



## WHITE PAPER

# Why Inline Data Reduction Is Required for Enterprise Flash Arrays

Sponsored by: Violin Memory

Eric Burgener  
September 2014

## IDC OPINION

---

With the explosion of mobile computing, social media, big data/analytics, and cloud computing technologies that are driving the buildout of the 3rd Platform of computing, data is expected to grow at a compound annual growth rate (CAGR) of 44% for at least the next five years. With heavily hard disk drive (HDD)-based storage infrastructures, enterprises are struggling to try and cope. Cost effectively and efficiently managing this data growth is *the* key challenge facing storage administrators today.

While data reduction technologies like compression and deduplication offer the promise of storing data much more cost effectively, on HDD-based systems, their use was generally limited to secondary storage environments like backup and archive. This is primarily because of the significant latencies these technologies introduced when used in conjunction with spinning disk, making them inappropriate to use inline in demanding application environments. Flash is a newer storage technology that is being widely deployed in 3rd Platform computing environments to meet evolving throughput and latency requirements, and the performance characteristics of flash are finally making the use of these types of data reduction technologies viable in primary storage environments.

Data reduction technologies like compression and deduplication can conservatively achieve average data reduction ratios of 6:1 against the types of mixed workloads that enterprises commonly run in 3rd Platform computing environments today. This reduction in the raw capacity required to store data represents significant savings for enterprises that are managing tens to hundreds of terabytes of capacity. The use of data reduction in combination with flash storage delivers a unique synergy: Flash performance enables the use of inline compression and deduplication, which in turn minimize storage capacity requirements and lower effective dollar-per-gigabyte costs while helping improve the write endurance of flash media.

While all-flash arrays (AFAs) were initially targeted at dedicated high-performance application environments, certain storage vendors have been enhancing their wares with enterprise-class features that are making these platforms appropriate for dense mixed-workload consolidation. With the primary obstacle to wider flash deployment in the datacenter being cost, the use of data reduction technologies significantly cuts the effective dollar-per-gigabyte costs of AFAs, making those platforms that support enterprise-class scalability, availability, reliability, and data management services that much more attractive. While enterprises are likely to keep an eye on dollar-per-gigabyte costs, IDC believes that a much better measure of the value that flash brings to the datacenter is total cost of ownership (TCO).

Enterprise-class AFAs that can be used for mixed-workload consolidation already offer a more compelling TCO than HDD-based arrays for 3rd Platform applications, and the addition of inline lossless data reduction only makes them that much more compelling.

## IN THIS WHITE PAPER

---

This white paper focuses on the importance of data reduction technologies like compression and deduplication for enterprise storage workloads, discussing both technology and business implications of their use with flash-based arrays. Then, it provides a brief overview of Violin Memory's implementation of these technologies in the company's flash-based enterprise storage solutions portfolio.

## SITUATION OVERVIEW

---

The 3rd Platform computing era is here. Driven by mobile computing, social media, and big data/analytics, the 3rd Platform of computing heavily leverages virtual infrastructure, flash, and cloud technologies to meet the performance, agility, and cost requirements demanded by enterprises today. Businesses are capturing and storing more data than ever before, leading to explosive data growth rates. IDC expects that between now and 2020, data that must be stored will grow at a CAGR of 44%. Efficiently managing data growth is *the* key challenge facing storage administrators today.

As the buildout of the 3rd Platform of computing continues, enterprises are deploying flash-based storage solutions to ensure that they can meet the performance requirements of today's new breed of real-time-oriented applications. These applications exhibit I/O patterns very different from client/server workloads that are difficult for HDD-based systems to cost effectively handle. Flash is at least an order of magnitude faster than spinning disk and draws significantly less power. With its much denser input/output operations per second (IOPS)-to-capacity ratio, it also requires far fewer flash devices to meet performance requirements, reducing floor space requirements and lowering the number of servers required to drive storage performance. IDC recommends that all datacenters should have at least some flash deployed to help them more cost effectively meet the performance requirements of 3rd Platform environments.

