

Final Report

DoE Grant: DE-FC02-12ER26106/DE-SC0008635

Institution: The University of Texas at Arlington

Project Title: Next Generation Workload Management and Analysis System for Big Data

PI of this report: Kaushik De

Lead PI: Dr. Alexei Klimentov

April 24, 2017

Period covered by this report: September 1, 2012 – August 31, 2016

Overview

We report on the activities and accomplishments of a four-year project (a three-year grant followed by a one-year no cost extension) to develop a next generation workload management system for Big Data. The new system is based on the highly successful PanDA software developed for High Energy Physics (HEP) in 2005. PanDA is used by the ATLAS experiment at the Large Hadron Collider (LHC), and the AMS experiment at the space station.

Initial PanDA development and subsequent operation and maintenance has been supported through funding from the HEP programs at the Department of Energy (DOE) and National Science Foundation (NSF). The next generation workload management and analysis system work described here was supported by ASCR/CompHEP through this project. The program of work described here was carried out by two teams of developers working collaboratively at Brookhaven National Laboratory (BNL) and the University of Texas at Arlington (UTA). These teams worked closely with the original PanDA team – for the sake of clarity the work of the next generation team will be referred to as the BigPanDA project. Their work has led to the adoption of BigPanDA by the COMPASS experiment at CERN, and many other experiments and science projects worldwide.

The main deliverables of this project include four work packages, as described in the proposal:

- WP1 (Factorizing the core): Factorizing the core components of PanDA to enable adoption by a wide range of exascale scientific communities.
- WP2 (Extending the scope): Evolving PanDA to support extreme scale computing clouds and Leadership Computing Facilities.

- WP3 (Leveraging intelligent networks): Integrating networking services and real-time data access to the PanDA workflow.
- WP4 (Usability and monitoring): Real time monitoring and visualization package for PanDA.

The teams at UTA and BNL were fully integrated. UTA contributed to three out of the four work packages. The UTA development team consisted of three computing specialists. Mikhail Titov was a student in the Computer Science Engineering (CSE) department at UTA, supported for three years under this project. He successfully completed his Ph.D. program under the supervision of Dr. Gergely Zaruba (CSE) and Dr. Kaushik De (HEP), both at UTA. Danila Oleynik is a senior software engineer with almost ten years of previous experience in developing software at CERN and DUBNA. He worked on this project for three years, primarily on pilot development and BigPanDA operations at Oakridge Laboratory Computing Facility (OLCF). Artem Petrosyan is a software engineer, also with previous development background acquired at CERN and DUBNA. Artem worked on the BigPanDA project for two years, primarily on the networking interface to PanDA. Mikhail, Danila and Artem formed the working team at UTA, along with PI De. The rest of the BigPanDA team was located at BNL. Weekly phone meetings and regular face-to-face meetings were held at BNL and UTA to coordinate the joint work.

Accomplishments

The UTA team had major accomplishments in three out of the four work packages during the duration of the project. The main areas of work at UTA included: network integration with PanDA, pilot code development for running PanDA at OLCF, monitoring of network information, data management and data location services, refactoring and packaging of pilot code, and opportunistic use of Titan at OLCF for ATLAS simulations. These work areas corresponds to work packages WP1, WP2 and WP3.

Simulations at Titan

Titan is one of the largest supercomputers in the world. Titan is a hybrid-architecture Cray XK7 system with a theoretical peak performance exceeding 27 petaflops. Titan features 18,688 compute nodes, (each with one 16-core AMD Opteron CPU and 1 NVIDIA Kepler K20X GPU), 299,008 x86 cores, a total system memory of 710 terabytes, and a high-performance proprietary network. The combination of these technologies allows Titan to achieve up to 10 times the speed of its predecessor, the Jaguar supercomputer, while consuming the same average power load and occupying the same physical footprint. The BigPanDA project provided the first important demonstration of the capabilities that a

workload management system (WMS) can have on improving the uptake and utilization of LCF from both application and systems points of view.

Through modifications to the pilot system in PanDA, Danila developed and implemented a new capability in BigPanDA to collect information about unused worker nodes on Titan, and based on that information to adjust workload parameters to fill free resources. Initial proof-of-concept tests of this mechanism achieved increased system utilization levels from 90% to 93% (14.3% of free cycles), as well as provided short wait times to ATLAS jobs submitted to Titan via PanDA; all with no negative impact on OLCF ability to schedule large, leadership-class jobs.

