

# Evaluation Report

## Data Domain Deduplication Storage Systems

### I. Executive Summary

#### Summary of Evaluation

It is the opinion of Kahn Consulting, Inc. that Data Domain deduplication storage systems provide secure and trustworthy environments for the long-term storage of electronic business records. The higher speeds and lower storage requirements enabled by Data Domain deduplication technology increase the ability for organizations to store data to meet business continuity and archiving needs. Data Domain systems also contain many controls designed to maintain the integrity and security of the data stored in the system. Data Domain systems can promote organizational efforts to conform to information management and compliance requirements.

*Not a legal opinion or legal advice. For all questions regarding compliance with specific laws and regulations seek legal counsel.*

WHERE LAW & TECHNOLOGY MEET



January 2009

## II. Evaluation Overview

Kahn Consulting, Inc. (“KCI”) was engaged by Data Domain, Inc. (“Data Domain”) to evaluate their inline deduplication storage systems.<sup>1</sup> The primary focus of this evaluation is the Data Domain capabilities that address the integrity, accessibility, security, and privacy of information managed by the system. In conducting this evaluation, KCI has assessed Data Domain deduplication storage systems using criteria derived from broad compliance requirements and best practices related to information management.

The explosion of information facilitated by the rapid technological advances in computing has created concerns about maintaining the integrity of this information. As data volumes increase, the standard model of tape backup is increasingly constrained by the limitations of the technology and business needs beyond backup. At the same time, events such as 9/11 and Hurricane Katrina have dramatically highlighted the need for effective, efficient and rapid disaster recovery capabilities.

Quick recovery is meaningless, however, if the enterprise cannot be assured of the integrity of the data—that the terabytes of data backed up will be retrieved in the same form in which the data was written to the storage device. This essential requirement is mandated by both business needs and legal and regulatory directives. For example, in litigation, the party seeking to introduce electronic evidence must be able to prove that it is authentic – that it is what it claims to be.

Furthermore, data integrity must be maintained over the long term. Various regulations require the retention of data over multi-year periods. For example 23 CFR § 203.60 requires drug companies to keep records relating to the distribution of drug samples for three years. In addition, once an organization becomes aware of the possibility of a lawsuit, it must preserve evidence relevant to that lawsuit until the lawsuit is concluded, which could be many years. Business needs may also dictate the retention of information for long periods. On the other hand, once information is no longer needed by the enterprise, it must be disposed of completely and securely. Government standards such as DoD 5220.22-M highlight the need to ensure that once data is gone, it is truly gone.

Business continuity and disaster recovery, archiving for records retention, and preservation of content for Legal Holds are all essential components of an overall information management strategy. The three areas, however, serve different business purposes. Business continuity and disaster recovery are aimed at restoring business operations as soon as possible in the unlikely event of system-wide failure. The purpose of archiving is to retain records in a trustworthy and accessible manner for the period of time taking into consideration business needs and legal requirements. Preservation of content for Legal Holds is focused on ensuring that all information responsive to a legal or regulatory matter is preserved and managed in a manner that protects it from deletion or material alteration.

Management of information throughout its entire life cycle has become critical in today’s information-centric society. As a result, the tools which help organizations manage their information assets are today as critical as the information itself.

## III. About Data Domain

Data Domain has developed a line of storage systems which retain data at an intermediate level (also known as “nearline” storage), between primary, online storage, and long-term storage

WHERE LAW & TECHNOLOGY MEET

**KAHN**  
CONSULTING INC.

on media such as tape. This tier or storage is traditionally used for backup and archive data. Data Domain deduplication storage systems perform inline data deduplication at high speeds.. Deduplication occurs in random access memory, so that only unique data is written to storage media. This same process facilitates the transfer of large volumes of data over low-bandwidth wide-area network (WAN) connections to increase remote disaster recovery capabilities. The material differences among Data Domain systems are storage capacities and performance. Consequently, the capabilities evaluated in this report are common to all Data Domain storage systems.

