

PAPER • OPEN ACCESS

## Networks in ATLAS

To cite this article: Shawn McKee and For the ATLAS Collaboration 2017 *J. Phys.: Conf. Ser.* **898**  
052006

View the [article online](#) for updates and enhancements.

### You may also like

- [The LHC Olympics 2020 a community challenge for anomaly detection in high energy physics](#)  
Gregor Kasieczka, Benjamin Nachman, David Shih et al.
- [Special issue on applied neurodynamics: from neural dynamics to neural engineering](#)  
Hillel J Chiel and Peter J Thomas
- [From biologically-inspired physics to physics-inspired biology](#)  
Alexei A Kornyshev

The advertisement features a green background on the left with the ECS logo and text. The right side has a dark teal background with white text and images of industrial robotics and a scientist.

**ECS**  
The  
Electrochemical  
Society  
Advancing solid state &  
electrochemical science & technology

**DISCOVER**  
how sustainability  
intersects with  
electrochemistry & solid  
state science research

# Networks in ATLAS

Shawn McKee<sup>1</sup>, For the ATLAS Collaboration

<sup>1</sup> Physics Department, University of Michigan, Ann Arbor, MI 48109-1040 USA

E-mail: smckee@umich.edu

**Abstract.** Networks have played a critical role in high-energy physics (HEP), enabling us to access and effectively utilize globally distributed resources to meet the needs of our physicists. Because of their importance in enabling our grid computing infrastructure many physicists have taken leading roles in research and education (R&E) networking, participating in, and even convening, network related meetings and research programs with the broader networking community worldwide. This has led to HEP benefiting from excellent global networking capabilities for little to no direct cost. However, as other science domains ramp-up their need for similar networking it becomes less clear that this situation will continue unchanged. What this means for ATLAS in particular needs to be understood. ATLAS has evolved its computing model since the LHC started based upon its experience with using globally distributed resources. The most significant theme of those changes has been increased reliance upon, and use of, its networks.

We will report on a number of networking initiatives in ATLAS including participation in the global *perFSNAR* network monitoring and measuring efforts of WLCG and OSG, the collaboration with the LHCOPN/LHCONE effort, the integration of network awareness into PanDA, the use of the evolving ATLAS analytics framework to better understand our networks and the changes in our DDM system to allow remote access to data.

We will also discuss new efforts underway that are exploring the inclusion and use of software defined networks (SDN) and how ATLAS might benefit from:

- Orchestration and optimization of distributed data access and data movement.
- Better control of workflows, end to end.
- Enabling prioritization of time-critical vs normal tasks
- Improvements in the efficiency of resource usage

## 1. Introduction

Innovation supporting science continues to increase requirements for the computing and networking infrastructures of the world. Instrumentation, storage, processing facilities and collaborative partners are often geographically and topologically separated, thus complicating the problems involved with data management. Global scientific collaborations, such as ATLAS, continue to push the network requirements envelope. Data movement in this collaboration is routinely including the regular exchange of many 10's of petabytes of datasets between the collection and analysis facilities in the coming years. This increased emphasis on the “network”, now a vital resource on par with the actual scientific process, implies that it **must** be a highly capable and reliable resource to ensure success; the lack thereof could mean critical delays in the overall scientific progress of distributed data-intensive experiments.



Content from this work may be used under the terms of the [Creative Commons Attribution 3.0 licence](#). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

We will report on the role of networking in supporting the scientific mission and goals of the ATLAS collaboration. Networks are fundamental to the distributed computing model ATLAS has developed and, as such, end-to-end network performance and network problems have a significant impact on the ability of ATLAS physicists to reach their scientific goals in a timely manner. In this paper we will discuss the ongoing efforts to monitor, measure and maintain our networks and exploratory work to integrate programmable networks into a future ATLAS global infrastructure.

The remainder of the paper will proceed as follows. Section 2 will discuss the ATLAS collaboration, as well as data movement requirements and expectations. Section 3 will discuss the work to monitor and measure our networks. Section 4 will discuss the ATLAS effort to analyze our network data. Section 5 will discuss PanDA and how it is evolving to better utilize the network. Section 6 will cover exploratory work to determine the impact of future networks on ATLAS. Section 7 concludes the paper.

