



# Solving the Storage Conundrum

The most difficult storage problems Internet Services face — and how to overcome them

## Abstract

There are two things that can be done to minimize damage to Internet Services Companies caused by unexpected downtime or outages. The first is providing customer transparency, giving them detailed explanations as to exactly what's going on. The second is to get rid of or reduce as many outages or performance degradations as possible by eliminating some of the most common root causes. Regrettably, too many Internet services organizations cut corners on their storage in a futile attempt to keep costs low. Far too often that approach ends up costing far more both in out-of-pocket costs, lost customers, lost reputation, and lost revenue. This paper will focus on the cure for common root causes directly attributed to the storage or file systems.

## Introduction

Internet services companies share a common goal of providing products and services over the Internet, a market that grows larger every day as more of the world's population goes online. The common economic value proposition for all Internet services comes from exceptional quality of service and performance at a reasonable price. It is the promise of this high quality of service and performance that causes so many customers to require some form of written guarantee. This guarantee takes the form of a service level agreement (SLA). The SLA is essentially a penalty clause that refunds a percentage of the customers' money when service and/or performance levels drop below the promised or agreed to level. The greater the disparity, the greater refund.

However, when service levels drop as the result of an outage or performance decrease, to the point when SLA refunds start to kick in, the real damage to both the service provider and the customer is far more than the refund and a lot more excruciating. Loss of reputation, users, customers, prospects, advertisers, and irretrievable revenue are all consequences.

Most Internet services have the normal human denial reaction to service outages and/or performance level drops. They recognize it can happen to them (which is why they provide SLAs), and they irrationally believe that it will not...until it does.

It happens to everyone, including those with outstanding service reputations. A simple Web search demonstrates that Rackspace (a service provider that prides itself on an impeccable service record), Amazon S3, Salesforce.com, Facebook, MySpace, Skype, NetFlix, YouTube, Mozilla, SourceForge (which includes the popular Slashdot Web site), CNN, Netcetera, GoDaddy, Joyent, and many others, have all have faced well-publicized, sizeable, outages and degraded performance at different times. Every single outage and/or performance degradation is painful while creating an exasperating fire drill for IT.

It is not an issue of "if" there will be an outage or performance degradation, the issues are when, how long, and how often they will occur. The root causes of an outage or service performance degradation are varied and numerous.

Unfortunately, no Internet service can provide 100% uptime. No matter what type of infrastructure that is put in place, there will always be some outages or performance degradations. There are two things that can be done to minimize damage. The first is providing customer transparency, giving them detailed explanations as to exactly what's going on. The second is to get rid of or reduce as many outages or performance degradations as possible by eliminating some of the most common root causes. Regrettably, too many Internet services organizations cut corners on their storage in a futile attempt to keep costs low. Far too often that approach ends up costing far more both in out-of-pocket costs, lost customers, lost reputation, and lost revenue. This paper will focus on the cure for common root causes directly attributed to the storage or file systems.

## Storage or File System Caused Service Outages and Performance Degradations

The primary storage or file system causes of service outages and performance degradations are:

- Inability to adapt to dynamically changing performance requirements in real time (a.k.a. spikes in performance demand);
- Online scalability;
- Operational complexity;
- Inadequate data availability.

### Inability of the storage or file system to adapt to dynamically changing performance requirements in real time

Non-predictive performance requirement is the most common and vexing storage induced service outage or performance degradation issue. It is not unusual for a Internet service to suddenly have an unexpected serious spike in IOPS or access requirements from a current customer. The customer could be running a new promotion that gets an unforeseen positive response. Or there is an unanticipated situation within a customer that puts unanticipated excessive strain on the storage



system. This was the root cause of a recent Amazon Web Services outage that took out their EC2, S3, and AWS and lasted for over 3 hours.

The outage was caused by a sudden unforeseen surge in a particular type of usage (PUTs and GETs of private files which require cryptographic credentials, rather than GETs of public files that require no credentials). One very large customer, plus several other customers, created the surge suddenly and without warning, dramatically increasing their usage. The storage system was not able to adapt and keep up. The strain forced it to shut down taking down all of Amazon's Web services with it.

The circumstances that shut down Amazon's online services are not atypical. It will happen to any service provider when their storage cannot adapt to sudden unexpected performance loads that exceed their storage systems' ability to keep up.

The Apple launch of the iPhone 3G is another example of a storage system not adapting to a temporary increased demand based on a onetime event. New, excited, Apple customers of the iPhone 3G had to wait from hours to days to activate their phones generating the wrong kind of press coverage. This is not the kind of customer experience Apple and AT&T want people to remember and can only negatively affect long-term sales and revenues.

