

BRAD JOHNS CONSULTING L.L.C

# Taming Big Data Storage with Crossroads Systems StrongBox

---

## Table of Contents

Table of Contents .....	2
Introduction .....	3
The Challenge.....	3
The Solution.....	4
Benefits .....	5
Data Preservation .....	6
Flexible Performance .....	6
Scalability.....	7
Future Proof .....	7
Summary.....	8
Trademarks and Special Notices .....	9

## Introduction

Big Data has burst onto the Information Technology scene. The confluence of advances in servers, analytic techniques and software has changed the way enterprises deal with computing infrastructures. The variety, volume and velocity of Big Data are also accelerating. Diverse applications such as simulation, visualization, modeling, seismic, video surveillance, and analytics are creating and processing unprecedented amounts of unstructured information. While these applications provide exciting new insights for business, they also place increasing requirements on the storage infrastructure and challenge storage management in a multitude of ways.

## The Challenge

So what exactly are the challenges that Big Data presents to storage management? First, Big Data demands retention of an unprecedented quantity of information—repositories of five, ten and even one hundred petabytes (PB) are common. Second, Big Data requires storage duration; much of this information may be kept for decades or longer. Third, Big Data presents a challenge to data preservation. The information in storage is often very valuable and difficult, if not impossible, to recreate; and for business, regulatory, or legal reasons, preservation in its original form may be essential. Fourth, while the data may lie dormant for long periods, when it is recalled, the performance requirements vary dramatically. For example, in high-performance computing massive amounts of data must be transmitted in a limited time frame. This process places considerable strain on the network and storage infrastructure, especially in comparison to online archive applications where the amount of data recalled is smaller, and the throughput requirements are much lower. Finally, not only are the Big Data repositories very large, they also grow at dramatic rates.

These drivers (volume, duration, variety, bandwidth and unstructured data growth) challenge traditional storage infrastructure and management techniques. The proliferation of storage devices makes management increasingly difficult. Service levels and recovery objectives are difficult to achieve. Large and growing quantities of

information also challenge resources such as budgets, floor space and energy requirements. Traditional storage solutions simply can't meet all of these requirements in a cost-effective manner.

In addition, an emerging Big Data storage management issue is data migration. Big Data will likely outlive the storage technology where it currently resides, and at some point in the future, the data will need to migrate to newer storage technology. This migration requirement raises some important questions:

- Is the software that created the file and understands the format still available? Is there other software that can understand the file format?
- What about the device interfaces that are used by the storage technology, will they be available and supported?
- How will the migration take place without impacting service levels during the migration period?

Given that the volume of data is so large, the need for data preservation so great, and the technical issues so varied, migration planning is essential to any Big Data storage management strategy.

### **The Solution**

The Crossroads StrongBox<sup>®</sup> addresses these concerns with an integrated, intelligent Big Data solution that combines StrongBox with IBM tape libraries and tape drives which support the open standard Linear Tape File System (LTFS). StrongBox is an easily installed appliance, which provides on-line file access to information stored on tape media using the LTFS format. Intelligent management of Big Data by StrongBox supports file server performance levels by retaining file stubs or whole files on disk cache. The size of the built-in disk cache used by StrongBox is flexible and tunable to meet the user's workflow requirements. Crossroads StrongBox stores every file on tape, leveraging the open standard LTFS to enable non-proprietary, self-describing, and fully portable storage for scalable, long-term data retention. The disk cache is scalable and flexible to support the users' performance needs while helping to eliminate latencies and optimizing the cost per gigabyte. StrongBox provides robust data protection via its tape media monitoring capabilities, file level hash codes, and policy-based second copy

and replication features. It includes the ability to transparently migrate<sup>1</sup> data from one tape storage technology to another.

StrongBox supports the IBM LTO Gen 5, Gen 6 and TS1140 tape drives. LTO Gen 6 is the latest generation of the Linear Tape Consortium tape drive. It features a 2.5 TB tape cartridge with 160 MB/Second native data rate and includes LTFS support. It supports a multi-vendor format specification for data interchange and multiple sources of media. LTFS-enabled drives and LTO tape are well suited for enterprises that prefer solutions from multiple suppliers and have significant data interchange requirements.