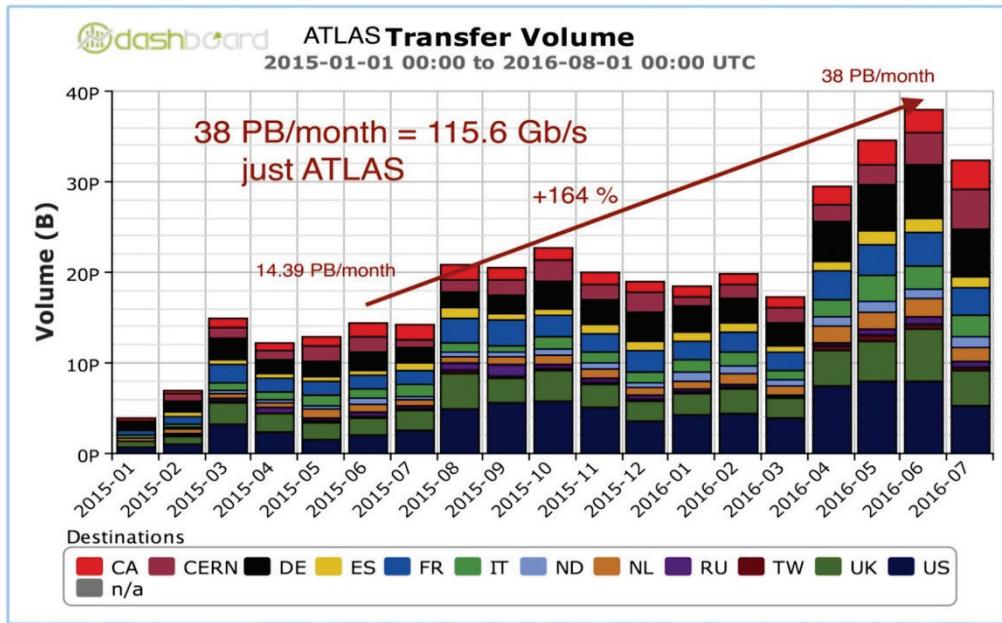
## 2. The ATLAS Collaboration

The ATLAS collaboration consists of over 3000 physicists and 1000 students from 38 countries and 178 Universities and Laboratories worldwide. This large group of scientists is working together at the Large Hadron Collider (LHC) [1] to learn about the basic forces that have shaped our Universe since the beginning of time and that will determine its fate. ATLAS physicists are exploring the frontiers of high-energy physics in a number of ways: explaining the origin of mass, exploring the range of validity of our standard model, searching for microscopic black holes, probing the existence of extra dimensions, and looking for evidence of an as-of-yet undiscovered particle that may explain the dark matter in our Universe.

To undertake these explorations, the collaboration has constructed the ATLAS detector [2] over a 15 year period and assembled it at Point 1 in the LHC ring. The detector is 45 meters long, 25 meters high and weighs about 7000 tons. The ATLAS detector employs a number of types of sub-detectors to measure attributes of the various particles resulting from the collision of counter-rotating beams of protons in the LHC. There are millions of electronic channels associated with the readout of the ATLAS detector. In effect, the set of all of these sub-detectors and associated readouts can be viewed as a very large 3-dimensional digital camera, capable of taking precise "pictures" 40 million times a second (proton beam bunches cross one another every 25 nanoseconds). The detailed information collected allows ATLAS physicists to reconstruct the underlying events and search for new physics.

If all the data ATLAS produces could be stored, it would fill more than 232,000 CDs per second, a rate (and corresponding data-volume) which is not feasible to support with current technology. Instead a set of hardware, firmware and software systems make fast decisions about what data is interesting to keep and results in a data rate of 400-1000 MBytes/sec into "offline" disk storage. Even so, this rate of data production results in many petabytes of data being produced by ATLAS each year. In addition, detailed simulations also produce Petabytes of data required to understand how the ATLAS detector responds to various types of events and validate that the ATLAS software works as expected. It is important to note that these large data volumes are common to all the LHC experiments and not just ATLAS.

