

Deployment of a WLCG network monitoring infrastructure based on the perfSONAR-PS technology

S Campana¹, A Brown², D Bonacorsi³, V Capone⁴, D De Girolamo⁵, A F Casani⁶, J Flix^{7,11}, A Forti⁸, I Gable⁹, O Gutsche¹⁰, A Hesnaux¹, S Liu¹⁰, F Lopez Munoz¹¹, N Magini¹, S McKee¹², K Mohammed¹³, D Rand¹⁴, M Reale¹⁵, S Roiser¹, M Zielinski¹⁶ and J Zurawski¹⁷

¹ CERN, Geneva, CH

² Internet2, Ann Arbor, MI 48104

³ Università di Bologna, IT

⁴ Università di Napoli and INFN, IT

⁵ INFN CNAF, Bologna, IT

⁶ Universidad de Valencia, ES

⁷ Centro de Investigaciones Energeticas, Medioambientales y Tecnologicas, CIEMAT, Madrid, Spain

⁸ The University of Manchester, Oxford Road, Manchester M13 9PL, UK

⁹ University of Victoria, CA

¹⁰ Fermi National Accelerator Laboratory, US

¹¹ Port d'Informacio Cientifica (PIC), Universitat Autonoma de Barcelona, Bellaterra (Barcelona), ES

¹² Physics Department, University of Michigan, Ann Arbor, MI 48109-1040, USA

¹³ University of Oxford, Wellington Square, Oxford, OX1 2JD, UK

¹⁴ Imperial College London, Blackett Lab, Physics Dept, Prince Consort Rd, London SW7 2BW, UK

¹⁵ GARR, IT

¹⁶ University of Rochester, US

¹⁷ ESnet, Lawrence Berkeley National Laboratory, Berkeley, CA 94720

E-mail: simone.campana@cern.ch

Abstract. The WLCG infrastructure moved from a very rigid network topology, based on the MONARC model, to a more relaxed system, where data movement between regions or countries does not necessarily need to involve T1 centres. While this evolution brought obvious advantages, especially in terms of flexibility for the LHC experiment's data management systems, it also opened the question of how to monitor the increasing number of possible network paths, in order to provide a global reliable network service. The perfSONAR network monitoring system has been evaluated and agreed as a proper solution to cover the WLCG network monitoring use cases: it allows WLCG to plan and execute latency and bandwidth tests between any instrumented endpoint through a central scheduling configuration, it allows archiving of the metrics in a local database, it provides a programmatic and a web based interface exposing the tests results; it also provides a graphical interface for remote management operations. In this contribution we will present our activity to deploy a perfSONAR based network monitoring infrastructure, in the scope of the WLCG Operations



Coordination initiative: we will motivate the main choices we agreed in terms of configuration and management, describe the additional tools we developed to complement the standard packages and present the status of the deployment, together with the possible future evolution.

1. Introduction

The WLCG infrastructure has moved from a very restrictive network topology, based on the MONARC model, to a more interconnected system, where data movement between regions or countries does not necessarily need to involve T1 centers. While this evolution brought obvious advantages, especially in terms of flexibility for the LHC experiment's data management systems, it also opened the question of how to monitor the increasing number of possible network paths, in order to provide a global, reliable network service. The perfSONAR [1] network monitoring framework has been evaluated and agreed to as a proper solution to cover the WLCG network monitoring use cases: it allows WLCG to plan and execute latency, traceroute and bandwidth tests between any instrumented endpoint through a central scheduling configuration, it allows archiving of the metrics in a local database, it provides a programmatic and a web based interface exposing the tests results; it also provides a graphical interface for remote management operations.

In this paper we will describe the relevant details of the perfSONAR implementation we choose and present our activity to deploy a perfSONAR based network monitoring infrastructure, in the scope of the WLCG Operations Coordination initiative. We will motivate the main choices we agreed in terms of configuration and management, describe the additional tools we developed to complement the standard packages and present the status of the deployment, together with the possible future evolution.