Data Domain systems are designed to support backup, archiving and disaster recovery. The deduplication of data provides greater storage and network efficiency and lowers storage costs whether the business purpose of storing information is to support business continuity or to support the long-term archiving of content for records retention purposes.

## Data Domain Architecture

**Inline deduplication.** Fast, inline deduplication is accomplished using the Data Domain SISL™ (Stream-Informed Segment Layout) scalability architecture. With SISL, incoming data is broken down into small data segments in RAM memory. A unique fingerprint is created for each segment, which allows the segment to be compared against others. The system seeks to store unique segments, along with reference information indicating where the segment appears in the data stream, on disk. The system also applies standard compression to this information before it is written to disk. According to Data Domain, the comparison process is able to identify 99% of duplicate segments in RAM, so that nearly all data written to disk is unique. By designing the system so that the majority of the deduplication processing occurs in RAM, the overall system functionality is enhanced. Data Domain estimates that, on average, data storage requirements are 5% of what they would be without deduplication for backup data sets, and 20% of what they would be for non-backup data. Data Domain states that the deduplication process can be applied to any type of data that is written to the system. Deduplication rates will vary depending upon the nature of the data.

**Data protection.** Data Domain's Data Invulnerability Architecture is designed to protect system data using a number of techniques, including:

- 1) **Data verification.** As data is written to the system, the Data Domain operating system will compute checksums, or unique values calculated for purposes of determining the accuracy and completeness of the data. Once the deduplicated data has been written to disk, checksums are calculated again by reading the data back from disk, and the two groups of checksums are compared. If problems are detected which cannot be corrected, the archiving process is repeated.
- 2) **Fault avoidance and containment.** A Data Domain system is designed to write only to new areas of the disk which reduces the likelihood that data is incorrectly overwritten. According to Data Domain, this technique also results in less overhead for the system, since it is not required to keep track of data eligible to be overwritten, which in turn reduces the chances of software errors which could lead to data corruption. Before actually writing data to the hard drives, it is written into a non-volatile RAM buffer, which is a type of random access memory not affected by power loss. In the event of a system failure and after a restart, the file system will verify the integrity of the data in this buffer before it is written to the disks, which reduces the likelihood that data will be lost because of the failure.

WHERE LAW & TECHNOLOGY MEET

**KAHN**  
CONSULTING INC.

In addition, file system checks are performed during initial verification and after each backup. Because the system only writes to new areas of disk, the time to perform file system checks are significantly reduced, as it is not required to rebuild the overhead data structures associated with keeping track of data which can be overwritten.

- 3) **Continuous fault detection and repair.** Data Domain systems use RAID technology, or Redundant Array of Independent Disks. RAID systems distribute data redundantly over a series of disks using a technique called striping, so that the loss of an individual disk will not result in a loss of data. Data Domain utilizes a level of RAID designated as RAID-6, which is designed to protect against the loss of two disks at the same time. Lower levels of RAID are designed to protect against the loss of one disk at a time. The Data Domain system is also designed to verify the integrity of the RAID stripes after the stripes have been written to the disk system. In the checksum process described above, if faults are detected, the system will request that RAID-6 utilize its redundancy to correct the data error. This checksum verification process is proprietary to Data Domain. The operating system also re-verifies the integrity of all data in a scheduled background process to protect against defects on the disk which may creep in over time.
- 4) **File System Recoverability.** All data is stored with metadata (information which describes the underlying content). If the metadata becomes corrupted, it can be rebuilt by scanning the data on disk. Data Domain utilizes snapshots for file recovery, which are designed to minimize both performance and storage requirements. Data snapshots are taken of the same pool of deduplicated, compressed sequences used by live data. The normal data deduplication process is then applied to any new data or updates to the data sets, and therefore only small increments of new or changed data are stored once the snapshot is taken.

