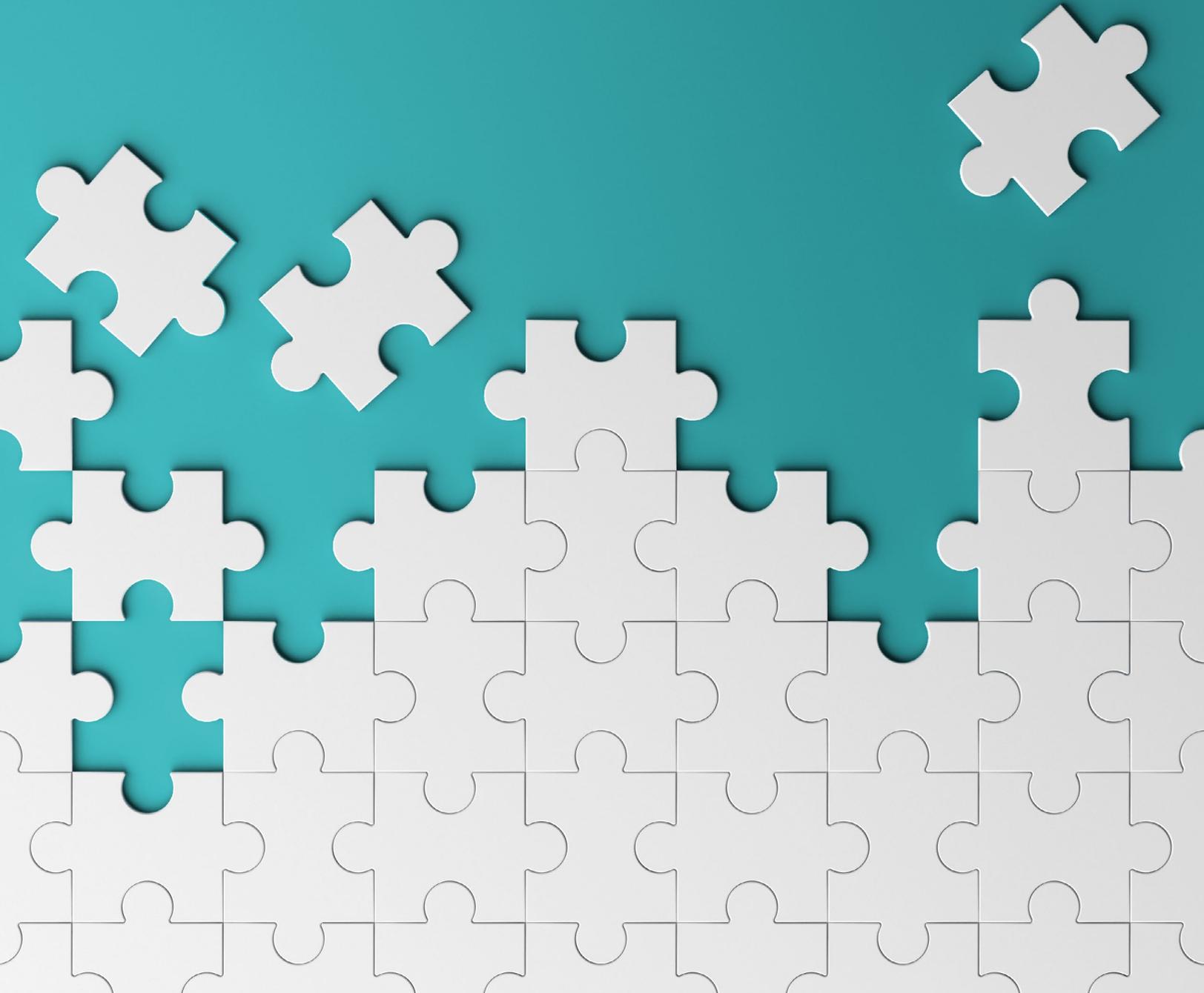


WHITE PAPER

# Consolidate Storage to Improve Availability and Lower Costs



# Table of Contents

Abstract .....	3
Introduction .....	4
Consolidation Works .....	5
Potential Storage Consolidation Benefits .....	6
Downtime Causes .....	6
Frequency of Repair Activities .....	7
Performance with Hardware Failures .....	8
Application RTOs and RPOs .....	8
Storage Array Availability Fundamentals .....	9
Real vs. Perceived Risks .....	9
Availability .....	9
High Availability System Design Guidelines .....	9
Mean Time Between Data Losses .....	11
Acquisition and Ownership Costs .....	11
Storage Consolidation Best Practices .....	12
Conclusion .....	12