2. Network Monitoring for WLCG Using perfSONAR-PS

WLCG sites are globally distributed and intrinsically rely upon high-performance networks to function effectively. Experience has shown that end-to-end network issues can be difficult to identify, localize and debug because of insufficient tools and the typically many administrative domains involved in WLCG network paths between sites. One prominent example was the BNL-CNAF network issue documented in https://ggus.eu/ws/ticket_info.php?ticket=61440 (transfer problem between CNAF and BNL). Resolving the issue took over 7 months and logged 72 entries in the ticket. This and many other examples highlight the requirement for WLCG to have a pervasive network monitoring infrastructure deployed able to identify issues when they occur and aid in debugging and localizing problems for quick resolution.

Obviously, some Research and Education (R&E) networks had long ago recognized the difficulties in finding and diagnosing wide-area network (WAN) issues. A consortium of R&E networks and Universities in North/South America and Europe developed the perfSONAR infrastructure to help address such issues. Within WLCG the USATLAS sites began testing and deploying a perfSONAR implementation called perfSONAR-PS in 2007 [2]. By 2010, with its usefulness demonstrated, the LHCOPN [3] choose to also deploy the perfSONAR-PS Toolkit at the CERN Tier-0 and the 10 Tier-1's worldwide. With this experience in place it was logical that WLCG adopt the perfSONAR-PS Toolkit for its network monitoring needs. In fall 2012 a dedicated operations task-force was formed to coordinate and manage deploying perfSONAR-PS at all WLCG sites worldwide.

3. Description of the perfSONAR-PS Toolkit

The perfSONAR-PS Toolkit is the result of an open source development effort based upon perfSONAR [4]. It targets creating an easy-to-deploy and easy-to-use set of perfSONAR services and is available as a bootable all-in-one system (ISO on CD or USB) or can be deployed onto the target host systems local disk via a 'netinstall' process. The current release (v3.x) is based upon CentOS6 and is available in 32 and 64-bit architectures.

Installation of the toolkit is recommended to be on dedicated physical hardware to ensure other applications or services don't interfere with network metrics and diagnostics provide by perfSONAR-

PS. For ease of maintenance the recommend deployment method is “netinstall” to local disk which utilizes a YUM repository for perfSONAR-PS components. Updates and security fixes are then easily available via ‘yum update’. Because the components are available via RPM it is also possible to integrate installing perfSONAR-PS with a site’s local provisioning and configuration management systems, though this obviously requires significantly more work for the person deploying. It should be noted that while it is not recommended, some sites have deployed the perfSONAR-PS toolkit on virtualized machines (VMs).

3.1. Details of the perfSONAR-PS Toolkit

Once the perfSONAR-PS Toolkit is deployed on a host it provides a number of features:

- a web based GUI allowing the administrator to configure the host and services running, define and schedule tests and display measurement results.
- a set of services that provide throughput tests (BWCTL), ping tests (PingER), latency/pack-loss tests (OWAMP), traceroute and two on-demand diagnostics (NDT,NPAD).
- a measurement archive which stores all test results and makes them available through an API.
- an integrated Cacti instance that can monitor and graph host metrics and could be configured to query, monitor and graph local network equipment via SNMP.

A perfSONAR-PS toolkit instance is typically configured to be either a “bandwidth” or a “latency” instance. It is not recommended to run both types of services on the same node because the bandwidth tests interfere with the latency tests and can cause false-positives indications of network problems.

4. perfSONAR-PS deployment in WLCG

4.1. Early phase of perfSONAR deployment

The first deployment of perfSONAR in the WLCG infrastructure started in the scope of monitoring the OPN performance. All the 11 WLCG T1s and CERN have been instrumented with two distinct perfSONAR hosts (one for bandwidth and the other for latency/packet loss tests) and configured to execute periodic measurements among all sites. At this stage there was no tool for automatic configuration: each service manager had to specify the hostnames of target services and configure the test parameters. All information was documented in static and text-based web pages, and kept up-to-date by a responsible person. This deployment model worked very well for a small and rather immutable system such as the LHCOPN and successfully provided important matrices of test results that were used to improve the stability of the system. The model was originated in USATLAS, which started a campaign of deployment of perfSONAR-PS across T2 sites in 2007. As more sites and regional clouds joined in deploying perfSONAR-PS the limitations of the original model become obvious: a trivial change at some sites (installation of a new host, change of alias or hostname) would have to be reflected in a manual reconfiguration change at all other sites. The effort become almost unmanageable when it was agreed to utilize perfSONAR to monitor the performance of the LHCONE [5] infrastructure, with many new sites joining every month. At the same time, most experiments understood the importance of a network monitoring infrastructure and asked for a broad deployment at all sites. In Fall 2012 it was therefore agreed to consolidate the perfSONAR deployment effort in WLCG and a dedicated task force was initiated under the scope of the WLCG Operations Coordination activity.

