



SECURITY AND COMPLIANCE FOR SCALE-OUT HADOOP DATA LAKES

ABSTRACT

This paper describes how the Dell EMC Isilon scale-out NAS platform protects the confidentiality, availability, and integrity of Hadoop data to help meet compliance regulations. By implementing the HDFS protocol natively, the Isilon storage system provides a multiprotocol data lake that secures Hadoop data with identity management, authentication, access control, file-level permissions, WORM, data-at-rest encryption, and auditing.

October 2016

The information in this publication is provided “as is.” DELL EMC Corporation makes no representations or warranties of any kind with respect to the information in this publication, and specifically disclaims implied warranties of merchantability or fitness for a particular purpose.

Use, copying, and distribution of any DELL EMC software described in this publication requires an applicable software license.

DELL EMC², DELL EMC, the DELL EMC logo are registered trademarks or trademarks of DELL EMC Corporation in the United States and other countries. All other trademarks used herein are the property of their respective owners. © Copyright 2016 DELL EMC Corporation. All rights reserved. Published in the USA. <10/16> <white paper> < H13354>

DELL EMC believes the information in this document is accurate as of its publication date. The information is subject to change without notice.

DELL EMC is now part of the Dell group of companies.

TABLE OF CONTENTS

EXECUTIVE SUMMARY	1
INTRODUCTION	1
Information Security and Regulatory Compliance	2
The Costs of Noncompliance	4
Compliance Problems with Hadoop	4
STORING HADOOP DATA ON ISILON SCALE-OUT NAS	5
Use Cases	5
Elasticity, Simplicity, Flexibility	6
Separating Data from Compute	6
The OneFS HDFS Implementation	6
SECURING HADOOP DATA	7
Role-Based Access Control for Administration	8
Compliance Mode, WORM, and the Root Account	8
Controlling Access to Hadoop Data	9
Access Zones	10
Identity Management	10
Kerberos Authentication.....	11
ID Mapping	12
User Mapping	12
Securing Data at Rest with Self-Encrypting Drives	13
Protecting Data In Transit with Partner Solutions.....	13
Supporting SEC Rule 17a-4.....	13
THE ONEFS UNIFIED PERMISSIONS MODEL.....	14
ACL Policies for Mixed Environments	14
The OneFS Permissions Model and Compliance	15
AVAILABILITY AND INTEGRITY	16
Isilon Architecture for Availability	16
OneFS Availability and Integrity Solutions	17

AUDITING AND MONITORING FOR COMPLIANCE	17
Auditing Mechanisms.....	18
Integrating with SIEM Tools	18
CONCLUSION.....	18

EXECUTIVE SUMMARY

Organizations are turning to Hadoop to implement centralized storage systems for all their enterprise data—in effect, a *data lake*. However a data lake built on Hadoop can present considerable challenges with respect to efficiency, simplicity, and security.

The Hadoop Distributed File System (HDFS) lacks some enterprise capabilities that facilitate efficient data storage. Storage architects, recognizing the limitations of HDFS, often implement alternative storage systems to hold enterprise data used for analytics. But implementing systems alongside of HDFS can end up creating an inefficient workflow, because data has to be moved into HDFS for analysis and then exported to obtain results.

In addition, because Hadoop is a distributed system designed for dual purposes (data storage and data-intensive computational analysis), it is difficult to secure. As Hadoop is often implemented as a multi-tenant service without client-server interactions, there is no single point of access where the system can be secured. The distributed nature of the system, coupled with many jobs running on many clients, makes Hadoop's native security capabilities difficult to implement and time-consuming to manage.

Ensuring compliance with a regulation or corporate data policies can require as many as 20 additional layers of security software, all of which must smoothly interoperate. Even with a layered approach to securing a native Hadoop system, however, compliance problems can linger. Connecting Hadoop to Active Directory with Apache Knox, for instance, controls access only to the system, not to directories or files, meaning personnel without a business need-to-know can access sensitive data. If a Hadoop data lake fails to implement adequate measures for information security, it increases the risk of security incidents.

The Dell EMC Isilon scale-out network-attached storage (NAS) platform delivers a multiprotocol data lake that helps secure Hadoop data with the following capabilities:

- Compliance mode
- Role-based access control for system administration
- Identity management
- Authentication
- Fine-grained access control to the file system
- Cross-protocol permissions and ACL policies
- User and ID mapping to associate one user with one ID
- Write-once, read-many storage (WORM)
- Encryption of data at rest
- Auditing of SMB events
- Auditing of RBAC administrative changes
- Integration with third-party tools to monitor security events and to encrypt data in transit

Hadoop compute clients gain access to stored data through a Hadoop Distributed File System (HDFS) interface. The distributed Isilon® OneFS® operating system implements the server-side operations of the HDFS protocol on every node in an Isilon storage cluster, and each node functions as both a namenode and a datanode.

The result is a highly efficient, scalable storage platform that safeguards the confidentiality, availability, and integrity of Hadoop data to help meet such compliance regulations as PCI DSS, FISMA, and HIPAA.

INTRODUCTION

Forward-thinking organizations are implementing technologies to store all their data in a centralized system—today frequently referred to as a *data lake*. As a central repository for disparate sources of information, a data lake enables organizations to transform their business with big data analytics. To store and analyze their data, many organizations select Hadoop. However a data lake built on Hadoop can present considerable challenges with respect to efficiency, simplicity, and security.

First, the Hadoop Distributed File System (HDFS) lacks some enterprise capabilities that facilitate efficient data storage. Storage architects recognizing the limitations of HDFS often implement alternative storage systems to hold enterprise data used for analytics.

But implementing systems alongside of HDFS can end up creating an inefficient workflow, because data has to be moved into HDFS for analysis and then exported for the results.

Second, a Hadoop data lake can become a highly complex undertaking, and with complexity, organizations see a rise in security risks and operating expenses. Ironically, most organizations address the security risks with a redundant, layered approach that further increases the system's complexity. Or, if an organization cannot adequately address security for the whole data lake, the organization will split the lake into segments, or ponds, to isolate the access to each segment by business role or need to know—in effect re-creating the information silos that the data lake was intended to eliminate.

Third, Hadoop lacks mature security capabilities to protect the growing amounts of data that it stores, either as a data hub itself or as a temporary analytics storage system. The analysis of enterprise data, even when data masking is used to obfuscate personal identifiable information, inevitably includes some degree of sensitive data.

Data that contains sensitive business information requires protection to meet internal data security policies and external compliance regulations like Sarbanes-Oxley. Data that contains sensitive information about customers, accounts, finances, health, credit cards, and so forth requires security controls to meet the compliance regulations of the organization's industry, such as the Payment Card Industry Data Security Standard (PCI DSS), the Health Insurance Portability and Accountability Act (HIPAA), and Security and Exchange Commission (SEC) Rule 17a-4.

Although regulatory requirements vary by industry, as the following table shows, the requirements to implement information security are often similar. The next section discusses these similarities.

Table 1 Regulatory requirements by industry

INDUSTRY	COMPLIANCE REGULATION
Credit card processing	Payment Card Industry Data Security Standard (PCI DSS)
Healthcare	Health Insurance Portability and Accountability Act (HIPAA)
Life sciences	Genetic Information Non-Discrimination Act (GINA)
Financial services	Sarbanes-Oxley Act (SOX), Dodd-Frank Act, Security and Exchange Commission (SEC) Rule 17a-4
Media and entertainment	The Motion Picture Association of America's security requirements for content movement
Government	Federal Information Security Management Act (FISMA) for U.S. government agencies

Information Security and Regulatory Compliance

Many compliance regulations are founded on the three points of what some auditors call the information security triangle: integrity, confidentiality, and availability. For example, the federal information processing standard known as FIPS 200 — titled Minimum Security Requirements for Federal Information and Information Systems — defines information security as "the protection of information and information systems from unauthorized access, use, disclosure, disruption, modification, or destruction in order to provide confidentiality, integrity, and availability."¹

U.S. federal information policy defines confidentiality, integrity, and availability as follows:

- Confidentiality means "means preserving authorized restrictions on access and disclosure, including means for protecting personal privacy and proprietary information"
- Integrity "means guarding against improper information modification or destruction, and includes ensuring information nonrepudiation and authenticity"
- Availability "means ensuring timely and reliable access to and use of information"²