## Abstract

This white paper shows that, in most environments, the concerns of CIOs, Storage Architects, and IT Directors about how storage consolidation increases the size of fault domains are unfounded because modern multi-controller scale enterprise arrays that have been market validated can effectively deliver 100% availability<sup>1</sup>. They have more usable availability, higher levels of fault tolerance, and more reliable non-disruptive software update capabilities than dual controller arrays<sup>2</sup>. Storage consolidation simplifies infrastructure topologies, by reducing the number of arrays being managed and opportunities for misconfigurations both of which contribute to downtime budgets. Storage consolidation also creates larger pools of free space that improve usable availability and operational efficiency relative to collections of smaller dual controller storage systems. It also reduces the frequency of repair activities as well as power and cooling requirements by decreasing the number of controllers and ports in the SAN infrastructure.

<sup>1</sup> "Market validated" arrays are arrays that have been in general available for a minimum of 9 to 12 months and have 10s of thousands of machine years of production experience.

<sup>2</sup> "Usable availability" is the ability to meet service level objectives in the presence of hardware failures.

## Introduction

Existing storage infrastructures that have “evolved” frequently suffer from avoidable inefficiencies such as: management complexity, a lack of agility, availability concerns, performance problems, capacity shortfalls, skills shortages, backup/recovery issues, and wasted money.

Storage consolidation projects solve many of these problems by providing infrastructure architects with an opportunity to re-evaluate past decisions and design a storage infrastructure tailored to current and future workloads. With the creation of self-managing storage arrays that make intelligent data placement decisions and access to a plethora of migration tools, technical risks, concerns about data availability, the size of fault domains, and skills shortages are no longer reasons to fear storage consolidation.

This whitepaper demonstrates that consolidation is a successful strategy for improving availability and operational efficiency. It also examines the primary causes of downtime and repair activities; provides a high-level tutorial on failure mathematics, describes the benefits of storage consolidation, and defines the best practices on executing a storage consolidation project.

## Consolidation Works

For skeptics of the “putting all your eggs in one basket” strategy, history has shown us that building sturdier baskets is a strategy that works, and you only have one basket to watch. Three great examples that highlight the success of this strategy are:

- ▶ **Air Transportation** - From the very beginning airplane manufacturers have focused on building bigger, safer planes that can fly farther, faster, and higher. The result is that air travel is safer per seat/mile than driving and plane crashes are a rarity. Why? Because bigger planes have economies of scale that make the addition of system redundancy and advanced safety features inherently more affordable relative to smaller planes, and there are fewer opportunities for air traffic control to misroute planes because there are fewer planes in the air.
  - ▶ More Flights
  - ▶ Limited Capacity
  - ▶ Less Redundancy
- ▶ **Crude Oil Transportation** - The growth in global crude oil shipments saw the application of the same “bigger is better” strategy that aircraft manufacturers have employed: fewer bigger, safer, double-hulled tankers equipped with GPS navigation and lots of automation instead of many smaller tankers. Since the cost of building a ship is tied tightly to how much they weigh (more low-tech than high-tech), and volumes go up faster than surface area (i.e. weight), building big double-hulled ships is inherently more cost effective than building smaller double-hulled ships. The result is faster, more reliable, accident-free delivery of oil around the world.
  - ▶ Fewer Flights
  - ▶ Lower \$/seat
  - ▶ More Redundancy
  - ▶ More Voyages/Risk
  - ▶ Higher Crew Costs
  - ▶ More Accidents
  - ▶ Faster to Market
  - ▶ Higher Resiliency
  - ▶ Better Safety
- ▶ **Storage Consolidation** - Market statistics from both IDC and Gartner indicate that users are pursuing a similar consolidation strategy as they attempt to improve storage infrastructure availability and operational efficiency while lowering costs. More specifically, statistics show annual PB shipments increasing even as the number of storage arrays shipping per year are in decline. In other words, average capacity configurations of storage arrays are increasing.