**Retention Lock.** Data Domain Retention Lock software gives users the option to securely “lock down” data so that it cannot be easily changed. Once the user establishes a retention date, the file cannot be altered or deleted until the date is reached. The retention lock for a specific file can be removed through a separate administrative interface, but the file cannot be deleted through that interface. Once retention rules controlling the information have been satisfied, the file can be deleted or changed. Retention parameters can be set on a file-by-file basis, while minimum and maximum retention periods can be set globally.

**Replication utilizing deduplication.** Data Domain Replicator software allows deduplicated data to be passed over a wide area network. According to Data Domain, its approach to replication allows network traffic to be reduced by up to 99%, thus facilitating the movement of large amounts of backup and archive data to secure offsite locations.

## IV. Data Domain Capabilities

This part of the Evaluation is divided into sections that describe the information integrity, information protection, retention, destruction, preservation, security, and disaster recovery capabilities that are desired in storage systems, explain why each capability is desired, and evaluate Data Domain’s compliance with each capability.

### Information Trustworthiness

WHERE LAW & TECHNOLOGY MEET



**Desired Capability.** Data should be stored in manner that ensures that it is trustworthy, and that it remains trustworthy over the long term, even for a period of years, until the information is no longer required.

**Information Management Principle.** Information is trustworthy if it can be demonstrated that it has not been altered and remains accurate since it was created or archived. Business best practices and many laws and regulations require digital information to have integrity. In order to be admissible in court, the party seeking to introduce electronically stored information must be able to demonstrate that information is authentic. The court in *Lorraine v. Markel American Ins. Co.*, 241 F.R.D. 534, 546 (D. Md. 2007) stated that “[o]ne method of authenticating electronic evidence...is the use of “hash values” or “hash marks” when making documents.” The Lorraine opinion uses “hash values” as applied to entire files. The Data Domain system applies hashes at the segment level. Because the system combines segments to form the files when read by the end user, a hash value applied to a file would be the same both before being read into the system and after restoration.

**Data Domain Capability.** The Data Domain system utilizes checksums to ensure the integrity of the data segments stored and the associated metadata when the data is written to disk. The system also periodically rechecks data integrity via background processes.

## Information Protection

**Desired Capability.** The storage system should protect information from the accidental or deliberate deletion, alteration or overwriting of that information.

**Information Management Principle.** Information management principles extending to information integrity from a system standpoint also apply to the maintenance of integrity from human errors or any malicious activity.

**Data Domain Capability.** The Retention Lock feature is designed to prevent data stored on Data Domain deduplication storage from being erased or changed while the retention period is in effect. Removing the lock on a given file to allow deletion is a multiple step process. The administrator must remove the lock using a separate administrative interface which does not in and of itself allow deletion. Once the retention date is reached, the file can be deleted. The user cannot make any changes to a file once a retention period has been established for that file. This functionality protects against inadvertent alteration or deletion of data. Deliberate deletion of data must take place with the cooperation of the administrator.

## Business Records Retention

**Desired Capability.** The storage system should support the management and retention of business records through the coding of records, assignment of retention periods, and protection of records from alteration or destruction through the retention period.

**Information Management Principle.** Many laws and regulations require that information be maintained for designated time periods. Best practices, as well as internal management policies, also require retention of information for certain periods. The information must be protected from any changes, internally or externally, over the established time period. These periods should be set at the record level.

WHERE LAW & TECHNOLOGY MEET

**KAHN**  
CONSULTING INC.

**Data Domain Capability.** Enforcement of retention periods is enabled by Retention Lock software. The Data Domain system allows retention periods to be set at the file level. Minimum and maximum retention periods may be set globally. Once the retention period is established, the system is designed to prevent the files from being overwritten or erased. Several file archiving applications can be leveraged to set the retention periods for archived files automatically. Once the retention date is reached, the file can be deleted. There is a separate administrative interface through which the retention lock on a file can be reset. The file cannot be deleted or altered through this interface. The user can delete or alter the file after its retention lock has been reset. The retention lock reset by the administrator is logged by the system.<sup>2</sup> This two step process for deletion or alteration of files prior to the expiration of the retention period is designed to guard against accidental deletion or alteration of data.