The compliance regulations that are not founded on the information security triangle incorporate at least data integrity and confidentiality as governing principles, both of which can be seen in the following high-level requirements of PCI DSS:³

Build and Maintain a Secure Network and Systems	Install and maintain a firewall configuration to protect cardholder data Do not use vendor-supplied defaults for system passwords and other security parameters
Protect Cardholder Data	Protect stored cardholder data Encrypt transmission of cardholder data across open, public networks
Maintain a Vulnerability Management Program	Protect all systems against malware and regularly update anti-virus software or programs Develop and maintain secure systems and applications
Implement Strong Access Control Measures	Restrict access to cardholder data by business need to know Identify and authenticate access to system components Restrict physical access to cardholder data
Regularly Monitor and Test Networks	Track and monitor all access to network resources and cardholder data Regularly test security systems and processes
Maintain an Information Security Policy	12. Maintain a policy that addresses information security for all personnel

The high-level requirements of PCI DSS elucidate common technical components of storing data in an enterprise storage hub to comply with various regulatory mandates:

- One user, one ID
- Authentication
- Access control, with access limited by role and need
- Auditing and monitoring
- Retention
- Encryption

One User, One ID: The requirement to identify each user with a unique name or number plays a role in many compliance regulations. HIPAA, for instance, sets forth the following requirement in Section 164.312(a)(2)(i): "Assign a unique name and/or number for identifying and tracking user identity." This HIPAA technical safeguard maps to several interrelated security controls:⁴

¹ The FIPS 200 standard specifies the minimum security requirements that U.S. federal agencies must meet through selecting security controls from NIST Special Publication 800-53, Security and Privacy Controls for Federal Information Systems and Organizations. The selection of the security controls depends on risk and other factors. The risk-management process begins by using FIPS Publication 199, Standards for Security Categorization of Federal Information and Information Systems, to categorize information and information systems so the appropriate security controls from NIST Special Publication 800-53 can be selected. Healthcare organizations governed by HIPAA must also select the appropriate security controls from NIST SP 800-53; see NIST Special Publication 800-66, An Introductory Resource Guide for Implementing the Health Insurance Portability and Accountability Act (HIPAA) Security Rule.

² Source: 44 U.S.C., Sec 3542. See <http://www.gpo.gov/fdsys/pkg/USCODE-2011-title44/html/USCODE-2011-title44-chap35-subchapIII-sec3542.htm>.

³ The table is from Payment Card Industry Data Security Standard Requirements and Security Assessment Procedures, Version 3.0, November 2013.

⁴ NIST Special Publication 800-66, An Introductory Resource Guide for Implementing the Health Insurance Portability and Accountability Act (HIPAA) Security Rule, maps the HIPAA Technical Safeguards to the security controls in NIST 800-53.

- Account management
- Identifier management
- Identification and authentication
- Access enforcement

Account management and identifier management provide the foundation for identification, authentication, and access enforcement. Authentication and access control are the technical core of such compliance regulations as HIPAA, FISMA, and PCI DSS. In its most general form, access control seeks to limit access to personal information or business records to only those with a legitimate need to obtain it. In other cases, access control requirements can be quite specific, such as in the case of PCI DSS Requirement 7.2.3, which stipulates that the default access control setting must be deny-all.

Another cornerstone of compliance regulations is auditing and monitoring. The rapidly growing unstructured data that populates Hadoop data lakes often contains sensitive information like intellectual property, confidential customer data, and company records. Auditing can detect fraud, inappropriate entitlements, unauthorized access attempts, and other anomalies. Government agencies as well as businesses in financial services, health care, life sciences, and media and entertainment must meet regulatory requirements developed to protect against data breaches, fraud, tampering, and data loss. PCI DSS, for instance, includes imperatives that specify how to monitor, track, and audit access to data that contains credit card account information.

Retention of records is another common denominator. With such compliance regulations as Sarbanes-Oxley and SEC 17a-4, certain records must be retained under tight security conditions to guard their integrity and reliability. SEC 17a-4, for instance, stipulates that critical data must be protected for a given period of time from malicious, accidental, or premature alteration or deletion.

Several compliance regulations cite encryption of data at rest as either a requirement or as an appropriate internal security control. For example, NIST Special Publication 800-53, titled Security and Privacy Controls for Federal Information Systems and Organizations, defines a set of baseline security controls for organizations governed by FISMA or HIPAA. The document includes a security control for the protection of information at rest, and the security control's supplemental guidance says, "Organizations have the flexibility to either encrypt all information on storage devices (i.e., full disk encryption) or encrypt specific data structures (e.g., files, records, or fields)."⁵

Requirements such as these are a sampling of the compliance regulations that a Hadoop data lake must implement to mitigate security risks and adequately protect the confidentiality, integrity, and availability of data.

The Costs of Noncompliance

If an organization fails to implement adequate measures for information security, it increases the risk of a security incident, and the consequences can be costly. Breaches, in particular, are a recurring problem—more than 510 million records with sensitive information have been broken into since January 2005, according to PrivacyRights.org.⁶ The Ponemon Institute estimates the cost of a data breach for a U.S. company in 2012 at \$188 *per record*.⁷

The costs of a mega breach, such as the one that compromised the records of Target, are much higher: Trade associations estimate that the breach at Target cost financial institutions more than \$200 million.⁸ Target itself reported \$61 million of expenses related to the breach in the fourth quarter of 2013.⁹ Sony estimates that the breaches it suffered in 2011 cost about \$171 million.¹⁰

eBay is the most recent case in point. In May 2014, it announced what may turn out to be the biggest-ever cyber-attack—hackers broke into a system holding the personal data of 233 millions eBay customers. "The scope for damage is absolutely huge and could be the biggest hack of all time, given the number of users eBay has," Rik Ferguson, global vice president of security research at security software firm Trend Micro, told *The Guardian* newspaper of London.¹¹

⁵ NIST Special Publication 800-53, Revision 4, Security and Privacy Controls for Federal Information Systems and Organizations, <http://csrc.nist.gov/publications/PubsSPs.html>.

⁶ As cited by PCI SSC Quick Reference Guide.

⁷ https://www4.symantec.com/mktginfo/whitepaper/053013_GL_NA_WP_Ponemon-2013-Cost-of-a-Data-Breach-Report_daiNA_cta72382.pdf

⁸ http://www.huffingtonpost.com/2014/02/18/target-data-breach-cost_n_4810787.html

⁹ http://www.nytimes.com/2014/02/27/business/target-reports-on-fourth-quarter-earnings.html?_r=0

¹⁰ [http://www.darkreading.com/attacks-and-breaches/sony-data-breach-cleanup-to-cost-\\$171-million/d/d-id/1097898?cid=rssfeed_iwkw_all](http://www.darkreading.com/attacks-and-breaches/sony-data-breach-cleanup-to-cost-$171-million/d/d-id/1097898?cid=rssfeed_iwkw_all)

¹¹ eBay urges users to reset passwords after cyber attack, *The Guardian*, Wednesday 21 May 2014, <http://www.theguardian.com/technology/2014/may/21/eBay-urges-users-to-reset-passwords-after-cyberattack>.

The personal information was apparently unencrypted. “It is inexcusable for a company the size of eBay with the amount of data it holds to not encrypt all personal information held and to not constantly be at the forefront of security technology,” said Alan Woodward, a professor from the department of computing at the University of Surrey.¹²

When an enterprise implements a data lake with Hadoop, it increases the likelihood of a security incident and its associated costs because Hadoop has not matured enough to protect the confidentiality and integrity of data.

Compliance Problems with Hadoop

In its native implementation the HDFS file system fulfills few compliance requirements for information security. Without separating Hadoop compute clients from the Hadoop Distributed File System, Hadoop is difficult to secure because it is a complex, distributed system of many client computers with a dual purpose—data storage and data-intensive computational analysis.

With a typical Hadoop implementation, the many clients mean that there is no single point of access for the system; instead, every node or client is a point of access. The distributed nature of the system, coupled with many jobs running on many clients, makes Hadoop's native security capabilities difficult to implement and time-consuming to manage. Most native Hadoop clusters implement container-level security in which Kerberos authenticates users, clients, and services and tokens authorize jobs, tasks, and file access.

Recent approaches to improving native Hadoop security, such as Apache Knox, center on adding perimeter security by integrating Hadoop with an identity management system like Microsoft Active Directory or OpenLDAP. With each subsequent add-on security component or layer, however, complexity increases. While security risks are reduced, the overhead of managing layers of security

increases exponentially with each layer. Approximating compliance with a regulation can require as many as 20 additional layers of security software, all of which must interoperate seamlessly.