### **Online Scalability**

What frequently works well in the small or test environment rarely works well when scaled up. The larger a system grows, the more complicated it becomes. Many Internet services don't discover this until after they start to grow. At that point, they are already experiencing the cost and complexity issues of attempting to make do with limited scalable storage.

The most common forms of storage found in Internet services are the do it yourself (DIY) bricks (servers with embedded direct attached storage), open source NAS (LINUX Samba), low cost NAS (Microsoft Storage Server), and the more expensive industry leader, NetApp. The conventional wisdom being that these systems are inexpensive, or in the case of NetApp, relatively inexpensive. Unfortunately, conventional wisdom is wrong, as it usually is.

DIY bricks, open source NAS, low cost NAS, even NetApp architectures are constrained in their ability to scale. They were designed for the small to medium organization requirements, not for the requirements of the service provider. Per Echo Labs, it is guaranteed that typical NAS or brick systems will run out of gas. As a result, these storage systems rarely allow online expansion of capacity or performance and even when they do (such as in the case of NetApp), expansion is still limited.

Common Internet service scaling requirements dictate that the storage system's performance and capacity must grow on-demand, just-in-time, online, and application non-disruptive, with nominal over-bought capacity. To meet provider scaling requirements with the DIY or NAS storage systems entails adding more systems. This is when the trouble starts.

More systems necessitates load balancing of data and more mounts for every server or client accessing their data. Every time a new system is required to increase capacity or performance, application disruptions increase and the probability of a serious outage or performance degradation increases exponentially. This type of scaling is by definition, offline. It only takes one system not balanced properly for SLA performance to degrade performance. Unfortunately, many DIY bricks and NAS systems have a file limit where performance does not simply degrade, it falls off a cliff. If the file system exceeds certain performance or file management thresholds, it is likely to cause a complete outage.

Scaling these storage systems begets operational complexity. Operational complexity begets human error. Human error begets system outages or performance degradation. As storage scales linearly, operational complexity increases by orders of magnitude.

### Operational Complexity

To manage the rapid increase in numbers of storage systems means constant data migration between systems. Deciding how to distribute the performance and data load as storage systems are added requires a detailed planning exercise. As the number of systems increases, so do the data movement combinations. What data is migrated to what system and how much is based on current performance load and number of files under management in each system. As straightforward as it seems, it is anything but. Different applications have different performance and capacity growth rates. Some applications are more mission critical than others and the higher priority injects a complication that makes ongoing scaling a nightmare of biblical proportions.

Applications go down hard when they run out of storage or when performance degrades so severely that the application times out as if the storage were disconnected. The probability of an outage becomes increasingly acute with each new system and further compounded by massive increase in management tasks. Most Internet services attempt to run a lean efficient IT administrative organization. The rapid upsurge in storage management operations runs against the grain forcing an alarming unexpected rise in operational budgets.

### Inadequate Data Availability

Another storage-induced service outage comes from inadequate data availability. Rarely are there adequate measures in place to provide sufficient data availability to prevent a service outage from a system failure; site failure; accidental data deletion; administrator error; malware attack, or something as simple as scheduled maintenance.

To keep data available when there is:

- A disk fault or failure requires RAID 1, 10, 5, or 6.
- Dual disk faults or failures require RAID 6 or volume mirroring.
- A storage system failure requires duplicate storage systems and software licenses with replication between them, or a clustered storage system with data accessible by all nodes.
- A site failure requires duplicate storage systems, software licenses, and server systems plus wide area replication.
- Accidental deletion, admin errors, and malware attacks require mountable time based snapshots (a.k.a. clones), or Continuous Data Protection (CDP).
- Scheduled maintenance (microcode upgrades patches, disk replacements or additions, etc.) requires that a copy of the data on the system being maintained be made available to its applications without application disruption. Doing this requires a mountable snapshot or clone that is accessible by the applications.

The typical DIY storage bricks or NAS simply do not have this level of data availability. Then there is the other cause of significant service outages, speed of recovery (a.k.a. recovery time objective or RTO).

When there is adequate data protection, it is just as likely that the data RTO is not fast enough to prevent a serious service outage. To do so requires implementing data protection capabilities that may not be available for the DIY storage bricks and NAS. That which is available calls for multiple software applications with completely different management schemes, consoles, and tasks. At the end of the day, data availability is still inadequate while complexity and operational budgets mount ever higher.

### Finding the Cure

The DIY and standard NAS approaches are short-sighted. This bootstrap approach may appear to minimize upfront costs but it does so by brutally raising operational expenses, escalating total cost of ownership, hamstringing growth, while ultimately risking long-term provider success.