## Destruction of Information

**Desired Capability.** The storage system should ensure that unwanted digital information is destroyed in a manner which ensures that it cannot be recovered.

**Information Management Principle.** When a file is deleted by a user, a computer system normally will only mark the file as deleted, which allows the file to be overwritten by the system. A significant period of time may elapse before the disk space containing the file is actually physically overwritten by the system. Therefore, additional functionality is often required to ensure permanent physical deletion of the file, to prevent the file from being recovered. In addition, the requirement to properly destroy certain types of private information is a requirement of existing and emerging privacy laws and regulations in the United States and abroad, including the Federal Trade Commission rules regarding the proper disposal of consumer information.<sup>3</sup>

**Data Domain Capability.** As individual data segments may be shared by more than one file, an individual segment will not be deleted until it is no longer used by any other file. The System Sanitization feature (included with Retention Lock software) is designed to overwrite segments in accordance with Department of Defense standard 5220.22-M for permanently deleting digital information, once those segments have been identified as no longer used by any file.

## Preservation of Information for Legal Holds

**Desired Capability.** The onset of an investigation, audit or lawsuit triggers the obligation to preserve information for the duration of such investigation, audit, or lawsuit. Any retention period applicable to this information must be suspended, and the storage system must ensure that the information is not deleted or otherwise disposed of.

**Information Management Principle.** The Federal Rules of Civil Procedure (and many corresponding state rules of civil procedure) provide a “safe harbor,” pursuant to which a party will not be subject to sanctions for the destruction of data due to the routine, good faith operation of an information system. However, the Advisory Committee notes to the Federal Rules define “good faith” to include a party’s obligation to intervene to ensure that any processes leading to destruction of data are suspended if the party is under an obligation to preserve information. Similarly, regulatory authorities will penalize a party if information relevant to an audit or investigation is destroyed when the party is under an obligation to preserve the information.

**Data Domain Capability.** Retention Lock software allows the retention period for information at the file level to be set for a virtually indefinite period. For the purposes of ensuring that information is preserved for a lawsuit, audit or investigation, Retention Lock can be used to

WHERE LAW & TECHNOLOGY MEET

**KAHN**  
CONSULTING INC.

prevent the file from being deleted or altered during that preservation period. Once all Legal Holds on a file are lifted, the administrator can then reset the retention lock to allow the file to be deleted if the retention period has otherwise expired.

## System Security

**Desired Capability.** The storage system must contain controls preventing unauthorized access to information.

**Information Management Principle.** An enterprise's information assets have value and must be protected. In addition, an organization can be exposed to liability if sensitive information it maintains on other parties is lost or exposed.

**Data Domain Capability.** Data Domain systems utilize the full range of UNIX permission attributes. In addition, standard Microsoft Windows permission controls are utilized in Windows environments. Retention Lock software also protects against alteration or deletion of files during their retention period. Data Domain replication capabilities also provide an important safeguard for data. Since data is replicated from site to site over secure networks, the risk of shipping tape media between sites and the cost of expensive tape encryption techniques are avoided.

## Disaster Recovery

**Desired Capability.** Storage systems should have the ability to replicate system data to a remote site to permit restoration of system functions in the event of a disaster.

**Information Management Principle.** Best practices require that mission critical information should be stored in at least two physically separate locations a sufficient distance from one another. These capabilities are often required by regulation. For example, organizations which possess protected health information under the Health Insurance Portability and Accountability Act (HIPAA) must “[e]stablish and implement procedures to create and maintain retrievable exact copies of electronic protected health information....Establish (and implement as needed) procedures to restore any loss of data....[and] Establish (and implement as needed) procedures to enable continuation of critical business processes for protection of the security of electronic protected health information while operating in emergency mode.”<sup>24</sup>

**Data Domain Capability.** Data Domain systems offer several features which protect against data loss:

- 1) **WAN based replication.** When Data Domain appliances are located at both the primary and the remote site, Data Domain Replicator software allows deduplicated data to be sent over the data line, resulting in a stated 99% reduction in network traffic. Accordingly, network-based replication of backup and archive data is feasible, permitting both backup and restoration of data over a data line. The alternative to a network-based backup is physically sending removable storage media off site, with associated risks involved in physical transfers of data (i.e. misplacement of media, possibility of transfer vehicles being involved in accidents, lost or stolen tape, etc.).