Even with a layered-approach to securing a native Hadoop system, compliance problems linger. Connecting Hadoop to Active Directory with Apache Knox, for instance, controls access only to the system, not to directories or files.

Another liability is that every Hadoop node retains a root account and password that could provide the wrong person with access to data stored in HDFS. System administrators might have access to sensitive stored data that they do not have a business need to access, and they might be able to delete, whether accidentally or intentionally, information that is supposed to be kept.

The information in a Hadoop data set can also increase complexity. A data lake at a hospital, for example, might contain sensitive health information about patients as well as their credit card information. The hospital personnel who have a business need to access the credit card information, however, might not have a business need to access the healthcare data. The traditional solution has been to create a separate data silo for each set of users and then provision access to each data silo for only those qualified to access the data.

An Dell EMC Isilon storage cluster not only provides a scalable NAS system for a Hadoop data lake but also secures Hadoop data with identity management, authentication, access control, file-system permissions, WORM, encryption for data at rest, and some auditing capabilities. Although neither an Isilon cluster nor any other system can fulfill all the requirements of a compliance regulation, an Isilon cluster can implement a range of security controls to help comply with most regulatory mandates. Third-party tools coupled with other technologies, such as encryption of data in transit, can further harden a data lake to improve security and compliance.

STORING HADOOP DATA ON ISILON SCALE-OUT NAS

Powered by OneFS operating system, the Dell EMC Isilon scale-out network-attached storage (NAS) platform delivers a scalable pool of storage with a global namespace and a native implementation of the HDFS protocol to provide an enterprise-grade, scale-out data lake. Isilon scale-out NAS is a fully distributed system that consists of nodes of modular hardware arranged in a cluster. The distributed Isilon OneFS operating system combines the memory, I/O, CPUs, and disks of the nodes into a cohesive storage unit to present a global namespace as a single file system.

The nodes work together as peers in a shared-nothing hardware architecture with no single point of failure. Every node adds capacity, performance, and resiliency to the cluster, and every node can act as a Hadoop namenode and datanode.

¹² eBay urges users to reset passwords after cyber attack, The Guardian, Wednesday 21 May 2014.

Hadoop compute clients gain access to stored data through a Hadoop Distributed File System (HDFS) interface. The OneFS namenode daemon is a distributed process that runs on all the nodes in the cluster. A compute client can connect to any node in the cluster over HDFS to access namenode services. For Hadoop analytics, the Isilon scale-out distributed architecture minimizes bottlenecks, rapidly serves big data, and optimizes performance for MapReduce jobs.

As nodes are added, the file system expands dynamically and redistributes data, eliminating the work of partitioning disks and creating volumes. The result is an efficient and resilient storage architecture that brings the security capabilities of an enterprise scale-out NAS system to storing data for analysis.

Use Cases

An Isilon cluster simplifies data management while reducing the time to gain insights from data. Although high-performance computing with Hadoop has traditionally stored data locally in compute clients' HDFS file system, the following use cases make a compelling case for coupling MapReduce with Isilon scale-out NAS:

- Store data in a POSIX-compliant file system with SMB, HTTP, FTP, and NFS workflows and then access it through HDFS for MapReduce
- Scale storage independently of compute as your data sets grow
- Protect data more reliably and efficiently instead of replicating it
- Eliminate HDFS copy operations to ingest data and Hadoop fs commands to manage data
- Implement distributed fault-tolerant namenode services
- Manage data with enterprise storage features such as deduplication, snapshots, and compliance mode
- Secure Hadoop data to help fulfill the requirements of compliance regulations

Elasticity, Simplicity, Flexibility

To handle the growth of big data, an Isilon cluster scales out dynamically, optimizes data protection, supports existing workflows with standard network protocols like SMB and NFS, and manages data intelligently with enterprise features like deduplication, automated tiering, and monitoring.

Hadoop's ratio of CPU, RAM, and disk space depends on the workload—factors that make it difficult to size a Hadoop cluster before you have had a chance to measure your MapReduce workload. Expanding data sets also makes sizing decisions upfront problematic. Isilon scale-out NAS lends itself perfectly to this scenario: Isilon scale-out NAS lets you increase CPUs, RAM, and disk space by adding nodes to dynamically match storage capacity and performance with the demands of a dynamic Hadoop workload.

An Isilon cluster fosters data analytics without ingesting data into an HDFS file system. With an Dell EMC Isilon cluster, you can store data on an enterprise storage platform with your existing workflows and standard protocols, including SMB, HTTP, FTP, REST, and NFS as well as HDFS—and secure the data by using NTFS/NFSv4 access control lists, access zones, self-encrypting drives, and other security capabilities.

Regardless of whether you store the data with SMB or NFS, however, you can analyze it with a Hadoop compute cluster through HDFS. There is no need to set up an HDFS file system and then load data into it with tedious HDFS copy commands or specialized Hadoop connectors. Combining multiprotocol data ingestion with Hadoop analytics produces a data lake that lays the foundation to transform your organization into an information-driven enterprise:

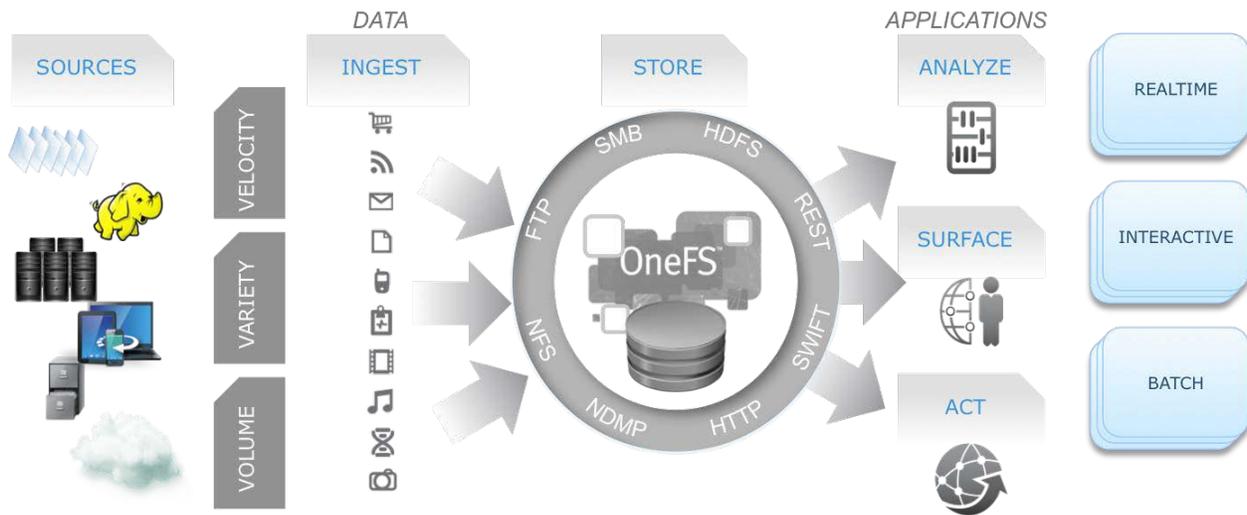


Figure 1. Combining multiprotocol data ingestion with Hadoop analytics

Separating Data from Compute

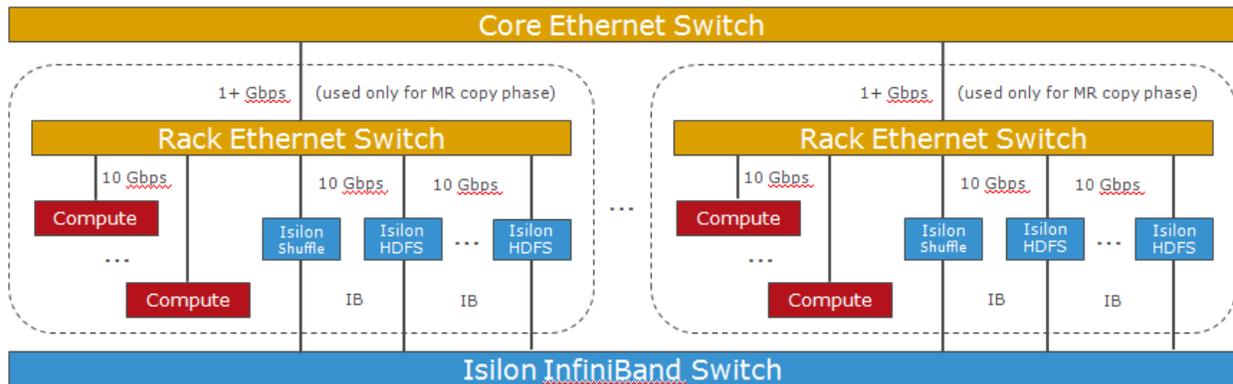
An Isilon cluster separates data from compute. As Hadoop nodes run MapReduce jobs, the machines access the data stored on an Isilon cluster over HDFS. OneFS becomes the HDFS file system for MapReduce nodes and other compute clients.