Because of the data-intensive nature of the ATLAS scientific program, the ATLAS collaboration implicitly relies upon having a ubiquitous, high-performing, global network to enable its distributed grid-computing infrastructure. Providing effective access to petabytes of data for thousands of physicists all over the world just wouldn't be possible without the corresponding set of research and education (R&E) networks that provide 1 to 10 to 100 Gigabits per second of bandwidth to enable ATLAS data to flow to where it is needed. As can be seen in Figure 1, recent ATLAS wide-area network use is rapidly increasing, continuing an exponential increase that has been observed since startup. This increasing use exemplifies the importance of networking for ATLAS and its globally distributed computing model.



**Figure 1** Recent ATLAS wide-area network use with a trend-line showing a 164% increase. ATLAS WAN use has grown to almost 38 Petabytes per month.

Typical network paths that ATLAS data traverses consist of multiple administrative domains (local area networks at each end and possibly multiple campus, regional, national and international networks along the path). The ability of the Internet to allow these separate domains to transparently inter-operate is one of its greatest strengths. However, when a problem involving the network arises, that same transparency can make it very difficult to find the cause and location of the problem.

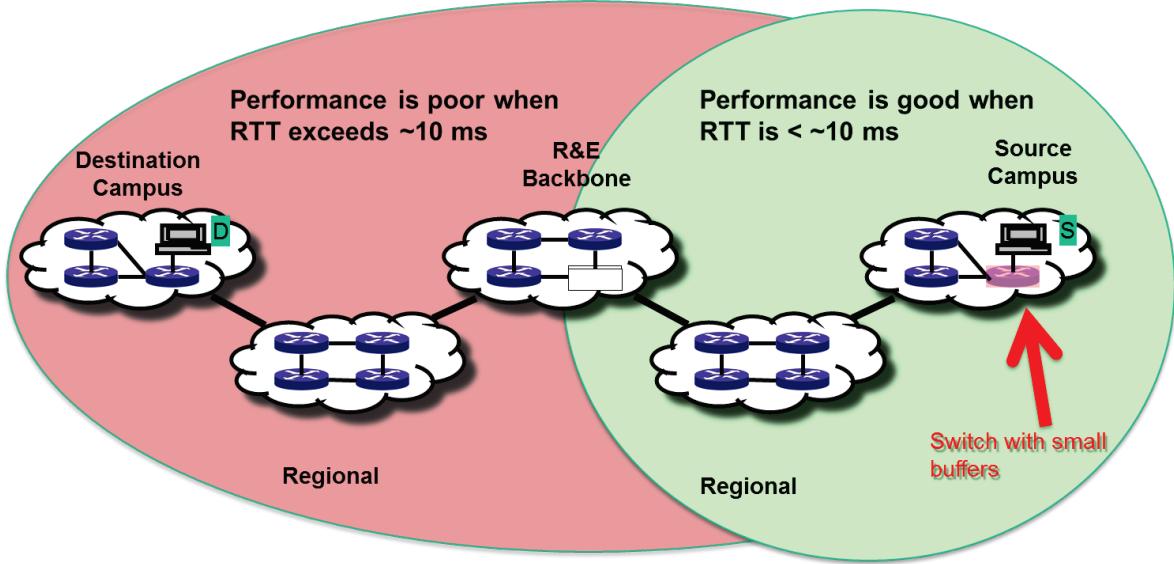
Because of both the criticality of the network for ATLAS normal operations and the difficulty in identifying and locating the source of network problems when they occur, the US ATLAS facility began deploying and configuring perfSONAR-PS (now referred to simply as perfSONAR) in 2008. Our goal was to provide our sites with a set of tools and measurements that would allow them to differentiate network issues from end-site issues and to help localize and identify network specific problems to expedite their resolution. This effort evolved into first an ATLAS and eventually a WLCG [3] and OSG [4] effort to monitor and measure our networks.

### 3. Monitoring and Measurement of our Networks

Network problems can severely impact ATLAS's workflows and have taken weeks or months to get addressed. End-to-end network issues are difficult to spot and localize because they are multi-domain (multiple independent administrators) and involve many components (end-systems, software, firmware, routers, switches, cables, etc). Standardizing on specific tools and methods allows ATLAS (and HEP in general) to focus resources more effectively and better self-support its collaborators. Thus we have chosen to use *perfSONAR*. *perfSONAR* is a framework that enables network performance information to be gathered and exchanged in a multi-domain, federated environment and its use in HEP was described in detail in a previous CHEP paper [5].