After these tests, Titan was fully integrated with the ATLAS PanDA based production system by Danila. The ATLAS experiment routinely runs Monte-Carlo simulation tasks on Titan, in support of dozens of ongoing data analysis efforts. All operations, including data transfers to and from Titan, are transparent to the ATLAS Computing Operations team and physicists. **Figure 1** shows an example from the ATLAS monitoring dashboard of running ATLAS production jobs on Titan in October of 2015. Even in these early stages of integration with PanDA, we could ramp up production quickly on Titan.

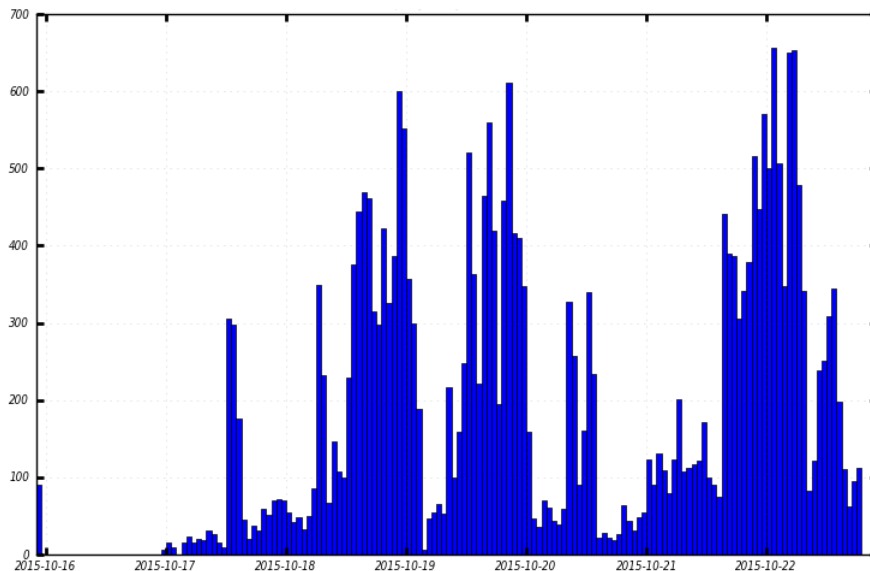


Figure 1: ATLAS production jobs on Titan

In **Figure 2**, we show the calendar year 2015 usage of ATLAS using the PanDA workflow management system on the Titan supercomputer, in units of Titan-core-hours, showing that this simulation production workflow has grown to consume an average of 2.7 million hours per month over the last three months of the year.

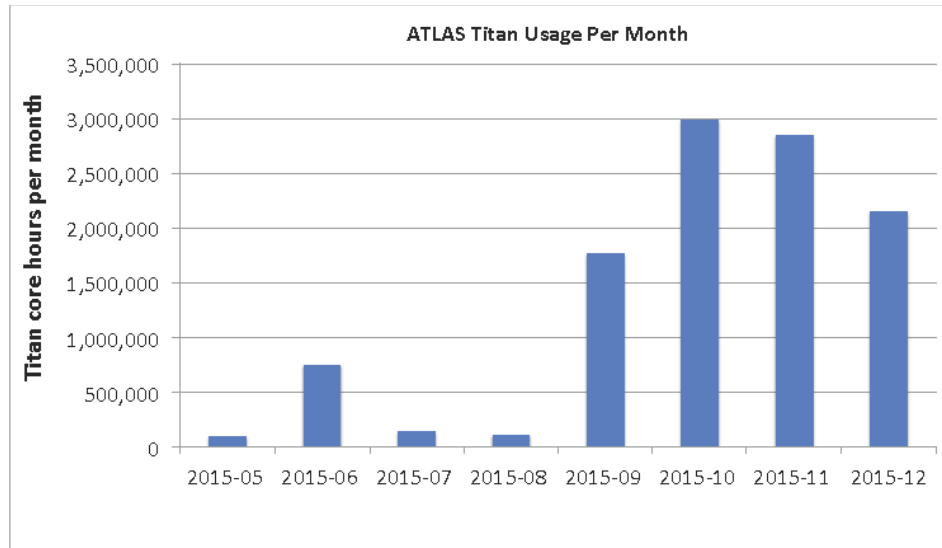


Figure 2: Growth in Titan-core-hour usage per month by ATLAS, 1 April 2015 - 31 December 2015

The trend of increasing production use of Titan by BigPanDA continued strongly in 2016. Various improvements and optimizations were made to the PanDA pilot. Data transfers were optimized. The number of pilot jobs running concurrently were increased. Through these and other enhancements, both the utilization of Titan, as well as the efficiency of backfill usage grew during 2016. In **Figure 3**, we show the Titan core hours used by ATLAS. The number of unusable Titan hours per month are showed for comparison. ATLAS used a total of 42.5 million core-hours on Titan in 2016 through August 31.