Storing data in a Isilon scale-out NAS cluster instead of Hadoop nodes streamlines the entire analytics workflow. Isilon's HDFS interface eliminates extracting the data from a storage system and loading it into an HDFS file system. Isilon's multiprotocol data access with SMB and NFS eliminates exporting the data after you analyze it. The result is that you cannot only increase the ease and flexibility with which you analyze data but also reduce capital expenditures and operating expenses. For more information, see [Dell EMC Isilon Scale-Out NAS for In-Place Hadoop Data Analytics](#).

The OneFS HDFS Implementation

OneFS implements the server-side operations of the HDFS protocol on every node, and each node functions as both a namenode and a datanode. The job tracker and task tracker functions remain the purview of Hadoop compute clients. OneFS contains no concept of a secondary namenode: Since every Isilon node functions as a namenode, the function of the secondary namenode—checkpointing the internal namenode transaction log—is unnecessary.

The cluster automatically load balances HDFS connections across all the nodes in the cluster. Because OneFS stripes Hadoop data across the cluster and protects it with parity blocks at the file level, any node can simultaneously serve datanode traffic as well as namenode requests for file blocks. Isilon provides rack-locality on its front-end network. A rack of Isilon nodes can assign compute clients to the Isilon nodes closest to a client's network switch to work with a network topology or to optimize performance. Client computers can access any node in the cluster through dual 1 GigE or dual 10 GigE network connections, as the following network diagram demonstrates.



A SmartConnect license adds additional network resilience with IP address pools that support multiple DNS zones in a subnet as well as IP failover. For more information, see [Dell EMC Isilon Best Practices for Hadoop Data Storage](#).

OneFS supports most major Hadoop distributions and projects, including Apache Hadoop, Cloudera, Hortonworks, Pivotal, HAWQ, Apache HBase, Apache Hive, Cloudera Impala, Cloudera Manager, and Apache Pig.

All the data protection, availability, integrity, and security capabilities of an Isilon cluster are available to safeguard Hadoop data. Because using an Isilon cluster to store HDFS data establishes a client-server relationship between Hadoop compute clients and the storage system, securing the data is easier and more efficient than with a Hadoop compute cluster alone: The data can be protected with security controls at the point of HDFS access and the storage system as a whole can be locked down.

The next sections describe the security solutions that an Isilon cluster can apply to Hadoop data to help meet compliance regulations and internal security policies.

SECURING HADOOP DATA

Although an Isilon cluster cannot fulfill all the requirements of a compliance regulation, an Isilon cluster implements a range of security controls to help comply with regulatory mandates. As an enterprise storage system, an Isilon cluster helps secure Hadoop data for compliance with the following capabilities:

- Compliance mode
- Role-based access control for system administration
- Identity management
- Authentication
- Fine-grained access control to the file system
- Cross-protocol permissions and ACL policies
- User and ID mapping to associate one user with one ID
- WORM
- Encryption of data at rest
- Auditing of SMB events
- Auditing of RBAC administrative changes

Third-party tools coupled with other technologies, such as encryption of data in transit, can further harden an Isilon Hadoop data lake to improve security and compliance.

Given the complexity of compliance regulations and the nuances of how they govern your systems in the context of your architecture, workflows, policies, security posture, and other factors, you should obtain an independent assessment by a third-party auditor to confirm that the implementation of a OneFS technological capability satisfies a compliance requirement.

Role-Based Access Control for Administration

OneFS includes role-based access control (RBAC) for administration. RBAC lets you manage administrative access by role. You can create separate administrator roles for security, auditing, storage, and backup. Then you can further tailor the RBAC role by assigning privileges to execute administrative commands.

By default, only the root and admin users in OneFS can log in to the command-line interface through SSH or the web administration interface through HTTPS. The OneFS root or admin user can then add other users to roles with privileges to perform administrative functions.

Assigning users to roles that contain the minimum set of necessary privileges can help fulfill such compliance regulations as PCI DSS Requirement 7.1.2: "Restrict access to privileged user IDs to least privileges necessary to perform job responsibilities."

RBAC can restrict the access and privileges of administrators so that, for example, the backup administrator is not assigned the same privileges as the overall system administrator. Restricting the access and privileges of administrators helps meet the guidance of PCI DSS Requirement 7.1.2 to insulate users without sufficient knowledge about an aspect of the storage system from making accidental changes or from modifying important security settings. Enforcing least privilege with RBAC also helps limit damage if an unauthorized person steals an administrative ID.

Administrators gain access to an Isilon cluster by using SSH for the command-line interface or HTTPS with SSL for the web administration interface. Encrypted administrative access supports PCI DSS Requirement 2.3: "Encrypt all non-console administrative access using strong cryptography. Use technologies such as SSH, VPN, or SSL/TLS for web-based management and other non-console administrative access."

In addition, all the commands executed by OneFS RBAC accounts are logged so that they can be monitored and audited; see the section on monitoring and auditing later in this paper.

Compliance Mode, WORM, and the Root Account

SmartLock is a OneFS feature that, after it is activated with a license, can lock down directories with write-once, read-many storage, commonly known as WORM. There are two ways to apply WORM:

- By putting the entire cluster into compliance mode and then applying SmartLock to specific directories
- By applying SmartLock to specific directories without placing the cluster in compliance mode

Compliance mode protects critical data from malicious, accidental, or premature alteration or deletion to help you comply with SEC 17a-4 regulations. Complying with SEC 17a-4 regulations is the intended use case for compliance mode, but it can help meet the requirements of other stringent compliance regulations.

Compliance mode imposes four constraints on the cluster:

1. Eliminates the root account and replaces it with a compliance administrator account that executes commands with *sudo*.
2. Activates a tamper-proof compliance clock to protect data in a WORM state.
3. Permanently disables privileged delete.
4. Permanently places the cluster in compliance mode: The cluster cannot be reverted to its previous state, and the root account cannot be restored.

For Hadoop analytics, the elimination of the root account might add managerial overhead and complicate system administration and troubleshooting. In compliance mode, most of the commands associated with a privilege can be performed through the *sudo* program. The system automatically generates a *sudoers* file of users from their roles. *sudo* activity is logged in a *sudo* log file, which is located in **/var/log/messages**.

In compliance mode, permissions must be managed deliberately, precisely, and consistently. Compliance mode should thus be reserved for complying with SEC 17a-4 regulations while analyzing data with Hadoop or a similar use case that requires you to eliminate the root account and monitor changes to the system's configuration. For more information, see the section later in this paper on SEC 17a-4 regulations.

Another option is to apply WORM to some directories without putting the cluster into compliance mode. You can restrict, but not eliminate, the use of the root account and manage system administrators with role-based access control and the *sudoers* file. The changes made by RBAC accounts can then be tracked to help fulfill compliance regulations for monitoring and auditing.

For FISMA and HIPAA, WORM can help put in place security controls that protect data at rest. "Integrity protection can be achieved, for example, by implementing Write-Once-Read-Many (WORM) technologies," SC-28 in NIST SP 800-53 says.

For more information on SmartLock and WORM, see [Automated Data Retention with Dell EMC Isilon SmartLock](#). For more information on auditing and monitoring, see [File System Auditing with Dell EMC Isilon, Dell EMC Common Event Enabler, and Varonis DatAdvantage](#).

Controlling Access to Hadoop Data

At the core of many compliance regulations is protecting data with strong access control measures. OneFS implements enterprise mechanisms to closely manage and strictly control access to the directories and files in the file system as well as the administrative interfaces.

