

ESG Lab Review

VCE Vblock Systems with EMC Isilon for Enterprise Hadoop

Date: November 2014 **Author:** Tony Palmer, Senior ESG Lab Analyst, and Mike Leone, ESG Lab Analyst

Abstract: This ESG Lab review documents hands-on testing of the ability of the VCE technology extension for EMC Isilon storage to drive value from enterprise Hadoop as part of a business analytics strategy. Testing focused on the functionality, security, simplicity, and performance of the converged VCE solution, which provides unstructured data analytics with Hadoop.

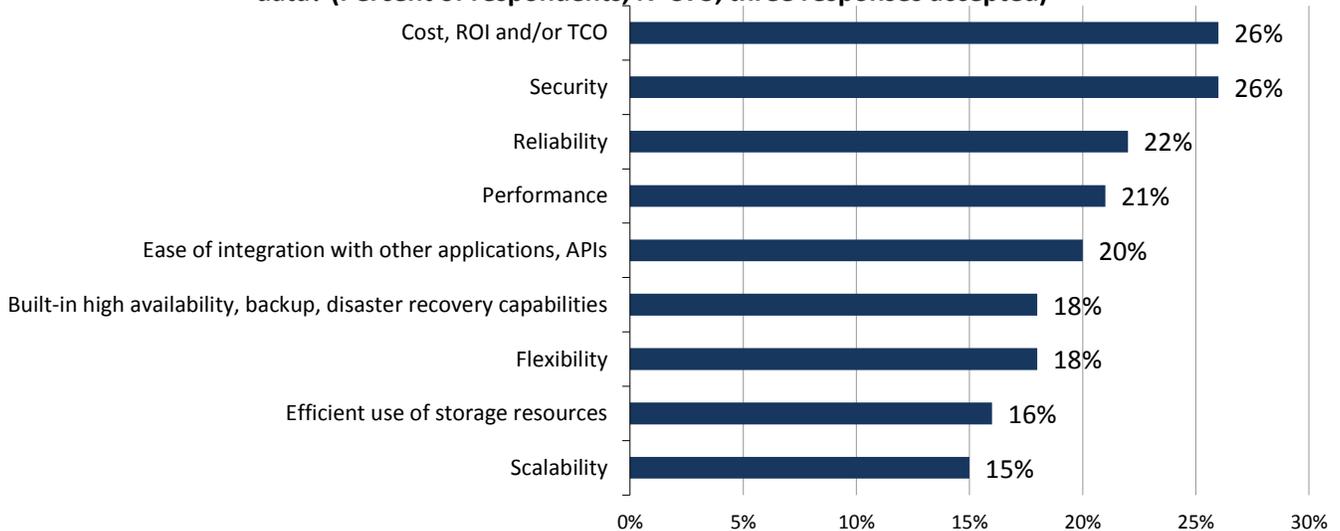
Background

Integrated computing platforms (ICPs) have transformed the way organizations invest in IT infrastructures, and the value of these turnkey solutions is building momentum within the IT community. ICPs, which consist of servers, storage, network, and virtualization, are garnering increased interest among businesses, especially in virtualized and hybrid cloud environments. This is due in part to the fact that ICPs can help ensure that infrastructure deployments are not a bottleneck, and can allow them to be enablers when rolling out or consolidating new or existing applications. Because of the pre-engineered and pre-validated nature of ICPs, organizations experience numerous business benefits. According to ESG research, these benefits include faster deployment times, improved service and support, ease of management, improved scalability, increased agility and reliability, and improved total cost of ownership (TCO).¹

The aforementioned business benefits of ICPs coincide with and complement what organizations are looking for when considering a business intelligence, analytics, and big data solution. In fact, in a recent ESG survey, respondents were asked what their most important attributes were when considering such a solution. As shown in Figure 1, top responses included cost/ROI/TCO, security, reliability, performance, and flexibility, and scalability.² It just seems natural for organizations to extend an ICP to drive value from their big data where they already run their mission-critical workloads.

Figure 1. Top Nine Most Important Considerations for Business Intelligence, Analytics, and Big Data Solutions

Which of the following attributes are most important to your organization when considering technology solutions in the area of business intelligence, analytics, and big data? (Percent of respondents, N=375, three responses accepted)



Source: Enterprise Strategy Group, 2014.

¹ Source: ESG Research Brief, [Integrated Computing Platform Trends](#), August 2014.

² Source: ESG Research Report, [Enterprise Data Analytics Trends: Market Drivers, Organizational Dynamics, and Customer Expectations](#), May 2014.

The Solution: VCE technology extension for EMC Isilon storage for Enterprise Hadoop