Flash is a must-have storage technology in the datacenter for reasons other than just pure performance. As flash is deployed in higher capacities within each datacenter, the secondary economic benefits of flash deployment include a very aggressive TCO for enterprise storage solutions that HDD-based systems just cannot touch. These benefits include far fewer devices needed to meet performance requirements, reduced energy and floor space costs, fewer servers needed to drive storage performance, and lower software licensing costs (due to the need for fewer servers). And flash, with its extremely low latencies, makes the use of inline data reduction technologies like compression and deduplication extremely viable for primary storage environments in a way that HDDs never could. The ability to deploy data reduction technologies in primary storage environments can easily reduce storage capacity requirements by 50-85% depending on the application workloads. For environments with hundreds of terabytes or more of data, this potentially represents huge savings.

## Enterprise Data Reduction Requirements

On their face, data reduction technologies make more efficient use of available storage capacity, increasing the usable capacity of devices over and above their rated raw storage capacities. The term *data reduction* generally refers to two types of technologies: compression and deduplication. Data compression involves encoding information using fewer bits than the original representation, whereas data deduplication eliminates duplicate copies of repeating data, replacing them with more space-efficient pointers to a single logical copy of the redundant data. Compression is generally used to identify short, repeated substrings inside individual files, whereas deduplication is intended to inspect large volumes of data, identifying large files or sections of files that are identical.

When it comes to data reduction, there are a variety of different algorithms. Algorithms vary in their data reduction ratios, their performance, the amount of overhead they generate, and whether they are a "lossless" or a "lossy" method. Lossless methods will be able to reconstitute the exact data that was originally operated on and are required where the exact accuracy of data is critical, but they may generate more overhead. Examples where lossless methods are required include numerical and statistical data as well as many types of text files. Lossy methods may result in higher data reduction ratios, higher performance, and less overhead, but they generally cannot reconstitute the exact data that was originally operated on but rather a "good enough" version for the relevant application. Examples where lossy methods are acceptable include some types of audio, image, and video files. Enterprise data requires lossless methods.

There are two general approaches to data reduction: inline and postprocess. Inline methods compress and/or deduplicate the data before writing it to primary storage, whereas postprocessing methods first write the full data sets to primary storage and then use an asynchronous background process to compress and/or deduplicate the data before writing it back. Inline methods can increase application latencies, but postprocess methods require more storage capacity up front.

For the most efficient capacity utilization, inline data reduction methods are preferable. With HDD-based systems, however, the storage devices were not fast enough to perform inline data reduction without impacting application performance. For this reason, data reduction with HDD-based systems was relegated primarily to secondary storage applications like backup and archive where low write latencies were not as critical. The use of flash media for primary storage, however, completely changes the situation. Given flash latencies that are at least an order of magnitude lower than HDDs, it becomes very viable to use data reduction technologies in primary application environments while still consistently providing sub-millisecond write latencies. This means that for most application workloads, the additional latencies introduced do not unduly impact application response times, allowing the use of data reduction technologies in high-performance primary storage application environments.

When flash-based systems are being deployed, having the data reduction performed inline provides several important benefits. First, it increases flash endurance. Flash array vendors have implemented a number of features, such as wear leveling, write minimization, and flash media overprovisioning, to increase the endurance of flash media. When data reduction is performed inline, only compressed or deduplicated data is actually written to flash, resulting in not only less data being written per I/O but also lower overall storage capacity requirements. Since write I/O is the primary cause of wear against flash media, anything that can be done to minimize it increases flash endurance.

For these reasons, data reduction in flash-based systems should be performed inline, maximizing flash endurance and keeping storage capacity requirements to a minimum.

Certain data sets will not benefit at all from either compression or deduplication. When these types of data sets are present, it is advantageous to be able to turn off data reduction. Data reduction algorithms use CPU cycles, and using them up without producing any benefit is not efficient. Also, there are certain data sets where pure performance is critical, and when these applications need latencies in the sub-200-microsecond range, it's unlikely that this could be achieved when data reduction is being actively performed. For these two reasons, customers may want to be able to turn off data reduction.