In general, to securely support HDFS, NFS, and SMB clients, OneFS does three main things:

- Connects to directory services, such as Active Directory, NIS, and LDAP, which are also known as identity management systems.
- Authenticates users and groups. Authentication verifies a user's identity and triggers the creation of an access token that contains information about a user's identity.
- Controls access to directories and files at the level of the file system. OneFS compares the information in an access token with the permissions associated with a directory or a file to allow or deny access at a granular level.
- OneFS has features that can help comply with such regulations as the Federal Information Security Management Act (FISMA), the Health Insurance Portability and Accountability Act (HIPAA), Sarbanes-Oxley, SEC 17a-4, and the Payment Card Industry Data Security Standard (PCI DSS). An Isilon cluster, for instance, includes the following general capabilities:
- Identifies and authenticates users and groups by using a directory service. By integrating OneFS with a directory service, you can also use the directory service for account management.
- Provides rules to map identities from multiple external directory services to a single, unique user ID.
- Authorizes users and groups and controls access across different protocols by using POSIX mode bits, NTFS ACLs, or an optimal merging of them.
- Implements a consistent, predictable permissions model across all file-sharing protocols to preserve the intended security settings for files, directories, and other objects in the file system. The ACLs defined on OneFS are enforced when files are accessed through HDFS.
- Includes ACL policies that are, by default, set to help ensure compliance. Such policies include preserving ACEs that explicitly deny access to specific users and groups. The policies also let you tune the cluster to meet your access control objectives.
- The following diagram summarizes how directory services (which are listed across the top of the diagram in the dark gray boxes), identity mapping, policies, and permissions play a role in the OneFS system of authentication and access control.

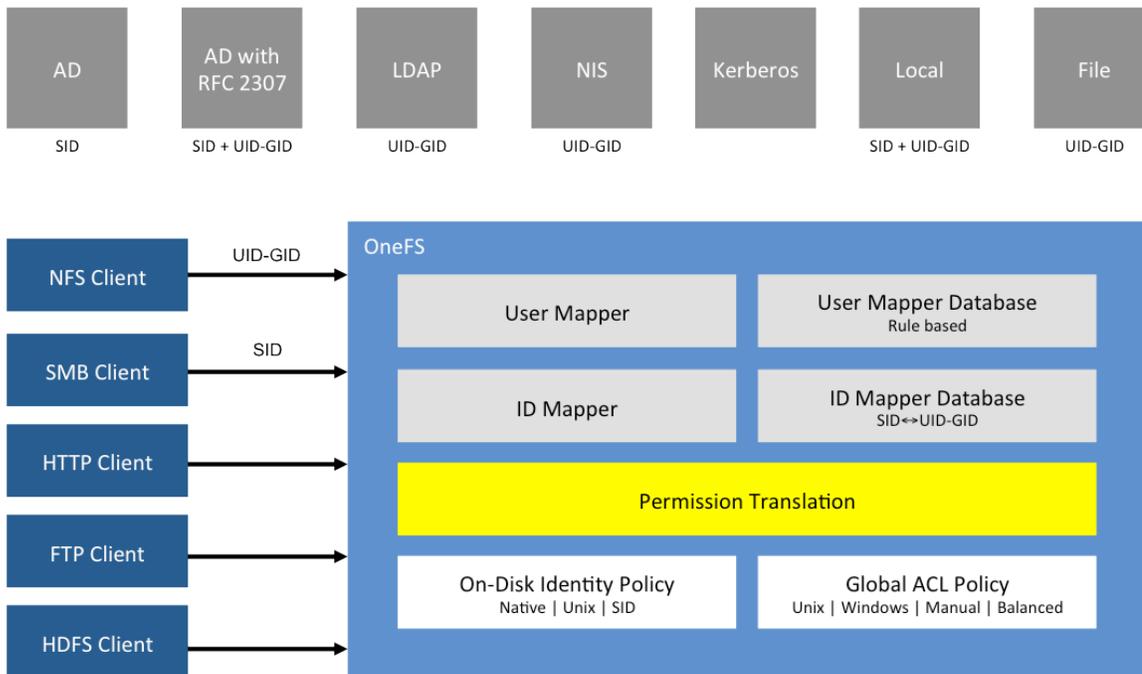


Figure 2 How directory services play a role in the OneFS system of authentication and access control

Access Zones

All of these authentication and authorization functions take place in an access zone—a virtual security context in which OneFS connects to directory services, authenticates users, and controls access to a segment of the file system. By default, a cluster has a single access zone for the entire file system. You may create additional access zones to allow users from different identity management systems, such as two untrusted Active Directory domains, to access different OneFS resources based on a destination IP address on the cluster. Access zones let you set up a cluster to work with multiple identity management systems, SMB namespaces, and HDFS namespaces.

The purpose of an access zone is to define a list of identity management systems that apply only in the context of a zone that contains SMB shares or different HDFS roots. As such, a key use case for an access zone is consolidating two or more Microsoft Windows file servers into a OneFS cluster and then analyzing the data with Hadoop. Another key use case is consolidating multiple Hadoop data sets into a single storage system but continuing to expose each data set with a unique root directory and then limiting access to only those who need to access a data set.

When a Hadoop user connects to an Isilon cluster, OneFS checks the directory services to which the user's access zone is connected for an account for the user. If OneFS finds an account that matches the user's login name, OneFS verifies the user's identity—that is, it authenticates the user. During authentication, OneFS creates an access token for the user. The token contains the user's full identity, including group memberships, and OneFS uses the token later to check access to directories and files.

When OneFS authenticates users with different directory services, OneFS maps a user's account from one directory service to the user's accounts in other directory services within an access zone—a process known as user mapping. A Windows user account managed in Active Directory, for example, is mapped by default to a corresponding UNIX account with the same name in NIS or LDAP. With a single token, a user can, if permitted, access files that were stored by a Windows computer over SMB and files that were stored by a UNIX computer over NFS or HDFS.

Similarly, to provide multiprotocol access to files with NFS, SMB, and HDFS, OneFS translates the permissions of Linux and Unix files to the access control lists of Windows files.

Identity Management

OneFS works with the following directory services to authenticate users and control access to files—functions that help satisfy compliance regulations for a unique ID for each user as well as authentication and access control:

- Active Directory. To work with UNIX and Linux systems, Active Directory includes optional support for UNIX attributes with an implementation of RFC 2307 Lightweight Directory Access Protocol (LDAP)
- Network Information Service (NIS)
- Local users and local groups
- File provider for accounts in `/etc/spwd.db` and `/etc/group` files. The file provider can add UNIX user and group account information from other systems

OneFS authenticates all the connections from any protocol with the directory service that you set up. For example, if you connect the cluster to Active Directory, the users in Active Directory can access the cluster through NFS, SMB, FTP, HTTP, and HDFS. For more information, see [OneFS Multiprotocol Security Untangled](#).

You can also use Microsoft Active Directory with Windows Services for UNIX and RFC 2307 attributes to manage Linux, UNIX, and Windows systems.¹³ Integrating UNIX and Linux systems with Active Directory centralizes identity management for Hadoop to address security controls that appear in several compliance regulations: account management, identification and authentication, and identifier management.

For example, HIPAA Section 164.312(a)(2)(i), titled Unique User Identification Implementation Specification, says that organizations should "Assign a unique name and/or number for identifying and tracking user identity." The import of PCI DSS Requirement 8.1.1 is the same: "Assign all users a unique ID before allowing them to access system components or cardholder data."

The HIPAA requirement maps to several security controls in NIST SP 800-53: For account management (AC-2), Isilon connects to Active Directory or another account management system, such as LDAP, for centralized account management. For identification and authentication (IA-2), OneFS performs identification and authentication of users by using Active Directory or another user directory, such as LDAP. For identifier management (IA-4), Isilon heeds the unique identifiers selected for individuals in an identity management system.

Active Directory can also implement *explicit deny*, and a compliance regulation that exemplifies its use relates to the International Traffic in Arms Regulations, or ITAR, which dictates that in the United States, information and material pertaining to defense and military technologies may only be shared with U.S. persons unless exempted or specially authorized by the Department of State. Organizations can face substantial legal fees and heavy fines if a foreign employee views ITAR-protected information.

If at a U.S. defense contractor, for instance, there are both foreign and domestic employees, *implicit deny* ensures that OneFS can bar a foreigner's access to sensitive data even when the foreign employee belongs to an Active Directory security group that would otherwise give the foreign employee access to the information.