4.2. Full scale WLCG perfSONAR deployment

The task force took a pragmatic approach and decided to start by focusing on the sites which were of primary interest from experiments. A list of 86 sites has been initially collected from ATLAS, CMS and LHCb. The sites have been organized in regions, mostly determined by their geographical location and the research network to which they belong, with few exceptions. All sites have been asked to deploy a first perfSONAR-PS instance to treat bandwidth tests and a second one to treat all other tests

(ping, packet loss, traceroute). Both instances had to be registered in the WLCG service registries (GOCDB or OIM for OSG sites) following given instructions.

The perfSONAR developers worked with USATLAS and later with WLCG to provide the ability to centrally define sets of related perfSONAR-PS Toolkits via the use of the so-called “mesh-configuration”. The idea is shown in Figure 1 and illustrates in perfSONAR-PS Toolkit instance reading from specific URLs that provide host and test configuration details in JSON format. In addition, some services were not as robust as needed for a required WLCG service and we worked closely with the Toolkit developers to improve the overall resiliency through a combination of bug-fixes and configuration changes.

Thanks to the mesh-configuration functionality, the deployment model has then been centralized: benefitting from the improvements in the perfSONAR-PS toolkit, which from release 3.3.1 do allow to import the configuration from a web based URL, all configurations have been stored in a single location available through HTTP. The configurations would include the definition of hosts, regions, test parameters and test scenarios. The test scenarios are based on the concept of “mesh” meaning a group of sites to which a test definition can be assigned: all sites of the mesh would trigger the given test with the parameters specified in the definition. This powerful scheme allowed defining different test scenarios for sites within the same region and across different regions. In particular it was agreed to configure:

- bandwidth tests (with a length of 30 seconds) every 6 hours for sites within the same region, every 12 hours between a T2-T1 pair of different regions and every week for any other combination;
- latency tests (packet loss), sending 10Hz of packets across all WLCG sites and measuring the loss based on batches of 600 packets (every minute);
- traceroute tests between all WLCG sites once per hour;
- ping tests between all sites every 20 minutes;

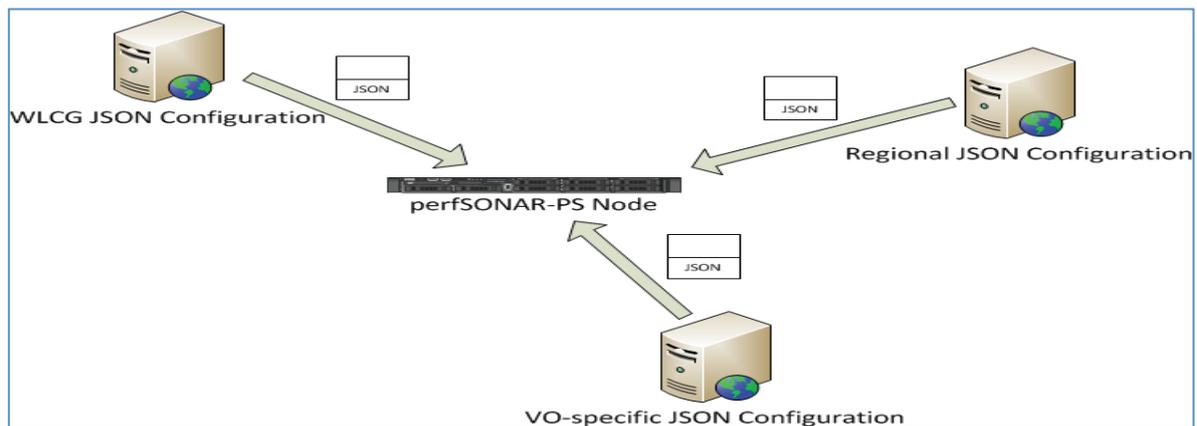


Figure 1: an illustration of the perfSONAR mesh configuration concept

Such configuration was agreed in order to achieve statistically meaningful network measurements without disrupting the overall production network activities.

