
Virtual Elephant

Thoughts on Enterprise and Cloud Native Architectures

Virtualized Hadoop + Isilon HDFS Benchmark Testing

During the VMworld EMEA presentation (Tuesday October 14, 2014) , the question around performance was asked again with regards to using Isilon as the data warehouse layer and what positives and negatives are associated with leveraging Isilon as that HDFS layer. As with any benchmark or performance testing, results will vary based on the data set you have, the hardware you are leveraging and how you have the clusters configured. However, there are some things that I've learned over the last year and a half that are applicable on a broad scale that can show the advantages to leveraging Isilon as the HDFS layer, especially when you have very large data sets (10+ Petabytes).

There are two benchmarking tests I want to focus on for this post. The tests themselves demonstrate the necessity for understanding the workload (Hadoop job), the size of the data set, and the individual configuration settings (YARN, MapReduce, and Java) for the compute worker nodes.

For each of these tests, we ran the virtualized Hadoop clusters on the very same x86 hardware, shared storage and Isilon arrays. The only two parameters that were modified between each test run was the size of the Hadoop cluster (worker count) and the size of each worker node. As the tests were repeated, it was possible for us to begin to understand the impact of the different configuration settings that can be made within the YARN and MapReduce config files in relation to the size of the worker nodes.

Compute hardware:

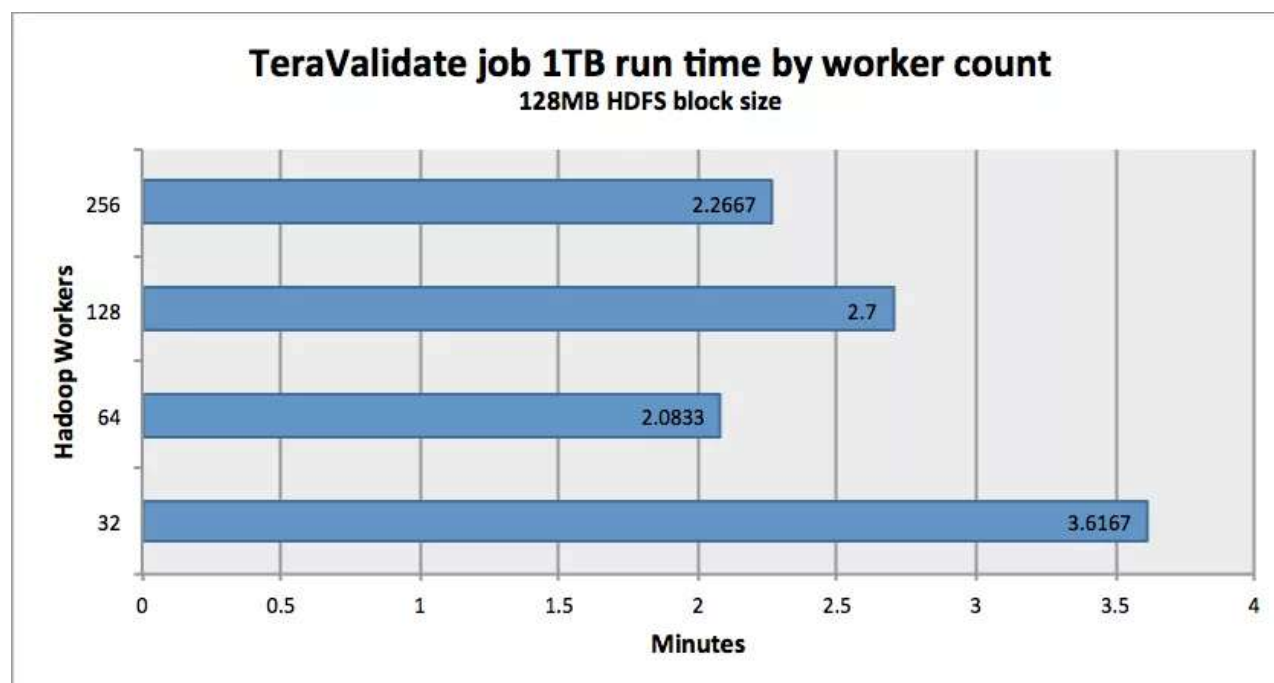
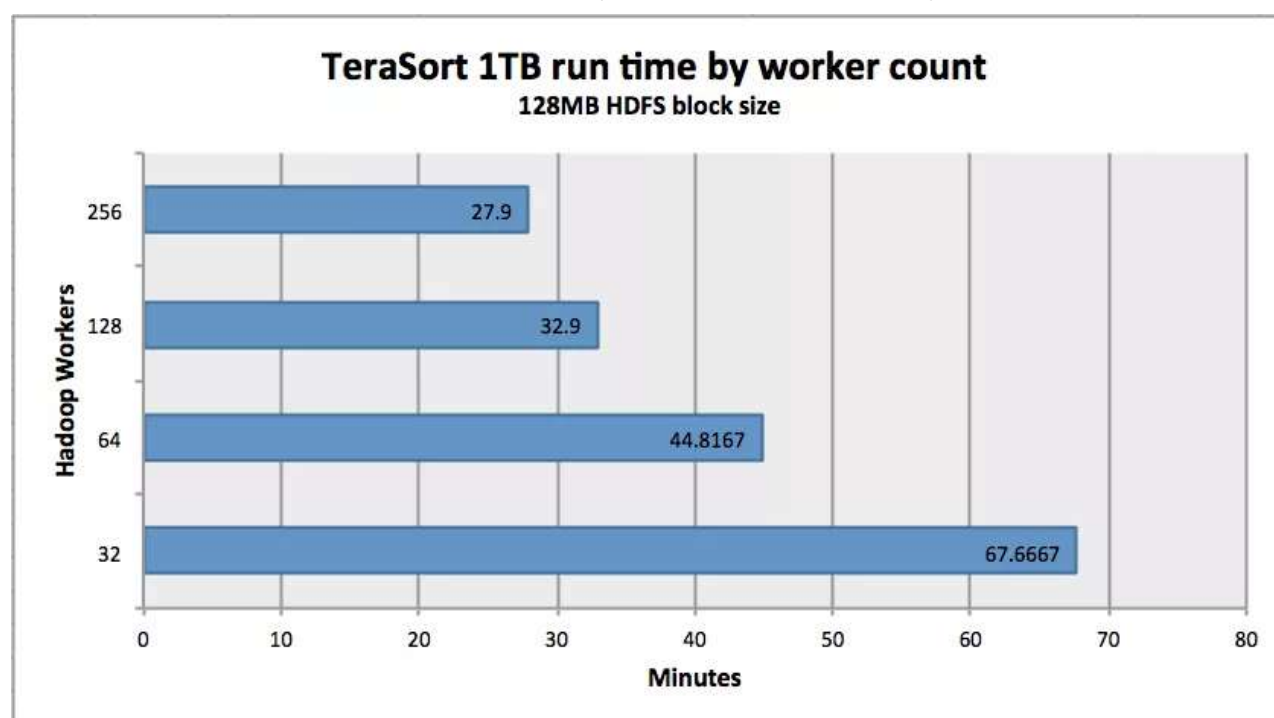
- 2 Cisco UCS 5108 Chassis
- 32 Cisco UCSB-B200-M3 Blade servers (Dual E5-2680v2 CPU, 128GB RAM)