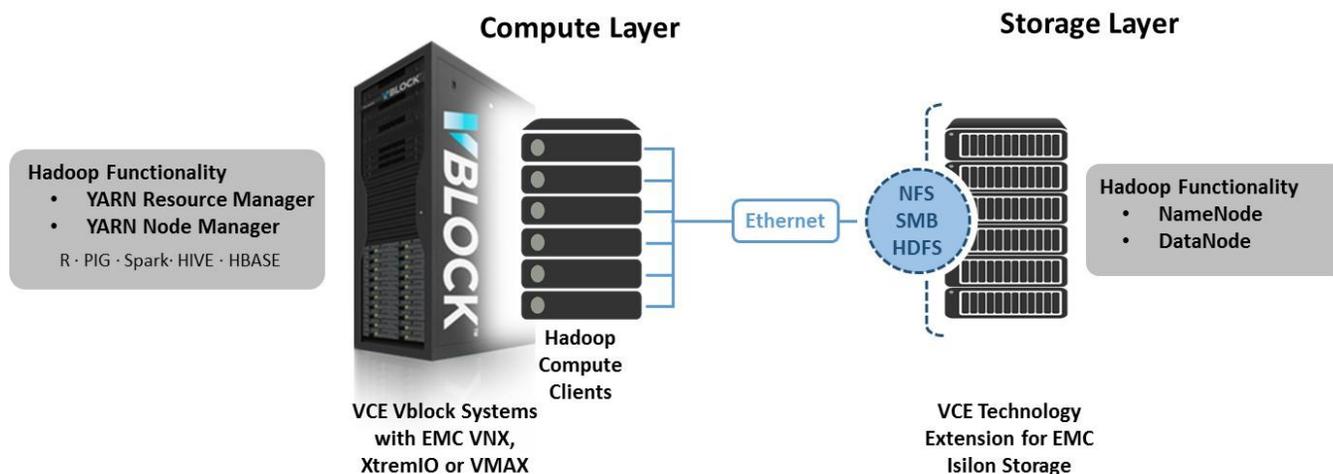
The VCE Vblock System is a well-known and extensively deployed integrated computing platform (ICP) that combines best-of-breed technologies from industry-leading vendors. VCE offers a range of Vblock Systems that businesses can choose from to suit their workload needs. With Cisco compute and networking, EMC storage and data protection, and VMware virtualization and management, Vblock Systems are designed to make it simpler and quicker for organizations to deploy a complete IT platform that has been pre-integrated, pre-validated, and pre-tested with release certification matrices (RCM) serving as a basis for VCE seamless support. Building off the preexisting success of VCE's integrated computing platform, VCE released the VCE technology extension for EMC Isilon storage to extend the business value of VCE. With analytics being a component of the dominant use cases for VCE, organizations can now deploy Hadoop distributions and tools from the Hadoop ecosystem along with dependent applications and databases to provide end-to-end analytics on Vblock Systems with EMC Isilon. Their benefits include

- **Flexible scale-out capacity and performance** from Vblock Systems with EMC Isilon to avoid traditional scale-up limitations in big data analytic environments. This scale-out approach provides improved data protection, data access, resiliency, high availability, manageability, and cost savings.
- **Augmenting the unstructured data analytics** storage environment with native HDFS support from EMC Isilon where unstructured data can be analyzed in-place, eliminating the need for ingest or staging. This augments the already highly prevalent structured data technology built into existing Vblock Systems with EMC Symmetrix VMAX, EMC VNX, and EMC XtremIO storage systems while deploying other analytics such as SAS or Splunk.
- **Predictable high availability and reliability** for extremely large data sets. As data sets reach petabyte scale, the pre-engineered Vblock Systems with scale-out Isilon storage offers high levels of component redundancy, helping to eliminate single points of failure.
- **Data protection** Vblock Systems take advantage of technologies from EMC, including EMC Avamar, EMC Data Domain, EMC RecoverPoint, and EMC VPLEX. EMC Isilon takes the integrated data protection of the complete solution a step further, by offering additional capabilities such as data-at-rest encryption for files, snapshotting, remote replication, and NDMP backups.
- **Networking Hadoop** workloads can be network-heavy, and the ability to choose and architect the leading networking technologies are key to success—contrasting with servers and appliances. Vblock Systems can employ Cisco Nexus 9000 Series Application Centric Infrastructure switches to provide application-driven automation, open software support, and hardware-based multi-tenancy.

As shown in Figure 2, the Hadoop compute resources are decoupled from the EMC Isilon shared storage to deliver a streamlined analytic workflow. There is no longer a need to create a separate environment to ingest data into a Hadoop cluster because the data can be written directly to Isilon using NFS, SMB, HTTP, or FTP and read by the Hadoop cluster using HDFS. This not only eliminates the need for a staging area in a traditional storage system, but also allows for analytics to be done on data that is in-place, while improving storage efficiency by eliminating the 3x replication required with traditional direct attached storage (DAS).

The multiprotocol approach of Vblock Systems 200, Vblock Systems 300, Vblock Systems 500, and Vblock Systems 700 with the VCE technology extension for EMC Isilon storage also helps extend Isilon's usefulness beyond Hadoop by including home directories, backup and archival, surveillance, and high-performance computing, while the range of block storage options available in Vblock Systems—EMC VNX, XtremIO, or VMAX—enable support for structured data management using traditional databases like Oracle and SQL Server.

Figure 2. Enterprise Hadoop with VCE Vblock Systems and VCE technology extension for EMC Isilon storage



EMC Isilon

Isilon scale-out NAS works as a distributed system that consists of multiple nodes in a cluster. The EMC Isilon OneFS operating system combines all hardware resources from each node, including memory, CPU, and disk, into a global namespace managed by a single file system. This creates a shared infrastructure with no single points of failure. As nodes are added to the cluster, the file system dynamically adjusts and redistributes data accordingly—as more nodes are added to a cluster, the performance and resiliency of that cluster increases. Similarly, increasing performance, capacity, and resources through upgrades and expansion of compute, storage, and network on Vblock Systems is simpler to manage as organizations aggregate workloads to manage large analytics and applications.