## Potential Storage Consolidation Benefits

The benefits that are most effective in building support for storage consolidation are those that address highest impact operational pain points. The following pain points are perennial favorites: lack of agility, availability concerns, performance problems, capacity shortfalls, skills shortages, backup/recovery issues, and budget constraints.

Following are sample user benefits that address these pain points:

- ▶ Consolidating storage makes stepping up in class from dual controller mid-range to multi-controller high-end arrays financially affordable.
- ▶ Modern multi-controller arrays have more scale and deliver more consistent performance in the presence of hardware failures and/or software updates.
- ▶ Lowering the \$/TB/month prices enables the purchasing of more storage.
- ▶ Self-managing storage that simplifies management keeps headcount flat even as capacity increases; this also applies to improvements in performance and D/R.
- ▶ Modern multi-controller hybrid arrays that implement RAID 6 or erasure coding increase mean time between data losses (MTBDLs) by orders of magnitude relative to RAID 1 and RAID 5 configurations.
- ▶ Using storage efficiency features such as data compression and deduplication lowers part counts, frequency of repair activities, and \$/PB costs.

## Downtime Causes

Table 1 shows, in descending order of frequency, the primary causes of storage system downtime. While one could debate the specific order of these causes, we can agree that storage and SAN related hardware failures are not primary sources of downtime events. Experience and root cause analysis have shown that when storage hardware failures cause downtime it is almost always a software bug that was exposed by a hardware failure that caused the outage. Were it otherwise, cloud-based monitoring and analytics would have no impact on storage system availability.

TABLE 1

Cause	Comment
Human errors	Proportional to complexity and frequency of interactions between storage admins and FEs with the storage arrays
Software bugs	Inversely proportional to code maturity and proportional to software size and complexity
Poor software change control	Improved by cloud-based analytics
Infrastructure misconfigurations	Proportional to infrastructure complexity and end-to-end configuration validation which is frequently lacking
Defective D/R testing	Failover/failback fail to work

# Frequency of Repair Activities

The frequency of repair activities is proportional to the number of components in a storage system and inversely proportional to the components' MTBFs. More parts means more failures which means more repair activities. Equation 1 defines the MTBF as the reciprocal of the failure rate, which means that a higher MTBF translates into a lower frequency of repair activities.

## Equation 1 – MTBF

**MTBF = 1/failure rate**

Since PB scale data centers have hundreds to thousands of HDDs and SSDs running 24x7, they account for almost all hardware repair activities. Putting theory into practice, Figure 1 shows that Seagate nearline HDDs rated at 1.2M MTBF have an annualized failure rate of 0.73%. Figure 1 shows that an array built with 480 of these Seagate nearline HDDs, approximately 4 PB of dual parity protected capacity with 12 TB HDDs, should experience no more than 3.5 HDD related repair activities per year, or no more than 1 HDD repair activity/PB/quarter.

**FIGURE 1**

MTBF (hours).....	1,200,000
Hours/year (24 x 365).....	8,760
Annual Failure Rate (AFR) (Hours per year / MTBF).....	0.73%
# of array HDDs.....	480
Annual frequency of HDD related repair activities (AFR x # of array HDDs).....	3.504

These relatively frequent repair activities coupled with disk capacities that are growing faster than data transfer rates, have made using more resilient data protection schemes and shrinking rebuild times critical design objectives in modern arrays because fast rebuild times reduce the “window of vulnerability” that occurs every time an HDD or SSD fails<sup>3</sup>. Shrinking the window of vulnerability improves data durability or mean times between data losses.

Two techniques that have proven their worth in reducing rebuild times are:

- ▶ Replacing the concept of spare disks with spare capacity which speeds data rebuilds by parallelizing the rebuild process.
- ▶ Only rebuilding data instead of the whole device’s capacity further reduces data rebuild times.