Kerberos Authentication

OneFS lets you set up Kerberos authentication by using Active Directory or a stand-alone MIT Kerberos 5 key distribution center. Verifying the identity of all Hadoop users and services with the highly secure Kerberos protocol helps fulfill compliance requirements for authentication.

Setting up a stand-alone KDC to process HDFS authentication requests instead of sending them straight to Active Directory is an approach that can reduce the load on the domain controllers when a large Hadoop cluster boots up and all the services attempt to authenticate at the same time. The KDC authenticates the services while maintaining tight security. After you set up the KDC, you can establish a trust between the KDC and Active Directory so that you can centrally manage all the users, groups, and principals in Active Directory.

The Kerberos security protocol is a strong method of addressing compliance requirements for authentication, such as HIPAA's "Person or Entity Authentication Standard" in Section 164.312(d): "Implement procedures to verify that a person or entity seeking access to electronic protected health information is the one claimed."

This HIPAA requirement maps to security control number IA-2 on identification and authentication in NIST SP 800-53. For IA-2, OneFS uses an identity management system such as Active Directory or LDAP to perform identification and authentication. In the context of multiple identity management systems, OneFS can uniquely identify and authenticate users by using its built-in user identifier and identity mapping rules.

¹³ In some versions of Microsoft Windows, Windows Services for UNIX is known as Identity Management for Unix (IDMU).

For PCI DSS, using Kerberos for authentication helps satisfy Requirement 8.2.1 on the storage system for Hadoop: "Using strong cryptography, render all authentication credentials (such as passwords/phrases) unreadable during transmission and storage on all system components."

ID Mapping

OneFS includes two mapping services—ID mapping and user mapping—to combine a user's identity and identifiers from several identity management systems into a single access token with a unique identity and identifier to help meet regulatory requirements for one user with one ID.

Although their names are similar, the ID mapping service differs from the user mapping service. The goal of the ID mapping service is to map Windows SIDs to UNIX UIDs and GIDs and vice versa in order to provide consistent, secure access across two or more protocols, such as HDFS, NFS, and SMB.

During authentication, the ID mapping service associates Windows identifiers with UNIX identifiers. When a user connects to an Isilon cluster over SMB, the ID mapping service maps the user's SIDs to UIDs and GIDs for access to files that were stored over NFS. By default, the ID mapping service matches accounts with the same name.

The ID mapping service can help implement security controls, such as the FISMA security control in NIST SP 800-53 for identifier management (IA-4), to manage identifiers and to establish a single identifier for each user.

User Mapping

While the ID mapper links account identifiers like UIDs and SIDs across directory services, the user mapping service combines access tokens from different directory services into a single token.

When the names of an account in different directory services match exactly, OneFS automatically combines their access tokens into a single token. For example, the user mapping service maps, by default, a Windows account named **YORK\jane** from Active Directory to a UNIX account named **jane** from LDAP and generates an access token that combines the group membership information from the two accounts. OneFS also automatically maps two group accounts with exactly the same name.

The user mapper lets you combine and control a user's identities. After the user mapper collects a user's identities from the directory services, the user mapper can modify the token by applying rules that you create. By manipulating tokens with rules, you can address use cases common to environments with several directory services:

- Merging several identities into a single token that works for access to files stored over SMB, NFS, and HDFS. The token can include supplemental groups from both Active Directory and LDAP
- Selecting a primary group when there are competing choices from Windows and UNIX
- Managing identities when Active Directory and LDAP serve as directory services; for example, you can authenticate with Active Directory but use a UNIX identity for Hadoop file access
- The user mapping service is another component of OneFS that can help implement security controls to meet requirements such as HIPAA Section 164.312(a)(2)(i), which requires a unique name and number to identify each user identity.

To address PCI DSS Requirement 8.1—"Define and implement policies and procedures to ensure proper user identification management for non-consumer users and administrators on all system components"—OneFS works with external identity management systems like Active Directory and provides methods such as user and ID mapping to manage, at a fine-grained level, user identification. The section later in this paper on multiprotocol permissions shows how these methods interact with the OneFS ACL model to establish fine-grained access control for Hadoop directories and files.

The user mapping service and the ID mapping service can also help address PCI DSS Requirement 8.5: "Do not use group, shared, or generic IDs, passwords, or other authentication methods as follows:

- Generic user IDs are disabled or removed.
- Shared user IDs do not exist for system administration and other critical functions.
- Shared and generic user IDs are not used to administer any system components."

For more information on the user mapping service and how it merges identities, see [Identities, Access Tokens, and the OneFS User Mapping Service](#).

Securing Data at Rest with Self-Encrypting Drives

The findings of a Ponemon Institute report demonstrate the relationship between encryption and a strong security posture. The main driver for using encryption, the report says, is lessening the impact of data breaches. Improving compliance with privacy requirements is a secondary concern.¹⁴

An Dell EMC Isilon cluster with self-encrypting drives secures data at rest to address the following use cases:

- Protect data and drives against theft
- Return failed drives holding sensitive data to their vendor while safeguarding the data and reducing drive retirement costs
- Manage keys automatically with an internal key manager to remove key management as a barrier to deployment and to eliminate the overhead tied to managing keys

The Isilon self-encrypting drives are FIPS 140-2 Level 3 validated. The drives automatically apply AES-256 encryption to all the data stored in the drives without requiring equipment other than the drives.¹⁵ The drives are certified for U.S. government use, and they follow the TCG Enterprise SED standard. If a self-encrypting drive is lost or stolen, the data on the disk is rendered unreadable.¹⁶

While other encryption solutions commonly used with Hadoop can degrade performance by as much as 30 percent, Isilon's self-encrypting drives distribute the workload across every drive to produce only a nominal effect on performance: The performance of read and write operations is estimated to be less than 5 percent slower than the performance of comparable unencrypted drives. Memory usage for key management is estimated to consume less than 0.001 percent of a node's total RAM. Key management is estimated to have no impact on a cluster's total available capacity. CPU utilization for authenticating access to SEDs and for managing keys is estimated to be less than 1 percent of a cluster's CPUs.

Section 404 of Sarbanes-Oxley requires companies to assess risks to their financial reporting and to implement internal controls to mitigate those risks. Self-encrypting drives help safeguard the integrity of financial information at rest and help mitigate the risk of theft of financial records.

Protecting Data In Transit with Partner Solutions

To protect data in transit, an Dell EMC Isilon cluster works with the Vormetric Encryption Agent to encrypt data on a Microsoft Windows, Unix, or Linux client before the data is transmitted over the network to an Isilon cluster. The Vormetric Data Security Manager integrates key management, data security policies, and audit logs with a centrally managed FIPS 140-2 certified appliance. When you combine the Vormetric agent with Intel's hardware-based encryption instruction set, the effect on performance is limited. For more information on how an Isilon cluster works with Vormetric to secure data in transit, see Security Solutions for Dell EMC Isilon Scale- Out NAS or contact your Dell EMC Isilon representative.

Dell EMC Isilon also works with AFORE CloudLink SecureVSA. It provides multi-tenant software-defined storage encryption to secure cloud workloads. CloudLink SecureVSA is a virtual storage appliance that establishes an encryption layer between virtualized applications and an Dell EMC Isilon storage cluster to encrypt data on a per-application or per-tenant basis. For more information, see the Dell EMC overview of the Afore solution.

Supporting SEC Rule 17a-4

With a SmartLock license, OneFS 7.1 or later can operate in SmartLock compliance mode. Compliance mode protects critical data from malicious, accidental, or premature alteration or deletion to help you comply with SEC 17a-4 regulations.

¹⁴ 2013 Global Encryption Trends Study, Ponemon Institute, February 2014, <https://www.thales-esecurity.com/cpn/global-encryption-trends-study>.

¹⁵ For more information, see the Specification for the Advanced Encryption Standard (AES), FIPS Publication 197, at <http://csrc.nist.gov/publications/PubsFIPS.html>.

¹⁶ For more information on cryptographic standards, key strengths, and algorithms, see NIST Special Publication 800-57, Part 1, at <http://csrc.nist.gov/publications/>.

SEC Rule 17a-4(f) sets forth requirements to preserve the records of exchange members, brokers, and dealers of financial securities. Electronic records must be stored in a non-rewriteable, non-erasable format that is frequently referred to as read-many, write-once storage, commonly known as WORM.