By using EMC Isilon Storage as the back-end of a Hadoop cluster, existing workflows that utilize SMB and NFS protocols make it easier for both IT and end-users to collect and manage the data they want to analyze. EMC Isilon also natively integrates the HDFS protocol to help organizations streamline the development and deployment of proven Hadoop workflows while reducing staging, mirroring, rewriting, and deleting processes. EMC Isilon fully supports Apache Hadoop along with all leading Hadoop distributions, including portability between distributions, and is fully supported by VCE. Also, with support for running simultaneous instances of different Hadoop distributions that access the same or different data sets on EMC Isilon, organizations can gain new levels of deployment flexibility and support to provide analytics as a service in the Vblock System with the technology extension.

How It Works

Data stored on the Isilon cluster is accessed over HDFS by Hadoop compute clients running MapReduce jobs on the VCE Vblock System. The Isilon OneFS operating system implements the HDFS protocol on every node in the cluster, enabling it to function as the native HDFS storage for the Hadoop compute clients. Each node in the Isilon cluster acts as a Hadoop NameNode and DataNode, with the NameNode working as a distributed process running on every node in the cluster, reinforcing high levels of availability and resiliency because no secondary NameNode is required. Hadoop manages job tracking and task tracking on compute nodes.

As HDFS connections come into the Isilon cluster, they are automatically load-balanced across each node in the cluster. The Isilon OneFS operating system stripes all data across the cluster and uses parity at a file level for data protection. This means that any node in the cluster can simultaneously serve traffic as a DataNode or NameNode. Rack locality can also be provided by assigning compute clients to specific Isilon nodes in the cluster.

VMware vSphere Big Data Extensions with VCE Vblock Systems and EMC Isilon

VMware first introduced project Serengeti in 2012, as an open-source project that makes it easier to deploy and manage Hadoop using virtualization. As Serengeti gained traction, VMware created a commercially supported version called VMware vSphere Big Data Extensions (BDE) as part of the vSphere Enterprise and Enterprise Plus Editions. BDE provides a simple, efficient, and flexible way to run and scale Hadoop applications by automating the deployment and management of Hadoop provisioning by the owner of the Hadoop cluster. While virtualizing the Hadoop cluster, the compute and storage functions of Hadoop may be separated, thereby creating an elastic compute environment. This can then be easily scaled to meet the demands of any Hadoop application. This can be taken a step further to offer on-demand scalability with the ability to add virtual nodes to existing bare-metal Hadoop deployments that already exist in an organization's infrastructure.

Being part of vSphere Enterprise and Enterprise Plus Editions, BDE is itself deployed as a pair of virtual machines: a management virtual machine and a template virtual machine. The management virtual machine serves as the brain of the system, while the template virtual machine is cloned as many times as necessary to meet the deployment requirements of a newly created virtualized Hadoop cluster. When using BDE with VCE and EMC Isilon, the Vblock System handles the storage and hosting of the Hadoop compute virtual machines and roles, including the Compute Master, Worker, and Client node groups. Because Isilon already decouples Hadoop storage from the compute roles, The BDE user simply creates the compute-only parts of the Hadoop cluster in the new virtual machines and provides an address to the HDFS interface on the Isilon cluster. This makes the deployment of the virtualized Hadoop cluster even faster.

Hadoop Performance

By separating compute and storage with VCE Vblock Systems and EMC Isilon's HDFS-enabled NAS, organizations can benefit from advancements in interconnected network and Ethernet connectivity technology to deliver high levels of performance and low latency processing. Resources can scale independently of one another, depending on requirements, to not only improve overall performance, but also achieve the additional levels of data protection, high availability, and scalability of a shared storage model.

ESG Lab used the Hadoop TeraSort suite to validate the HDFS and MapReduce layers of a VCE Vblock Systems and EMC Isilon joint-solution. The suite consists of three different tests that measure the performance of different areas in the Hadoop clusters:

- **TeraGen** generates a random data set, which MapReduce jobs can be run against by TeraSort.
- **TeraSort** combines testing the HDFS and MapReduce layers of a Hadoop cluster using the input data from TeraGen and MapReduce to sort the data.
- **TeraValidate** validates that TeraSort properly sorted the data.

The first phase of performance testing leveraged a Vblock System 340 with EMC VNX storage and EMC Isilon with SSDs to complete the series of tests on a growing data set. The test bed consisted of a cluster running Hadoop 2.3 with 16 Hadoop worker nodes running on Cisco B200 M3 blade servers with 256GB RAM and dual 10-core 2.80 GHz E5-2680 processors for each, with two virtual machines per blade, and no map output compression. An EMC Isilon cluster consisting of eight Isilon X400 nodes was employed for HDFS data space³, and an EMC VNX 5400 with 90 x 600GB drives was used for shuffle space.⁴ The Cisco Network configuration included Nexus 5548Up and UCS Fiber Interconnect 6248. When optimizing for performance, ESG encourages organizations to vary compute, network, and storage along with the virtual machine consolidation ratio, compression, and other YARN parameters to meet the desired price-performance target. Other factors such as security, availability, and other mission-critical features must be balanced with pricing considerations.

The data set size was scaled from 100GB to 1TB and job completion time was monitored in each test case. The results are shown in Figure 3 and Table 1. The key takeaways come in terms of speed and predictability. All three tests not only

³ Note: Testing was executed with previous generation X400 nodes due to lab availability, current Isilon systems ship with X410 nodes, which should offer higher baseline performance.