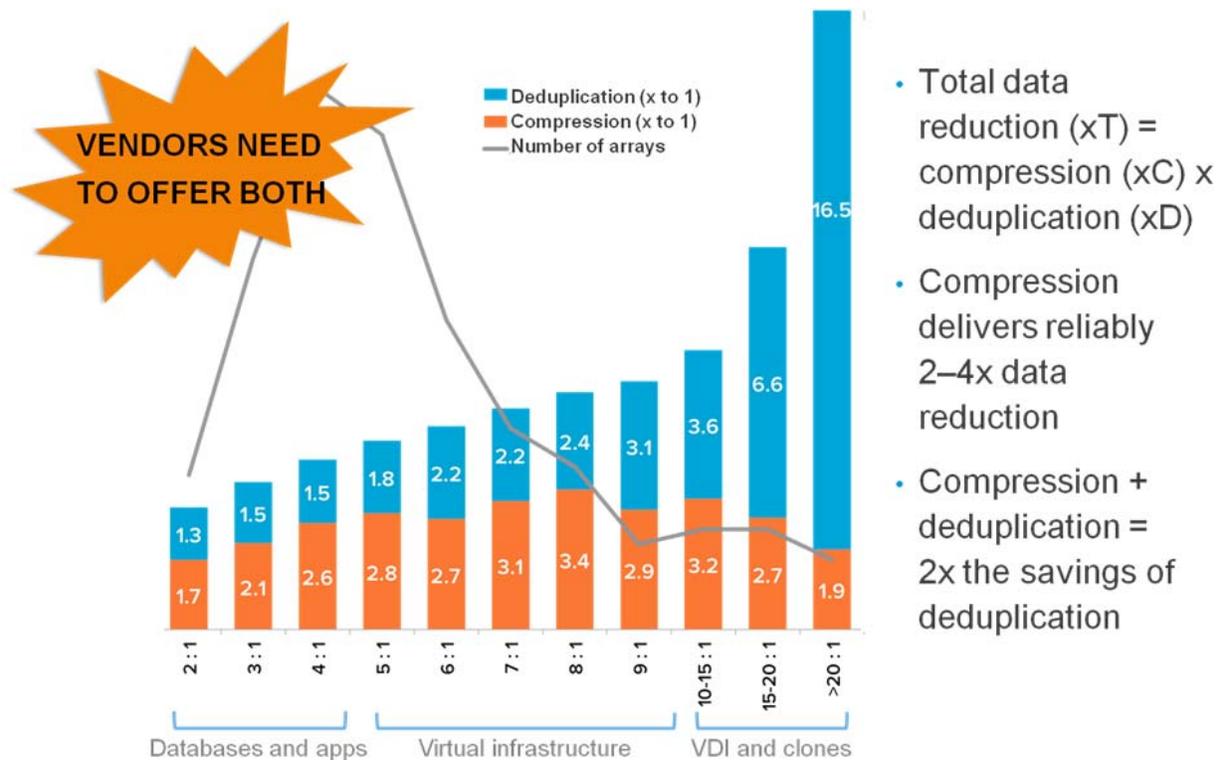
The level of granularity at which data reduction is enabled or disabled could also be important. If a flash-based array is being deployed just for a single application, then the ability to enable/disable data reduction at the system level is fine. But if customers intend to use an array to host a mix of storage workloads, then the ability to enable/disable data reduction at the workload level becomes extremely important. The ability to control data reduction more granularly, say at the level of the volume (for block-based storage) or file (for file-based storage), provides the added flexibility necessary to use a storage solution most efficiently for mixed-workload consolidation.

