



# Standalone Mesh-config (MCA)

- ❄ Central configuration and control of our meshes has been via OSG's MyOSG/OIM application...not easy to share access / management
- ❄ Soichi was pulled back into efforts to finish up the MCA (Mesh-Configuration Admin) GUI
  - ❑ 20% effort from December-March, 10% April-May funded by perfSONAR/IU
- ❄ Documentation at <http://docs.perfsonar.net/mca> and at <https://github.com/soichih/meshconfig-admin>
- ❄ Issues tracked at <https://github.com/soichih/meshconfig-admin/issues>
- ❄ OSG instance running at <https://meshconfig-itb.grid.iu.edu/> (create an account to play with this)
  - ❑ Now 258 hosts imported from OIM/GOCDB
  - ❑ New API available <https://meshconfig-itb.grid.iu.edu/apidoc/>
  - ❑ **Pushing to production this week:** <https://meshconfig.grid.iu.edu/>

# MCA Screenshot

Secure | <https://meshconfig-itb.grid.iu.edu/#!/configs/58d07d2931480800201aae24>

## MESH CONFIGS

Showing registred mesh configs

18 Mesh Configs

/ps-testbed	perSONAR Testbed	Jan 22, 2017
/soichi-test	test	Jan 23, 2017
/test	Soichi's Test Config	Jan 23, 2017
/test-dynamic	Soichi's Dynamic HostGroup Test	Jan 30, 2017
/us-atlas	USATLAS Mesh Config	Jan 30, 2017
/us-cms	USCMS Mesh Config	Mar 20, 2017
/UK	UK Meshconfig	Mar 20, 2017
/lhcone-all	LHCONE Mesh Config	Mar 20, 2017
/wlcg-lhcb-latency	WLCG LHCB Latency Mesh	Mar 20, 2017
/opn-all	OPN Config	Mar 20, 2017
/wlcg-cms-latency	WLCG CMS Latency Mesh	Mar 20, 2017
/wlcg-lhcb-bandwidth	WLCG LHCB Bandwidth Mesh	Mar 20, 2017
/wlcg-cms-bandwidth	WLCG CMS Bandwidth Mesh	Mar 20, 2017

### AUTO MESH CONFIG

Enter hostname of perSONAR toolkit instance to generate a mesh config URL containing tests for that instance as test endpoints.

Enter hostname

MeshConfig URL: <http://meshconfig-itb.grid.iu.edu/pub/config/opn-all> [Open MeshConfig](#)

**Name \*** OPN Config

**Description**

**Admins** Soichi Hayashi <hayashis@iu.edu> Shawn McKee <smckee@umich.edu> Marian Babik <marian.babik@cern.ch>


Users who can update this configuration

**Tests**

Enabled (include in mesh config) [Remove Test](#)

**Test Name** TCP BWCTL Test Between OPN Bandwidth Hosts

**Service Type** Bandwidth (bwctl) **Mesh Type** Mesh



# Primary Features of MCA

- ❄ Gathers and organizes information on hosts from a combination of sources
  - ❑ Imports from perfSONAR global lookup service entries; **we can discover perfSONARs deployed anywhere if they join the right community**
  - ❑ Able to gather information from registration databases like GOCDB / OIM
- ❄ Auto-completion / entry of values
- ❄ Context dependent user interface (see Testspecs)
- ❄ Can be easily installed outside of OSG
- ❄ Provides a RESTful interface to allow easy monitoring and software-controlled config
  - ❑ <https://meshconfig-itb.grid.iu.edu/apidoc/>
- ❄ Supports filtering and dynamic host groups
  - ❑ Can now build dynamic meshes, e.g., **all CentOS6 hosts who are members of the ATLAS community**
- ❄ Able to support both perfSONAR 3.x and 4.x

# Experiments Test Framework (ETF)

- ❄ WLCG is working with the LHC experiments to develop and provide a service monitoring capability called ETF
  - ❑ Experiments Test Framework(ETF) builds upon Nagios and Check\_mk to provide a flexible monitoring application
  - ❑ See details at <http://etf.cern.ch/docs/latest/>
- ❄ The WLCG Throughput WG is leveraging this to update our monitoring of perfSONAR using ETF\_ps
  - ❑ There is a docker version of ETF at <https://gitlab.cern.ch/etf/docker>
  - ❑ We have an etf\_ps version in docker: mbabik/etf:1.2.8p20
  - ❑ See [https://etf.aglt2.org/etf/check\\_mk/](https://etf.aglt2.org/etf/check_mk/) for a running example
- ❄ Our goal is to replace the OSG psomd.grid.iu.edu and perfsonar-itb.grid.iu.edu check\_mk instances with this

# ETF Example: psum06

❄ Shown below is an example of monitoring a v4.0 pS host

- ❑ pS specific checks including new info from JSON publication

Check MK Raw 1.2.8p176

Services of Host psum06.aglt2.org 17 rows /DC=ch/DC=cern/OU=Organic Units/OU=Users/CN=mckee/CN=500323/CN=Sha

Tactical Overview

Hosts	Problems	Unhandled
286	14	14

Services

Services	Problems	Unhandled
4627	1995	1995

Quicksearch: psum0

Views: Overview, Hosts, Host Groups, Services, Service Groups, Metrics, Business Intelligence, Problems, Inventory, Other