4.3 Data aggregation and data mining

The raw measurements provided by the various tests are stored locally at each perfSONAR-PS

US cloud throughput measurement											
	---	0	1	2	3	4	5	6	7	8	9
0:BNLBNL-Test (lhcmon.bnl.gov)	---	2.29	2.45	1.54	2.38	2.94	3.84	1.31	0.04	1.17	---
1:AGLT2 (psmsu02.aglt2.org)	1.01	---	0.98	0.31	1.14	0.94	0.41	0.72	0.02	0.12	---
2:AGLT2 (psum02.aglt2.org)	1.07	1.24	---	1.31	2.80	2.04	1.21	1.90	0.02	0.00	---
3:MWT2 (iut2-net2.iu.edu)	1.03	0.03	0.04	---	6.50	6.00	0.92	0.93	0.05	0.12	---
4:MWT2 (mwt2-ps02.campuscluster.illinois.edu)	1.98	0.77	0.59	0.27	---	5.32	0.92	0.93	0.01	0.12	---
5:MWT2MWT2(UC) (uct2-net2.uchicago.edu)	2.88	1.03	1.84	3.34	5.19	---	0.92	0.84	0.02	0.14	---
6:NET2 (atlas-npt2.bu.edu)	1.03	1.85	1.23	0.72	1.82	2.31	---	1.57	0.01	0.14	---
7:SWT2 (ps2.occhep.ou.edu)	1.42	1.37	0.76	0.20	0.42	2.43	1.71	---	0.03	0.27	---
8:SWT2 (psuta02.atlas-swt2.org)	1.44	0.01	0.00	0.09	0.11	0.06	0.22	0.05	---	0.01	---
9:WT2 (psnr-bw01.slac.stanford.edu)	1.62	1.50	0.51	0.40	0.00	0.54	0.51	1.71	0.00	---	---

Figure 2: The perfSONAR Modular Dashboard throughput measurements for the USATLAS region

Topology						FTS transfer (per file)				perfSONAR					WAN data access (WN-SE)			
SecSite	SecCloud	SecTier	DetSite	DetCloud	DetTier	DDM Sonar				perfSONAR								
						AvgBRS (MB/s)	EvS	AvgBRM (MB/s)	EvM	AvgBRL (MB/s)	EvL	MinThr (MB/s)	AvgThr (MB/s)	MaxThr (MB/s)		MinPL	AvgPL	MaxPL
Talco-LOG2	TW	T1	RAL-LOG2	UK	T1	0.51+/-0.63	285	7.25+/-5.02	336	7.88+/-5.47	649	0.3	0.3	0.3	0.0	65.7	329.0	n/a
Talco-LOG2	TW	T1	INDP-CC	FR	T1	0.52+/-0.66	55886	6.34+/-2.94	6121	16.10+/-6.07	1617	0.5	0.5	0.5	600.0	600.0	600.0	n/a
TRUMF-LOG2	CA	T1	Talco-LOG2	TW	T1	0.41+/-0.41	400	1.25+/-0.24	38	2.88+/-1.30	5	0.4	0.5	0.6	0.0	0.0	1.0	n/a
pc	ES	T1	Talco-LOG2	TW	T1	0.04+/-0.09	162	0.00+/-0.00	0	0.00+/-0.00	0	0.3	0.6	0.8	0.0	0.0	0.0	n/a
FZK-LOG2	DE	T1	Talco-LOG2	TW	T1	0.17+/-0.24	1178	1.01+/-0.23	505	16.93+/-11.49	5	0.5	1.3	2.2	0.0	0.0	0.0	n/a
BNL-ATLAS	US	T1	RAL-LOG2	UK	T1	0.29+/-0.51	45183	3.71+/-1.71	2697	21.06+/-15.41	879	1.5	1.7	1.9	0.0	18.2	229.0	n/a
Talco-LOG2	TW	T1	FZK-LOG2	DE	T1	0.83+/-1.08	280	4.70+/-2.82	36	16.37+/-9.10	125	1.9	2.0	2.3	0.0	0.1	2.0	n/a
INFN-T1	IT	T1	Talco-LOG2	TW	T1	0.29+/-0.46	540	1.87+/-0.76	6	0.00+/-0.00	0	1.7	2.0	2.3	0.0	0.0	0.0	n/a
pc	ES	T1	RAL-LOG2	UK	T1	0.61+/-0.31	5202	6.32+/-2.22	216	20.80+/-9.55	4	1.5	2.4	2.5	0.0	57.5	357.0	n/a
BNL-ATLAS	US	T1	INDP-CC	FR	T1	1.63+/-2.05	101375	15.30+/-6.78	28627	39.26+/-12.94	5481	2.5	3.3	4.4	0.0	0.0	0.0	n/a
NOGF-T1	NO	T1	Talco-LOG2	TW	T1	0.09+/-0.13	4488	1.40+/-0.62	67	19.33+/-0.81	5	3.7	3.8	4.3	0.0	0.0	0.0	n/a
INDP-CC	FR	T1	Talco-LOG2	TW	T1	0.36+/-0.57	4641	3.58+/-2.00	3840	9.12+/-6.52	1067	3.3	4.2	5.3	0.0	0.0	0.0	n/a
FZK-LOG2	DE	T1	RAL-LOG2	UK	T1	0.47+/-0.74	70705	7.44+/-6.32	7598	14.03+/-16.74	6770	2.8	4.4	9.9	0.0	24.2	193.0	n/a
RAL-LOG2	UK	T1	Talco-LOG2	TW	T1	0.06+/-0.19	13355	0.96+/-0.34	528	0.00+/-0.00	0	6.7	6.7	6.7	0.0	0.6	5.0	n/a