The IBM TS1140 tape drive has unique features that are particularly valuable in high-performance and capacity Big Data applications. The TS1140 also boasts the lowest hard error rate in the industry at 1 bit in  $10^{20}$  bits<sup>i</sup>. The hard error rate for the TS1140 tape drive is three orders of magnitude higher than LTO for improved reliability. In very large environments, the investment in tape media is significant. The IBM TS1100 family of tape drives has a long history of protecting and leveraging this investment. For example, the TS1140 can read the 10 year old original 3592 tape 300 GB JA tape cartridge. While there is no guarantee that this same level of compatibility will be afforded in future, media reuse is clearly an important objective of the TS1140 tape drive.

The enterprise-class IBM System Storage TS3500 tape library supports both the TS1140 and LTO Gen 6 tape drives. This library scales in a modular fashion to a massive 15 library tape complex supporting 2.7 Exabytes (2700 PB) of data, providing dual accessors for high availability, and supporting high-density expansion frames. The expansion frames are well suited for long-term archiving. For example, they can support 1,000 TS1140 tape cartridges in a single frame, providing 4 PB of storage, or 1,320 LTO Gen 6 tape cartridges, providing 3.3 PB of storage in a 10 square feet footprint.

## **Benefits**

The StrongBox solution for Big Data provides vital storage management benefits:

---

<sup>1</sup> Transparent migration is a planned StrongBox feature.

- Robust data preservation
- Easy, non-disruptive scalability
- A flexible architecture that supports a wide range of throughput performance
- Future proofing with non-proprietary standards

### **Data Preservation**

The Crossroads StrongBox Big Data solution incorporates robust data preservation features. StrongBox supports enterprises with critical data integrity needs with advanced features including tamper checking, HASH codes and no-delete. These capabilities combine to ensure that data is not accidentally or intentionally altered.

The reliability of physical storage media is an important aspect of any Big Data solution, and StrongBox's incorporation of tape storage addresses this concern. With new tape media formulations rated for a thirty year life, modern tape media is the most reliable storage available. In contrast, numerous studies reveal that hard disk drives used in traditional disk storage systems typically have annual failure rates of approximately 3%.<sup>2</sup> One such study by Instrumental Inc.<sup>3</sup> demonstrated the practical impact of different storage technologies. Reading as little as 11 terabytes on a consumer disk drive results in an error based on typical error rate specifications. This event would occur in just 32 hours of reading the disk at the full data rate. In contrast, the first error on an IBM TS1140 tape would occur in 1,512 years!

Data preservation is also enhanced by StrongBox's ability to create one or more copies of the data, based on user policy. The copy tape can be exported from the StrongBox and physically transported, or the data can be electronically transmitted to a StrongBox at a secondary location.

### **Flexible performance**

Throughput performance requirements in Big Data environments vary significantly, and the StrongBox solution can be architected to meet widely diverse requirements. In high-bandwidth environments, StrongBox may be configured with large amounts of external

---

<sup>2</sup> <http://static.usenix.org/event/fast07/tech/schroeder/schroeder.html/>

<sup>3</sup> Instrumental Inc. , Tape: Comparison of LTO and Enterprise, April 2013

disk storage to optimize recall performance. Users can define how long files reside on the cache to provide full disk performance of file reads. In addition, StrongBox includes the ability to “pre-fetch” data from tape to the disk cache to improve performance. For example, if the user knows a set of files will be needed, StrongBox can seamlessly move those files from tape to the disk cache. Pre-fetching effectively eliminates any latency associated with recalling information from the tape libraries.

StrongBox also supports the high sequential read bandwidth requirements of many Big Data applications. For example, 100 TB of data could reside on both the disk cache and tape. The cache provides rapid accessibility for high-performance needs, while data is always protected on tape.

However, if the primary requirement is for long-term preservation, while remaining online and cost efficient, the StrongBox solution can be architected with a smaller amount of disk cache relative to the amount of tape storage. For example, StrongBox could be configured with one petabyte of native capacity on tape and require only eight terabytes of disk cache. While the data appears “online” to the user, the solution offers dramatic cost and energy savings by storing the primary data on tape storage.