⁴ See Tables 3 and 4 in the Appendix for the details of the test environment for Vblock Systems with EMC Isilon.

completed jobs quickly, but also did so faster than ESG Lab expected in some cases, based on past experience with the TeraSort suite on similarly sized clusters. As the data set size increased, job duration met or exceeded ESG Lab expectations in all test cases.

Figure 3. Running the TeraSort Test Suite on a Growing Data Set

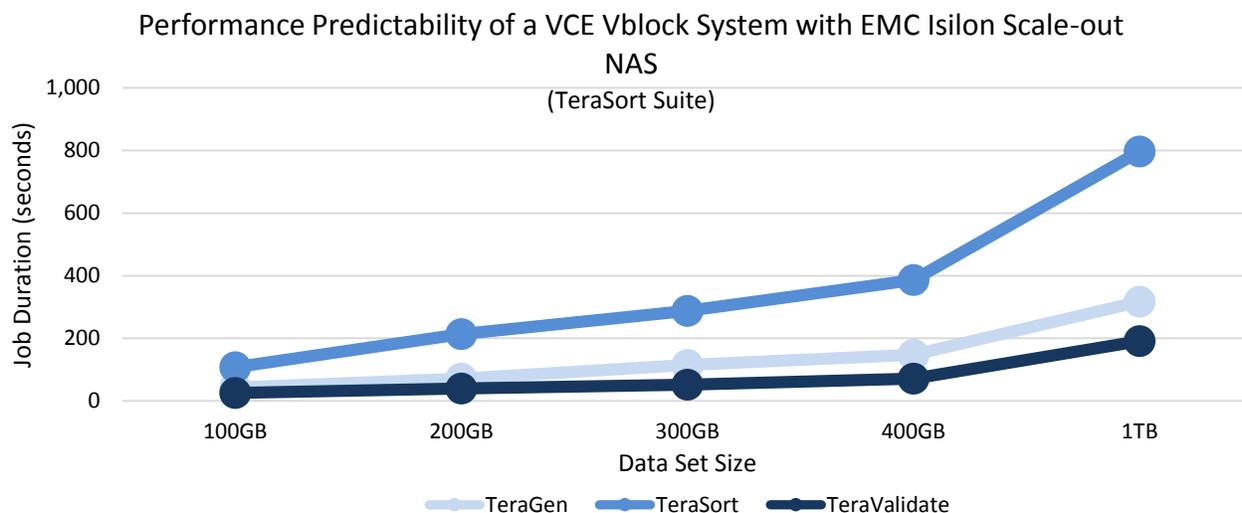


Table 1. Running the TeraSort Test Suite on a Growing Data Set

Data Set Size	Time for a Job to Complete (seconds)		
	TeraGen	TeraSort	TeraValidate
100GB	43	108	25
200GB	71	214	40
300GB	115	288	52
400GB	147	386	71
1TB	317	797	191

Next, ESG Lab compared the performance of a traditional Hadoop cluster consisting of commodity servers and DAS to a VCE Vblock Systems and EMC Isilon implementation of Hadoop. The goal was to understand the performance benefits that can be achieved when leveraging a platform that separates the storage from the compute and eliminates the need to ingest data for analytic jobs.

The traditional Hadoop cluster⁵ consisted of 16 nodes, while the directly compared test bed for VCE and Isilon consisted of 16 compute nodes and eight storage nodes. The compute nodes in the traditional Hadoop cluster leveraged servers that had two E5-2650 v2 CPUs, 64GB of RAM, eight 2.5" 300GB 10K SAS 6Gbps internal drives for HDFS, temp files, and Hadoop shuffle space, and a single 10GbE connection; while VCE Vblock Systems with EMC Isilon leveraged Cisco UCS B200 M3 blades with 256GB RAM, 2x E2680 v2 CPUs, and no internal storage. A single worker virtual machine within each compute node was configured with 32 vCPUs, and 56GB of RAM.

In all test cases, the total time for each job to complete was monitored and the results are shown in Figure 4. For each test, the VCE and Isilon solution outperformed the traditional Hadoop environment. The VCE and EMC Isilon configuration completed the TeraGen job in less than half the time it took traditional Hadoop. VCE and EMC Isilon reduced the time required for the TeraSort job by a third, and TeraValidate completed 26 percent faster.

⁵ See Table 5 in the Appendix for the details of the traditional Hadoop cluster configuration.

Figure 4. Comparing Performance of Traditional Hadoop to a VCE Vblock Systems with EMC Isilon Scale-out NAS