Data reduction technologies (compression and deduplication) do not reduce rebuild times, but they do reduce the number of devices needed to hold a given amount of data, thereby reducing failure rates: intelligent rebuilds. High-end multi-controller architectures have the compute power and bandwidth needed to create performant implementations of these techniques, even in the presence of hardware failures.

For comparative purposes let us assume that a user has distributed these 480 HDDs across four dual-controller midrange arrays to contain the size of fault domains. Let us further assume an AFR of 1% because server vendors do not publish server AFRs and failure rates above 1% can create customer

<sup>3</sup> Common examples of more resilient data protection schemes include: dual parity, erasure coding, and Reed Solomon

satisfaction and business problems. The math yields an expected annual frequency of repair activities for non-media related hardware failures of 0.08 per year. Even if we increased the AFR to 10% or an MTBF of only 87,600 hours, the expected AFRA would still be less than once per year. In other words, array controller failures are essentially a “don't care” when deciding to consolidate storage.

## Performance with Hardware Failures

Since controller failures, absent software bugs or botched repairs, are non-critical events, the visible impact of a controller failure in a dual or multi-controller array generally manifests as a reduction in performance (IOPS), throughput (GB/s), and/or latency (milliseconds), not a loss of data accessibility or data integrity. Failure mathematics and the concept of usable availability inherently favor multi-controller arrays because the performance impact of controller failures is inversely proportional to the number of active controllers in an array. A single controller failure in an active/active dual-controller array could decrease performance by up to 50% or more if it forces the surviving controller to switch from write-in to write-through cache mode; an active/active/active three controller array by up to 1/3; a four-controller array by up to 25%, et al. The use of the phrase “up to” is not a subtle “get out of jail free card,” but an acknowledgment that sharing the increased workload across multiple surviving controllers not running at maximum capacity further minimizes the impact of controller failures.

Users that leave themselves with 25% to 30% of excess performance headroom will rarely notice a failure aside from proactively sent warning messages from the storage array vendor. Headroom is valuable because it enables arrays to gracefully tolerate hardware failures. It also improves usable availability by hiding software bugs that are only exposed when an array is under extreme stress. Maintaining 25% to 30% headroom also provides users not using COD or consumption-based pricing models with enough time to negotiate cost effective upgrades.

Caveat - organizations that are unable to reserve enough headroom to accommodate unforeseen performance problems or organic growth may benefit from a strategy of non-consolidation because of this topology's inefficiencies, but at the cost of lower operational and financial efficiency.

## Application RTOs and RPOs

Application RTOs and RPOs both depend upon the frequency of snapshots and the amount of data being protected. Increasing the frequency of snapshots shortens RPOs by shrinking the time between the last snapshot being taken and a recovery initiation. It also shortens RTOs by shrinking the number of transactions that need to be rolled forward or the amount of data that needs to be restored. These insights inevitably culminate in the concept of continuous data protection or the snapshotting of every change made to data. While the concept is alluring, it suffers from two inherent problems: creating snapshots adds software overhead (i.e., updating metadata) to every write operation, and snapshots consume capacity with every write operation.

Hence snapshot frequency is determined by performance considerations and limited by economics. Once again, upgrading to multi-controller architectures enables frequent snapshots, per the previous performance discussion. It also makes recovering from a data corruption event faster than recovery on a dual-controller array, bottlenecked by CPU cycles or bandwidth. Consolidating many storage arrays into

fewer multi-controller arrays with fewer and larger pools of storage may wash away any potential multi-controller array \$/TB price premiums by reducing stranded capacity and complexity, improving staff productivity, and lowering downtime costs.

## Storage Array Availability Fundamentals

### REAL VS. PERCEIVED RISKS

Equation 2 relates risk, fault domains, and availability in a manner consistent with common sense. Increasing the size of the blast radius or fault domain increases risk and improving availability decreases risk. Since anything less than 100% availability creates a real risk of downtime, storage arrays cannot have any single points of failure (SPOFs), must have fault tolerance, and provide non-disruptive everything: software updates, repair activities, and capacity upgrades.