An Isilon cluster with a SmartLock license protects Hadoop data subject to the requirements of SEC Rule 17a-4 by using compliance mode. A compliance assessment by an independent, third-party auditor found that the Isilon solution, when configured and implemented properly, fulfills the requirements of SEC Rule 17a-4(f); for more information about the third-party compliance assessment or the proper implementation and configuration of compliance mode, contact a Dell EMC Isilon representative.

Isilon supports meeting the SEC Rule 17a-4(f) requirements that are directly related to the recording, storage and retention by performing the following functions:

- Preserves the records in a non-erasable, non-rewriteable format with integrated control codes and features to retain records. The retention period can be extended for a legal hold or regulatory investigation.
- Automatically verifies the accuracy and quality of the recording process with a built-in verification process that includes creating block-level checksums to detect errors and to check integrity.
- Uniquely identifies and serializes each stored record.
- Replicates record files and associated retention metadata either locally or remotely.

Isilon compliance mode also works with RainStor's data compression, immutable data retention, and auditing. RainStor provides full data lifecycle management, including expiration, purge, and record tagging, to help address aspects of Sarbanes-Oxley, SEC Rule 17a-4, the Dodd-Frank Act, and the Communications EU Data Protection regulations; for more information, see [Dell EMC Isilon Scale-Out NAS and RainStor Hadoop Solution](#).

THE ONEFS UNIFIED PERMISSIONS MODEL

A standard Hadoop implementation provides only basic Unix-type permissions. Each file or directory is assigned an owner and a group; read-write permissions can be assigned to the owner, the group, and everyone else. Security and compliance problems arise, however, when for a file or a directory you need to assign different combinations of read and write access to different groups. Such problems are compounded because a standard Hadoop implementation does not maintain the ACLs of Microsoft Windows files when they are copied over from Windows shares.

In contrast, OneFS controls access to directories and files with POSIX mode bits and NTFS access control lists. To foster multiprotocol data access, OneFS maps the POSIX mode bits of a file from a Linux or Unix system to the permissions model of the Windows system, and vice versa. The result is that the permissions on directories and files remain intact for HDFS users and applications. The ACLs defined on OneFS are enforced when files are accessed through HDFS.

The OneFS permissions model helps satisfy compliance regulations for storing unstructured data by maintaining the intended security levels of directories and files across such protocols as NFS, SMB, and HDFS. An Isilon cluster includes the following capabilities to manage permissions:

- Authorizes users and groups and controls access across different protocols by using POSIX mode bits, NTFS ACLs, or an optimal merging of them.
- Implements a consistent, predictable permissions model across file-sharing protocols to preserve the intended security settings for files and directories in the file system.
- Includes ACL policies that are, by default, set to help ensure compliance from the start. Such policies include preserving ACEs that explicitly deny access to specific users and groups. The policies can also manage how permissions are initially set or modified to meet your access control objectives.
- For a discussion of how OneFS maps permissions between the security models of Unix and Windows systems, see [Dell EMC Isilon Multiprotocol Data Access with a Unified Security Model](#).

ACL Policies for Mixed Environments

An Isilon cluster includes ACL policies that control how permissions are processed and managed. By default, the cluster is set to merge the new permissions from a *chmod* command with the file's ACL. Merging permissions is a powerful method of preserving intended security settings while meeting the expectations of users. In addition, managing ACL policies manually gives you the following options to address compliance requirements in environments that mix NFS, SMB, and HDFS:

- ACL creation over SMB
- *Chmod* on files with ACLs
- Inheritance of ACLs created on directories by the *chmod* command from a Unix client
- *Chown* and *chgrp* on files with ACLs
- Who is allowed to run the *chmod* and *chown* commands
- Treatment of rwx permissions
- Group owner inheritance
- Removing ACLs from a UNIX client
- Owner permissions
- Group permissions
- Deny ACEs
- Changing interpretation of *utimes*
- Read-only DOS attribute
- The display of mode bits

For a description of the policies and a discussion of their usage, see [Dell EMC Isilon Multiprotocol Data Access with a Unified Security Model](#).

The OneFS Permissions Model and Compliance

The OneFS unified permissions model, identity mapping, and ACL policies address several compliance requirements from FISMA, HIPAA, and PCI DSS related to identification and access limitations.

Among the categories of security controls in SP 800-53 is access control. Within the access control category, a prominent control is access enforcement. "The information system," the document says, "enforces approved authorizations for logical access to the system in accordance with applicable policy."

The control includes supplemental guidance and control enhancements that detail how access enforcement mechanisms like ACLs are employed when necessary to control access between users and objects in the target information system. The supplemental guidance is to be applied as appropriate to implement security controls. The control enhancements, which add functionality to a control or increase the strength of a control, are to be applied when an information system requires greater protection to address a heightened potential impact of loss or to address the results of a risk assessment.

The unified security model of OneFS helps you conform to some of the supplemental guidance and control enhancements of the access enforcement security control. Most of the supplemental guidance and many of the control enhancements can be implemented by using systems such as Active Directory or LDAP to authenticate users and groups and authorize access to resources. OneFS works with both Active Directory and LDAP to help implement the access enforcement security control and much of its supplemental guidance for Hadoop assets stored on an Isilon cluster.

The following control enhancements directly apply to how a multiprotocol file system secures directories and files across dissimilar permissions models:

"The information system enforces a Discretionary Access Control (DAC) policy that: (a) Allows users to specify and control sharing by named individuals or groups of individuals, or by both; (b) Limits propagation of access rights; and (c) Includes or excludes access to the granularity of a single user."

With its default handling of changes to ACLs and its default ACL policies, the unified security model of OneFS helps you conform to this control enhancement without additional configuration, complexity, or managerial overhead.

First, because OneFS preserves ACLs across protocol boundaries, it allows users to specify and control sharing by naming individuals and groups of individuals in ACEs. The entries are maintained regardless of whether files are accessed over NFS, SMB, or HDFS. Even when OneFS creates ACLs for files from a Unix system for Windows users, the permissions of named individuals and groups are preserved unless you select an ACL policy that overrides them.

Second, you can limit propagation of access rights by using the policy that manages ACLs created on directories by the Unix `chmod` command. To limit propagation, make sure the policy is set to not make them inheritable.

Third, the ACEs used by OneFS can include or exclude access down to the granularity of a single user.

The default configuration of OneFS simply and effectively delivers access enforcement across protocols to help you enforce approved authorizations to stored electronic resources, including not only the supplemental guidance but also the more rigorous enhancement for a discretionary access control policy.

The granular permissions model of OneFS also helps implement aspects of the PCI DSS requirements for access control. PCI DSS requirement 7.2.3, for example, is that the default of the access control setting is to deny-all, which OneFS does by default.

Permissions at the level of the file system help the enterprise data hubs of healthcare organizations deal with the HIPAA workplace security standard (HIPAA Section 164.308(a)(3)(i)):

"Implement policies and procedures to ensure that all members of its workforce have appropriate access to electronic protected health information, as provided under paragraph (a)(4) of this section, and to prevent those workforce members who do not have access under paragraph (a)(4) of this section from obtaining access to electronic protected health information."

HIPAA Section 164.308(a)(3)(i) maps to the several relevant security controls in NIST SP 800-53: AC-1, AC-5, AC-6. Isilon supports AC-5, the separation of duties, by honoring the group memberships specified in Active Directory or another directory service. As such, an Isilon cluster can help implement separation of duties for users based on the settings in the directory service. For AC-6, least privilege, Isilon supports the access rights that you set for your Hadoop data.

But implicit in the first part of HIPAA Section 164.308(a)(3)(i) is the notion of availability. Availability means that, in the context of an enterprise data hub, the data is available to those who need it when they need it. The next section discusses how an Isilon cluster addresses availability and integrity.

AVAILABILITY AND INTEGRITY

Compliance regulations like FISMA and HIPAA combine the three points of the information security triangle—availability, confidentiality, and integrity—to form the basis for information security. To address the availability and integrity aspects of compliance for Hadoop data, Isilon includes many features, but a key component upon which the features rest is the scale-out, distributed architecture of an Isilon cluster.