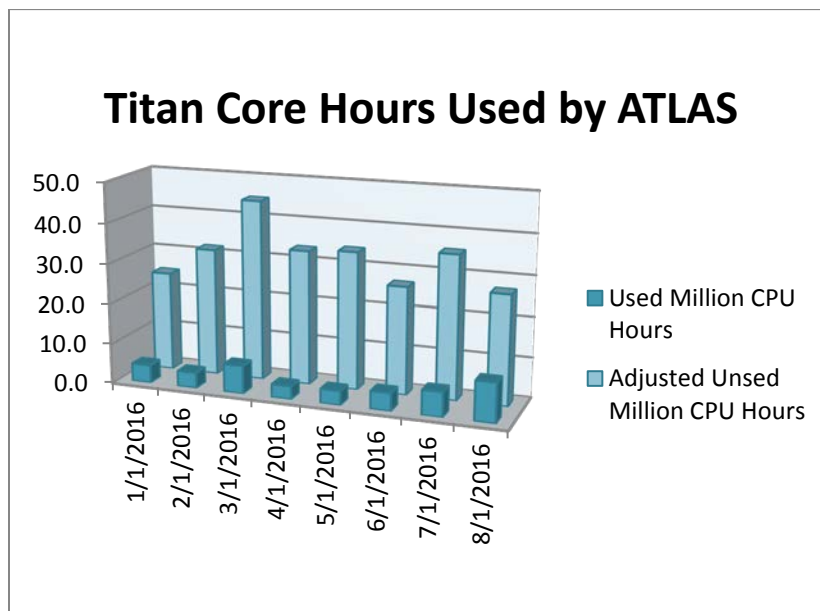
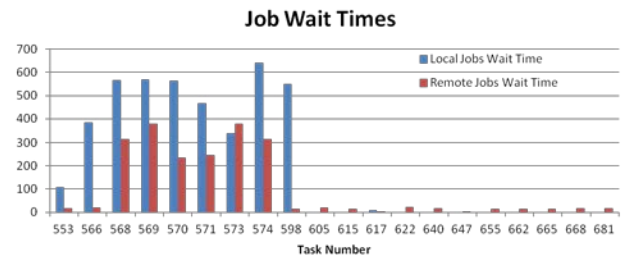


Figure 3: Growing utilization of Titan-core-hour per month by ATLAS in 2016

Network Integration

Progress in integrating PanDA with network resources was presented at the ISGC 2014 conference in Taiwan. This is the first example of a workload management system using network information to optimize the scheduling of distributed work. By using past data access rates as a component of job brokering, in conjunction with FAX (Federated data access using the XRootD protocol), we showed that users could reduce data access congestions, and get their work started faster. The job wait times are shown here on the y axis of a plot, where the blue bars denote normal jobs, while the red bars show jobs using network information. The x axis shows arbitrary sets of similar jobs. This work was done in collaboration with the NSF funded ANSE project. This work was carried out by Artem at UTA.



Network monitoring

A new framework was developed to store and provide access to various network statistics through PanDA. Three different data stores were implemented, along with multiple access protocols. Information from perfSonar, file transfers and direct read access are now available and reliably updated for PanDA and other systems. A screen shot of network data monitoring is showed here. This work was done by Artem at UTA.

