## Myth #1: Concern that virtualization will add significant performance overhead to a Hadoop cluster.

This is a common question from users in the early stages of considering virtualizing their Hadoop clusters. Engineers at VMware (and some of its customers) have done several iterations over multiple years of performance testing of Hadoop on VMware vSphere® with various hardware configurations. These tests have consistently shown that virtualized Hadoop performance is comparable to, and in some cases better than that of a native equivalent.

In 2015, a lengthy set of tests conducted on vSphere 6 with 32 host servers and 128 virtual machines, and more, showed that a MapReduce task finished in 12% shorter time on vSphere than the equivalent non-virtualized or native system with four virtual machines per server.

As with any platform on which Hadoop runs, the details of the setup matter. The disk storage, virtual machine placement and networking in particular need to be organized in keeping with the known best practices in order to get the highest performance from the system. The same principles apply to the native world. VMware has documented those best practices and built them into the Hadoop cluster-provisioning tool, VMware vSphere Big Data Extensions™. You can read more on this at
http://www.vmware.com/resources/techresources/10452

## Myth #2: Virtualization requires the use of shared storage.

This is a misunderstanding of the features of virtualization. VMware vSphere works very well with non-shared direct-attached storage (DAS) and hosts HDFS data on that storage. The Hadoop distribution vendors frequently recommend DAS-type storage for cost and performance reasons. With vSphere, each physical disk/spindle in DAS may be presented as a unique datastore to the hypervisor. Virtual disk files (VMDKs) would then be placed onto those datastores. This is a well-understood and tried-and-trusted mechanism. There are large virtualized Hadoop clusters running today that are entirely DAS-based and have no shared storage present at all.

## Myth #3 Hadoop cannot work with Shared Storage.

This is not true, and we are in fact now at a point where a number of users of Hadoop are requesting shared storage to back their clusters. Shared storage comes in many forms, such as SANs, VMware Virtual SAN™ or software-defined storage, NFS devices and HDFS-aware NAS storage mechanisms. SANs and NFS have been deployed in many VMware installations before Hadoop became popular, so they have become associated with vSphere generally, but they are not a pre-requisite, as seen in Myth #2.

The important factor to bear in mind with your choice of storage is the effective bandwidth that is available in terms of Mbytes/second. One can measure the effective bandwidth by using a loading tool such as IOmeter to mimic the traffic seen in Hadoop (long sequential I/Os of 64MB, 128 MB or higher). This is a different measurement to the classic IOPS (I/Os per second) that is used for measuring suitability for an RDBMS or an older style of data storage. Provided the required bandwidth is available to be shared across the number of servers that will be attached to the SAN or NAS, then Hadoop will be deployable there. In general, we see adopters of virtualized Hadoop placing smaller clusters (10 physical servers) for trial purposes on their SAN-based storage, if they intend to place significant performance load on those clusters.

With HDFS-aware NAS type storage, we have seen several deployments already of virtualized Hadoop clusters where the HDFS data is contained solely on the NAS device. The virtual machines in that case contain the compute nodes of Hadoop, such as the Resourcemanager, Nodemanager and Container processes. This has also been shown to scale up to over 100 connected servers running the compute-side virtual machines

### Myth #4: The Hadoop Distribution Vendors do not support virtualized Hadoop.

This is not true. The major vendors of Hadoop software have engaged with VMware to test and validate the behavior of their products on vSphere. These are documented in solution briefs, reference architectures and validation guides that are available from the vendors. VMware's policy is to work with the distro vendor should any problem arise to solve the customer's problem.

### Myth #5: The latest versions of Hadoop products are not available on a virtual platform.

This is a timing question. The vendors' products really depend on the operating system and tools such as Java runtime. Provided the latest Hadoop technology is tested and supported by the Hadoop vendor on the newer versions of these, and VMware supports the particular guest OS, then they will work on vSphere. The installation automation technology implemented in BDE might be a short time behind in its capability to deal with the latest Hadoop version, but this does not mean that the latest version of the Hadoop technology is not capable of being virtualized.

### Myth #6: No one else is doing virtualization of Hadoop. So why should I?

The opening statement here is not true. VMware is aware of many organizations that are in various stages of testing/deploying/running Hadoop on vSphere. A number of these customers have given public talks on the subject and others have had their deployments documented in case studies. Some of these deployments are in the hundreds of servers with multiple hundreds of virtual machines on those servers. There are many good reasons to virtualize Hadoop. One of the main reasons is the agility and simplicity of management you get from abstracting your Hadoop cluster away from the hardware, allowing you to share resources among clusters of different types and versions. For more information on customers who have deployed, please visit http://www.vmware.com/products/big-data-extensions.

## Myth #7: Hadoop and vSphere need have a specialized version to run together.

This is incorrect. The vSphere environment and the Hadoop software can be combined and run out-of-the box. VMware has donated a set of features that make the Hadoop topology aware of virtualization and those are now built into the distro vendors' products (called the Hadoop Virtualization Extensions). These ensure, for example, that all replicas of an HDFS datablock do not live on a group of virtual machines that reside on the same host server. This is a parameter that is expressed at Hadoop cluster creation time and is now a part of the standard Hadoop distribution code. No special versions of vSphere or Hadoop software are required for them to run well together.

## Myth #8: I can deploy Big Data using Containers as a better option to virtualizng Hadoop.

We understand that containers are a big area of interest to many users today; indeed they represent virtualization at a different level in the software stack. However, containers are not that suitable for encompassing all facets of Big Data – though they can play a part. It would be a mighty container indeed that could contain 30-40TB of data. I would not want to fire up such a container or tear it down. Or package it and move it around from one developer to another.

Containers are really for a different purpose than holding the actual HDFS data. They are useful for wrapping the compute-oriented components of Hadoop, such as hosting the Nodemanager and the executable containers/JVMs that run in the same (virtual) machine's operating system. Ideally these would be stateless components in a container landscape. This sort of container may be run inside a virtual machine. But that bulky datastore that makes up the core of Big Data does not belong in a container or set of containers, quite yet. That big data store may be linked to one or more containers, using any number of mechanisms – but it is not wrapped by one. So think carefully about the usage of containers for big data!

## Learn More

To take a deep dive, read Virtualized Hadoop Performance with VMware vSphere 6 on High Performance Servers http://www.vmware.com/resources/techresources/10452.

To read more blogs on Big Data, visit https://blogs.vmware.com/vsphere/2015/12/eight-myths-about-virtualizing-hadoop-dispelled.html

For more information, visit http://www.vmware.com/products/big-data-extensions.