# perfSONAR: On-board Diagnostics for Big Data

J. Zurawski*, S. Balasubramanian*, A. Brown‡, E. Kissel†, A. Lake*, M. Swany†, B. Tierney* and M. Zekauskas‡

*Energy Sciences Network (ESnet)
Lawrence Berkeley National Laboratory, Berkeley, CA, USA
Email: {zurawski, sowmya, andy, bltierney}@es.net
†Indiana University
School of Informatics and Computing, Bloomington, IN, USA
Email: {ezkissel, swany}@indiana.edu
‡Internet2
Ann Arbor, MI, USA
Email: {aaron, matt}@internet2.edu

*Abstract*—Big science data necessitates the requirement to incorporate state-of-the-art technologies and processes into science workflows. When transferring "big data", the network infrastructure connects sites for storage, analysis and data transfer. A component that is often overlooked within the network is a robust measurement and testing infrastructure that verifies all network components are functioning correctly. Many researchers at various sites use perfSONAR[1], a network performance measurement toolkit to isolate many types of network problems that reduce performance. perfSONAR is an essential tool that ensures scientists can rely on networks to get their data from end-to-end as quickly as possible.

## I. INTRODUCTION

Innovation can often be disruptive to "business as usual". Many scientific disciplines are beginning to develop innovative processes to modify traditional operational workflows, in an effort to adopt data-intensive methodologies. As an example, the field of genomics has experienced a monumental decrease in the size and cost of sequencing technology, along with a subsequent increase in data accuracy. This trend is illustrated by the graph shown in Figure 1. Older sequencing technology was prohibitively expensive, large in size, and incapable of producing finely detailed results; the emerging genomics technologies have facilitated a move toward sequencer deployment in smaller facilities, with fewer researchers required, yet are still capable of producing large data sets. While this has created economic incentive to purchase new technology, it does neglect another crucial component in the scientific workflow: the ability to analyze and store results that are produced.

Computational components, in the form of cluster or supercomputing resources require power, cooling, and access
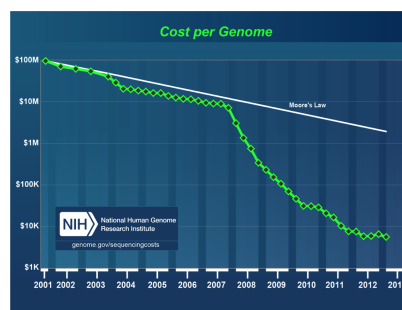


Figure 1: Cost required to sequence a genome in relation to Moore's law [14].

to fast, efficient networks to be most effective. Shared resources, such as those provided by Grids, Clouds, and dedicated facilities like computing centers funded by the NSF[2] and DOE[3], remain popular with domain researchers of all types who find it infeasible to operate private infrastructure.

With the advent of high-speed networks, and accompanying software designed to efficiently broker the migration of research data, it is possible for users of all levels of sophistication to integrate remote analysis and storage into their scientific workflow. The U.S. Department of Energy's Energy Sciences Network (ESnet)[4] has studied scientific network patterns for a number of years. A plot[5] of historic network traffic illustrates a need for an effective "conveyor belt" for science; researchers will be buried in the deluge of data that will arrive in the near future as they turn observational data into analyzed results at an accelerated pace, over great distances, and involving numerous collaborators (see Figure 2). Networks are indeed a critical cog in this machinery, and must be working at peak efficiency with adequate capacity to ensure success.

[1] perfSONAR-PS. http://psps.perfsonar.net

[2] National Science Foundation. http://www.nsf.gov

[3] Department of Energy Office of Science. http://science.energy.gov

[4] ESnet. http://www.es.net

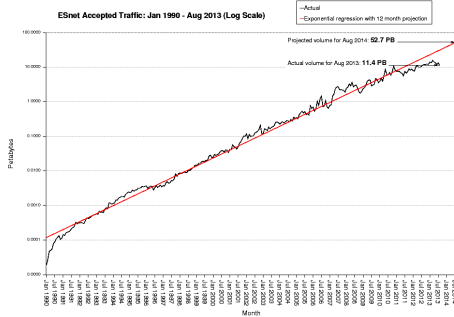[5] ESnet Statistics. http://stats.es.net/top.html

Figure 2: ESnet historic network traffic. $11PB$ of aggregate network traffic was observed in 2013. The one year projection estimates the need to handle four times this.