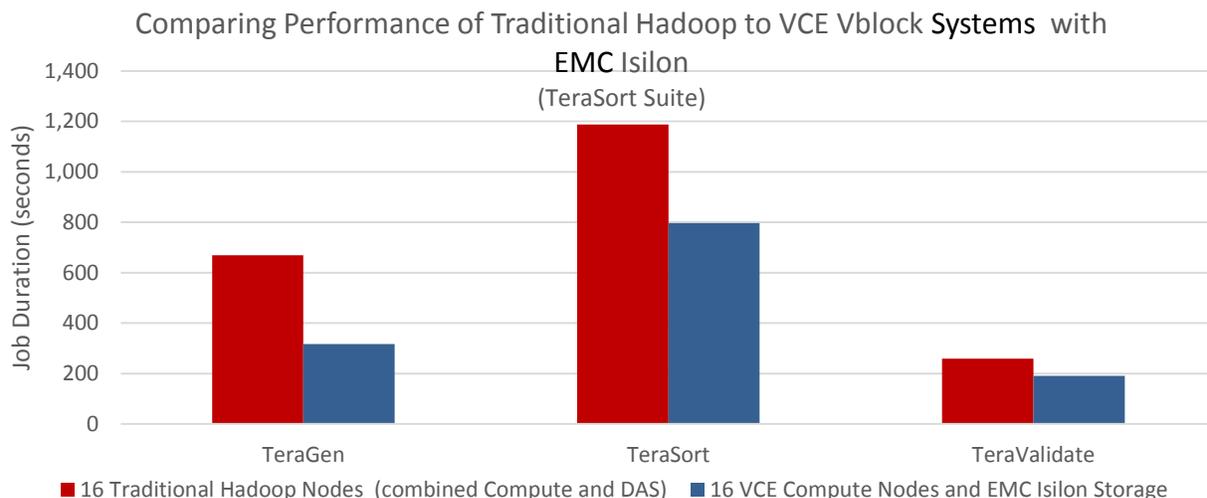


Table 2 lists the detailed results of the performance comparison tests.

Table 2. Running the TeraSort Test Suite on a Growing Data Set

Configuration	Time for a Job to Complete (seconds)		
	TeraGen	TeraSort	TeraValidate
16 Traditional Hadoop Nodes	669	1,187	259
16 VCE Compute Nodes and EMC Isilon Storage	317	797	191

Why This Matters

In traditional Hadoop implementations, organizations have leveraged commodity servers and DAS as combined compute and storage nodes. To meet the ever-increasing requirements of scalability and performance, organizations must add servers with DAS, growing compute and storage together. Though this approach is widely accepted, it does not provide the ultimate flexibility organizations need to avoid potential bottlenecks or overprovisioning of resources, and it doesn't address growing security and protection requirements. The big data problem is, at its most basic, similar to the structured database problem. Finding and correcting one bottleneck typically reveals another, and this is particularly hard to solve in a DAS or appliance environment where compute and storage are tied together.

VCE Vblock Systems with EMC Isilon addresses this issue because the validation and testing is designed to ensure that implementation, use, and management deliver maximum ease of use, protection, value, reliability, and consistency of performance.

ESG Lab validated that VCE Vblock Systems with EMC Isilon are well suited to deliver levels of virtualized Hadoop performance comparable to bare-metal installations in a scalable, flexible package. EMC Isilon's single file system architecture provides an ideal foundation to independently scale Hadoop storage and compute resources to meet performance and capacity requirements as needed. When comparing TeraSort Suite test results on a traditional Hadoop configuration with a VCE Vblock Systems and EMC Isilon, the latter yielded significant performance benefits, completing Hadoop Jobs in as little as half the time.

Flexibility, Efficiency, Security, and Compliance

In selecting technology solutions in support of business intelligence and other analytics projects, there are many considerations to take into account, and decision criteria aren't as simple as "better, cheaper, and faster"—they span a broad range of attributes. While financial concerns—such as cost and predicted return on investment—top the list, other top evaluation criteria include operational requirements like security, reliability, and flexibility.

ESG Lab looked at multiple areas of functionality with a goal of exploring the true differentiation VCE with EMC Isilon brings to Hadoop and big data analytics. The test bed for this phase of testing consisted of two Hadoop "mini" clusters running Hadoop 2.3 with three worker nodes each for compute using the same EMC Isilon cluster for storage with a distinct storage zone set up for each cluster: mini1-master-0 and mini2-master-0.

First, the use of multiple zones in a global namespace to provide secure and shared capacity was examined. ESG Lab created a file on the mini2 cluster, and was unable to view or access the file from the mini1 cluster.

Next, ESG Lab looked at authentication and permissions. A file was created over SMB, and the ESG Lab analyst viewed it using the same user account over both NFS and HDFS. Both the contents and permissions of the file were verified to be correct and identical. Then a file was created using HDFS and read using NFS and SMB with the same result. Finally, NFS was used to create a third file and the test was repeated, also successfully. This test was performed with local accounts, but Active Directory and LDAP can be used for user authentication as well. Translation between protocols occurred in real time and was completely transparent. The simultaneous access from SMB and HDFS is unique to EMC Isilon as of this writing, and enables data in place processing with no need to copy or ingest data.

This capability also enables secure, transparent sharing of data between clusters, either using permissions to allow another cluster to view and access data across zones, or by creating a separate zone for sharing and giving multiple clusters access to the shared zone. ESG Lab tested all of these capabilities, and they worked precisely as expected.