The typical network problem involves packet-loss or packet reordering along a wide area network path. To illustrate the impact of packet loss on long network paths we can use the example shown in Figure 2. Assuming the links shown are 10 Gbps, even a small loss can significantly impair the throughput. A 0.0046% loss (1 out of 22k packets) on 10G link results in very different throughput, depending upon the round-trip time (RTT):

- with **1ms RTT: 7.3 Gbps**
- with **51ms RTT: 122Mbps**
- with **88ms RTT: 60 Mbps** (factor 120 decrease relative to 1ms RTT)



**Figure 2** A wide area network path from a source of network traffic to a destination through many routers and switches. The behaviour of TCP in the presence of packet loss degrades significantly with round-trip time (RTT) and packet loss “close” to the source can mask network problems when traffic is to nearby destinations.

The impact of packet reordering and jitter can also be significant. Referring to Figure 3 we note the impact of packet reordering when there is significant jitter. At 70ms RTT on 10 Gbps link, a 60 second test results in significantly different throughput with only 1% packet reordering depending upon jitter:

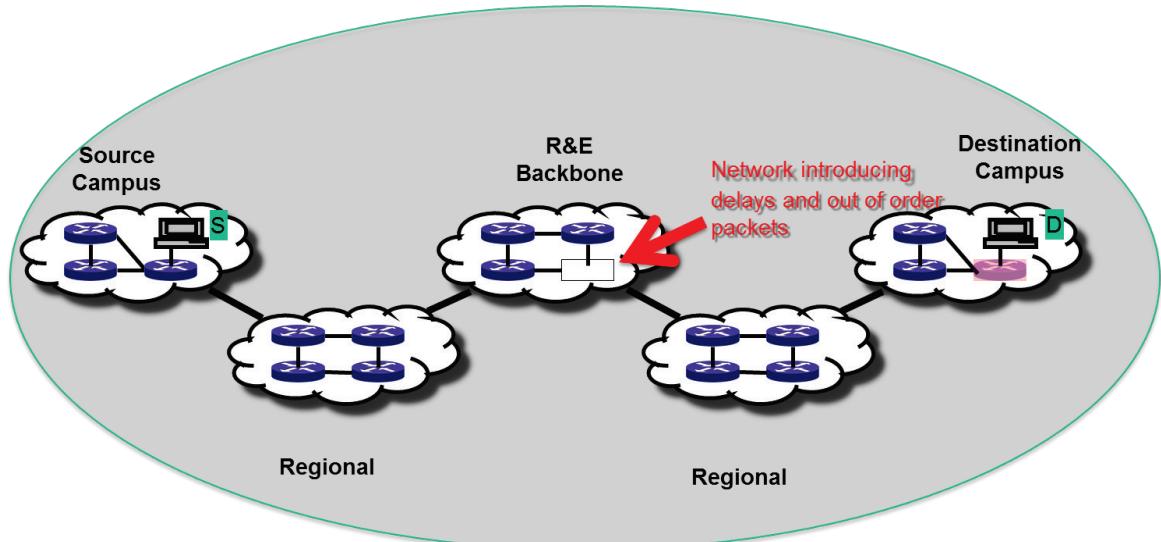
- with 1% re-ordering, **0.2 ms** jitter: **8.45 Gbps**
- with 1% re-ordering, **1.0 ms** jitter: **1.10 Gbps**

As we have seen, it is critical to understand when problems arise in the network adversely impacting ATLAS’s ability to use the network effectively. We rely upon *perfSONAR* to monitor and measure our networks. *perfSONAR* provides a number of standard metrics we 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
  - Measurements are only useful when we know the exact path they are taking through the network.
  - Tracepath additionally measures the Maximum Transmission Unit (MTU) on the end-to-end path but is frequently blocked from operating correctly because of incorrectly or over-zealously configured firewalls along the path.

### 3.1. Organizing and Maintaining ATLAS Networking

The ATLAS collaboration is benefiting-from and participating-in a number of efforts to instrument, measure, monitor, understand and control our networks. Since 2012 the Open Science Grid (OSG) has had a networking area whose goal is to provide network information and support to its members and collaborators. Since 2014 the WLCG has operated the Network and Transfer Metrics Working Group (NTMWG) which is responsible for instrumenting, measuring and reliably gathering *perfSONAR* network metrics from our networks. All ATLAS Tier-1 and Tier-2 centers are mandated to deploy *perfSONAR* Toolkit instances co-located with their storage resources. The deployment campaign was led by the WLCG *perfSONAR* Deployment Task-force [6] which completed its work in 2014.