The complexity of computer networks can make troubleshooting problems difficult. A misconfigured device or a physical abnormality can introduce loss or corruption that looks like loss that will lead to performance degradation anywhere in this shared infrastructure. Devices inserted to protect the network can also limit performance. Performance limits are one kind of "friction" to effective use of the network. Users that perceive the network as unreliable, whatever the actual reason, will learn to mistrust the resource. This perception has caused many collaborations to feel that bulk shipment of storage via the postal system can deliver more throughput than a modern network.

Traditional science applications, including those that migrate data from acquisition site to analysis facilities, are known to rely on the transmission control protocol (TCP) [16] for numerous use cases. TCP is robust in many respects, however, the very mechanisms that make TCP so reliable also make it perform poorly when network conditions are not ideal [3], [11], [13]. In particular, TCP interprets packet loss as network congestion, and reduces the "sending" rate when loss is detected: even a tiny amount of packet loss is enough to dramatically reduce TCP performance and draw out network use from minutes to hours to potentially days. Thus, all the networks in paths that support data-intensive science must strive to provide TCP-based applications with loss-free service, if these applications are to perform well in the general case.

Operational soundness is a high priority for the maintainers of these networks, particularly when there are science drivers pushing the overall network design. The Science DMZ[6], a design paradigm developed by ESnet, has been adopted by numerous institutions as a method to reduce the overall friction that is known to exist in modern converged network designs [6]. Architectural and technical choices will lead to performance gains for network users. This paradigm is featured as a simple block diagram in Figure 3; complexity

---

[6]Science DMZ. http://fasterdata.es.net/science-dmz/

---

has been reduced from typical network deployment choices. Along with simplified design and operation, the paradigm features a rich set of diagnostic abilities to ensure proper operation over time.

Traditional monitoring tools have not scaled beyond the administrative boundary of a domain for a variety of reasons: it may not have been a requirement in the original design, or policy constraints outside of the control of the tool may limit desired functionality. A loose coupling between deployed tools is often required: there must be enough control to enable sharing of policy and measurement. perfSONAR is an innovative federated monitoring tool designed with multi-domain operation as the core principle [10]. This framework inserts a layer of middleware between the measurement tools, and user facing products such as graphical interfaces and alarming systems. Policy (e.g. who can view measurements, who can request them), location and discovery, and a data abstraction layer to normalize the output of diverse tools so they can be consumed and analyzed in a coherent and useful manner [18], [20] are all provided via perfSONAR. perfSONAR is unique in that the combined sum of functionally is only possible via the contributions of individual tools. perfSONAR is a powerful component in identifying and removing "friction" from networks.
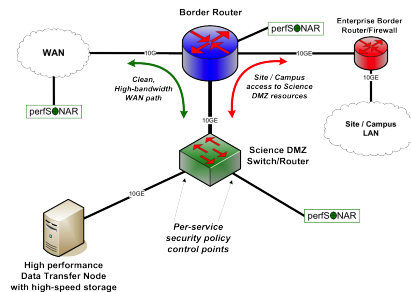


Figure 3: ESnet's Science DMZ architectural design pattern.

The remainder of this paper will proceed as follows: Section II will introduce perfSONAR as a network monitoring solution that has a broad appeal to operations staff as well as scientific end users. Section III discusses suggested use of the monitoring tools. Section IV will present specific use cases of the perfSONAR framework, related to scientific operations, and used in conjunction with related approaches to modify network architectures. Finally, Section V discusses related work, including work that leverages the perfSONAR framework.

## II. perfSONAR Software

Performance monitoring is critical to the discovery and elimination of so-called "soft failures" in the network. Soft failures are problems that do not cause a complete failure that prevents data from flowing, but that cause perceived poor performance. Examples of soft failures include packet

loss due to failing components, routers forwarding packets using the management CPU, or inadequate configuration of network devices. Soft failures often go undetected for many months or longer, as most network management and error reporting systems are designed for reporting "hard failure", such as loss of a link or device.

perfSONAR is a service oriented approach to performance monitoring. Functionality is divided into individual components; many of which work on their own but are also designed to work in harmony with each other and with remote instantiations controlled by others. Federation is a crucial design pattern, and facilitates the software as being an "end-to-end" way to monitor, diagnose, and solve network performance issues [9].
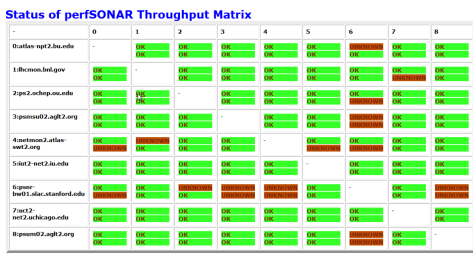


Figure 4: Network performance dashboard based on perfSONAR as used by the ATLAS collaboration.