## Setting Expectations About Data Reduction Ratios

Data reduction ratios will vary by workload, and enterprises are advised to test vendor implementations against their specific data sets before making any savings calculations. But empirical data reviewed by IDC indicates that certain workloads benefit more from compression while others benefit more from deduplication. Compression seems to work better for database and other application workloads (e.g., Microsoft Exchange, SharePoint, SQL, Oracle), while deduplication is a clear winner for virtual desktop infrastructure (VDI) and other environments that heavily use clones. What is clear from Figure 1, however, is that the ability to leverage *both* technologies against mixed workloads produces overall higher data reduction ratios than using just one technology. To maximize the ability for an enterprise storage vendor to support different types of workloads with its storage solutions and maximize space savings, it is critical that the company offer both technologies in its data management services portfolio.

FIGURE 1

## Inline Compression *and* Deduplication



Source: IDC, 2014

## Flash Will Transform Enterprise Storage Solutions over Time

While flash brings huge performance benefits to virtual environments, the secondary economic benefits of flash deployment will not start to kick in until a datacenter has deployed flash in sufficient capacities. All flash-based arrays offer high performance, but not all of them offer the mature set of scalable data management services that are necessary to support mixed-workload consolidation. Without those data management services, flash-based arrays will continue to be deployed primarily for dedicated application environments with lower storage capacities.

For flash-based arrays to become viable platforms for mixed-workload consolidation, they must offer the same set of data management services found on today's proven HDD-based enterprise storage workhorses, including snapshots, cloning, thin provisioning, replication, and an ability to integrate seamlessly into preexisting datacenter workflows. These services must be provided in addition to several other features: a resilient platform that provides online expansion and device replacement, reconfiguration, maintenance, and firmware upgrades; an ability to scale to hundreds of terabytes of usable capacity while consistently providing predictable performance on an application-by-application basis; and a usable life consistent with the depreciation life cycles that are in use for enterprise storage today.

For flash-based arrays, this latter requirement needs to be met with flash management software that delivers reliability and endurance that meet or exceed the level of reliability and endurance of HDD-based solutions.

Once AFAs deliver on this feature set, there is one potential remaining obstacle: the dollar-per-gigabyte cost associated with flash media. Although IDC believes that dollar per gigabyte is not the right metric to evaluate the true value that flash brings to the datacenter environment – TCO at the system level is a much better measure – enterprises are likely to continue to keep an eye on it. The use of data reduction technologies can cut the effective dollar-per-gigabyte costs of flash significantly, making these technologies an excellent fit for a flash-based enterprise-class storage platform.

IDC believes that ultimately, AFAs will become the preferred enterprise storage workhorses for primary application environments. The performance they bring to the table is an absolute requirement for 3rd Platform computing that HDD-based systems just cannot cost effectively provide, the secondary economic benefits of flash deployed in high capacities will outrun the TCO that HDD-based systems can bring to the table, and lossless data reduction technologies just make the flash array TCO story that much more compelling. None of this can happen, however, unless the AFA can support dense mixed-workload consolidation while continuing to deliver consistently predictable performance as configuration scales. IDC expects that this transition will take five to seven years at the industry level, although there are AFA vendors today that are already enabling dense mixed-workload consolidation, making the all-flash datacenter a very viable option for primary storage.

## THE VIOLIN MEMORY SOLUTION

Violin Memory was a pioneer in the flash space, shipping its first product, a flash memory-based caching appliance, in 2007. The original founders of the company were primarily chip architects focused on designing a "flash-optimized storage architecture" to get the most I/O performance, in terms of high throughput and low latency, out of a storage system. The choice of a flash-based architecture was critical in meeting this objective. Violin's solutions portfolio includes the Violin Maestro All Flash Caching Appliance, the Violin 6000 Series AFA, the Violin Windows Flash Array, and the Violin Concerto 7000 Series AFA – all leveraging the company's Flash Fabric Architecture for high storage performance, resilience, and density. Violin sells flash-based storage solutions through a direct sales force and a worldwide network of resellers, including Dell, NEC, and Fujitsu.

Violin's customer base includes close to 500 companies (mostly Fortune 2000 customers) worldwide in financial services, manufacturing, media and entertainment, healthcare, government, technology, retail, and other markets. The company has had considerable success selling to market leaders in their industries, with three of the top 10 largest corporations, six of the top 10 largest telcos, three of the top 10 largest retailers, the top 5 largest software companies, and nine of the top 15 largest IT companies in the world as its customers.