DDM Sonar						perfSONAR						FAX	
AvgThr (MB/s)	EvL	AvgThr (MB/s)	EvL	AvgThr (MB/s)	EvL	MinThr (MB/s)	AvgThr (MB/s)	MaxThr (MB/s)	MinPL	AvgPL	MaxPL	FAX	ratio
1.05±0.19	10	7.48±1.45	11	12.54±0.72	918	12.4	54.7	55.9	0.0	0.0	2.0	n/a	
0.85±0.04	10	0.97±0.20	682	26.48±13.48	10	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
0.42±0.06	10	0.09±0.11	10	0.00±0.00	0	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.39±0.06	10	1.02±0.04	10	0.00±0.00	0	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.50±0.07	10	2.91±0.02	10	0.00±0.00	0	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.40±0.06	10	2.45±0.05	10	3.13±0.79	10	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.12±0.28	405	4.13±1.44	1575	4.59±0.90	2053	164.2	172.3	180.3	0.0	0.0	0.0	0.0	
2.10±1.08	4920	8.78±0.32	10075	14.06±23.95	4008	0.0	0.0	0.0	0.0	0.0	0.0	0.72	
0.47±0.11	5	1.22±0.39	9	0.00±0.00	0	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.37±0.11	10	1.14±0.20	5	2.53±0.15	10	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.97±0.04	10	7.53±3.91	10	0.00±0.00	0	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.16±0.36	10	0.90±2.04	10	00.02±0.11	10	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.94±0.06	10	0.41±1.33	10	0.00±0.00	0	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.50±0.25	10	4.95±1.63	10	21.09±0.01	10	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
1.13±0.11	10	7.17±1.44	510	0.00±0.00	0	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
0.82±0.33	10	6.90±1.82	10	30.36±11.35	10	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
1.14±0.09	10	6.50±2.41	10	0.00±0.00	0	n/a	n/a	n/a	n/a	n/a	n/a	n/a	

Data management

Mikhail worked on data movement and data management aspects of PanDA. Scientific computing has advanced in its ability to deal with massive amounts of data, since the data production capacities have increased significantly over the last decades. Most large science experiments require vast computing and data storage resources to provide results or predictions based on the data obtained. For scientific distributed computing systems with hundreds of petabytes of data and thousands of users it is

important to keep track not just of how data is distributed in the system, but also of individual user's interests in the distributed data (implicit interconnection between user and data objects). This however requires the collection and use of specific statistics such as correlations between data distribution, user preferences and the mechanics of data distribution.

Mikhail's Ph.D. thesis work focused on user activities and interests in such a distributed computing system, namely BigPanDA. He investigated whether data that was gathered in the past in PanDA shows any trends indicating that users could have mutual interests that would be repeated for the next data usages, using data mining techniques such as association analysis, sequential pattern mining, and basics of the recommender system approach. He showed that such common interests between users indeed exists and thus could be used to provide recommendations (in terms of collaborative filtering) to help users with their data selection process.

Publications and Publicity

- The public-facing project web page is available at <http://pandawms.org/>.
- The source code is available through the public GIT repository <https://github.com/PanDAWMS>
- A.Klimentov et al. Next generation workload management system for big data. In 16th International Workshop on Advanced Computing and Analysis Techniques in Physics Research. plenary talk. September 2014. Prague, Czech Rep.
- Sergey Panitkin, Danila Oleynik, Kaushik De, Alexei Klimentov, Alexandre Vaniachine, Artem Petrosyan, Torre Wenaus, and Jaroslava Schovancova. Integration of PanDA workload management system with Titan supercomputer at OLCF. Technical report, ATL-COM-SOFT-2015-021, 2015.
- A. Klimentov et al. Next Generation Workload Management System for Big Data on Heterogeneous Distributed Computing. J. Phys. Conf. Ser., 608(1):012040, 2015.
- World's Most Powerful Accelerator Comes to Titan with a High-Tech Scheduler. BNL Newsroom, May 07, 2014. <http://www.bnl.gov/newsroom/news.php?a=24864>.
- Update on ATLAS and CMS latest results. CERN Press Office, Dec. 16, 2015. <http://press.web.cern.ch/update/2015/12/update-atlas-and-cms-latest-results>.
- D.Oleynik et al. K.De, A.Klimentov. Integration of Tier-1 Grid Center with High Performance Computers at NRC-KI for LHC experiments and beyond HENP. In Computing in High Energy and Nuclear Physics Conference. April 2015, Okinawa, Japan.
- M Borodin, K De, D Golubkov, T Maeno, D Oleynik, A Petrosyan, A Klimentov, T Wenaus, P Nilsson, S Jha, et al. PanDA Beyond ATLAS: A Scalable Workload Management System for Data Intensive Science. Technical report, ATL-COM-SOFT-2014-010, 2014.