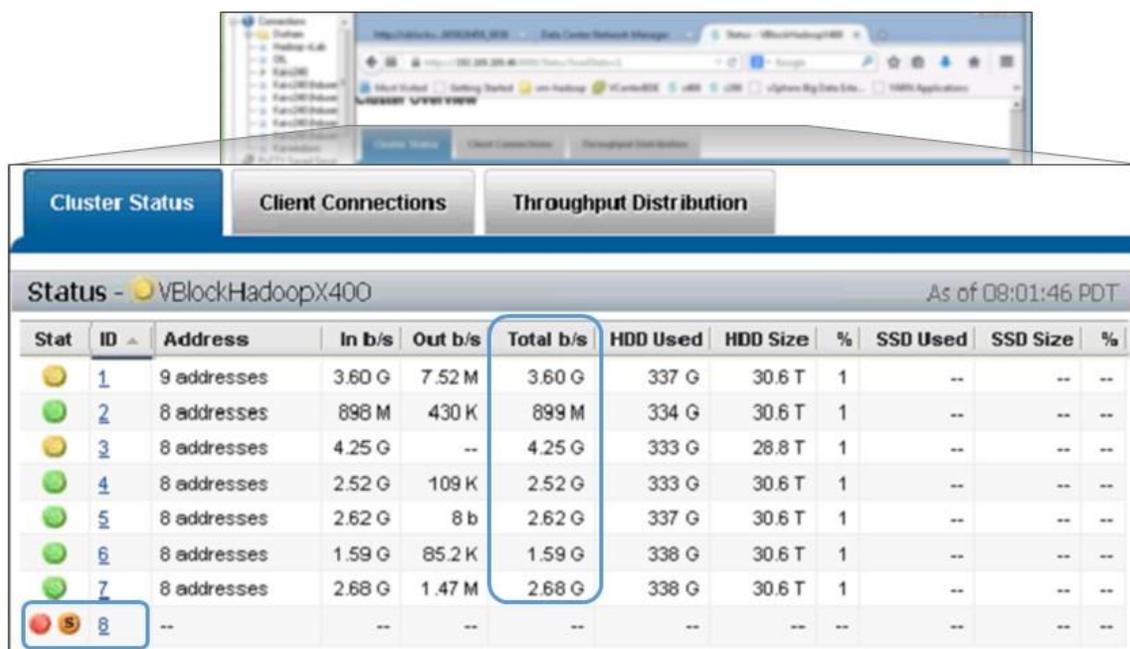
EMC Isilon SmartLock provides write once, read many (WORM) functionality designed to support corporate governance or full compliance. When SmartLock is enabled, administrators can set files and folders to WORM using the REST API, HDFS, NFS, or SMB. ESG Lab tested SmartLock and confirmed that files cannot be deleted before their expiration date when they're locked, but are still available for analytics. When compliance mode is enabled, the root user is removed and files may not be deleted.

Next, ESG Lab examined running MapReduce on a growing, open file. This scenario would be useful if analytics needed to be run on server log files without having to periodically copy the log files to a staging area and ingest them into HDFS. A large file was written to the Isilon cluster using a python script to write words to a file continuously. While this script was running and the file was growing, a MapReduce job was run to count the words and output the total number of words in the file. The job was run multiple times, showing a growing number of words until the word generator script completed. This test was run from a Linux system over NFS and from a Windows system over SMB.

Finally, in order to test resiliency of the Vblock Systems and EMC Isilon-based Hadoop cluster, ESG Lab intentionally powered down one of the eight Isilon nodes, removed it from the cluster and observed the results. Before removing the node, TeraSort was run using the healthy system with eight nodes against a 1TB data set. When the node was removed, EMC Isilon immediately began restriping data across the remaining seven nodes.

While the data from the removed node was building across the remaining seven, TeraSort was re-run against the 1TB data set. Where TeraSort completed in 873 seconds with the healthy eight-node cluster, the seven-node cluster was able to complete the TeraSort job in 999 seconds, while the data was rebuilding. This represents a 12 percent difference in performance, or almost precisely seven-eighths of the performance of the healthy eight-node cluster.

Figure 5. Running TeraSort During Data Rebuild



Why This Matters

A growing number of organizations are deploying big data and analytics platforms to improve the efficiency and profitability of their businesses. ESG asked mid-sized organizations and enterprises which business initiatives would drive the most technology spending in their organizations, and improved data analytics, business intelligence, and customer insight was the third most-cited response.⁶ In the same survey, one in five organizations reported a problematic shortage of existing skills in business intelligence and data analytics. This situation demands a modern infrastructure where organizations can deploy precious resources to productively and securely perform their analytic tasks as an extension of their current environments.

ESG Lab validated that VCE Vblock Systems with EMC Isilon provide real and tangible advantages for organizations deploying Hadoop for big data and analytics. Vblock Systems with EMC Isilon simplified the deployment and management of Hadoop clusters. By eliminating the need for staging storage and data ingest, Vblock Systems with EMC Isilon facilitated the removal of duplicated data, reduced overall job time, and provided the ability to transparently share data across and between clusters with robust authentication and permissions.

Security and compliance were also robust, enabling multi-tenancy and read-only access to data when needed. ESG Lab also validated that vSphere Big Data Extensions (BDE) allow the automatic provisioning of Hadoop nodes on demand, as needed, for both virtual Hadoop clusters and as virtualized node additions to existing bare-metal clusters.

Businesses can easily extend new or existing Vblock Systems running business intelligence or related tasks with EMC Isilon to support Hadoop projects. Since Hadoop is often run in silos owned by business units, and requires extra steps to perform business analytics, bringing all these workloads into the VCE environment and running them as a service should help decrease time to insight while increasing the utilization, operational efficiency, visibility, and control needed in Hadoop environments in the enterprise.

⁶ Source: ESG Research Report, [2014 IT Spending Intentions Survey](#), February 2014.