The perfSONAR system can run continuous checks for latency changes and packet loss, as well as periodic "throughput" tests (a measure of available network bandwidth). An example of this periodic probing is shown in Figure 4, and utilized by the ATLAS physics collaboration[7] to monitor network performance between participating facilities in their collaboration. If a problem arises that requires a network engineer to troubleshoot the network infrastructure, the tools necessary to work the problem are already deployed [5], [12].

To illustrate the effectiveness of perfSONAR, consider the case where a network device is experiencing a small amount of packet loss. The problem has gone unnoticed by the local staff, and is really only manifested for use cases that require large amounts of capacity or via use cases that span great latencies. Now let's assume a remote user, one that is located several domains away from problem area, wishes to access a scientific resource in the form of a long-lived file transfer: because of the size and longevity of the task the user will be impacted by this performance abnormality, and they will be left with many questions:

- Is his data movement software working correctly?
- Are the hosts involved (e.g. both his local resources, and those at the scientific repository) "tuned" for the task at hand?

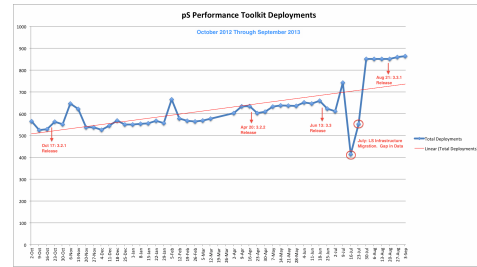[7]USATLAS Dashboard. https://perfsonar.racf.bnl.gov:8443/exda/?page=25&cloudName=USATLAS



Figure 5: perfSONAR deployement growth since October 2012

- Is the network functioning correctly? If not, how can we figure out which network has a problem when the path involves several domains, and possible hundreds of devices?

perfSONAR can address these three questions with a variety of techniques. perfSONAR contains measurement tools such as BWCTL[8], OWAMP[9], and NDT[10], which are designed to emulate the behavior of common network activities such as bulk data movement and video transfer. The results from such measurement tests can be extrapolated to the use case of a typical scientist. If the measurement tools behave poorly along with the scientific application, it is an indication that the host or network may be malfunctioning.

A constellation of deployed perfSONAR instances located at key network junctures can be used to tests the users path on an end-to-end basis. In the previous example, the user could deploy the tools directly to their resources via easy-to-install packages. The remote site could do the same, or could allocate an entire "Performance Node"[11] to be used for long term monitoring functionality. Networks in the middle may have similar test nodes available. Debugging becomes an exercise in path verification for end-to-end and end-to-middle paths until the data loss that is impacting network performance can be found, and corrected. As shown in Figure 5, the number of deployments has steadily grown over the past year, and trends suggest this will continue.

## III. DEPLOYMENT STRATEGIES

perfSONAR works best when it is available along network paths. A robust deployment strategy depends on the type of network that is involved. For instance, having a performance tester located near the scientific collaborators is the most sensible deployment strategy. Equally, locations where traffic intermingles, e.g., exchange points or the core of a university campus, is also an important region to monitor. Backbone providers with several points of presence (PoPs) could

[8]BWCTL. http://www.internet2.edu/performance/bwctl/

[9]OWAMP. http://www.internet2.edu/performance/owamp/

[10]NDT. http://www.internet2.edu/performance/ndt/

[11]pS Performance Toolkit. http://psps.perfsonar.net/toolkit

make testing resources at each, as a service to downstream customers [2].

Collaborations are often well formed and feature regular traffic patterns that relate to the workflow. For example, if data is captured at one facility, but must be processed at others, a regular pattern of data exchange will exist between these actors, and thus there is a need for measurement activities to ensure proper operation. Many operations groups, such as XSEDE[12] and members of large collaborations, such as the LHC, recommend a measurement schedule that is a "full mesh", e.g., all sites test to all other sites several times a day. This builds a history of performance results, and allows for easy correlation against expected values.