Violin provides enterprise-class AFAs with high-performance, scalable, and resilient platforms that support online expansion and device replacement, reconfiguration, maintenance, and firmware upgrades. The Concerto 7000 AFA, Violin's current flagship, consistently delivers sub-500-microsecond latencies and scales to support up to 500,000 IOPS. Violin supports a broad set of data management services on its AFAs, including snapshots, cloning, thin provisioning, and replication. To ease integration into preexisting datacenter workflows, Violin supports APIs from VMware (including VADP, VAAI, and

VASA), Microsoft (including VSS and ODX), and popular application vendors like Oracle (RMAN). Based on this mature enterprise-class feature set, Violin has developed a strong group of bellwether reference accounts that have already extensively used Violin's solutions for mixed-workload consolidation. Violin also has some very large customers that are already very close to running an all-flash datacenter on Violin AFAs for their primary application environments.

The one key data management feature that was missing from Violin's portfolio was data reduction. With the introduction of the Concerto 2200 Data Reduction Appliance in August 2014, Violin has addressed this gap. The Concerto 2200 is an appliance that sits in front of up to four Violin AFAs, providing inline lossless compression and deduplication that can be enabled/disabled at a very granular level. This first release supports NFS ingest, giving customers the ability to turn data reduction on or off at the file, share, or share group level, and can be used with Violin 6000 or 7000 Series AFAs. Data reduction for block storage configurations is expected to follow in early 2015. Concerto 2200-based solutions are targeted at VDI and virtual server infrastructure (VSI) environments.

The Concerto 2200 dashboard allows customers to see the data reduction ratios they are achieving in real time against their workloads, providing critical information to help optimize the use of these technologies to reduce effective dollar-per-gigabyte costs. Data reduction ratios will vary by workload, but across typical mixed virtual workloads, customers can expect to achieve average ratios of 6:1, a conservative assumption that would extend the usable capacity of a single Violin 7000 Series AFA to 672TB. Under this assumption, a Concerto 2200-based VDI solution would deliver a street cost of under \$2/GB for a 2,500 persistent virtual desktop configuration.

The Concerto 2200 solution consists of two appliances working together to provide high-availability data reduction services across up to four connected Violin 6000 or 7000 Series AFAs. The connections from the hosts to the appliances are Ethernet, but the appliances are connected to Violin AFAs across Fibre Channel (FC). Compressed and deduplicated LUNs behind the appliances can be on one array or spread across up to four arrays. The appliance architecture allows customers to apply data reduction features to installed Violin AFAs or new Violin AFAs – or some combination thereof.

As with all inline data reduction, this approach will introduce some additional latencies. With high-performance network infrastructure (10GbE on the front end and 8Gb FC on the back end), Concerto 2200-based storage solutions will still be able to consistently deliver sub-millisecond latencies for mixed workloads. This type of performance represents a noticeable improvement at the application level for most applications, and for those applications that require even lower latencies, Violin offers the option to turn data reduction off. This approach gives customers the flexibility to define the feature profiles they need to meet various application requirements – a valuable capability for AFAs targeted for mixed-workload consolidation.

## FUTURE OUTLOOK

---

Flash-optimized true enterprise-class AFAs will ultimately become the enterprise storage workhorses in most datacenters, and support for granularly controlled lossless inline data reduction will be a baseline feature requirement in those platforms. Vendors like Violin that are targeting their platforms for dense mixed-workload consolidation will be forced to offer AFAs to stay competitive. Those vendors that offer granular enablement/disablement of data reduction will provide the ability to use these storage solutions most efficiently in mixed-workload environments.

The performance benefits that flash brings to datacenters are undeniable, but a good number of organizations – over 40% – have not yet deployed flash in production. These organizations cited the cost as the single biggest obstacle to deployment. The availability of inline data reduction features on flash-based arrays further narrows the dollar-per-gigabyte disparity between HDD and flash on acquisition and further tips the TCO balance in favor of those AFAs that can viably support dense mixed-workload consolidation.