The Cisco servers were connected up to the SAN fabric through a pair of UCS 6296 Fabric Interconnects. Each blade was setup to boot from a dedicated SAN LUN for ESXi. The VMDKs for each Hadoop worker node was attached to the same SAN device providing the boot LUNs. The "scratch" space for the Hadoop jobs was run within each VMDK for the specific worker node, this was not setup to be kept on the Isilon — which is an option.

Cluster size and node configuration:

The tests were ran against four different cluster configurations, limited to the same amount of physical hardware resources to show the differences between the cluster and node sizes. Performing the tests in this manner allows you to see the effectiveness of scaling out the number of nodes within a Hadoop cluster and what effect the node size has within each cluster deployment.

- 32-node Hadoop cluster: 8 vCPU, 58GB RAM per node
- 64-node Hadoop cluster: 4 vCPU, 29GB RAM per node
- 128-node Hadoop cluster: 2 vCPU, 14.5GB RAM per node
- 256-node Hadoop cluster: 1 vCPU, 7.25GB RAM per node



As you can see, there are some improvements you would expect to see and there are areas (64 nodes vs 128 nodes) where additional investigation is required. The numbers themselves are interesting, however beyond saying that "Yes, Isilon can effectively provide an HDFS layer to a Hadoop compute-only cluster", I believe there is still some investigation that can and will take place.

Specifically, the next test cases are three fold using the same physical hardware that we are deploying in our production private cloud environment and the same dataset used in the above tests:

1. Traditional Hadoop clusters without virtualization.
2. Virtualized HDFS data-only cluster (in lieu of Isilon-backed HDFS) and separate compute-only virtualized Hadoop nodes.
3. Isilon-backed HDFS with separate compute-only virtualized Hadoop nodes.

I am of the opinion completing the above tests and comparing the results will help us determine what strategy is best and provide us with a firm understanding of all the advantages and disadvantages to any of the IaaS solutions for Hadoop.

If you are interested in learning more about the above tests and environment that was used to run them, there will be a white paper coming out from EMC soon and I will make that available when it is published.

I encourage you to take time and investigate how to leverage Isilon storage arrays to take advantage of the HDFS protocol and determine for yourself if it makes sense in your environment. As with any technology shift, there are positives and negatives and it is up to us to determine for ourselves what works best for *our environments*.

Share this:





Virtual Elephant. All rights reserved. Copyright 2016.
The opinions expressed on this site are solely my own and do not represent the opinions of my employer.
Nucleus by GalussoThemes.com
Powered by WordPress