Network metrics vary, and can tell different characteristics of behavior. For instance, "achievable bandwidth" is a measurement of how a well behaved application could expect to perform on a given network segment when current conditions are considered, including the network capacity, congestion, and factors on the host and operating system. Tools such as iperf[13] and nuttcp[14] are designed to exercise this particular metric. Latency, a lighter weight yet still important measurement of the time required to traverse network links, is useful for applications that have "real-time" sensitivities. Latency can be measured in terms of a round-trip time (e.g., through the "ping" tool) or on a one-way basis (e.g. by using OWAMP). Finally, a measurement of packet-loss, as seen by either applications or the network devices themselves can be provided by passive measurement mechanisms like SNMP or active tools like OWAMP. These metrics tell an important story individually about the realities of a network (end-to-end or individual segments), but are most useful when interpreted together.

## IV. Scientific Case Studies

To highlight the utilitarian nature of perfSONAR when used in a diverse networking environment, we present two use cases that demonstrate the necessity of regular monitoring when handling data-intensive science requirements. While these use cases show problems discovered by manual examination of data, the diagnostic information delivered via this framework forms the basis for future advancements that could be used fo fully automate diagnostics. Analysis frameworks, such as On-Time-Detect [4], are capable of consuming raw perfSONAR data from distributed sources and are a closer step to machine guided network repair. The following examples illustrate that scientific use cases can be fragile, and require stable and reliable networks to function properly.
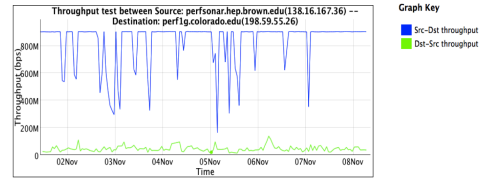


Figure 6: Network performance via the BWCTL tool. The placement of this test server mimicked that of researchers traversing the site firewall.

### A. Brown University Physics Department

Brown University[15] is the home to numerous research groups. Their high energy physics group[16] participates in the Large Hadron Collider (LHC) Compact Muon Solenoid (CMS)[17] experiment. Physicists from the university routinely download data sets from remote locations (e.g. Fermilab[18] in the United States, or directly from the LHC site at The European Organization for Nuclear Research, CERN). Most of the data sets being downloaded can range in size from hundreds of gigabytes to several terabytes.

The physics group at Brown requires a stable network, and observed through perfSONAR monitoring, shown in Figure 6, that performance into the university from remote sites was more than an order of magnitude below the performance outbound. Additional testing and analysis of the network found that the campus security devices were incapable of handling the needs of data-intensive science occurring on the campus. An open question for the campus emerged: how can security be implemented in a sensible manner, and yet not impact the requirements of the scientific community by impeding network performance?

The campus adopted the approaches recommended by the Science DMZ design pattern in an effort to remove the friction from the physics departments network; additional paths were created and dedicated to researchers along with the implementation of sensible security policies that were able to deliver the same overall goals as a general purpose firewall, without harming the sensitive science flows.

### B. National Energy Research Scientific Computing Center

The National Energy Research Scientific Computing Center (NERSC)[19] is a Department of Energy computing facility. This center houses numerous computing and storage resources for many research disciplines. It is common for

---

[12]XSEDE Dashboard. http://ps.ncar.xsede.org/maddash-webui/

[13]iperf. http://dast.nlanr.net/Projects/Iperf/

[14]nuttcp. http://wcisd.hpc.mil/nuttcp/Nuttcp-HOWTO.html

[15]Brown University. http://www.brown.edu

[16]Brown University HEP. http://www.het.brown.edu

[17]CMS, http://home.web.cern.ch/about/experiments/cms

[18]Fermilab, http://www.fnal.gov
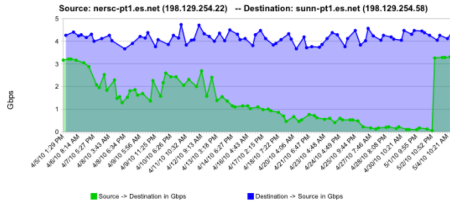
[19]NERSC. http://www.nersc.gov

Figure 7: Observed performance of the BWCTL measurement tool, through a failing network device on the ESnet network. This failure impacted all users at the NERSC computing facility.

researchers located at national labs and universities to maintain arrangements with NERSC as a part of their science workflows; namely as the destination for data analysis or the long term storage of data or results.