Figure 3: The ATLAS Site Status Board network measurements view

instance in a MySQL database. The information are then collected, aggregated and stored in the perfSONAR Modular Dashboard [6]. This dashboard provides a GUI for quick visualization of the test results, with aggregation by site and/or region. For example Figure 2 shows the matrix of perfSONAR bandwidth tests between USATLAS sites, as displayed in the perfSONAR dashboard. Additionally, the dashboard exposes the results programmatically via HTTP. Various applications collect and consume the JSON output and consume the numbers for different use cases. For example Figure 3 shows the ATLAS Site Status Board [7], where perfSONAR measurements of throughput and packet

loss are utilized to complement the network measurement information coming from file transfer statistics between two storages and between storage and Worker Node.

5. Conclusions

The status of the perfSONAR-PS deployment in WLCG is graphically summarized in Figure 4: 81 sites have deployed perfSONAR-PS instances and published in GOCDB/OIM (3 of them need attention concerning the publication), while 31 site are still missing. The deployment issues are clustered in regions where the WLCG perfSONAR deployment misses a strong contact person. In the near term we plan to work in parallel on three different areas:

- The perfSONAR deployment will be completed at the missing sites. Additionally, we need to verify that existing instances are running the most recent version and import correctly the central configuration.
- We will instrument service monitoring and, in future, alarming in case an instance fails to schedule tests or in general starts malfunctioning
- We will work with the perfSONAR Modular Dashboard developers to rationalize and improve the various views and configuration of the tool.



Figure 4: The perfSONAR deployment status in WLCG

References

- [1] Hanemann A., Boote J., Boyd E., Durand J., Kudarimoti L., Lapacz R., Swany M., Trocha S., and Zurawski J., “PerfSONAR: A Service-Oriented Architecture for Multi-Domain Network Monitoring”, International Conference on Service Oriented Computing (ICSOC 2005), Amsterdam, The Netherlands, 2005
- [2] Tierney B., Metzger J., Boote J., Brown A., Zekauskas M., Zurawski J., Swany M., Grigoriev M., “perfSONAR: Instantiating a Global Network Measurement Framework”, 4th Workshop on Real Overlays and Distributed Systems (ROADS’09) Co-located with the 22nd ACM Symposium on Operating Systems Principles (SOSP), January 1, 2009
- [3] <http://lhcopn.web.cern.ch/lhcopn/>
- [4] Zurawski J., Balasubramanian S., Brown A., Kissel E., Lake A., Swany M., Tierney B., Zekauskas M., “perfSONAR: On-board Diagnostics for Big Data”, 1st Workshop on Big Data and Science: Infrastructure and Services Co-located with IEEE International Conference on Big Data 2013 (IEEE BigData 2013), October 6, 2013
- [5] <http://lhcone.net/>

- [6] The perfSONAR Modular Dashboard projec: <https://github.com/PerfModDash/PerfModDash>
- [7] The ATLAS SSB at <http://dashb-atlas-ssb.cern.ch/dashboard/request.py/siteview>