Bookmarks: Add Bookmark, EDIT

WATO - Configuration: Main Menu, Hosts, Host Tags, Global Settings

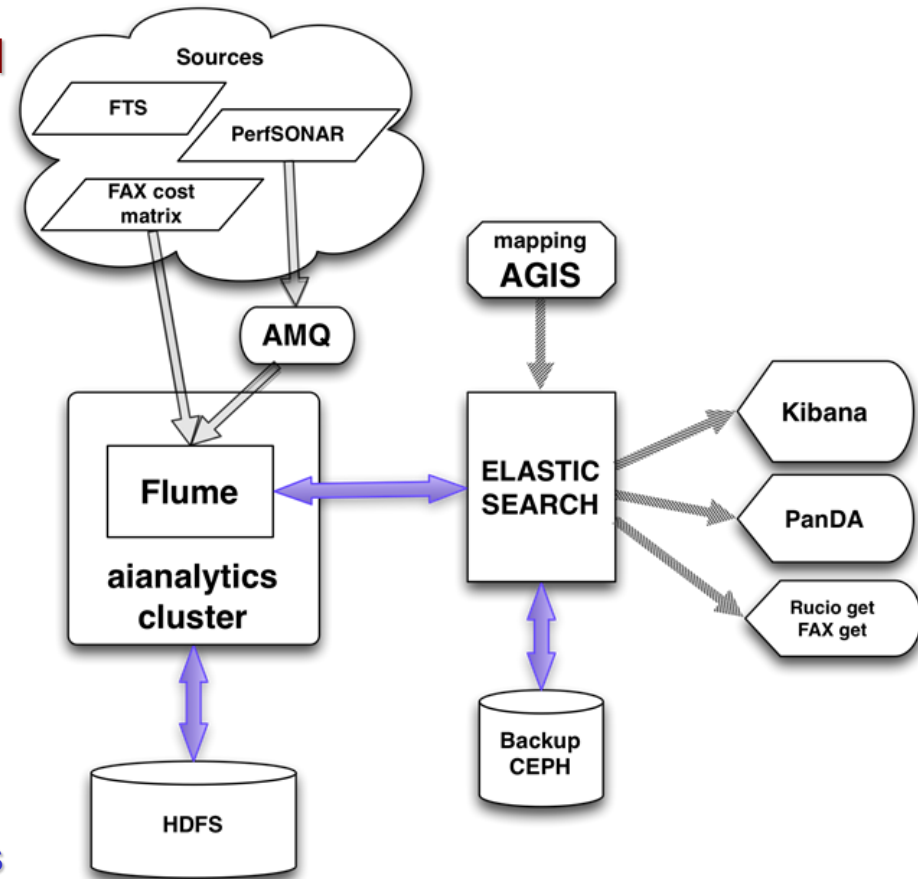
State	Service	Icons	Status detail	Age	Checked	Perf-O-Meter
OK	OSG datastore freshness: bwctl		OK - esmond is complete	22 hrs	21 min	
OK	OSG datastore freshness: trace		OK - esmond is complete	22 hrs	16 min	
OK	perfSONAR configuration: contacts		OK - Contact and organization found	21 hrs	59 min	
OK	perfSONAR configuration: location		OK - Location: -84.4762249/42.7240755	21 hrs	59 min	
OK	perfSONAR configuration: meshes		OK - Meshconfig auto URL configured	21 hrs	59 min	
OK	perfSONAR esmond freshness: bwctl		OK - esmond is complete	22 hrs	12 min	
OK	perfSONAR esmond freshness: trace		OK - esmond is complete	22 hrs	8 min	
OK	perfSONAR esmond freshness: trace rev		OK - esmond is complete	22 hrs	3 min	
OK	perfSONAR hardware check		OK - CPU:2/4cores/1993.710Mhz RAM:4GB NIC:1Gbps/9000MTU/IPv6 enabled	21 hrs	59 min	
OK	perfSONAR json summary		OK - Toolkit metadata successfully retrieved	21 hrs	59 min	
OK	perfSONAR services: bwctl		TCP OK - 0.005 second response time on psum06.aglt2.org port 4823	21 hrs	8 min	
OK	perfSONAR services: http/https		OK - Toolkit homepage reachable	21 hrs	59 min	
OK	perfSONAR services: ndt/npad disabled		OK - NDT/NPAD disabled and not running	21 hrs	59 min	
OK	perfSONAR services: ntp		OK - NTP synchronized	21 hrs	59 min	
OK	perfSONAR services: pscheduler		OK - pScheduler stats retrieved	22 hrs	25 min	
OK	perfSONAR services: regular testing/pscheduler		OK - pScheduler is enabled and running	21 hrs	59 min	
OK	perfSONAR services: versions		OK - Toolkit version found 4.0-0.17.rc3.el6	21 hrs	59 min	

# ATLAS Network Analytics

- ❄ Ilija Vukotic/U Chicago is leading the effort to get network metrics into an analytics platform.
- ❄ This analytics service indexes historical network related data while providing predictive capabilities for network throughput.