**Figure 3 A wide area network path from a source of traffic to a destination through links that introduce delays and packet reordering. In this case TCP throughput, even without packet loss, can be significantly degraded because of packet reordering and jitter (the variation in inter-packet arrival timing).**

Recently ATLAS has taken the lead in trying to analyze and better understand the network metrics being gathered by OSG and this work will be described in section 5. The ATLAS work is well aligned with the OSG goal of providing effective alarming and alerting for network problems.

One of the important activities of the NTMWG was to create a support unit [7] to coordinate responses to potential network issues. Tickets opened in the support group can be triaged to the right destination by networking experts from ATLAS or the other experiments. Many issues are potentially resolvable within the working group because of the information available from perfSONAR. More complex network issues can at least be identified and directed to the appropriate network support centers along with any additional supporting information. This has resulted in significantly decreasing the time it has taken to resolve network issues.

Lastly we should mention the LHC Optical Private Network (LHCOPN) and LHC Open Network Environment (LHCONE) [8] efforts. The LHCOPN was created in 2006 to implement, manage and maintain dedicated network circuits between CERN (the Tier-0) and the set of Tier-1 centers worldwide. A group of physicists, network engineers and members of global Research and Education (R&E) networks have met 2-3 times per year to manage and develop the LHCOPN since its inception. Because of the success of the LHCOPN in meeting the needs of ATLAS and the other LHC experiments, discussions were started concerning how this effort might be expanded to incorporate the needs of the Tier-2s and even Tier-3 computing sites worldwide. This discussion led to the formation of the LHCONE effort in 2012 and subsequently joint meetings with LHCOPN. The LHCOPN and LHCONE efforts are providing high-energy particle physics experiments like ATLAS with customized planning, services and development to support their global network requirements.

#### 4. ATLAS Analytics and Network Data Analysis

The volume and complexity of the network metrics that are being gathered globally by OSG and WLCG have created a pressing need to get this data into a location suitable for filtering and analysis. The metrics are gathered along many paths across our global R&E networks and measure various characteristics of those paths which change with time. To be able to more fully understand how our networks are operating and especially to be able to identify and localize network problems, we need to apply more complex analysis to our metrics than we can do with our existing *perfSONAR* toolkit capability.

In late 2015 the OSG networking group began working with Ilija Vukotic, University of Chicago, to feed network metrics into a new ELK (Elasticsearch, Logstash and Kibana) instance that was already capturing many useful ATLAS metrics. Data being gathered by OSG from the global set of WLCG-related *perfSONAR* instances was published to an ActiveMQ message bus instance hosted at CERN and then sent to the Chicago ELK analytics instance via a customized Flume instance. This analytics service indexes historical network related data while providing predictive capabilities for network throughput. Further details are provided in these conference proceedings [9].

Using this analytics platform for network metrics was immediately valuable. By having all this data query-able we were able to ask questions like: “Which sites have more than 2% packet loss to more than 80% of their testing partners for the last 12 hours?” Being able to quickly find and localize network problems is critical for our infrastructure performance. We are working on defining standard alarm tables that continually update as new data is gathered that will serve as the basis of an eventual alerting system.

One of the conclusions we were able to reach by having this network analytics capability is that much of our ATLAS 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, firmware and hardware settings that lead to inefficient use of our networks. This wealth of data now available and analyzable can identify bottlenecks and poor performance. We are now working to automatically and consistently find and fix such problems in ATLAS resources.

## 5. PanDA and ATLAS Workflow and the Network

ATLAS relies upon the Production and Distributed Analysis (PanDA) workload management system [10] to coordinate and optimize the collaboration’s set of tasks across its global resources. PanDA is responsible for selecting the job execution site and it does this via a multi-level decision tree involving task brokerage, job brokerage and a dispatcher. It also includes predictive workflows like the PD2P (PanDA Dynamic Data Placement). Site selection only used processing and storage requirements.