The Bigger Truth

The big data market is continuing to evolve as organizations are facing increased pressure to drive value from their big data investments. Hadoop offers great advantages to many organizations, especially by being able to leverage commodity hardware to handle massive amounts of data, but challenges quickly arise as those data sets continue to grow. Ingest rates and analytic service-level agreements can tax infrastructures and create bottlenecks, while limited flexibility makes it hard for organizations to reliably and cost-effectively scale resources independently of one another. More importantly, security and reliability risks come to the fore as businesses mature their Hadoop deployments into larger scale production as part of the end-to-end analytic value chain including traditional business intelligence or other popular analytics such as SAS or Splunk.

By leveraging an industry-proven ICP in VCE Vblock Systems and combining it with EMC Isilon and VMware vSphere Big Data Extensions, organizations get a fully integrated platform that meets and grows with their big data and analytics requirements. BDE makes deploying Hadoop clusters quick and easy, while Vblock Systems and Isilon deliver a scale-out infrastructure that achieves higher levels of security, reliability, availability, and performance for unstructured data analytics in production environments. With native HDFS support, the flexible, single file system architecture of Isilon eliminates the need to move data into a Hadoop cluster and instead brings Hadoop to the data. And by separating Hadoop compute functionality from Hadoop storage functionality, resources can easily scale to meet big data capacity or performance demands without having to overprovision.

ESG Lab validated that VCE Vblock Systems with EMC Isilon delivered Hadoop performance favorably comparable to a bare-metal installation in a scalable, flexible package. EMC Isilon's multi-protocol architecture provides a foundation to independently scale Hadoop storage while VCE Vblock Systems scale compute resources to meet performance and capacity requirements as needed. VCE Vblock Systems with EMC Isilon completed jobs as much as twice as fast as a traditional Hadoop configuration. While this testing used EMC VNX for Hadoop shuffle, users are not limited to VNX for shuffle space. EMC VNX, VMAX, and XtremIO are all available in Vblock Systems. Organizations can choose the platform based on workload and budgetary demands.

ESG Lab testing revealed that VCE Vblock Systems with EMC Isilon provide real and tangible benefits for organizations deploying Hadoop for big data and analytics. Vblock Systems with EMC Isilon simplified the deployment and management of Hadoop clusters. Vblock Systems with EMC Isilon eliminated the need for staging storage and data ingest, which facilitated the removal of duplicated data, reduced overall job time, and provided the ability to transparently share data across and between clusters with robust authentication and permissions.

Security and compliance support were also robust, enabling multi-tenancy and read-only access to data when needed. ESG Lab also validated that VMware vSphere Big Data Extensions (BDE) allows the automatic provisioning of Hadoop nodes on demand for both virtual Hadoop clusters and as virtualized node additions to existing bare-metal clusters.

Cisco, EMC, and VMware have invested in the differentiating features and functionalities integrated into VCE Vblock Systems to ensure that load is balanced across the infrastructure, ensuring maximum utilization. This enables jobs to complete faster, which has a direct, positive impact on the bottom line. Whether deploying a traditional Hadoop distribution, expanding an already deployed Hadoop implementation, or starting fresh with a new Hadoop implementation that you seek to combine with your existing IT environment, ESG Lab suggests checking out the VCE Vblock Systems with EMC Isilon for enterprise Hadoop to advance your analytic strategy.

Appendix

Table 3. Test Environment—Key Components and Configuration for Vblock Systems with EMC Isilon

Element	Configuration
VCE Vblock System 340	<p>VCE Management</p> <ul style="list-style-type: none"> • Advanced Management Platform (AMP-2) <p>Fabric Interconnect</p> <ul style="list-style-type: none"> • Cisco UCS 6248UP Fabric Interconnect <p>Servers</p> <ul style="list-style-type: none"> • 16 X Cisco UCS B200 M3 <p>Server Details</p> <ul style="list-style-type: none"> • 2 X Xeon Intel E5-2680V2 (2.8 GHz) • 256GB memory (16 X 16 GB) • Host local storage: None (diskless; SAN Boot) • 2 X Cisco UCS VIC-1240 • Cisco UCS 5108 Blade Server Chassis <p>Storage – EMC VNX 5400</p> <ul style="list-style-type: none"> • Connectivity: Fibre Channel • Drive Count: 92 X 600GB 10K SAS (HUC10906) • Single RAID5 Storage Pool Configuration <ul style="list-style-type: none"> ○ 90 Drives in Storage Pool—all LUNs listed below provisioned from this pool ○ SAN Boot devices for UCS blades (16 LUNs X 20GB) ○ Storage for ESXi and virtual machines (10 LUNs X 2TB) <ul style="list-style-type: none"> ▪ Virtual machines include BDE server, master and workers ○ Hadoop shuffle space (10 LUNs X 2TB) • 8Gb Fibre Channel Networking—8 lanes from switch to VNX <p>Networking</p> <ul style="list-style-type: none"> • Switches: 48 ports Cisco • A pair of N5K-C5548UP capable of 10 Gigabit Ethernet, Fibre Channel, and FCoE switch
EMC Isilon (offered as VCE technology extension for EMC Isilon storage)	<ul style="list-style-type: none"> • HDFS configuration: <ul style="list-style-type: none"> ○ 512MB Hadoop block size • Nodes: 8 • Isilon X400-4U-Dual-96GB-2x1GE-2x10GE SFP+-34TB-800GB SSD • OneFS 7.1.1.0 with patch-130611 • Dual 10 Gig Ethernet connectivity per Isilon node • Additional pair of N5K-C5548UP

Table 4. Test Environment—Key Components and Configuration for Vblock System with EMC Isilon—Continued