NERSC was an early adopter of the Science DMZ paradigm, and has maintained perfSONAR testers for a number of years. In particular they participate in regular testing activities with their upstream provider (ESnet) along with other major research labs around the country.

Figure 7 shows a graph of performance measurements captured over a number of months at NERSC. This graph illustrates a common problem involving the gradual failure of an optical networking component. The measurement, an achievable bandwidth metric delivered by the iperf tool, captured the slow decline in available bandwidth until an alarm was finally triggered that prompted engineers to investigate further. This problem is particularly challenging and related to the old fable of "frog boiling" since it occurred slowly and did not raise other alarms related to packet loss metrics or passive observations from the network device itself.

## V. Related Work

End to end monitoring and network measurement is an often researched and published topic. Services such as NWS [17], [19] and MDS [8] provided early monitoring for distributed applications on the grid. Scientific collaborations, including the Large Hadron Collider (LHC)[20] Virtual Organization from the High Energy Physics Space, created their own software to meet mission demands [15]. Commercial offerings, including Solarwinds[21] and Cisco Prime[22] have introduced performance monitoring tools over the years to address the issues of health and performance, but often require a fully homogeneous deployment. The IETF has also tried to standardize architecture and protocols in recent efforts, many of which relate to governmental sponsored surveys of a countries network capabilities[23]. Many of these

---

[20]LHC. http://home.web.cern.ch
[21]Solarwinds. http://www.solarwinds.com
[22]Cisco Prime. http://www.cisco.com/en/US/prod/netmgtsw/prime.html
[23]A Reference Path and Measurement Points for LMAP. https://datatracker.ietf.org/doc/draft-ietf-ippm-lmap-path/

efforts have a multi-domain component, and they have tried to unify the tasks of measurement, storage, processing, and visualization to ease the deployment burden on operators and deliver much needed functionality to end users.

perfSONAR is unique in that the development team had an early realization to the measurement problem; many tools have solved key problems in the ecosystem, but lacked a cohesive mechanism to "glue" the final results into all-encompassing solution. perfSONAR focuses on this "middleware" aspect to facilitate sharing, discovery, and access, without attempting to recreate seminal work related to actual measurements and analysis. A related project from the GENI [1] collaboration is Periscope, and includes the Unified Network Information Service (UNIS) [7]. A holistic view of the network is key to the successful operation of distributed computing architectures. Supporting network-aware applications and application driven networks requires a detailed understanding of network resources from multi-layer topologies, associated measurement data, and in-the-network service location and availability information. The perfSONAR system unifies a suite of monitoring services and tools with a common data model and protocols in order to measure network performance on various devices and across end-to-end paths. Periscope builds on, and uses, existing perfSONAR service deployments and implements enhanced versions of the perfSONAR protocols to provide new functionality for pervasive, scalable monitoring, and to improve the usability of the system for environments such as the GENI testbed.

## VI. Conclusion

Scientific innovation will continue to adopt data-intensive strategies in the years to come. Addressing "big data" requirements calls for a system wide approach: computational components, storage, and networks must all work in harmony to ensure success. Networks in particular are prone to complications due to their design and usage patterns, there is a requirement that performance monitoring should ensure both local and end-to-end success scenarios.

perfSONAR is a framework designed to federate testing on a global scale, and offers the ability to capture performance metrics of diverse types in an automated and seamless fashion. These metrics, when delivered through analysis tools, can directly benefit the network operations and scientific research communities by ensuring that all components are working at peak efficiency.

## VII. Disclaimer

This document was prepared as an account of work sponsored by the United States Government. While this document is believed to contain correct information, neither the United States Government nor any agency thereof, nor the Regents of the University of California, nor any of their employees, makes any warranty, express or implied, or assumes any

legal responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by its trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or any agency thereof, or the Regents of the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or any agency thereof or the Regents of the University of California.

REFERENCES

[1] Global Environment for Network Innovation. http://geni.net.

[2] Internet2 Network Observatory. http://www.internet2.edu/observatory/.

[3] C. Barakat, E. Altman, and W. Dabbous. On TCP Performance in a Heterogeneous Network: A Survey. *IEEE Communications Magazine*, 38:40–46, 2000.

[4] P. Calyam, J. Pu, W. Mandrawa, and A. Krishnamurthy. Ontimedetect: Dynamic network anomaly notification in perfsonar deployments. In *IProc. of IEEE Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems (MASCOTS)*, 2010.