#### *Equation 2 – Risk*

**Risk = Blast radius X (1- Availability)**

Staying with the transportation analogies, when a supertanker breaks up, it's a catastrophe; when a fuel truck crashes, it's bad but it's rarely a catastrophe. This makes it easier to invest more in a supertanker's safety features than a fuel truck, and protecting a fuel truck is inherently harder than protecting the supertanker.

### AVAILABILITY

Equation 3 defines the relationship between availability, MTBF, and MTTR. It also highlights that achieving 100% data availability requires fault tolerance and that all repair activities are non-disruptive. It also shows that low MTBFs result in frequent repair activities rather than downtime..

#### *Equation 3 – Availability*

**Availability = MTBF/(MTBF + MTTR)**

For those evaluating the technical risks of storage consolidation, but not interested in more detailed failure mathematics analysis, consolidation into fewer multi-controller high-end storage arrays increases MTBDLs by reducing rebuild times, usable availability by reducing the impact of hardware failures on a percentage basis and lowering the frequency of repair activities by reducing your storage infrastructure parts count.

## High Availability System Design Guidelines

Storage array hardware availability is determined by the factors listed in Table 2. However, it is important to note that 100% storage array availability does not guarantee that data will never be lost. Protection against data loss is provided by RAID and erasure coding schemes with the degree of protection limited by performance and cost considerations which will be explored later.

TABLE 2

Factor	Comment
The number of components in the system	If it's not there it can't break. This and cost considerations are the major drivers toward simplicity and reducing the system's parts count.
The Mean Time Between Failure of every component	High quality components fail less frequently than consumer grade components <sup>4</sup>
The number of critical failure modes in the array <sup>5</sup>	No SPOFs is the aspirational goal because software will always remain a SPOF even after critical hardware failures are eliminated.
The mean time to repair	100% availability demands that all repairs be non-disruptive.

There are no inherent availability advantages to scale-up versus scale-out architectures because the fault tolerance, recovery, and non-disruptive repair capabilities built into Infinidat and many other storage arrays, have practically decoupled availability from hardware failure and software updates. This shifts the focus on improving storage availability to improving software quality, recovery capabilities, and reducing the number of parts needed to build the array.

Reducing the parts count, in addition to lowering costs, reduces the frequency of hardware repair activities, which reduces the opportunities for Field Engineers to make mistakes that bring the storage system down. Analyzing software quality and recovery capabilities is an inherently subjective analysis because it involves so many unknowns. Hence the use of the frequency of software updates, excluding functional enhancements, as an analog of code quality.

Since the MTBF of controller electronics is unaffected by the IOPS pushed through them, the most obvious way to reduce the parts count is to push more IOPS through each controller. The number of IOPS that a controller can support is determined by CPU performance and software efficiency. Comparing annual CPU performance improvements to HDD data transfer rate improvements shows microprocessor improvements outstripping media improvements: approximately 40% per year for microprocessors vs. 10% to 15% per year for HDDs. This comparative advantage drives down the CPU to capacity ratio needed to avoid performance bottlenecks. The difference between CPU performance improvements and data transfer rate improvements favors the continued building of scale-up arrays that also possess scale-out capabilities, if and when they are ultimately needed.

These trends bestow frequency of repair activity and cost advantages on scale-up arrays because they generally have fewer controllers and supporting electronics - HBAs, NICs, power supplies, fans, etc. than scale-out arrays of equivalent capacity. This advantage generally increases with capacity because scale-out arrays frequently add capacity by adding nodes that include controllers and their supporting electronics. Leveraging fewer components into a lower cost of goods and the environmental advantages of scale-up versus scale-out arrays helps explain the continued success of scale-up arrays in the marketplace.

<sup>4</sup> Manufacturing defects reduce component MTBFs and are therefore not included in Table 1

<sup>5</sup> Critical failures are failures that bring the array down or that require disruptive repairs.