Element	Configuration												
Hadoop	<ul style="list-style-type: none"> • Hadoop Version 2.3.0 • Cluster Topology: Host_AS_RACK • Non-default parameters <ul style="list-style-type: none"> ○ Tasks <table border="1" data-bbox="727 520 1507 730"> <thead> <tr> <th>Benchmark</th> <th>Map</th> <th>Reduce</th> </tr> </thead> <tbody> <tr> <td>TeraGen</td> <td>1024</td> <td>0</td> </tr> <tr> <td>TeraSort</td> <td>2048 (derived)</td> <td>1024 and 320</td> </tr> <tr> <td>TeraValidate</td> <td>320 (derived)</td> <td>1</td> </tr> </tbody> </table> ○ yarn.nodemanager.resource.memory-mb=80000 ○ mapreduce.map.memory.mb=8000 ○ mapreduce.task.io.sort.mb=1024 ○ mapreduce.reduce.memory.mb=8192 ○ mapreduce.map.java.opts=-Xmx6000m ○ mapreduce.reduce.java.opts=-Xmx6000m 	Benchmark	Map	Reduce	TeraGen	1024	0	TeraSort	2048 (derived)	1024 and 320	TeraValidate	320 (derived)	1
Benchmark	Map	Reduce											
TeraGen	1024	0											
TeraSort	2048 (derived)	1024 and 320											
TeraValidate	320 (derived)	1											
Hypervisor	<ul style="list-style-type: none"> • vSphere 5.5 – EMC PowerPath/VE enabled • vCenter 5.5 • Storage: HDFS for Isilon <ul style="list-style-type: none"> ○ Isilon HDFS protocol ○ VNX Fibre Channel for Shuffle • vSphere Distributed Switch with multiple VLANs 												
Virtual Machines	<ul style="list-style-type: none"> • Deploy via VMware vSphere Big Data Extension 2.0 <ul style="list-style-type: none"> ○ VM version 10 ○ VMware Tools installed ○ HDFS using Isilon HDFS protocol ○ VNX for VM and shuffle • 2 VMs per server <ul style="list-style-type: none"> ○ 96GB memory ○ 15 vCPUs 												
Linux	<ul style="list-style-type: none"> • Provisioned by BDE 2.0 details include <ul style="list-style-type: none"> ○ CentOS 6.4 x84_64 ○ /etc/security/limits.conf <ul style="list-style-type: none"> ▪ nofile=32768 ▪ nproc=32000 ○ SELINUX disabled (/etc/selinux/config) ○ Java version "1.7.0_51" ○ HDFS based on Isilon HDFS protocol • Shuffle partition formatted with EXT4 												

Table 5. Test Environment—Traditional Hadoop Cluster Configuration

Element	Configuration												
Hadoop	<ul style="list-style-type: none"> • Hadoop Version 2.0.5 • Non-default parameters <ul style="list-style-type: none"> ○ dfs.blocksize=512M ○ io.file.buffer.size=131072 ○ mapreduce.map.java.opts=-Xmx1536m ○ mapreduce.map.memory.mb=2048 ○ mapreduce.map.output.compress=false ○ mapreduce.reduce.java.opts=-Xmx1536m ○ mapreduce.reduce.memory.mb=2048 ○ mapreduce.task.io.sort.factor=3000 ○ mapreduce.task.io.sort.mb=768 ○ yarn.app.mapreduce.am.resource.mb=1024 ○ yarn.nodemanager.resource.memory-mb=45076 ○ Tasks <table border="1" data-bbox="699 905 1490 1115"> <thead> <tr> <th>Benchmark</th> <th>Map</th> <th>Reduce</th> </tr> </thead> <tbody> <tr> <td>TeraGen</td> <td>80</td> <td>0</td> </tr> <tr> <td>TeraSort</td> <td>2048 (derived)</td> <td>100</td> </tr> <tr> <td>TeraValidate</td> <td>100 (derived)</td> <td>1</td> </tr> </tbody> </table>	Benchmark	Map	Reduce	TeraGen	80	0	TeraSort	2048 (derived)	100	TeraValidate	100 (derived)	1
Benchmark	Map	Reduce											
TeraGen	80	0											
TeraSort	2048 (derived)	100											
TeraValidate	100 (derived)	1											

The goal of ESG Lab reports is to educate IT professionals about data center technology products for companies of all types and sizes. ESG Lab reports are not meant to replace the evaluation process that should be conducted before making purchasing decisions, but rather to provide insight into these emerging technologies. Our objective is to go over some of the more valuable feature/functions of products, show how they can be used to solve real customer problems and identify any areas needing improvement. ESG Lab’s expert third-party perspective is based on our own hands-on testing as well as on interviews with customers who use these products in production environments. This ESG Lab report was sponsored by VCE.

All trademark names are property of their respective companies. Information contained in this publication has been obtained by sources The Enterprise Strategy Group (ESG) considers to be reliable but is not warranted by ESG. This publication may contain opinions of ESG, which are subject to change from time to time. This publication is copyrighted by The Enterprise Strategy Group, Inc. Any reproduction or redistribution of this publication, in whole or in part, whether in hard-copy format, electronically, or otherwise to persons not authorized to receive it, without the express consent of The Enterprise Strategy Group, Inc., is in violation of U.S. copyright law and will be subject to an action for civil damages and, if applicable, criminal prosecution. Should you have any questions, please contact ESG Client Relations at 508.482.0188.