## 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
- ❄ Initial “Alarm” query prototyped and tested for Source-Dest paths with high packet-loss



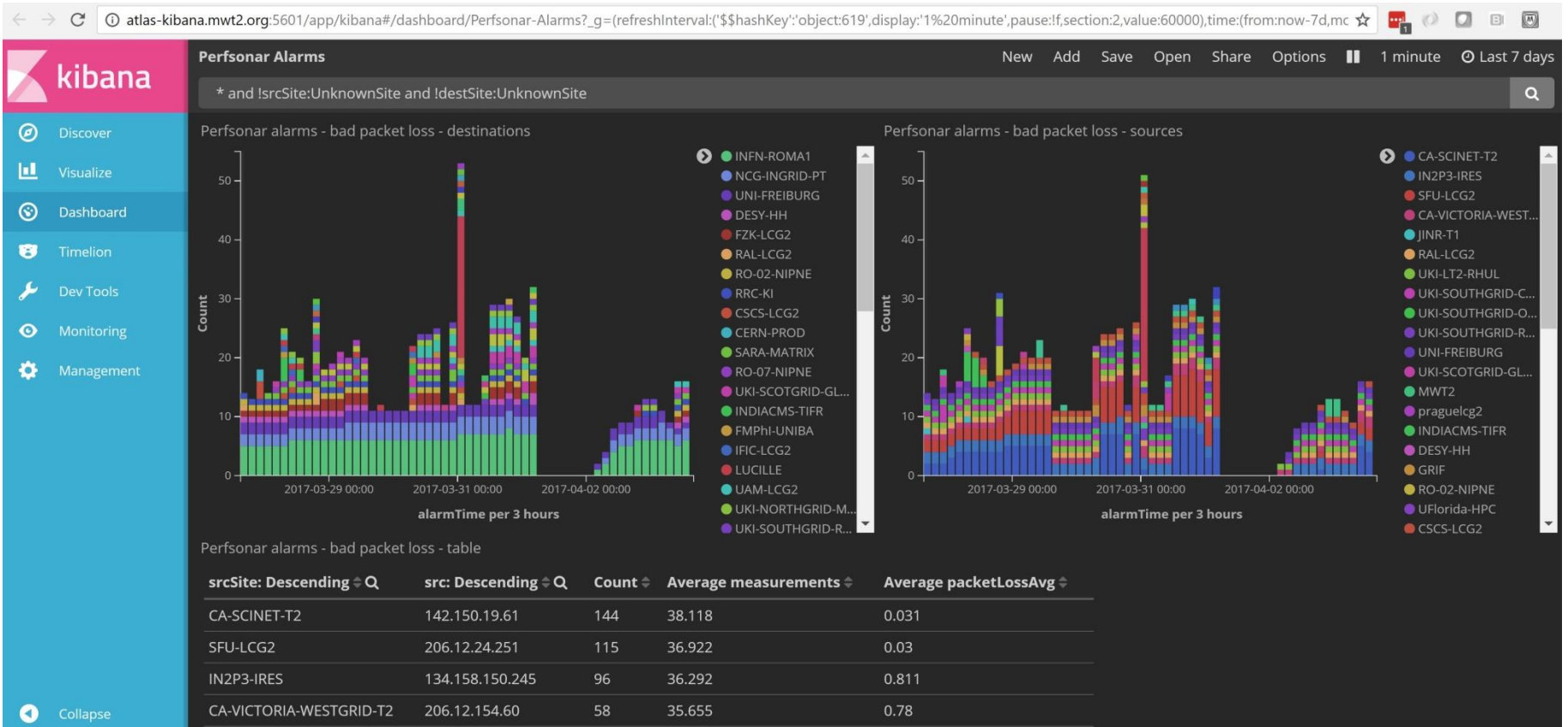
# Recruiting for Exploiting our Data

- ❄ With the current systems we have in place there are many opportunities to “extract value” from the data we collect
- ❄ We would like to encourage anyone with an interest in using the data or trying out new ways of exploiting the data to get involved.
- ❄ There are many possible areas to work in:
  - ❑ Network topology: visualization, analysis, user interfaces
  - ❑ Site analysis: identifying bottlenecks, misconfiguration, problems
  - ❑ Network issues: robust problem identification, localization, alarming
  - ❑ Notification systems: low-noise alerting, customization, user interfaces
- ❄ The next few slides highlight some medium/long term work