## Mean Time Between Data Losses

MTBDLs are at least as important to reliable data center operations as system availability because they can result in longer duration events than software bugs or hardware failures. Most vendors answer questions about reliability or availability with claims ranging from 99.99% to 100.00% availability but are reluctant to discuss rebuild times and anticipated HDD and SSD related repair activities. Rebuild times and MTBDLs are inseparable and rebuild times may become very long, hours to days in busy traditional architecture arrays.

Following are the factors that influence MTBDLs.

- ▶ The MTBF of the SSDs or HDDs that are storing data - a higher MTBF (i.e., the quality of the components) reduces the frequency of repair activities and data rebuild times
- ▶ The total number of SSDs and HDDs in the storage array - because it influences the probability of multiple SSD or HDD failures occurring within a RAID or erasure group.
- ▶ The number of SSDs or HDDs in a RAID or erasure group - more parts mean more failures, more repair activities, and more time rebuilding data because larger RAID or erasure groups hold more data than shorter ones.
- ▶ The number of failures that can be tolerated in a RAID group - RAID 1, 10, and 5 guarantee data integrity in the presence of single HDD or SSD failures; RAID 6 guarantees data integrity in the presence of two HDD or SSD failures. Erasure codes, because they store data as systems of equations, can protect against any number of failures if its impact on performance and cost are acceptable.
- ▶ Rebuild times - Intuitively we appreciate that shorter rebuild times provide higher MTBDL, but its actual impact on MTBDL is often greatly underestimated because it is a variable that is not memorable.

## Acquisition and Ownership Costs

Competition between on-premises storage vendors and cloud providers has significantly eroded vendors' ability to charge different prices for HDD and SSDs installed in high-end arrays vs. mid-range arrays. Hence differences in acquisition and ownership costs between dual-controller and multi-controller arrays are increasingly being influenced by differences in controller costs, hardware maintenance, and one-time/annual software licensing charges. Thus the \$/TB cost delta between dual and multi-controller arrays shrinks as configurations grow, especially at multi-petabyte scale.

Given that storage arrays historically have had high list prices, complex pricing models, and aggressively negotiated one-off discounts there are no simple rules of thumb (ROT) that can be used to estimate actual \$/TB/month pricing with a high degree of accuracy. InfiniBox with its sub-millisecond response times, multi-petabyte scale, all-inclusive software pricing model, and disruptive \$/TB pricing further complicates efforts at creating simple ROTs for \$/TB pricing as does the heavy use of replication software, management tools, and scripting that create strong dependencies even as they hide the architectural ugliness of arrays whose legacies may stretch back decades.

# Storage Consolidation Best Practices

Storage consolidation projects frequently fail because they challenge the status quo, and change creates winners, losers, and risks. Following are some of the most important best practices that maximize the probability of consolidation projects being successfully completed.

- ▶ Obtain the support of senior management before proceeding with any storage consolidation project.
- ▶ Build a team that includes important stakeholders: storage architects, operations, developers, finance, and legal.
- ▶ Include changes in infrastructure acquisition and ownership costs; downtime costs, migration costs, and productivity improvements in your risk/reward and ROI analysis.
- ▶ Give your vendors a vested interest in the success of your consolidation projects by having them share data migration costs and risks with you.
- ▶ Sell the benefits of storage consolidation to build support among users.

## Conclusion

At multi-petabyte scale, consolidating storage onto multi-controller arrays, especially arrays with all-inclusive software pricing, capacity on demand (COD), and consumption-based pricing models, makes storage consolidation the optimal business decision. This decision should take into account the impacts to availability, performance, staff productivity, and total operating cost.

Vendors that show a willingness to compete on price and are willing to take ownership of array configuration and data migration further increase the probability of a successful consolidation project. Infinidat is such a storage vendor.



**STANLEY ZAFFOS** Sr. VP, Product Marketing, Infinidat

Prior to joining Infinidat, Stanley Zaffos was a Research VP with Gartner focused on Infrastructure and Operations Management. His areas of expertise cover storage systems, emerging storage technologies, software-defined storage, hyper-converged infrastructure, and hybrid cloud infrastructure. He's worked with numerous clients to develop messaging and collateral that maximizes the impact of their product announcements and sales training and has helped them to define roadmaps to ensure ongoing competitive advantage.