To eliminate or mitigate storage system based outages and performance degradations necessitate a different way of thinking about storage and infrastructure. The storage systems must be capable of dealing with Internet services' requirements including and not limited to:

**Adaptability to dynamically changing performance requirements in real time**

- Unexpected spikes in IOPS or throughput performance are met on-demand
- Unexpected spikes in capacity are met on-demand

**Online scalability**

- Independent linear scaling of both performance and capacity
- Application non-disruptive

**Operationally simple regardless of the scale**

- From the very small to the incredibly large, management stays simple

**Data availability is optimum no matter the disaster or situation**

- Available and online if there is a single or dual disk fault in a RAID group
- Available and online if there is a storage system or controller failure
- Available and online if there is a site failure
- Available and online if there is an accidental deletion, admin error, or malware
- Available and online during scheduled maintenance
- In the event of a disaster data availability is incredibly fast

**The Answer – BlueArc Titan Series**

BlueArc designed its Titan series to eliminate each and every one of these backbreaking problems.

**Titan adapts dynamically to changing performance requirements in real time**

Most other unified NAS/SAN storage systems today cannot adapt to performance demand spikes as well as the BlueArc Titan. It was designed from the ground up for the best possible performance. Common performance-draining storage functions such as CIFS and NFS have been put into the silicon. Additionally, BlueArc is well known for ultra-low latency (measure of how quickly a server can respond to clients' requests) for server IO request handling. The lower the latency, the greater number of clients and IO that can be handled. Titan maintains this rapid response even under load as each part of the architecture is designed to handle the maximum load simultaneously with distributed memory, state engines and pipelines that enable this massive parallelism.

**BlueArc Titan provides industry-leading Online Scalability**

Titan is architected to scale both up – the scaling capacity within the NAS/SAN filer, and out – the clustering of more than two heads in a single image. Titan scales up by virtualizing backend RAID storage. Because the Titan can scale capacity by adding additional RAID controllers as well as hard disk drives, more processing power is available for the storage subsystem. This allows performance and capacity to scale linearly. Standard NAS filers do not scale the RAID processing meaning that at a give scale, performance actually declines as capacity increases.

Titan scales out by seamlessly clustering up to eight Titan nodes in a single system. When paired with its clever implementations of virtualized backend SAN storage, thin provisioning, and multi-head clustering, the scaling is online, application non-disruptive, and effortless.

Each Titan supports astonishing scalability including:

- 12,000 Windows (CIFS) users
- 60,000 UNIX (NFS v2 or NFS v3) users
- Up to 16 million files per directory
- Up to 8 billion total files
- Up to 200,000 IOPS (NFS v3)
- Up to 20Gbps aggregate throughput
- Up to 4 PB of total usable storage (per Titan or cluster)
- Up to 256 TB per file system

With up to 8 Titans in a single clustered system image, Titan sets the bar for scalability.

A major reason for Titan's uncommon scalability is that its smart hardware architecture and hardware-based file system provides task processing at extraordinarily high speed while its IO latency is minimal.

### **BlueArc Titan Simplifies Operations**

NAS is by definition the easiest form of storage to implement and utilize. However, its shortcomings in performance, spikes, online scalability, operational complexity, and data availability limit its usefulness for Internet services. Scaling traditional NAS storage can lead to a proliferation of storage servers creating more work for overburdened operational staff.

Titan goes further than other storage systems. With an intuitive graphical user interface, the Titan is easy to implement, easy to operate, easy to scale up, easy to scale out, easy to manage, and easy to use. It makes both common and complicated operations simple by enabling Internet service companies to massively scale NAS storage behind a single managed platform

### **BlueArc Titan provides best-in-class Data Availability**

The Titan's best in class data availability includes:

#### **More than 500,000 hours MTBF (> 57 years)**

- High MTBF means a lot less probability of unscheduled downtime. Less scheduled downtime equals higher data availability.

#### **RAID 6 protection**

- RAID 6 is the only RAID protection that protects against a second drive failure or fault in a RAID group while the first is being rebuilt. The probability of a 2nd drive failure or fault goes up exponentially when the 1st drive fails.
- Titan's RAID 6 means a much lower probability of data loss, which equates into much higher data availability.

#### **Snapshots & Quick Restore**

- BlueArc Titan can store up to 1,024 "rules-based" snapshot images per volume. Rules-based snapshots are tied to business rules allowing entity versus a volume-based snapshot management. Snapshots are configurable to allow users to automatically and rapidly recover their deleted or older versions of files from snapshots, without requiring the system administrator to restore it first from tape, while preserving all files permissions.
- This result is increased data availability by multiple orders of magnitude.

#### **Incremental Data Replication (IDR)**

- IDR provides incremental backup of a volume or a directory on the volume. IDR uses snapshots as the basis for replication, utilizing the last snapshot as the reference point for the next replication to occur. Snapshot based incremental copying enables backup of files in use that would otherwise be skipped while replicating only changed files, drastically reducing the amount of replication traffic. IDR supports inter-Titan replication (local or remote) as well as intra-Titan replication.