### **Scalability**

The StrongBox solution provides easy, massive capacity scaling with a transparent, intelligent storage layer between the server and tape storage. A single StrongBox can support up to 2 billion files of user data. With a single TS3500 tape library that can support up to 180 PB, a StrongBox can scale and manage multiple petabytes as a single file system. As capacity needs grow, tape media can be simply added to the tape library and made available for StrongBox use. By utilizing the expansion frames available with the TS3500 tape library, large increments of capacity can be added transparently to the environment; each expansion frame holds either 4 PB (TS1140) or 3.3 PB (LTO Gen 6) when fully populated with tape media.

### **Future Proof**

With the planned roadmap features<sup>ii</sup>, StrongBox will support the user or applications’ ability to migrate files from one tape technology to another. Transparently to the

applications, and based on user policy, future versions of StrongBox will start migrating data to new tape formats, reducing the complexity of migration, and providing an automated tool to lower the cost of managing the actual data migration. Of course, StrongBox utilizes intrinsic data verification capabilities to validate the authenticity of the migrated data.

To facilitate future migration efforts, StrongBox employs LTFS, an industry-standard format supported by the Storage Networking Industry Association (SNIA). This vendor-neutral format eliminates the potential dependency on proprietary software. Not only does the use of LTFS simplify future migration, but it also significantly reduces costs by removing the need for the long-term support and maintenance contracts associated with proprietary software.

### **Summary**

Big Data poses a myriad of storage management challenges. These challenges are large and growing, and addressing them is vital to the ongoing operations of the enterprise. Users of Big Data require new, innovative storage solutions to cost effectively manage this information while meeting demanding service levels and compliance needs. The StrongBox Big Data solution is a powerful new tool that integrates IBM tape libraries to simplify management of storage repositories and lower associated costs.

## **Trademarks and Special Notices**

This report was sponsored by Crossroads Systems.

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both.

Linear Tape-Open, LTO, the LTO Logo, Ultrium, and the Ultrium logo are trademarks of HP, IBM Corp. and Quantum in the U.S. and other countries

Other company, product, or service names may be trademarks or service marks of others.

---

### **<sup>i</sup> Disclaimers**

This information is provided for planning and information purposes only, and does not represent an express or implied warranty, guarantee or contractual commitment. THE SUBMITTAL OF THIS DATA CREATES NO WARRANTY, EXPRESSED OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR APPLICATION.

The reliability metrics set forth are based on the following statements, definitions and assumptions:

- Availability estimates are averages derived from measured field performance data and/or current engineering projection data. The following are excluded events from availability:
  - Vandalism, abuse, fire, flood damage or acts of God, or any other factor beyond IBM's reasonable control.
  - Scheduled outages of any nature, which include but are not limited to, remedial or preventative maintenance or engineering changes.
  - Operator errors, operational procedures or procedural changes that cause system outages.
  - Application or operating system errors outside of the storage subsystem which prevent access to data.
  - Corrupt data introduced into the solution.
  - Loss of or interruption to electrical power (IBM highly recommends use of an uninterruptible power system (UPS) to minimize power-related outages).
  - Network related error or application hang conditions
  - Installation of any new products, features or changes not defined as part of the proposed configuration.
  - Any outages or degraded performance that occurs due to delay or postponement by the customer in obtaining support or required services needed to preserve integrity and availability of the proposed solution.
- Field performance data reflects an average across installed subsystem configurations and features in typical application environments. The data assumes industry standard operating system recovery procedures and that environmental / operational characteristics and maintenance is performed according to IBM maintenance standards, procedures and schedules.
- Engineering projection data assumes the use of the industry standard operating system error recovery procedures, assumes maintenance is provided according to IBM maintenance

standards, procedures and schedules and assumes the units are operated in accordance with IBM environmental and operational specifications.

<sup>ii</sup> Plans represent goals and objectives only and are provided for planning and information purposes only. They are subject to change without notice.