

Implementaiton of Hadoop 2 on PBS-based HPC Systems

Joshua Hull, Mark Baker, Alex Berk
School of Computing
Clemson University
Clemson, SC 29632
Email:
{jhull, mnbaker, aberk}@clemson.edu

Abstract—Mark Baker, Alex Berk, and Joshua Hull implemented Hadoop 2 on a PBS-based HPC System on the Palmetto Cluster in a version update from Hadoop 1. Development saw challenges to creating a cluster environment on the Palmetto cluster that still would work with existing Hadoop 1 scripts and infrastructure and PBS script editing to support the additional features and components like YARN in the Hadoop 2 framework. Thorough testing with WordCount and TerraSort was used to measure the success of implementation and initialization of the Hadoop 2 cluster.

I. INTRODUCTION

As part of a group project in Distributed and Cluster Computing Mark Baker, Alex Berk, and Joshua Hull worked to implement the myHadoop2 framework and environment on a cluster of the Palmetto supercomputer. This environment will allow users to run the updated Hadoop 2 framework in an on demand cluster on Palmetto. Hadoop 2 introduced additional tools that need to be handled by myHadoop2, such as the YARN resource management system. These new tools help to reduce bottlenecks in Hadoop and generally increase performance. The purpose of myHadoop2 is to give the users of Palmetto and other PBS based systems a way to set up a dynamic Hadoop 2 cluster in order to take advantage of these new features while still allowing use of the original Hadoop 1 system and scripts.

II. MYHADOOP2 IMPLEMENTATION

A. Technical Implementation

The group working for several weeks modified the existing myHadoop scripts to work with Hadoop 2 from their original Hadoop 1 setup. The modifications involved adjusting the scripts to configure the new YARN resource manager system as well as other new features in Hadoop 2 framework like initializing the resource manager on the namenode of the cluster and changes to the PBS scripts to allow Hadoop 1 scripted jobs to work in the Hadoop 2 framework. The script will need to be able to modify the new components configurations based on the resources that the PBS system has allocated to the user.

In addition to implementing the myHadoop2 script the group also benchmarked the script against its predecessor

in similar cluster environments on Palmetto. The group implemented two standard Hadoop benchmarks: WordCount and TerraSort. These benchmarks were implemented in identical manners and only changed where the differences between Hadoop versions necessitated them to be.

B. Group Members

The group members that were responsible for implementing myHadoop2 successfully were as stated before: Joshua Hull, Alex Berk, and Mark Baker. Each team member brought along various pieces of experience and was responsible for components that reflect their strengths in regards to the projects requirements for success.

Joshua was responsible for benchmark development and testing along with collaborating with the partners on various parts of the project. His approximate four years in using the Linux operating system and installing and configuring software for it meant that he was able to quickly install and configure a system with both Hadoop and Hadoop 2 installed on it. Joshua was also responsible for developing the benchmarks that allowed the team to compare the performance of myHadoop2 to its predecessor using the testing situations in WordCount and TerraSort. The development of the benchmarks was aided by the number of years of Java experience Joshua has learned and worked with.

Alex was primarily responsible for the development of the myHadoop2 script with assistance from Joshua. Alex is currently working with CCIT to get Hadoop 2 working on the Palmetto cluster, and was able to bring this experience to the project and successfully implement that goal for both the success of the project and the goal of CCIT. He is directly involved in the implementation of Hadoop 2 over OrangeFS as well, which greatly benefited him in progressing the myHadoop scripting on the cluster efficiently and with the best performance.

Mark was responsible for documentation of the project and collaboration on development with Alex and Joshua with the implementation to complete this objective. Mark has worked in a functional analyst role over the past three years and has had experience in documentation of analyses and projects and worked closely with Joshua and Alex throughout the development and implementation process to bring a complete

and detailed understanding of the project to its documentation and presentation.

III. CHALLENGES AND DIFFICULTIES

A. Script Modification

During the script editing and testing of the project, Josh came across issues in configuring a Hadoop session scripts and resource system with the new implementation of Hadoop 2. Understand the configuration scripts were a particular challenge and took time to for Josh to educate himself in understanding the major components to the scripts and to find out the tools associated with Hadoop to configure them appropriately for our project. While working on the configuration elsewhere there were times when the configuration itself proved confusing due to their locations in Hadoop 2 deciding where to maintain them versus Hadoop 1 that was well known. With the new features prevalent in Hadoop 2 there was also the issue of implementing the new resource tracking system that wasn't initially set up when creating the system in our project, but with the work Alex had been working on for CCIT to set up their own configuration of Hadoop 2 we were able to borrow the OpenClemson source[1] for the resource system to assist us in getting ours working correctly. One of the last troubles we had to overcome in setting up Hadoop 2 was dealing with the configuration for situations that would be outliers to normal operation of the environment. Such a scenario as with InfiniBand nodes were a contingency Josh had to work through but was easily fixed because this configuration was only needed when running 10GB or more connections between nodes and since the Palmetto is not set up for this kind of implementation a change to the configuration to opt out of managing node names for it was a simple solution to fix any issues that might've arisen even when not using the system.

B. Hadoop 2 Environment

One of our earliest challenges was in setting up the new Hadoop environment on the Palmetto without disrupting the current Hadoop 1 framework and optimizing Hadoop 2 to run its configuration in compatibility with the current Hadoop. Fortunately Alex had been working with CCIT in Clemson on getting this done with OrangeFS and through research and contact found the configured implementation of Hadoop 2.2 from the San Diego Supercomputer Center, and with a few bug fixes to let certain conditions on Palmetto work correctly managed to incorporate the new Hadoop in for scripting and testing without having to greatly edit or modify an older version and now has a well-organized environment for Hadoop users.

IV. CONFIGURATION AND TESTING

As with all new systems, they must be tested in various configurations and with different test suites and the implementation of Hadoop 2 for this project was no exception. Once the initial system was up and running on the Palmetto we ran a battery of tests using WordCount and TerraSort programs for

easy measuring and performance testing due to their reliable and simple architecture that would allow us to test various configurations of jobs on Hadoop 2. At the same time we also ran these same configurations and programs through the Hadoop 1 system so we could then compare results back with Hadoop 2 to see if the performance differences were correctly evaluated and if any issues arose during actual testing and deployment of the new system.

A. Hadoop 2 Testing

B. Hadoop 1 Testing

C. Overall Performance

V. CONCLUSION

This project had many challenges and situations where it took us time to learn and try different ways in getting our system to run through trial-and-error. Now while the system may not be perfect for every use conceivable now, our group believes that the new Hadoop 2 system has proven itself to work and to have better performance and reliability than our current Hadoop 1 and has laid the groundwork for further testing and expansion to replace the old system with ours for everyday use on Palmetto. Through Josh, Alex, and my efforts we have managed to create a new system for our user on the cluster for better big data computation and suggest that our work should be continued to be improved on to make it the best system possible. We ask that our work, all the source and environment, be open to others to use and to develop with to create fixes and improvements to our implementation and to move the Hadoop 2 system to a full replacement of the Hadoop 1 and have a dedicated system in place with scalability for expansion and upgrade with Hadoop versions yet to be developed. We wish everyone to make full use of our PBS-based HPC system with Hadoop 2 and thank you for the opportunity to make it possible in the first place.

ACKNOWLEDGEMENT

REFERENCES

- [1] [Online]. Available: <https://github.com/OpenClemson/myhadoop>