WHERE LAW & TECHNOLOGY MEET

**KAHN**  
CONSULTING INC.

- 2) **Internal system redundancy.** Utilization of RAID-6 technology protects against two disk failures. In other RAID configurations, any other simultaneous disk error will cause data loss.

## Application Protection

**Desired Capability.** Many applications, such as Microsoft Exchange or Oracle database applications, recommend full backups for optimal system recovery performance. For example, Microsoft recommends that Exchange 2003 environments be backed up daily.

**Information Management Principle.** In many instances, full backups of mission-critical applications can permit an organization to resume operations quickly in the event of a disaster and remain in compliance with applicable disaster recovery provisions of laws such as HIPAA.

**Data Domain Capability.** According to Data Domain, its deduplication technology can be applied to database backups. A daily full database backup will have a large amount of redundancy when compared to a full backup made the previous day. Thus, the significant redundancies found in daily full database backups can result in up to between 20x and 50x average reduction in storage space. The high backup speeds of Data Domain systems can permit more frequent full backups of large databases in narrower operational downtime windows than might otherwise be possible.

## V. About Kahn Consulting

Kahn Consulting, Inc. (KCI) is a consulting firm specializing in the legal, compliance, and policy issues of information technology and information lifecycle management. Through a range of services including information and records management program development; electronic records and email policy development; Information Management Compliance audits; product assessments; legal and compliance research; and education and training, KCI helps its clients address today's critical issues in an ever-changing regulatory and technological environment. Based in Chicago, KCI provides its services to Fortune 500 companies and government agencies in North America and around the world. Kahn has advised a wide range of clients, including International Paper, Dole Foods, Sun Life Financial, Time Warner Cable, Kodak, McDonalds Corp., Hewlett-Packard, United Health Group, the Federal Reserve Banks, Ameritech/SBC Communications, Prudential Financial, Motorola, Altria Group, Starbucks, Mutual of Omaha, Sony Corporation, and the Environmental Protection Agency. More information about KCI, its services and its clients can be found online at: [www.KahnConsultingInc.com](http://www.KahnConsultingInc.com).

WHERE LAW & TECHNOLOGY MEET



## VIII. Endnotes

<sup>1</sup> In undertaking this engagement, KCI exclusively relied upon information supplied by Data Domain through internal and external documentation, and interviews with Data Domain representatives. KCI does not conduct independent laboratory testing of information technology products, and as such, did not evaluate the products in a laboratory setting or otherwise field-test any Data Domain products.

<sup>2</sup> As with any sensitive system, organizations using the Retention Lock capabilities should ensure that administrators with the authority to delete file systems and perform other significant activities in a Data Domain system have proper training and security clearance.

<sup>3</sup> "Disposal of Consumer Report Information and Records," 16 CFR Part 682.

<sup>4</sup> 45 CFR § 164.308(a)(7)(ii).

**Entire contents © 2009 Kahn Consulting, Inc. ("KCI"). Reproduction of this publication in any form without prior written permission is forbidden. KCI shall have no liability for errors, omissions or inadequacies in the information contained herein or for interpretations thereof. The reader assumes sole responsibility for the selection of these materials to achieve its intended results. The opinions expressed herein are subject to change without notice. All rights reserved. [www.KahnConsultingInc.com](http://www.KahnConsultingInc.com) [info@KahnConsultingInc.com](mailto:info@KahnConsultingInc.com) 847-266-0722**

WHERE LAW & TECHNOLOGY MEET

**KAHN**  
CONSULTING INC.