# Network Analytics: Packet Loss

❄ Metrics from perfSONAR are sent to Elastic Search

❄ Kibana example of packet loss <http://tinyurl.com/lfaj293>





# R&E Network Data: ESnet Interfaces

❄ Kibana interface info <http://tinyurl.com/lwhtb6r>



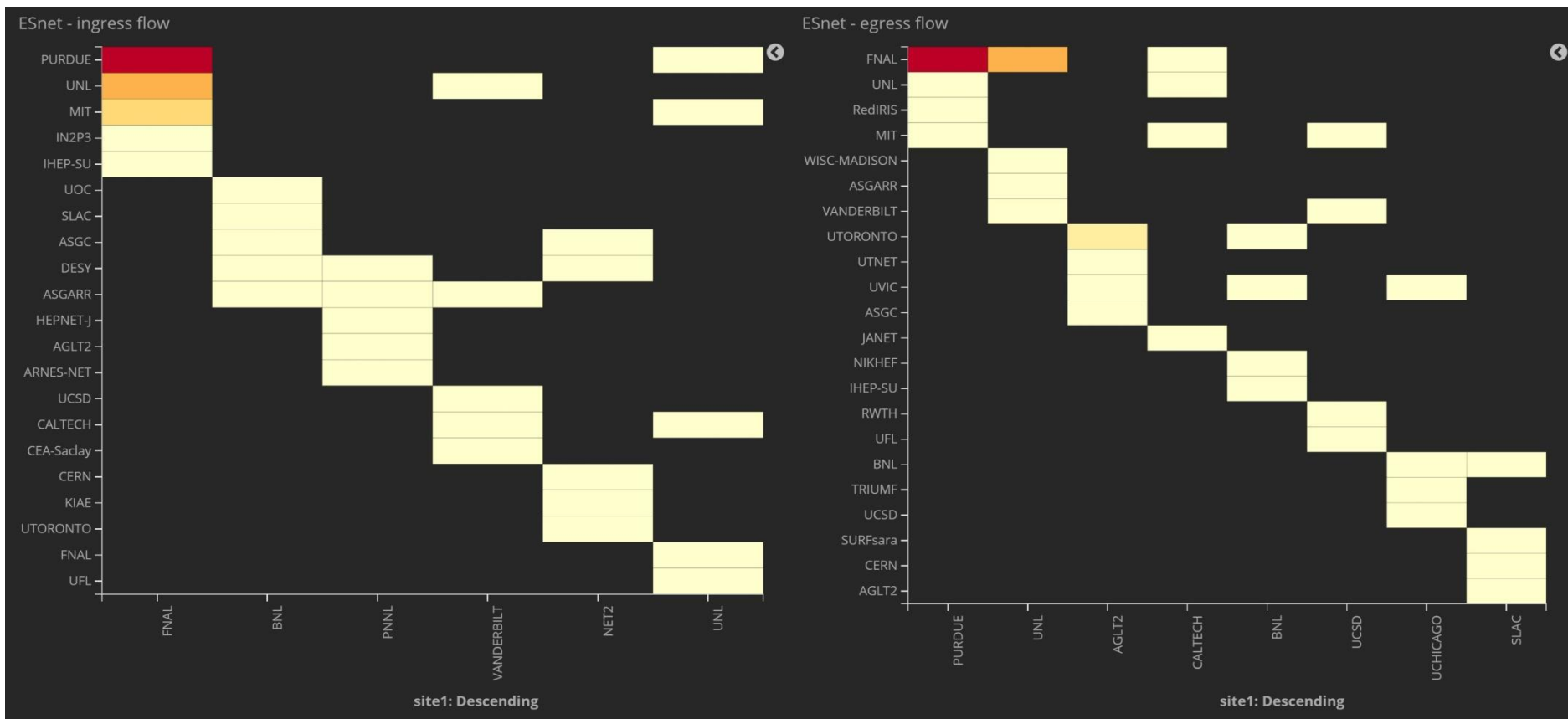
# R&E Network Data: ESnet Flows

❄ Kibana flow info <http://tinyurl.com/kf6ohhc>



# R&E Network Data: ESnet Flow Matrix

❄ Kibana flow matrix info <http://tinyurl.com/kf6ohhc>

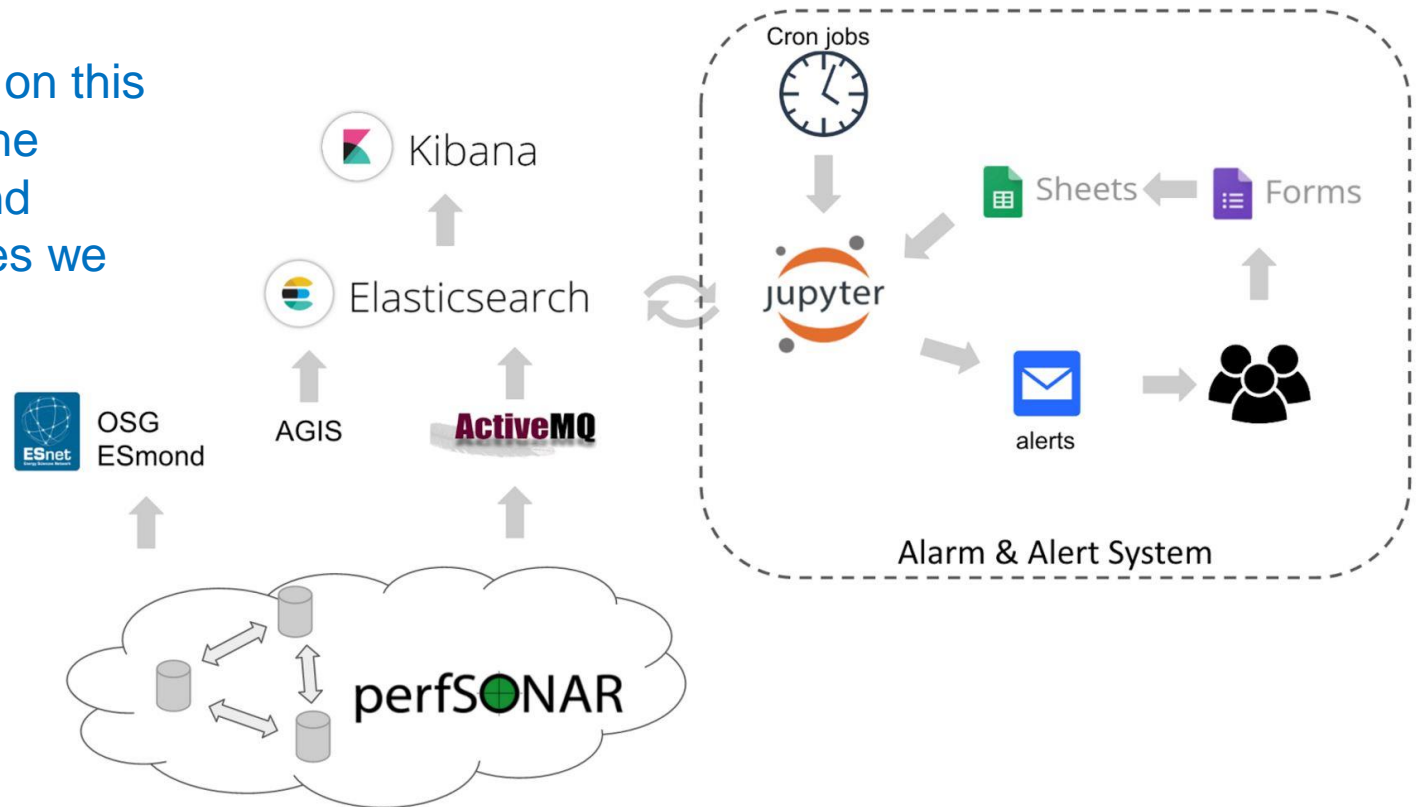


# Getting SNMP/Netflow from R&E Networks

- ❄ As you have seen we have ESnet data being regularly gathered into ElasticSearch.
  - What about additional networks?
- ❄ Within the next two weeks we should also be regularly gathering CERN networking data as well (Marian Babik will clone Ilija's ESnet code to use with CERN data)
- ❄ Can we plan to get GEANT and Internet2 data as well?
  - Let's talk over coffee or dinner!

# Prototype Network Alerting / Alarming

The system shown on this diagram is one of the possible alerting and alarming possibilities we are exploring



- Data from perfSONAR reaches Elastic Search which has a co-located Jupyter instance at the University of Chicago.
- Cron jobs run selected queries via Jupyter which find network problems, creating tables of **alarms**
- Higher-level analysis determines when **alerts** should be sent (and to whom)

# Enabling Alarming

- ❄ OSG has a medium term goal of alerting and alarming on network issues and LHCONE/LHCOPN may want to piggy-back on this effort.
- ❄ Milestone completed: technical design of a suitable analysis system based upon existing time-series technologies