Recently PanDA has evolved [11] to incorporate network information as another component for site selection because of the impact the network can have, both positively and negatively, on task completion times and failure rates. The ATLAS analytics platform mentioned in section 4 is used to summarize recent network metrics from FAX (federated storage access system) and *perfSONAR* and make it available for PanDA use. This data augments other information PanDA already uses such as job completion metrics, errors and timeouts per site.

The longer-term goal for PanDA is to go beyond network monitoring and treat the network as a managed resource. Can we incorporate network provisioning, orchestration and control via software defined network capabilities as part of PanDA? Initial simple tests have shown that network knowledge is useful and beneficial for all phases of the job cycle. In both the ANSE[12] and BigPanDA[13] projects we have added “hooks” into PanDA that could allow control of the network once production quality mechanisms are in place to support that across at least some of our networks. This would be a first and never attempted before for large scale automated workload management systems. To make this a reality for ATLAS will require new, production quality capabilities from future networks.

## 6. Exploring Future Networks for ATLAS

Future networks won’t just have increased capacity but are also enabling new interfaces and modes of operation that can allow end-users to control some aspects of how data flows across networks. This is referred to as software defined networking (SDN) and the primary impetus for these capabilities is driven by commercial entities. In fact, SDN originated with Google and its attempt to better orchestrate its data-center and wide-area networks.

An important question for ATLAS is whether or not SDN is something that can improve ATLAS’s ability to use its distributed resources. A group of people in the US from four of the ATLAS Tier-2 centers (AGLT2, MWT2, SWT2 and NET2) have begun to explore SDN for ATLAS. The idea is that

we need a way to compare and contrast the impact of controlling the network versus using the network as-is for real production ATLAS work. This means we need a non-disruptive way to incorporate ATLAS production systems into an SDN testbed.

To do this we began working with the LHCONE point-to-point effort which has been exploring the use of SDN to setup end-to-end network connections between LHCONE sites. The stumbling block has always been getting SDN capabilities all the way to the source or sink of data. The R&E networks may have different ways to setup circuits or control the “backbone” network but these never reach into the local area networks (LANs) nor to the servers hosting the data. For the US ATLAS sites, we proposed to solve this problem through the use of Open vSwitch [14].

With Open vSwitch (OVS) we have the ability to non-disruptively add SDN capability to existing ATLAS production data servers at a few of our Tier-2 sites. This will give us the ability to selectively test how ATLAS production behaves with and without SDN features in place.

Our plan is to deploy OVS on ATLAS production storage systems at all the participating Tier-2 sites as follows:

- Measure baseline performance on our systems for a few days.
- Install Open vSwitch v2.6.1 via RPM on all storage servers
- Reconfigure the network to move the server IP address from the Network Interface Card (NIC) to the Open vSwitch virtual switch
- Verify continued normal operation of the system now that all network traffic is passing through the vSwitch

Once we confirm that our systems continue to behave the same as before installing OVS, we can proceed to test features of SDN between our participating sites and compare and contrast the site performance when using these capabilities versus not using them. For example, one nice feature of OVS is the ability to shape traffic by pacing the rate at which network packets are inserted in the network. If we know storage systems are only capable of sourcing or sinking data at a certain maximal rate, we can shape the corresponding network traffic to match that rate. This not only helps reduce the load on the end-systems but also results in much better average performance across the network based upon tests that have been conducted using OVS [15]. In addition there is very little cost in terms of server resource use (roughly 1% additional CPU) to accurately shape traffic up to 100 Gbps. Finally, having OVS in place additionally allows various kinds of software defined network controllers (e.g., OpenDaylight, Ryu, Floodlight) to see and interact with our servers, giving us the possibility to orchestrate the network end-to-end for the first time.

We will need to do extensive testing once this capability is in place to understand the possible benefits for ATLAS. One of the challenges involved is getting complete instructions for our various end-sites on how to non-disruptively deploy OVS while those servers are in production. Assuming we can successfully get this functioning within the US, we have requests from the Tier-1 sites in the Netherlands and in Germany to also join in our testing. This will be important to test the possible impact using the long fat network pipes across the Atlantic.