Isilon Architecture for Availability

Scale-out NAS systems are different from traditional scale-up systems. The architecture of a Dell EMC Isilon scale-out NAS system contains no single master for the data and no concept of a high-availability (HA) pair. Instead, Isilon scale-out NAS is a fully distributed system that consists of nodes of modular hardware arranged in a cluster. The distributed Isilon OneFS operating system combines the memory, I/O, CPUs, and disks of the nodes into a cohesive storage unit to present a global namespace as a single file system. The nodes work together as peers in a shared-nothing hardware architecture with no single point of failure.

The result is a highly resilient storage architecture. The OneFS operating systems handles a failure by distributing the load of a failed node to the remaining nodes. The system keeps just enough redundant information to reconstruct the data on a node or disk that fails, and the amount of overhead needed to protect against failure decreases as nodes are added to the cluster.

Compared with traditional scale-up NAS systems, a scale-out architecture provides a more resilient foundation for data protection and data availability. In its 2013 report titled “Critical Capabilities for Scale-Out File System Storage,” Gartner rated Dell EMC Isilon highest among storage vendors for resiliency—the platform’s capabilities for provisioning a high level of system availability and uptime.¹⁷

The design of Isilon’s clustered architecture supports the following availability objectives:

- No single point of failure
- Tolerance for multi-failure scenarios
- Fully distributed single file system
- Pro-active failure detection and preemptive drive rebuilds
- Fast drive rebuilds
- Fully journaled file system
- Flexible, efficient data protection

For efficiency and reliability, OneFS stripes data to guard it with parity blocks at the file level instead of parity disks. At the same time, OneFS protects data with forward error correction, or FEC—a highly efficient method of reliably protecting data. FEC encodes a file’s data in a distributed set of symbols, adding space-efficient redundancy. With only a part of the symbol set, OneFS can recover the original file data. In a cluster with five or more nodes, forward error correction delivers as much as 80 percent efficiency. On larger clusters, FEC provides as much as four levels of redundancy.

OneFS Availability and Integrity Solutions

In addition to data protection with striping and forward error correction, OneFS includes the following software modules and features to help protect the integrity, availability, and confidentiality of data. Here is a partial listing of solutions that help protect Hadoop data; some of these modules require a separate license.

ANTIVIRUS. OneFS interoperates with Internet Content Adaptation Protocol (ICAP) servers to scan for viruses and other threats. This antivirus capability can help fulfill PCI DSS Requirement 5.1: “Deploy anti-virus software on all systems commonly affected by malicious software (particularly personal computers and servers).”

INTEGRITYSCAN. An IntegrityScan job examines the file system for inconsistencies by systematically reading every block and verifying its associated checksum. Unlike traditional ‘*fsck*’ style file system integrity checking tools, IntegrityScan runs while the cluster operates, eliminating the need for downtime. If IntegrityScan detects a checksum mismatch, OneFS generates an alert and automatically attempts to repair the block.

ISILON DATA INTEGRITY. Isilon Data Integrity (IDI) protects file system structures against corruption with 32-bit CRC checksums. All Isilon blocks use checksum verification. Metadata checksums are housed in the metadata blocks themselves, whereas file data checksums are stored as metadata, thereby providing referential integrity. All checksums are recomputed by the initiator, the node servicing a particular read, on every request. If a recomputed checksum does not match the stored checksum, OneFS generates a system alert, logs the event, attempts to repair the block, and returns the repaired block to the client if the block was successfully repaired.

PROTOCOL CHECKSUMS. In addition to blocks and metadata, OneFS also provides checksum verification for Remote Block Management (RBM) protocol data. RBM is a unicast, RPC-based protocol developed by Isilon for use over the internal InfiniBand network. Checksums on the RBM protocol are in addition to the InfiniBand hardware checksums provided at the network layer. The RBM checksums detect and isolate machines that contain faulty hardware components.

AUTOMATED CLUSTER REPLICATION AND FAILOVER. SyncIQ replicates data on another Isilon cluster and automates failover and failback operations between clusters. If a cluster becomes unusable, SyncIQ fails over to another Isilon cluster.

SNAPSHOTS. SnapshotIQ protects data with a snapshot—a logical copy of data stored on a cluster. A snapshot can be restored to its top-level directory. SnapshotIQ provides features to meet recovery point objectives.

¹⁷ "Critical Capabilities for Scale-Out File System Storage," Gartner, Inc., published Jan. 24, 2013, <http://www.gartner.com/technology/reprints.do?id=1-1DYP0VR&ct=130206&st=sb>.

NDMP BACKUP AND RESTORE. OneFS can back up data to tape and other devices through the Network Data Management Protocol.

ACCELERATED DRIVE REBUILDS. The time it takes a storage system to rebuild data from a failed disk drive is crucial to the data reliability of the system. With the advent of four terabyte drives and the creation of increasingly larger single volumes and file systems, typical recovery times for multi-terabyte drive failures can extend to a week or more. During this period, storage systems are vulnerable to additional drive failures and the data loss and downtime that can result.

Since OneFS is built upon a distributed architecture, OneFS can exploit the CPUs, memory, and spindles from many nodes to efficiently reconstruct data from failed drives in a parallel process. Because Isilon is not bound by the speed of any one drive, OneFS can recover from drive failures extremely quickly—an efficiency that increases with cluster size.

Isilon availability software and options contribute to ensuring that an enterprise can meet its recovery time objective (RTO), which is the allotted amount of time within a service level agreement to recover and restore data. For complete information about the data availability features of OneFS, see [High Availability and Data Protection with Dell EMC Isilon Scale-Out NAS](#).

AUDITING AND MONITORING FOR COMPLIANCE

The rapidly growing unstructured data that populates Hadoop data lakes often contains sensitive information like intellectual property, confidential customer data, and company records. Auditing can detect fraud, inappropriate entitlements, unauthorized access attempts, and other anomalies. Government agencies as well as businesses in financial services, health care, life sciences, and media and entertainment must meet regulatory requirements developed to protect against data breaches, fraud, tampering, and data loss.

With some compliance regulations, auditing file system operations like file creation or deletion is required to demonstrate compliance. In other scenarios, the goal of auditing is to track administrative changes. Another requirement is to track activities like logon events.

Auditing Mechanisms

OneFS provides several auditing mechanisms to ensure the availability, integrity, and confidentiality of the cluster and the data it stores:

- Support for SNMP versions 1, 2c, and 3 to remotely monitor hardware components, CPU usage, switches, and network interfaces for integrity
- A virtual appliance, called InsightIQ, to monitor and analyze the performance of an Isilon cluster to forecast capacity and maintain availability
- A RESTful application programming interface to automate monitoring and retrieve statistics
- Auditing of system configuration events to track changes by administrators
- SMB protocol monitoring to track user access and record file events such as opening files, deleting directories, viewing security settings, and modifying permissions

Integrating with SIEM Tools

The SMB event monitoring and auditing integrates with Varonis DatAdvantage, Symantec Data Insight, and other security information and event monitoring tools (SIEM). On OneFS, the events are logged on the node that an SMB client connects to and then stored in a file in **/ifs/ifsvar/audit/logs**. The logs automatically roll over to a new file once the size reaches 1 GB. The default data protection scheme for the audit log files is +3. To help meet regulatory requirements that require two years of audit logs, the audit log files are not deleted.

After an event is logged, a forwarding service sends the event to the Dell EMC Common Event Enabler with an HTTP PUT operation. The Dell EMC Common Event Enabler then forwards the event to an endpoint, such as Varonis DatAdvantage. The Varonis application coalesces the events to generate reports that contain the following information:

- An access summary that displays a log of daily events
- A sensitive access summary that displays a log of attempts to access files
- Directory access statistics
- User access statistics
- Tactical access statistics

For more information on auditing and monitoring, see *File System Auditing with Dell EMC Isilon*, *Dell EMC Common Event Enabler*, and *Varonis DatAdvantage*.

CONCLUSION

Powered by the distributed OneFS operating system, the Dell EMC Isilon scale-out network-attached storage (NAS) platform delivers a scalable, multiprotocol data lake to help secure Hadoop data with the following capabilities:

- Compliance mode
- Role-based access control for system administration
- Identity management
- Authentication
- Fine-grained access control to the file system
- Cross-protocol permissions and ACL policies
- User and ID mapping to associate one user with one ID
- WORM
- Encryption of data at rest
- Auditing of SMB events
- Auditing of RBAC administrative changes
- Integration with third-party tools to monitor security events and to encrypt data in transit
- Combining these capabilities with Isilon's high-availability solutions protects the integrity, confidentiality, and availability of Hadoop data to improve information security and compliance.