- ❄ Current operating implementation gathers all perfSONAR data OSG sends to CERN and puts it in ElasticSearch.
- ❄ Jupyter instance regularly runs cron tasks to analyze data
  - ❑ Near-term goal: anyone can subscribe to simple alert-emails.
  - ❑ Not “production” yet; needs further interface tweaking to make it easy for users to use
  - ❑ Still need to determine where this service (or its equivalent) should be homed (ATLAS specific resources, WLCG?, OSG?)
  - ❑ **Notify Ilija or Shawn if you want to test this**

# Debugging ASGC

- ✧ A ticket was opened [https://ggus.eu/index.php?mode=ticket\\_info&ticket\\_id=119820](https://ggus.eu/index.php?mode=ticket_info&ticket_id=119820) in March 2016:

Hello,

ATLAS recently suffers from various network related problems with access to resources in ASGC. Problems reported via proper channels are solved but soon they reappear. Examples of recent issues: stratum-1 at ASGC (GGUS ticket 119557) ; NDGF - ASGC transfer problems (GGUS ticket 119276). WLCG weekly meeting suggest to open this general ticket and assign it to WLCG Network Monitoring unit, which could use perfonar infrastucture to better understand these problems.

Jiri Chudoba as ATLAS Computing Run Coordinator

- ✧ ASGC connected by dedicated 10Gbit via StarLight and then via AMS to CERN-perfSONAR coverage in Asia quite limited
- ✧ Testing to ASGC from many places showed no more than **500Mbit**; StarLight to ASGC was not more than **200Mbit**, so moved focus to ASGC
- ✧ Asked ASGC for network map and asked them to re-connect their bwctl node directly to border router while keeping latency node on Cisco N7K