#### **Advanced NDMP (LAN-free backup)**

- The BlueArc Titan hardware implementation of this standard server-free backup keeps applications and their data online during the backup process. Traditional backup is an extremely server intensive commonly requiring the servers, applications, and their data to be offline during the backup. Titan's SiliconServer Architecture™ avoids this issue with a unique dedicated separate-state machine to data movement during the backup and restore processes.
- Because this silicon based NDMP handles all of the data movement during these operations, servers, applications, and data remains available and online.



### **Anti-Virus Support**

- Full integration with and support for the most popular anti-virus infrastructure from Symantec, McAfee, Sophos and Trend Micro prevents “Malware” from corrupting Titan’s storage pools. Files are scanned directly upon open or close procedures, configurable per share scanning, or scanning via inclusion/exclusion lists.
- Prevention and or spread of Malware infections, makes data availability higher and measurably simpler.

### **Data Migrator**

- Data Migrator automates data migration from one storage tier to another, based on pre-defined policies, rules and triggers. Data migration occurs transparently to both end users and applications.
- This permits data to remain available online longer, by moving aged or non-performance centric data to a lower cost storage tier transparently, without notifying users or changing applications.

### **Scalable N-Way Active-Active Clustering**

- BlueArc Titan’s scalable N-Way active-active clustering technology gives shared SAN storage resource access to all Titan nodes while providing fail-over for data protection of mission critical applications. Clusters configurations support start at 2 nodes and scale up to 8 node configurations today.
- Not only does this provide increased data availability if a Titan node were to fail, it also provides continuous data availability during upgrades and maintenance. No scheduled downtime is required as each node is upgraded or maintained on a sequential basis.

### **Cluster Name Space**

- The Titan Cluster Name Space (CNS) creates a Global Name Space (GNS) across multiple Titan storage pools. Multiple file systems appear as a single unified directory structure giving both CIFS and NFS clients global access to the data through any node in the cluster.
- This simplifies file system growth beyond 256TB reduces operational complexity (e.g. no moving data between file systems), and increases data availability by eliminating multiple mount points.

### **Remote Volume Mirroring (RVM)**

- Asynchronous or Synchronous volume mirroring is very useful for organizations like Internet services that cannot lose any data nor tolerate downtime, to protect the data.
- RVM means higher and faster data availability in the event of a system or site failure.

### **WORM file system**

- Robust, file system-based Write-Once, Read-Many (WORM) functionality prevents file system data from being modified or deleted until a specific retention date has passed.
- Permits data availability to be online longer with easier accessibility.

### **Cluster Read Caching**

- Cluster Read Caching allows read intensive workloads such as data mining, analytics and library to be accessible concurrently across all clustered nodes.
- This quantifiably increases aggregate read intensive workloads in both throughput and IOPS adjusting easily for performance spikes, preventing service outages, and increasing data availability.

Titan’s data availability capabilities mean that suffering an outage because of inadequate data availability is an exceedingly remote possibility.

## **Conclusion**

Internet services of all types are remarkably competitive. With quality of service at a reasonable price being king, it is foolish to implement storage systems that will ultimately cost more out-of-pocket, run off customers, decimate the business’ reputation, while reducing revenues and growth.

It doesn’t have to be that way. Any Internet services company seeking to minimize service outages while maximizing uptime, and maintaining a cost effective cost infrastructure that grows with them, should consider using the BlueArc Titan series as its primary storage. It is a sure-fire way to eliminate those Internet services storage-induced outage blues.

## About BlueArc

BlueArc is a leading provider of high performance unified network storage systems to enterprise markets, as well as data intensive markets, such as electronic discovery, entertainment, federal government, higher education, Internet services, oil and gas and life sciences. Our products support both network attached storage, or NAS, and storage area network, or SAN, services on a converged network storage platform.

We enable companies to expand the ways they explore, discover, research, create, process and innovate in data-intensive environments. Our products replace complex and performance-limited products with high performance, scalable and easy to use systems capable of handling the most data intensive applications and environments. Further, we believe that our energy efficient design and our products' ability to consolidate legacy storage infrastructures, dramatically increases storage utilization rates and reduces our customers' total cost of ownership.



**BlueArc Corporation**  
*Corporate Headquarters*  
50 Rio Robles Drive  
San Jose, CA 95134  
t 408 576 6600  
f 408 576 6601  
www.bluearc.com

**BlueArc UK Ltd.**  
*European Headquarters*  
Queensgate House  
Cookham Road  
Bracknell RG12 1RB, United Kingdom  
t +44 (0) 1344 408 200  
f +44 (0) 1344 408 202