It is IDC's contention that organizations are mistaken if they think flash is still too expensive to deploy for most primary application environments running on the 3rd Platform. Flash is not necessarily needed for all workloads, but in virtual environments running primary application workloads that demand performance, its use makes for more efficient and cost-effective storage configurations. This assumes, however, the availability of flash-based arrays that offer the enterprise feature set necessary to support mixed workloads. AFAs that deliver performance but lack those features will not allow the higher dollar-per-gigabyte costs of flash to be spread out across multiple applications to the point where they become more cost effective than HDD-based solutions.

## CHALLENGES/OPPORTUNITIES

---

With its new management team in place, Violin has come a long way from the issues surrounding its post-IPO performance in 2013. Clear and consistent marketing messages around its new focus on AFAs need to be combined with its success with mixed-workload consolidation and highlighted with case studies around those customers that are already well on the way to an all-flash datacenter built around Violin's solutions portfolio. This past year has seen a number of new products and features from Violin, including the Windows Flash Array; a new high-end array (the Concerto 7000); synchronous and asynchronous replication to provide better disaster recovery options; the integration of newer, more dense flash modules; and the ability to cluster storage together for improved performance, scalability, and resiliency.

The one area that was lacking in the Violin platform portfolio was inline data reduction, and Violin has now delivered that for NFS environments, with a stated intent to extend this to block-based storage in early 2015. Virtual environments today support a mix of block- and file-based applications, so AFAs targeted for mixed-workload consolidation need to support both. Violin chose to implement data reduction in a way that offers good flexibility to its customers – compression and deduplication can be enabled/disabled at a granular level, a choice that supports its stated intent to more aggressively pursue mixed-workload consolidation going forward, and the appliance-based approach allows the economic benefits of data reduction to be quickly and easily realized for installed base customers, preserving existing investment, as well as new customers.

## CONCLUSION

---

Data reduction capabilities are critical to improve the overall value proposition of storage, and flash makes their use extremely viable in high-performance primary storage environments – just one of the ways in which flash is transforming the datacenter. As low as any vendor can price its flash-based array to make it competitive on a dollar-per-gigabyte basis, the addition of data reduction capabilities to that same array will only improve its value proposition – as well as potentially providing other endurance benefits.

On flash-based arrays, there is a clear definition of how these features need to be implemented to maximize the additional value they bring to the table. Both compression and deduplication need to be available, they need to be implemented inline in a manner that is lossless and minimizes any associated latency hit, and they need to be granularly controlled at a file, file share, or share group level for NFS environments and at a LUN or VM level for block-based storage. When implemented in this manner, data reduction not only increases usable storage capacity, lowering the effective dollar-per-gigabyte cost of flash usage, but also improves flash endurance by minimizing the average amount of data written per I/O.

## About IDC

International Data Corporation (IDC) is the premier global provider of market intelligence, advisory services, and events for the information technology, telecommunications and consumer technology markets. IDC helps IT professionals, business executives, and the investment community make fact-based decisions on technology purchases and business strategy. More than 1,100 IDC analysts provide global, regional, and local expertise on technology and industry opportunities and trends in over 110 countries worldwide. For 50 years, IDC has provided strategic insights to help our clients achieve their key business objectives. IDC is a subsidiary of IDG, the world's leading technology media, research, and events company.

## Global Headquarters

5 Speen Street  
Framingham, MA 01701  
USA  
508.872.8200  
Twitter: @IDC  
[idc-insights-community.com](http://idc-insights-community.com)  
[www.idc.com](http://www.idc.com)

---

### Copyright Notice

External Publication of IDC Information and Data – Any IDC information that is to be used in advertising, press releases, or promotional materials requires prior written approval from the appropriate IDC Vice President or Country Manager. A draft of the proposed document should accompany any such request. IDC reserves the right to deny approval of external usage for any reason.

Copyright 2014 IDC. Reproduction without written permission is completely forbidden.