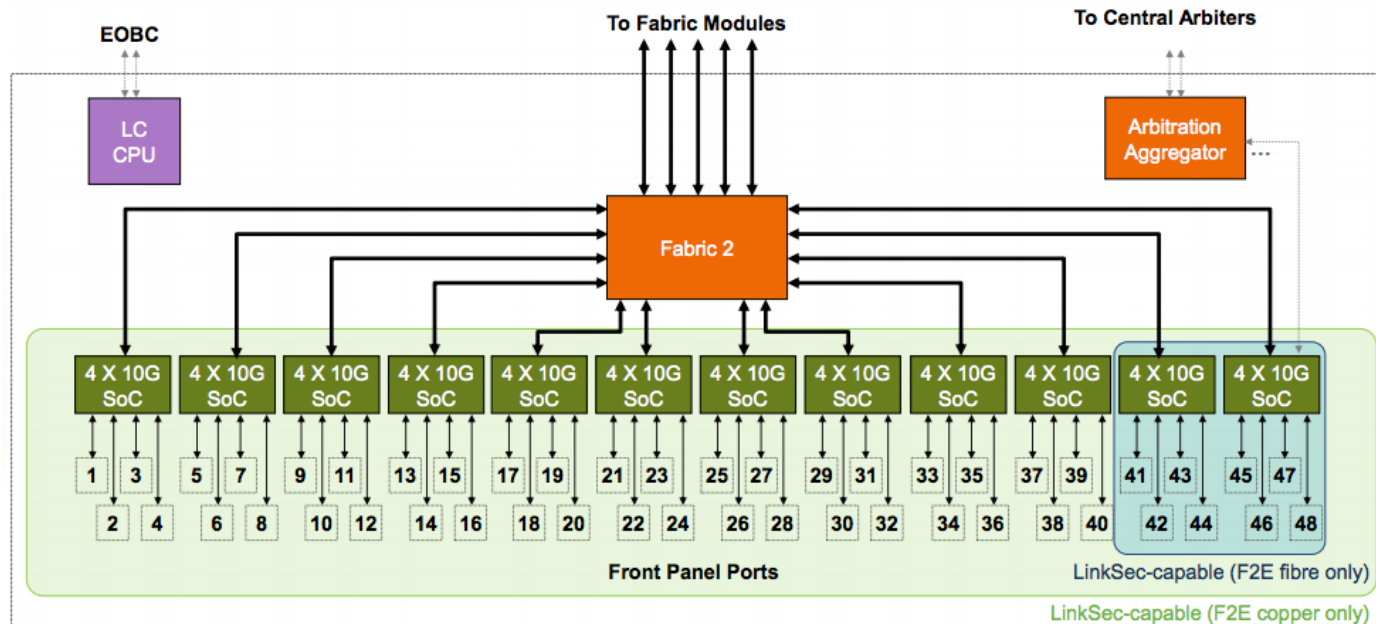
# ASGC Switch Issue Discovered

❄ Focus then moved to Cisco N7K (diagram below), which looks to have just 72MB of memory in total for all 12 ASICs.

- ❑ Testing from StarLight directly to ASGC border router revealed Cisco N7K to be the problem