[5] S. Campana, D. Bonacorsi, A. Brown, E. Capone, D. D. Girolamo, A. F. Casani, J. F. Molina, A. Forti, I. Gable, O. Gutsche, A. Hesnaux, L. Liu, L. L. Munoz, N. Magini, S. McKee, K. Mohammad, D. Rand, M. Reale, S. Roiser, M. Zielinski, and J. Zurawski. Deployment of a wlcg network monitoring infrastructure based on the perfsonar-ps technology. In *20th International Conference on Computing in High Energy and Nuclear Physics (CHEP 2013)*, 2013.

[6] E. Dart, L. Rotman, B. Tierney, M. Hester, and J. Zurawski. The science dmz: A network design pattern for data-intensive science. In *IEEE/ACM Annual SuperComputing Conference (SC13)*, Denver CO, USA, 2013.

[7] A. El-Hassany, E. Kissel, D. Gunter, and M. Swany. Design and implementation of a Unified Network Information Service. In *10th IEEE International Conference on Services Computing (SCC 2013)*, June, 2013.

[8] S. Fitzgerald. Grid information services for distributed resource sharing. In *Proceedings of the 10th IEEE International Symposium on High Performance Distributed Computing*, HPDC '01, pages 181–, Washington, DC, USA, 2001. IEEE Computer Society.

[9] M. Grigoriev, J. Boote, E. Boyd, A. Brown, J. Metzger, P. DeMar, M. Swany, B. Tierney, M. Zekauskas, and J. Zurawski. Deploying distributed network monitoring mesh for lhc tier-1 and tier-2 sites. In *17th International Conference on Computing in High Energy and Nuclear Physics (CHEP 2009)*, 2009.

[10] A. Hanemann, J. Boote, E. Boyd, J. Durand, L. Kudarimoti, R. Lapacz, M. Swany, S. Trocha, and J. Zurawski. Perfsonar: A service-oriented architecture for multi-domain network monitoring. In *International Conference on Service Oriented Computing (ICSOC 2005)*, Amsterdam, The Netherlands, 2005.

[11] M. Mathis, J. Semke, J. Mahdavi, and T. Ott. The macroscopic behavior of the tcp congestion avoidance algorithm. *SIGCOMM Comput. Commun. Rev.*, 27(3):67–82, July 1997.

[12] S. McKee, A. Lake, P. Laurens, H. Severini, T. Wlodek, S. Wolff, and J. Zurawski. Monitoring the us atlas network infrastructure with perfsonar-ps. In *19th International Conference on Computing in High Energy and Nuclear Physics (CHEP 2012)*, New York, NY, USA, 2012.

[13] S. Molnár, B. Sonkoly, and T. A. Trinh. A comprehensive TCP fairness analysis in high speed networks. *Comput. Commun.*, 32(13-14):1460–1484, 2009.

[14] G. E. Moore. Cramming more components onto integrated circuits. *Electronics*, 38(8), April 1965.

[15] H. Newman, I. Legrand, P.Galvez, R. Voicu, and C. Cirstoiu. Monalisa: A distributed monitoring service architecture. In *International Conference on Computing in High Energy and Nuclear Physics (CHEP 2003)*, 2003.

[16] J. Postel. Transmission Control Protocol. Request for Comments (Standard) 793, Internet Engineering Task Force, September 1981.

[17] M. Swany and R. Wolski. Representing dynamic performance information in grid environments with the network weather service. In *Cluster Computing and the Grid, 2002. 2nd IEEE/ACM International Symposium on*, pages 48–48, May.

[18] B. Tierney, J. Metzger, J. Boote, A. Brown, M. Zekauskas, J. Zurawski, M. Swany, and M. Grigoriev. perfsonar: Instantiating a global network measurement framework. In *4th Workshop on Real Overlays and Distributed Systems (ROADS09) Co-located with the 22nd ACM Symposium on Operating Systems Principles (SOSP)*, 2009.

[19] R. Wolski, N. T. Spring, and J. Hayes. The network weather service: A distributed resource performance forecasting service for metacomputing. *Journal of Future Generation Computing Systems*, 15:757–768, 1999.

[20] J. Zurawski, M. Swany, and D. Gunter. A scalable framework for representation and exchange of network measurements. In *2nd International IEEE/Create-Net Conference on Testbeds and Research Infrastructures for the Development of Networks and Communities (TridentCom 2006)*, Barcelona, Spain, 2006.