## 7. Conclusion

Networking is a critical component for ATLAS and underlies our distributed computing model. Problems in the network can cause significant degradation for ATLAS workflows and can be very hard to identify, locate and fix. To address this ATLAS has a working infrastructure in place to monitor and measure our networks using *perfSONAR* and is benefitting-from and contributing-to efforts around networking in OSG, WLCG and the LHCOPN and LHCONE communities. ATLAS is also leading the effort to make complex analysis of network data possible and working towards new capabilities in network notification and alerting, predictive network behavior and the identification of problematic sites and servers. The ATLAS PanDA system is also evolving to take better advantage of network knowledge and to prepare for future network capabilities that may allow control and orchestration of our networks. Lastly, a group in ATLAS is exploring the possible impact of SDN and testing how it may be able to benefit ATLAS.

### Acknowledgements

The author would like to acknowledge the contributions of the *perfSONAR* project, namely staff from ESnet, GEANT, Indiana University, Internet2, the University of Michigan and all the contributors shown at <http://www.perfsonar.net/about/who-is-involved/>.

We gratefully acknowledge the support of the Department of Energy (grant DE-SC0007859) and the National Science Foundation (grant PHY-1148698) which supported the work in this paper.

### References

- [1] Lyndon E and Philip B 2008 LHC Machine *Journal of Instrumentation* **3** S08001
- [2] Dunford M and Jenni P 2014 The ATLAS experiment *Scholarpedia* **9** 32147
- [3] Bird I 2011 Computing for the Large Hadron Collider *Annual Review of Nuclear and Particle Science* **61** 99-118
- [4] Pordes R, Altunay M, Avery P, Bejan A, Blackburn K, Blatecky A, Gardner R, Kramer B, Livny M, McGee J, Potekhin M, Quick R, Olson D, Roy A, Sehgal C, Wenaus T, Wilde M and Würthwein F 2008 New science on the Open Science Grid *Journal of Physics: Conference Series* **125** 012070
- [5] McKee S, Lake A, Laurens P, Severini H, Wlodek T, Wolff S and Zurawski J 2012 Monitoring the US ATLAS Network Infrastructure with perfSONAR-PS *Journal of Physics: Conference Series* **396** 042038
- [6] Campana S, Brown A, Bonacorsi D, Capone V, Girolamo D D, Casani A F, Flix J, Forti A, Gable I, Gutsche O, Hesnaux A, Liu S, Munoz F L, Magini N, McKee S, Mohammed K, Rand D, Reale M, Roiser S, Zielinski M and Zurawski J 2014 Deployment of a WLCG network monitoring infrastructure based on the perfSONAR-PS technology *Journal of Physics: Conference Series* **513** 062008
- [7] Babik M and McKee S 2014 WLCG network throughput support unit.
- [8] Martelli E and Stancu S 2015 LHCOPN and LHCONE: Status and Future Evolution *J.Phys.Conf.Ser.* **664** 052025
- [9] Vukotic I and Gardner R 2017 Big Data Analytics Tools as Applied to ATLAS Event Data. In: *22nd International Conference on Computing in High Energy and Nuclear Physics (CHEP2016)*, ( Journal of Physics Conference Series
- [10] Maeno T, De K, Wenaus T, Nilsson P, Stewart G A, Walker R, Stradling A, Caballero J, Potekhin M, Smith D and The Atlas C 2011 Overview of ATLAS PanDA Workload Management *Journal of Physics: Conference Series* **331** 072024
- [11] Maeno T, De K, Klimentov A, Nilsson P, Oleynik D, Panitkin S, Petrosyan A, Schovancova J, Vaniachine A, Wenaus T, Yu D and the Atlas C 2014 Evolution of the ATLAS PanDA workload management system for exascale computational science *Journal of Physics: Conference Series* **513** 032062
- [12] NSF 2012 CC-NIE Integration: ANSE (Advanced Network Services for Experiments).
- [13] De K 2015 The BigPanDA Project. pp An overview of the BigPanDA project, funded by DOE ASCR and DOE HEP.
- [14] Pfaff B, Pettit J, Kopenen T, Jackson E J, Zhou A, Rajahalme J, Gross J, Wang A, Stringer J, Shellar P, Amidon K and Casado M 2015 The Design and Implementation of Open vSwitch. In: *USENIX, (NSDI)*
- [15] Newman H, Mughal A, Kcira D, Legrand I, Voicu R and Bunn J 2015 High speed scientific data transfers using software defined networking. In: *Proceedings of the Second Workshop on Innovating the Network for Data-Intensive Science*, (Austin, Texas: ACM) pp 1-9