## 48-Port 1G/10G F2 / F2E I/O Module Architecture

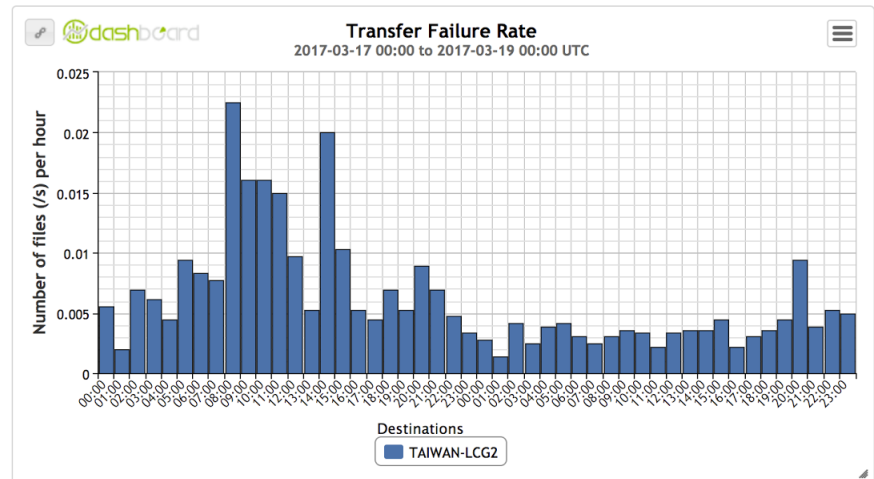
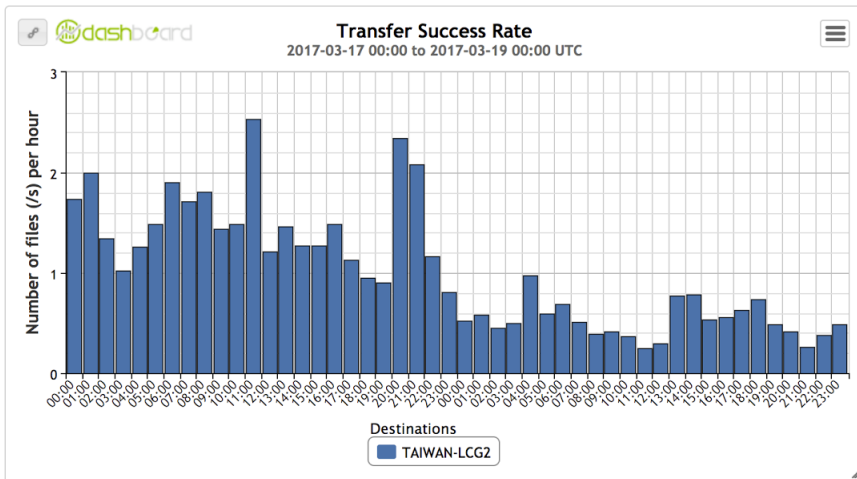
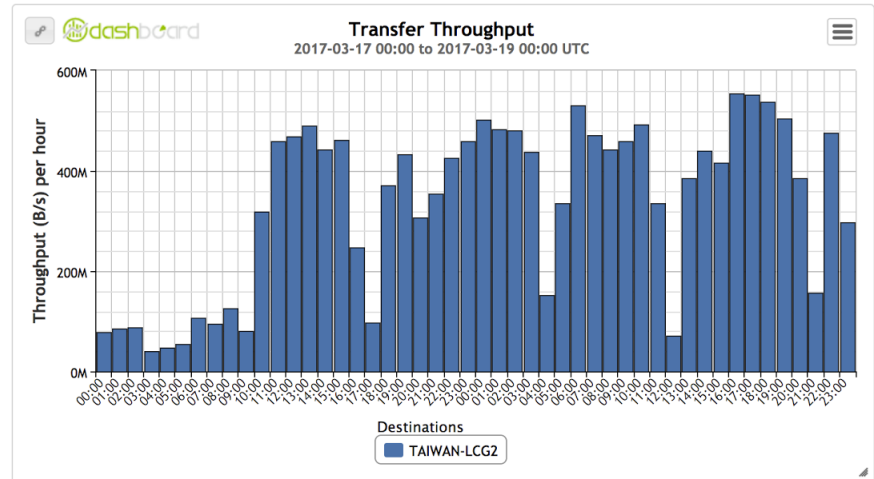
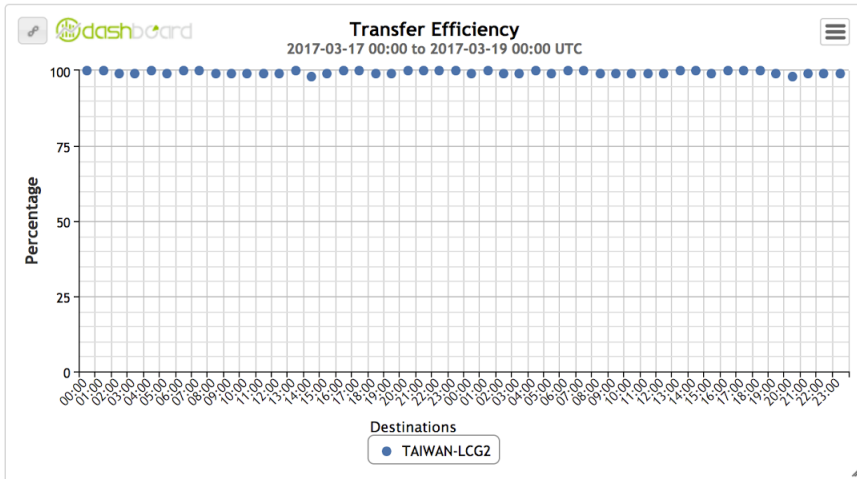
N7K-F248XP-25 / N7K-F248XP-25E / N7K-F248XT-25





# Resolving the ASGC Problem

- Asked ASGC to re-wire so that DTNs are hooked directly to border router
- Re-tested with perfSONAR and asked ATLAS to re-test using Rucio



# ASGC Summary

- ❄ **In summary, investigation took less than 30 days** (we posted recommendation to ASGC in May 2016 and there were some initial delays before we got going)
  - ❑ perfSONAR was critical to achieving this
  - ❑ We managed to gain factor 10, we're at ~ **4.4Gbps** now
  - ❑ Can we get better ?
- ❄ Major achievement for **WLCG** network throughput, shows that at least in Asia (but we have other cases showing that also elsewhere) there are benefits to be made.
- ❄ With respect to **LHCOPN/LHCONE** we're clearly missing guidance on high-latency (distant) sites network architectures and sizing the router buffers
  - ❑ Can we organize some kind of recommendation or point to something suitable for sites to refer to?

# REMINDER: WLCG Support Unit

- ❄ **Reminder:** We have a GGUS support unit (WLCG Network Throughput; [https://wiki.egi.eu/wiki/GGUS:WLCG\\_Network\\_Throughput](https://wiki.egi.eu/wiki/GGUS:WLCG_Network_Throughput)) used to report incidents (mailing list: wlcg-network-throughput at cern.ch)
- ❄ Experiments can report potential network performance incidents.
  - ❑ WLCG perfSONAR support investigates and confirms if this is network related issue.
  - ❑ Once confirmed, it will notify relevant sites and will try to assist in narrowing down the problem to particular link(s). Tracking of ongoing incidents will be via the WG page.
- ❄ Sites observing a network performance problem should follow their standard procedure, i.e. report to their network team and if necessary escalate to their network provider.
  - ❑ If confirmed to be WAN related, WLCG perfSONAR support unit can assist in further debugging. For the policy issues, sites escalate to the WLCG operations coordination.
  - ❑ [https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics#Network\\_Performance\\_Incidents](https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics#Network_Performance_Incidents).
  - ❑ [https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics#Network\\_Throughput\\_Support\\_Unit](https://twiki.cern.ch/twiki/bin/view/LCG/NetworkTransferMetrics#Network_Throughput_Support_Unit)
- ❄ **LHCOPN/LHCONE experts are very important in this coordinated activity.**

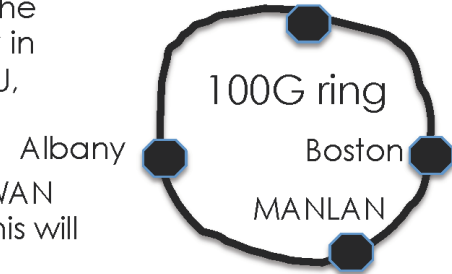
## Other Items

# NET2 Networking & LHCONE

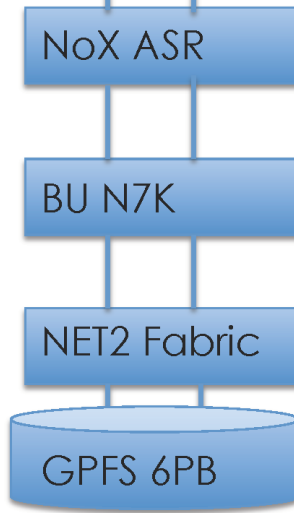
Both the Boston University and the Harvard University components of the U.S. ATLAS Northeast Tier 2 Center are located at the MGHPC facility in western Massachusetts, along with research computing facilities for BU, Harvard, MIT, Northeastern and the UMASS system.

We are in the process of creating a central DMZ for high throughput WAN access and high throughput internal access on the MGHPC floor. This will change how we get to LHCONE over time:

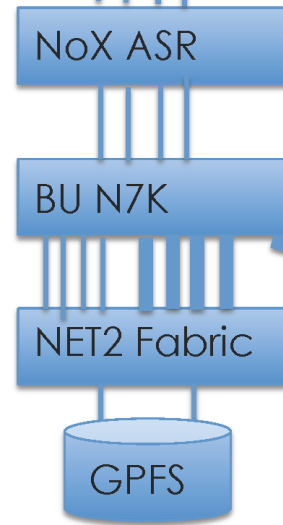
MGHPCC/Chicopee



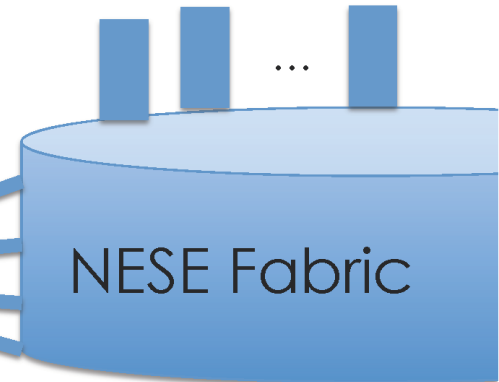
NOW  
LHCONE via tunnel to MANLAN



Coming Soon  
LHCONE via tunnel to MANLAN



LHCONE via tunnel to MANLAN  
N x 100G to NoX



Saul Youssef  
Boston University

**Goal: Federation of storage!**

# Playing with SDN in ATLAS

- ❄ Future networks won't just have larger capacity...
- ❄ A group 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 moved 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
  - ❑ **BENEFIT:** Traffic shaping can result in significantly improved use of the WAN for some paths
- ❄ 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
- ❄ **UPDATE:** Ben Mack-Crane and Galen Mack-Crane are now working with us to document OVS options and test non-disruptive production node migration paths.
  - ❑ *Testbed in-place with SL6.8 and CentOS 7.3 physical nodes and bonded 10G links*
  - ❑ *Results to be posted here: [https://www.aglt2.org/wiki/bin/view/Main/Open\\_vSwitch/InstallOpenvSwitch](https://www.aglt2.org/wiki/bin/view/Main/Open_vSwitch/InstallOpenvSwitch)*

# Summary and Action Items

- ❄️ **We need to plan for a campaign to clear up remaining LHCONE/LHCOPN problems and get v4.0 in place**
  - ❑ Need more instances in Asia in the regional R&E networks!!
  - ❑ New alarming service should help focus on network problems
- ❄️ **New monitoring, management and data analytics capabilities are coming soon.**
  - ❑ Need to setup LHCONE/LHCOPN to use these capabilities.
  - ❑ Need R&E network information from more network providers
- ❄️ **As we fix known issues and get perfSONAR instances reliably operating, we free up time to pursue possible issues in the network itself, rather than the framework that gets us network metrics.**
  - ❑ Looking for volunteers to get involved in network analytics, especially on LHCONE/LHCOPN issues

# Discussion/Questions/Comments?



# 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>

## ❄ New 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>

❑ New mesh config info: <http://meshconfig-itb.grid.iu.edu> Send feedback to Soichi

## ❄ Use-cases document for experiments and middleware

<https://docs.google.com/document/d/1ceiNITUJCwSuOuvbEHZnZp0XkWkwkPQTQic0VbH1mc/edit>